



**Universitatea “Transilvania” din Braşov**  
**FACULTATEA DE INGINERIE ELECTRICĂ ŞI ŞTIINŢA**  
**CALCULATOARELOR**

## **TEZĂ DE DOCTORAT**

**SISTEME ADAPTIVE CU REŢELE NEURONALE  
ARTIFICIALE NEUROMORFE. REALIZĂRI CU  
DISPOZITIVE HARDWARE RECONFIGURABILE**

Doctorand:  
**ing. László BAKÓ**

Conducător științific:  
**prof. dr. ing. Iuliu SZÉKELY**

**2009**



## Cuvânt înainte

Domeniul Inteligenței Artificiale și a mașinilor cu autoinstruire constă în sisteme artificiale și tehnici care manifestă un comportament „inteligent” și pot învăța din exemple. O gamă variată de algoritmi care sunt utilizați în acest domeniu de cercetare sunt inspirați de funcționarea cortexului mamiferelor. Creierul este format dintr-o rețea mare de celule, numite neuroni, care sunt conectate între ele într-un mod complicat. Rețele de funcții non-liniare simple, numite Rețele Neuronale, surprind unele concepte computaționale prezente și în creier. În ultimii ani, în aceste modele au fost incluse din ce în ce mai multe proprietăți de natură biologică. Aceste cercetări se concentrează pe Rețele Neuronale Pulsative (RNP), rețele de sisteme cu dinamică simplă, care surprind caracteristicile temporale simple ale neuronilor biologici, care comunică informația printr-o serie de impulsuri. Aceste serii de impulsuri sunt secvențe de evenimente discrete în timp, în care doar momentul impulsului reprezintă informație transmisă, nu și amplitudinea acestuia. Cercetările recente au arătat că aceste rețele sunt mai puternice și mai eficiente decât rețelele de tip analog folosite înainte, și recent au fost construite câteva aplicații foarte convingătoare cu ajutorul acestor rețele. Deoarece informația este reprezentată în timp, RNP-urile pot procesa informația temporală într-un mod intrinsec. Modelele mai exacte din punct de vedere biologic, folosite în RNP-uri, ne permit să tragem concluzii asupra funcționării creierului și a rețelelor biologice, pe baza rezultatelor obținute prin folosirea modelelor RNP. De asemenea, sugestii pentru construirea de sisteme artificiale pot fi obținute în mod direct din studiile neurobiologice.

Deși RNP-urile au multe proprietăți benefice, și acestea au fost demonstrate cu succes de câteva aplicații și implementări, crearea de aplicații de succes folosind RNP este încă foarte greoaie din punct de vedere ingineresc. În continuare subliniem câteva probleme de actualitate ale RNP: codificarea informației în serii de impulsuri, antrenarea supervizată a RNP, și implementările hardware. Scopul primordial al acestei teze este de a oferi tehnici și instrumente de implementare hardware care ajută RNP-urile să fie general aplicabile ca instrumente ingineresti.

Există ipoteze multiple asupra modului în care neuronii biologici codifică informația în serii de impulsuri. Avem două ipoteze principale bazate pe idei opozante: momentul exact al impulsurilor nu este important, deoarece toată informația este înmagazinată în rata de activare și la cealaltă extremă: toată informația este înmagazinată în momentul exact al impulsurilor. O schemă de codificare care apropie diferența de concept dintre aceste două scheme este codificarea cu filtre: informația codificată într-o serie de impulsuri poate fi extrasă cu succes folosind un filtru liniar. Din punct de vedere ingineresc însă, codificarea informației într-o serie de impulsuri, folosind această schemă, nu este deloc simplă.

O problemă foarte importantă a RNP-urilor este lipsa regulilor de învățare supervizate, general aplicabile, care să poată antrena RNP-urile în mod robust și rapid. Există o singură regulă de învățare supervizată, dar aceasta se poate aplica doar în cazul unei scheme de codificare specifice, care de fapt ignoră natura temporală intrinsecă a RNP. De asemenea nu are o convergență de antrenare robustă, este înceată și necesită o topologie de rețea specială, destul de mare.

Putem îmbunătăți viteza de convergență și cerințele de rețea ale acestei reguli prin antrenarea tuturor parametrilor neuronilor pulsativi. Celelalte probleme rămân nerezolvate însă. Din acest motiv vom introduce o nouă regulă de învățare, bazată pe aproximarea analogică a unui neuron pulsativ. Această regulă poate fi aplicată când folosim o versiune generalizată a schemei de codificare a momentelor de activare, converge în mod rapid și robust. Necesită doar topologii de rețea mici și este ușor implementabilă în hardware digital.

La prima vedere se pare că datorită naturii mai complexe a ecuațiilor necesare pentru calcularea RNP, acestea vor necesita mai mult timp de calcul decât rețelele tradiționale analogice. Vom demonstra că RNP-urile sunt bine adaptate pentru a fi implementate pe hardware digital. Mai multe proprietăți, cum ar fi natura binară, pulsativă a comunicației pot fi eficient implementate în hardware. S-a demonstrat deja, că se pot obține creșteri mari în viteză și în eficiența de consum prin implementarea de RNP în hardware analogic. Pe de altă parte, ne vom concentra asupra hardware-ului digital, deoarece folosind dispozitive digitale reconfigurabile, este posibilă implementarea designului RNP real în hardware, chiar și cu reconfigurare parțială în timpul funcționării când este nevoie. Vom arăta că există mai multe arhitecturi hardware posibile pentru a implementa RNP, fiecare având alt domeniu și alte caracteristici de viteză. Acest lucru permite proiectantului să aleagă o implementare optimă după ce s-au luat în considerare diverse constrângeri impuse de aplicația dată. Sunt posibile accelerări foarte mari sau implementări foarte mici, în funcție de ce necesită aplicația dată.

Pentru a demonstra tehnicile prezentate în această teză se vor prezenta mai multe aplicații ingineresti și aplicații de test (benchmark).

Aș dori să exprim mulțumiri pe această cale în primul rând d-lui Prof. dr. ing. Iuliu Székeley care mi-a fost alături cu sfaturi, m-a îndrumat, sprijinit și m-a ajutat pe toată durata elaborării tezei în calitate de conducător științific.

Mulțumesc domnilor Prof. dr. ing. Dávid László și Conf. dr. ing. Márton László pentru sprijinul acordat în rezolvarea problemelor ivite pe durata elaborării tezei.

De asemenea, mulțumesc colectivului Catedrei de Inginerie Electrică din cadrul Facultății de Științe Tehnice și Umaniste a Universității Sapientia, pentru sprijinul acordat și pentru atmosfera de lucru indispensabilă realizării tezei.

Mulțumesc Institutului Programelor de Cercetare al Fundației Sapientia și Fundațiilor Communitas și Eurotrans pentru bursele acordate, respectiv CNCSIS UEFISCSU, pentru acordarea grantului de cercetare TD în decursul finalizării stagiului doctoral.

În final, dar nu în ultimul rând, doresc să mulțumesc membrilor familiei mele în egală măsură, pentru încurajare, răbdare și pentru suportul fără de care nu aș fi reușit să parcurg acest drum.





# CUPRINS

<b>Cuprins .....</b>	<b>V</b>
<b>Lista figurilor .....</b>	<b>ix</b>
<b>1. Introducere.....</b>	<b>1</b>
1.1. Inteligență artificială și tehnici de machine learning .....	1
1.2. Rețele neuronale artificiale .....	3
1.3. Fundamente biologice.....	3
1.3.1. Neuronul pulsativ ideal .....	3
1.3.2. Impulsuri de activare.....	3
1.3.3. Sinapsele .....	4
1.3.4. Dinamismul neuronal.....	4
1.3.5. Pragul de activare sau aprindere și potențialul de activare .....	5
1.4. Sisteme de control neuronal.....	5
1.4.1. Control neuronal indirect .....	6
1.4.2. Model de rețea neuronală parțială sau parametrică.....	7
1.4.3. Control neuronal Direct .....	10
1.4.3.1. Neuro-control fără modelul sistemului.....	10
1.4.3.2. Control neuronal bazat pe model .....	11
1.4.3.3. Regulator neuronal robust pe bază de model .....	11
1.5. Organizarea tezei .....	12
<b>2. Modelarea și simularea rețelelor neuronale neuromorfe.....</b>	<b>13</b>
2.1. Alegerea modelului optim pentru neuronii pulsativi corticali .....	13
2.1.1. Introducere .....	13
2.1.2. Dispute curente în domeniul RNP .....	14
2.1.3. Caracteristici neuro-computaționale .....	15
2.2. Realizarea învățării Hebbiene competitive .....	19
2.2.1. Excitație echilibrată.....	19
2.2.2. Discuții.....	20
2.3. Metode de propagare înapoi a erorilor (back-propagation) în rețelele neuronale pulsative codificate temporal .....	21
2.3.1. Propagarea înapoi a erorii .....	22
2.3.2. Problema XOR.....	25
2.3.3. Gradientul erorii și rata de învățare.....	26
2.4. Dezvoltarea modelelor neuronale neuromorfe.....	27
2.4.1. Modelul unei rețele neuronale artificiale formale cu neuroni pulsativi .....	28
2.4.2. Modelul de reacție la impuls .....	29
2.5. Simularea rețelelor neuronale pulsative: o sinteză a instrumentelor și a strategiilor .....	29
2.5.1. Introducere .....	29
2.5.2. Strategii de simulare.....	30
2.5.2.1. Formalismul unui sistem hibrid.....	30
2.5.2.2. Folosirea liniarităților pentru simulări sinaptice rapide.....	31
2.5.2.3. Algoritmi sincroni sau bazați pe tact.....	32
2.5.2.4. Algoritmi asincroni sau bazați pe eveniment .....	33
2.5.3. Prezentare generală a mediilor de simulare .....	36
2.5.3.1. NEURON .....	36
2.5.3.2. GENESIS .....	37
2.5.3.3. NEST.....	38
2.5.3.4. SPLIT .....	39
2.5.3.5. Mvaspike.....	39
2.6. Implementări software proprii, validarea modelelor dezvoltate .....	40
2.6.1. Idei de bază .....	40
2.6.2. Algoritmul simulării.....	40
2.6.3. Încorporarea algoritmilor de învățare în procesul de simulare .....	41
2.6.4. Învățare bazată pe valori de prag .....	43

2.7.	Rezultate experimentale a simulărilor software.....	44
2.8.	Concluzii.....	45
<b>3.</b>	<b>Dispozitive FPGA și implementarea rețelelor neuronale artificiale .....</b>	<b>47</b>
3.1.	Evoluția dispozitivelor reconfigurabile.....	48
3.1.1.	Evaluarea performanței aplicațiilor bazate pe circuite FPGA.....	49
3.1.1.1.	Caracteristici FPGA .....	50
3.1.1.2.	Programarea circuitelor FPGA.....	50
3.1.1.3.	Lățimea de bandă și localizarea datelor .....	51
3.1.1.4.	Memorie suficientă?.....	52
3.1.2.	Clase de FPGA.....	52
3.1.3.	Elementele unui circuit FPGA din familia Xilinx Spartan3 .....	55
3.1.3.1.	Elementele Bloc RAM .....	56
3.1.3.2.	RAM distribuit în CLB .....	59
3.2.	Implementări RNA pe FPGA.....	61
3.2.1.	Moduri de cuplare și configurări pentru FPGA .....	61
3.2.2.	Maparea algoritmilor de rețele neuronale artificiale pe sisteme FPGA .....	62
3.2.3.	Implementarea funcțiilor de activare .....	68
3.2.4.	Aritmetică neurală/Reprezentarea datelor.....	72
3.3.	Concluzii.....	74
<b>4.</b>	<b>Posibilități de implementare total paralelă a modelelor neuronale pulsative cu circuite FPGA</b>	<b>77</b>
4.1.	Considerații inițiale.....	77
4.2.	Caracteristici generale ale sistemelor hardware pentru implementarea rețelelor neuronale. Motivații pro și contra.....	77
4.2.1.	Implementări analogice.....	77
4.2.1.1.	Avantajele implementărilor analogice: .....	77
4.2.1.2.	Dezavantajele implementărilor analogice: .....	77
4.2.2.	Implementări digitale .....	78
4.2.2.1.	Avantajele implementărilor digitale.....	78
4.2.2.2.	Dezavantajele implementărilor digitale.....	79
4.2.3.	Calificarea implementărilor .....	79
4.2.4.	Implementări RNA pe FPGA.....	79
4.3.	Sinteză asupra diferitelor implementări hardware de rețele neuronale artificiale.....	80
4.3.1.1.	Legea lui Amdahl.....	81
4.3.1.2.	Legea lui Gustafson.....	81
4.3.2.	Implementări hardware neuronale analogice .....	82
4.3.3.	Implementări hardware neuronale digitale.....	83
4.3.4.	Aplicații ale implementărilor hardware de rețele neuronale .....	84
4.3.4.1.	Accurate Automation Corp (AAC) .....	84
4.3.4.2.	Adaptive Solutions CNAPS .....	84
4.3.4.3.	IBM ZISC036 .....	85
4.3.4.4.	Intel 80170NX Electrically Trainable Analog Neural Network (ETANN).....	85
4.3.4.5.	Irvine Sensors 3DANN .....	85
4.3.4.6.	National Semiconductor NeuFuz/COP8 Microcontrollers.....	86
4.3.4.7.	Nestor NI1000.....	86
4.3.4.8.	Philips L-Neuro chips .....	87
4.3.4.9.	Sensory Circuits RSC-164 Speech Recognition Chip.....	87
4.3.4.10.	Siemens MA-16 chip, SYNAPSE-3 Neurocomputer.....	87
4.4.	Implementări software și hardware ale rețelelor neuronale bazate pe impulsuri.....	88
4.4.1.	Neuronii pulsativi.....	88
4.4.2.	Probleme de implementare.....	88
4.4.3.	Implementare pe calculator paralel .....	89
4.4.4.	SpikeNET.....	90
4.4.4.1.	Sistemul de simulare rețele neuronale neuromorfe SpikeNET.....	90
4.4.5.	Rețele neuronale pulsative, dezvoltări hardware.....	91
4.5.	Implementarea cu circuit FPGA a modelului neuronal propriu dezvoltat, validat prin simulare software .....	95
4.5.1.	Descrierea sumară a sistemului de dezvoltare FPGA .....	96

4.5.1.1.	Structura plăcii de dezvoltare.....	96
4.5.1.2.	Elementul central al plăcii XSA100, circuitul FPGA Xilinx XC2S100.....	97
4.5.1.3.	Mediul de dezvoltare software VHDL Xilinx ISE Webpack.....	97
4.5.2.	Procesul de realizare practică a modelului neuronal, implementarea sinapsei.....	98
4.5.2.1.	Fazele implementării sinapsei în FPGA.....	98
4.5.2.2.	Descrierea detaliată a funcționării sinapsei.....	100
4.5.2.3.	Îmbunătățirea sinapsei.....	100
4.5.2.4.	Unitatea serială de intrare date (USID).....	101
4.5.3.	Fazele proiectării somei (al corpului neuronal).....	102
4.5.3.1.	Descrierea detaliată a funcționării unității centrale a somei, adică a modului de calcul a potențialului de membrană.....	103
4.5.4.	Stadiul actual de dezvoltare a somei.....	105
4.6.	Prezentarea modelului neuronal implementat complet.....	106
4.6.1.	Programul de control al circuitului neuronal prin portul paralel al PC.....	106
4.6.2.	Câteva rezultate experimentale de testare a NP implementat, cu diferite configurații parametrice.....	107
4.7.	Implementarea unei rețele neuronale pulsative artificiale în circuitul FPGA. Rezultate experimentale.....	109
4.7.1.	Implementarea unei RNP în circuit FPGA.....	109
4.7.2.	Descrierea algoritmului de învățare supervizată implementat.....	110
4.7.3.	Măsurători și rezultate experimentale.....	111
4.8.	Analiza performanțelor rețelei neuronale neuromorfe implementate.....	113
4.8.1.	Rezultatele obținute de rețeaua neuronală implementată în hardware.....	113
4.8.2.	Comparație cu o rețea neuronală clasică realizată în software.....	114
4.9.	Concluzii.....	115
<b>5.</b>	<b>Implementări parțial paralele ale RNP utilizând microcontrolere încorporate în circuite FPGA</b> .....	<b>117</b>
5.1.	Prezentarea rezultatelor actuale ale domeniului din literatura de specialitate.....	117
5.2.	Clasificarea nesupervizată de mulțimi complexe (complex clusters) cu rețele de neuroni pulsativi.....	118
5.3.	Implementarea unei rețele de neuroni pulsativi cu întârziere.....	119
5.3.1.	Structura rețelei.....	119
5.3.2.	Variabile de intrare continue codificate în decalaje temporale.....	119
5.3.2.1.	Implementarea codării valorilor de intrare în decalaje temporale pe FPGA – Neuronii de intrare, de tip pseudo-RBF.....	120
5.4.	Realizarea unei rețele de neuroni pulsativi cu antrenare on-chip și aplicarea sa practică.....	121
5.4.1.	Rețeaua de neuroni pulsativi realizată în FPGA.....	123
5.4.1.1.	Modelul neuronal RBF pulsativ implementat.....	123
5.4.2.	Arhitectura rețelei implementate pentru clasificarea de fascicule cu domenii receptive.....	124
5.4.2.1.	Structura proiectului VHDL care implementează rețeaua neuronală.....	125
5.4.2.2.	Neuronul de intrare.....	126
5.4.2.3.	Modulul sinapsă.....	128
5.4.2.4.	Modulul de some al rețelei neuronale pseudo-RBF-pulsative.....	128
5.4.2.5.	Modulul de comunicație cu calculatorul.....	128
5.5.	Realizări de rețele neuronale neuromorfe utilizând microcontrolere soft-core înglobate în circuitele FPGA.....	128
5.5.1.	Microcontrolerul soft-core Xilinx PicoBlaze.....	129
5.5.2.	Implementare aplicativă: Detector de componente de frecvență.....	130
5.5.2.1.	Rețeaua neuronală pulsativă implementată pe FPGA.....	130
5.5.2.2.	Implementarea neuronilor de intrare – codificarea variabilelor de intrare în impulsuri decalate temporal.....	131
5.5.2.3.	Blocul neuronilor de intrare.....	133
5.5.2.4.	Modulul de implementare a sinapselor.....	135
5.5.2.5.	Modulul neuronilor de ieșire – implementarea somei cu multiple microcontrolere încorporate Xilinx PicoBlaze.....	136
5.5.2.6.	Dinamismul funcționării modului somă – Rezultate experimentale.....	137

5.5.2.7.	Descrierea sistemului test-bench dezvoltat. Prepararea setului de date de antrenare .....	139
5.5.2.8.	Testarea sistemului neuronal implementat .....	140
5.5.3.	Implementare test benchmark: Clasificarea setului de date Fisher IRIS.....	142
5.5.3.1.	Descrierea setului de date. Formularea problemei de clasificare .....	142
5.5.3.2.	Rețeaua neuronală pulsativă implementată pe FPGA .....	143
5.5.3.3.	Implementarea neuronilor de intrare – codificarea variabilelor de intrare în impulsuri decalate temporal .....	143
5.5.3.4.	Modulul de implementare a sinapselor .....	144
5.5.3.5.	Dinamismul funcționării modulului somă – Rezultate experimentale .....	144
5.5.3.6.	Descrierea sistemului testbench dezvoltat. Prepararea setului de date de antrenare .....	145
5.5.4.	Sistemul de dezvoltare utilizat pentru realizarea proiectelor cu procesoare încorporate .....	145
5.6.	Realizări de RNP cu microcontroler pe 32 de biți, încorporat în circuitele FPGA.....	146
5.6.1.1.	Microcontrolerul Xilinx MicroBlaze.....	146
5.6.2.	Implementare aplicativă: Recunoaștere de caractere cu o rețea neuronală neuromorfă multistrat .....	147
5.6.2.1.	Formularea problemei de recunoaștere de caractere .....	147
5.6.2.2.	Dezvoltarea componentelor hardware ale sistemului încorporat.....	147
5.6.2.3.	Structura rețelei neuronale pulsative dezvoltate.....	148
5.6.2.4.	Modelul neuronal utilizat .....	148
5.6.2.5.	Prezentarea algoritmului de învățare implementat în RNP .....	150
5.6.2.6.	Proiectarea componentelor software ale sistemului încorporat .....	151
5.6.2.7.	Rezultate experimentale .....	153
5.6.3.	Implementare benchmark: Clasificarea setului de date Wisconsin Breast Cancer ..	157
5.6.3.1.	Formularea problemei de clasificare .....	157
5.6.3.2.	Elementele hardware ale sistemului încorporat dezvoltat .....	157
5.6.3.3.	Principiul de funcționare perifericului de codificare a intrărilor .....	158
5.6.3.4.	Codificarea valorilor de intrare .....	159
5.6.3.5.	Structura rețelelor neuronale pulsative implementate pentru rezolvarea clasificării setului de date Wisconsin .....	159
5.6.3.6.	Algoritmul de învățare implementat .....	161
5.6.3.7.	Utilizarea proiecției lui Sammon pentru vizualizarea mulțimilor multidimensionale .....	161
5.6.3.1.	Rezultate experimentale .....	162
5.6.4.	Mediul de dezvoltare utilizat pentru sisteme încorporate .....	164
5.7.	Concluzii .....	166
<b>6.</b>	<b>Concluzii finale și contribuții originale.....</b>	<b>169</b>
6.1.	Concluzii finale.....	169
6.2.	Contribuții originale.....	173
<b>Bibliografie</b>	<b>.....</b>	<b>175</b>
ANEXA 1 – Distribuția componentelor utilizate pe circuitele FPGA (detector frecvențe vs. clasificare IRIS dataset) .....		183
ANEXA 2- Schema de implementare a unui neuron pulsativ .....		185
ANEXA 3- Schema de implementare a unei rețele neuronale pulsative.....		187
ANEXA 4- Schema de implementare a rețele neuronale pulsative pentru clasificare setului de date IRIS... ..		189

## LISTA FIGURILOR

1. Figura 1.1 Primele modele de neuroni și rețele neuronale .....	1
2. Figura 1.2 Reprezentarea schematică a conceptului de învățare în rețele neuronale artificiale .....	2
3. Figura 1.3 Structura unui neuron natural.....	3
4. Figura 1.4 <b>A.</b> Schema unei sinapse chimice; <b>B.</b> Fotografia electro-microscopică a unei sinapse (Trappenberg, 2002) .....	4
5. Figura 1.5 Neuronul post-sinaptic nr. $i$ recepționează impulsuri de activare de la neuronii pre-sinaptici $j=1,2,4$	
6. Figura 1.6 Rețea neuronală privită ca o cutie neagră .....	6
7. Figura 1.7 O rețea neuronală poate fi incorporată într-o schemă de control predictiv bazat pe model .....	7
8. Figura 1.8 Rețeaua neuronală ca estimator de parametrii.....	8
9. Figura 1.9 Control cu rețea neuronală bazată pe modelul invers al procesului .....	8
10. Figura 1.10 Regulator în buclă închisă cu autoacordare cu RNA .....	9
11. Figura 1.11 Modelare regulator .....	10
12. Figura 1.12 Neuro-control fără modelul sistemului .....	11
13. Figura 1.13 Control neuronal bazat pe model .....	11
14. Figura 1.14 Regulator neuronal robust pe bază de model .....	12
15. Figura 2.1 Structura de bază a unui neuron pulsativ .....	13
16. Figura 2.2 Tipuri de neuroni pulsativi .....	16
17. Figura 2.3 Comparatie a diferitelor modele de neuroni pulsativi .....	18
18. Figura 2.4 Corelația dintre încărcarea pre- și post-sinaptică .....	20
19. Figura 2.5 Relația dintre $\partial x_j$ și $\partial t_j$ pentru un spațiu $\mathcal{E}$ în jurul $t_j$ .....	23
20. Figura 2.6 Învățare XOR: numărul mediu necesar de iterații de învățare pentru a atinge suma propusă a erorilor pătratică de 1.0. ....	26
21. Figura 2.7 Învățare XOR: Numărul de cazuri din 10 care au convers.....	27
22. Figura 2.8 Numărul de cazuri din 10 care au convers pentru diferite valori ale $\tau$ .....	27
23. Figura 2.9 Forma funcției de prag al unui neuron natural .....	28
24. Figura 2.10 Forma funcției de răspuns a NP natural (sus PPSE și PPSI jos) .....	28
25. Figura 2.11 Un algoritm simplu bazat pe tact .....	32
26. Figura 2.12 Un algoritm simplu bazat pe eveniment cu interacțiuni sinaptice.....	34
27. Figura 2.13 Un algoritm simplu bazat pe eveniment cu interacțiuni sinaptice non-instantanee.....	35
28. Figura 2.14 Schema algoritmului de simulare.....	41
29. Figura 2.15 Regula de învățare bazată pe valori de prag.....	43
30. Figura 2.16 O rețea neuronală realizată în mediul de simulare dezvoltat Aplicația: recunoaștere de caractere dintr-o matrice de 7x5 puncte (albastru – neuroni de intrare, verde – neuroni de ieșire, roșu – sinapse).....	44
31. Figura 2.17 Aceeași problemă, rezolvată cu o rețea cu un strat ascuns.....	44
32. Figura 2.18 IA în timpul funcționării rețelei neuronale prezentate în Figura 2.16. ....	44
33. Figura 2.19 Rezultatul măsurătorii efectuate pe ieșirea RNP din Figura 2.17. Primii treizeci de pași temporali .....	45
34. Figura 3.1 Caracteristici FPGA: evoluția creșterii unităților logice și a vitezei de tact.....	50
35. Figura 3.2 Caracteristici FPGA: evoluția creșterii vitezei de calcul și a lățimii de bandă.....	50
36. Figura 3.3 Principalele clase de FPGA.....	52
37. Figura 3.4 Exemplu de rutare tip matrice simetrică .....	53
38. Figura 3.5 Exemplu Bloc de conectare.....	54
39. Figura 3.6 Exemplu de bloc comutator .....	54
40. Figura 3.7 Organizarea elementelor în circuitul FPGA.....	55
41. Figura 3.8 Ciclu de scriere în modul WRITE_FIRST .....	59
42. Figura 3.9 Fluxul de date în modul WRITE_FIRST .....	59
43. Figura 3.10 Ciclu de scriere în modul READ_FIRST .....	59
44. Figura 3.11 Fluxul de date în modul READ_FIRST .....	59
45. Figura 3.12 Ciclu de scriere în modul NO_CHANGE.....	59
46. Figura 3.13 Fluxul de date în modul NO_CHANGE .....	59
47. Figura 3.14 Fluxul de date Single-Port și Dual-Port RAM distribuit.....	60
48. Figura 3.15 Diagrama de timp a unei operații de scriere.....	60
49. Figura 3.16 Semnalele de control la RAM distribuit.....	61
50. Figura 3.17 Aproximare Taylor a funcției sigmoid.....	68
51. Figura 3.18 Aproximarea liniară pe porțiuni a funcției de activare sigmoid .....	69

52. Figura 3.19 Aproximarea funcției sigmoid prin metoda CRI.....	69
53. Figura 3.20 Aproximarea funcției de activare sigmoid prin polinom pătratic de gradul doi.....	70
54. Figura 3.21 Aproximarea funcției de activare prin două segmente neliniare.....	71
55. Figura 3.22 Aproximarea funcției sigmoid clasic printr-o funcție particulară pe bază de puterile lui 2.....	71
56. Figura 3.23 Reprezentarea valorii 7-16 prin flux de impulsuri.....	73
57. Figura 4.1 Pașii OCR, conform manualului de utilizare al Adaptive Solutions CNAPS.....	81
58. Figura 4.2 Legea lui Amdahl.....	81
59. Figura 4.3 Modelul neuronal conceput de Thorpe.....	89
60. Figura 4.4 Structura de bază a sistemului SpikeNET.....	90
61. Figura 4.5.- Structura sistemului din (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005) pe un circuit FPAG Xilinx Virtex XC2V8000.....	92
62. Figura 4.6 Pașii de execuție a sistemului cu elemente de rețea multiplexate prezentată în (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005).....	92
63. Figura 4.7 Diagrama bloc a topologiei procesorului neuronal SIMD.....	94
64. Figura 4.8 Structura rețelei din (Upegui, Peñ a-Reyes, & Sánchez, 2004).....	94
65. Figura 4.9 Diagramele bloc ale unui neuron și a unei sinapse din SIMD.....	95
66. Figura 4.10 Schema arhitecturală a plăcii de dezvoltare FPGA XSA100, produsă de firma XESS (SUA), cu un circuit FPGA Xilinx XC2S100.....	96
67. Figura 4.11 Structura internă a circuitelor FPGA din familia Spartan (stânga), părțile componente al unui bloc logic programabil (dreapta) (CLB).....	97
68. Figura 4.12 Fazele procesului de proiectare FPGA în mediul Xilinx Webpack.....	97
69. Figura 4.13 Schema structurală a sinapsei.....	98
70. Figura 4.14 Schema de conectare interioară a sinapsei.....	100
71. Figura 4.15 Schema de principiu a sinapsei.....	101
72. Figura 4.16 Schema de conectare a sinapsei modificate.....	101
73. Figura 4.17 Schema constructivă a USID.....	102
74. Figura 4.18 Transmisia serială de date către FPGA.....	102
75. Figura 4.19 Schema constructivă a circuitului ce realizează soma neuronului.....	103
76. Figura 4.20 Modelul somei implementate, într-o fază intermediară a proiectării.....	104
77. Figura 4.21 Circuitul de testare a modului central al somei, MPOT.....	104
78. Figura 4.22 Schema somei implementate în neuronul pulsativ complet.....	105
79. Figura 4.23 Valori ale potențialului de membrană, eșantionate în timpul testării.....	105
80. Figura 4.24 Schema de principiu a somei implementate hardware în FPGA.....	106
81. Figura 4.25– Impulsuri axonale într-un experiment cu potențial de prag scăzut.....	107
82. Figura 4.26 IA utilizând valori inițiale mari ale ponderilor și impulsuri de intrare distribuite – prima parte.....	108
83. Figura 4.27– IA utilizând valori inițiale mari ale ponderilor și impulsuri de intrare distribuite – partea a doua.....	108
84. Figura 4.28 Arhitectura rețelei neuronale pulsative implementate.....	109
85. Figura 4.29 Perechile de modele de intrare, utilizate în cursul învățării (notați ordinea și valorile biților.....	110
86. Figura 4.30 Comunicația dintre diferitele echipamente ale sistemului.....	111
87. Figura 4.31 Transmisii seriale ale valorilor de intrare cu răspunsurile axonale date de neuroni la acestea ...	112
88. Figura 4.32 Citirea valorilor de pondere sinaptică (hexadecimal) între două faze de transmitere a perechii de modele de intrare.....	112
89. Figura 4.33 Variația ponderilor sinaptice a primului neuron în timpul antrenării.....	113
90. Figura 4.34 Variația ponderilor sinaptice a celui de-al doilea neuron în timpul antrenării.....	113
91. Figura 4.35 Curba de variație a erorii rețelei neuronale feed-forward, realizate cu scop de comparație în Matlab, în timpul procesului de antrenare.....	114
92. Figura 4.36 Eroarea medie pătratică a primului neuron al RNP implementate hardware.....	114
93. Figura 4.37 Eroarea medie pătratică a celui de-al doilea neuron al RNP implementate hardware.....	114
94. Figura 5.1 Structura de bază a rețelei implementate.....	119
95. Figura 5.2 Conceptul de codificare a variabilelor de intrare în impulsuri decalate temporal.....	120
96. Figura 5.3 Parametrii principali al circuitelor Xilinx FPGA din familia Spartan 3.....	120
97. Figura 5.4 Modul BRAM, 1024 cuvinte de 4 biți.....	120
98. Figura 5.5 Funcțiile triunghiulare ale domeniilor receptive utilizate la implementare.....	121
99. Figura 5.6 Valorile funcțiilor triunghiulare ale domeniilor receptive cu care s-au inițializat memoriile BRAM.....	121
100. Figura 5.7 Principiul de funcționare al modelului neuronal implementat.....	123
101. Figura 5.8 Structura logică a rețelei neuronale implementate.....	124
102. Figura 5.10 Schema de conectare a modulelor principale, în vederea punerii în funcțiune a rețelei neuronale.....	125

103.Figura 5.9 Schema bloc a structurii rețelei neuronale implementate.....	125
104.Figura 5.11 Schema de conexiune a unui neuron de intrare.....	126
105.Figura 5.12 Modulul de control funcțional.....	126
106.Figura 5.13 Stimul de intrare pentru simularea modulului BRAM.....	127
107.Figura 5.14 Rezultat de simulare, valori ale funcțiilor domeniilor receptive locale (semnalul <b>dataout</b> ) pentru valori crescătoare ale variabilei de intrare (semnalul <b>xval</b> ).....	127
108.Figura 5.15 Stimul pentru testarea codării intrării în impulsuri decalate temporal.....	127
109.Figura 5.16 Rezultat de simulare al codării intrării în impulsuri. De exemplu, pentru valoarea de intrare 25, se vor activa neuronii de intrare 0 și 1, cu decalaje de 5 respectiv 12 pași (cicluri de tact), iar ceilalți neuroni emit impuls doar la sfârșitul ciclului de funcționare (după 16 pași), generat de un modul dedicat pentru acest scop.....	127
110.Figura 5.17 Modulul sinapsă.....	128
111.Figura 5.18 Modulul de some.....	128
112.Figura 5.19 Modulul de comunicație cu calculatorul.....	128
113.Figura 5.20 Evoluția funcțiilor integrate în circuitele FPGA.....	129
114.Figura 5.21 Diagrama bloc a microcontrolerului soft-core Xilinx PicoBlaze încorporat în RNP implementate.....	130
115.Figura 5.22 Schema bloc a rețelei pseudoRBF-spiking implementată.....	130
116.Figura 5.23 Spațiul intrărilor pentru problema propusă, cu exemple de centre de focus și vecinătățile acestora care vor fi învățate de RNP din FPGA.....	131
117.Figura 5.24 Metoda de codificare a valorilor de intrare în impulsuri decalate temporal.....	131
118.Figura 5.25 Funcțiile triunghiulare ale domeniilor receptive a neuronilor de intrare, cu scalarea decalajelor temporale.....	131
119.Figura 5.26 Modulul BlockRAM utilizat la stocarea decalajelor temporale.....	132
120.Figura 5.27 Schema de conexiune unui bloc BRAM al unui neuron de intrare.....	132
121.Figura 5.28 Stimul de intrare pentru simularea funcționării unui modul de codificare a intrărilor RNP.....	133
122.Figura 5.29 Rezultatul simulării unui modul de codificare a valorilor de intrare a RNP.....	133
123.Figura 5.30 Structura blocului neuronilor de intrare.....	133
124.Figura 5.31 Schema de conexiune blocului neuronilor de intrare.....	134
125.Figura 5.32 Stimul de intrare pentru blocul de neuroni de intrare.....	134
126.Figura 5.33 Rezultatul simulării blocului de neuroni de intrare ce codifică una din variabilele de intrare a RNP.....	134
127.Figura 5.34 Parametrii modelului neuronal.....	135
128.Figura 5.35 Intrările și ieșirile modulului sinapsă.....	135
129.Figura 5.36 Notății în descrierea funcționării somei.....	136
130.Figura 5.37 Divizarea unui pas temporal în procesarea RNP implementate hardware.....	136
131.Figura 5.38 Diagrama de execuție a programului în limbaj de asamblare rulat de procesoarele PicoBlaze care implementează modulele de somă.....	137
132.Figura 5.39 Dinamismul somei cu inhibare lateral.....	138
133.Figura 5.40 Dinamismul somei fără inhibare laterală.....	138
134.Figura 5.41 Spațiul valorilor de intrare utilizate pentru antrenarea RNP implementate în FPGA.....	139
135.Figura 5.42 Transformata FFT a semnalului analog din Figura 5.43.....	139
136.Figura 5.43 Semnalul analog prelucrat.....	139
137.Figura 5.44 Valori de antrenare concentrare în jurul punctelor de focus prescrise, și altele aleator alese ca zgomot.....	139
138.Figura 5.45 Diagrama bloc a sistemului testbench realizat.....	140
139.Figura 5.46 Diagrama de execuție a programului în limbaj de asamblare rulat de procesorul PicoBlaze care implementează Unitatea de Monitorizare și Control.....	141
140.Figura 5.47 Reprezentarea valorilor bazei de date Fisher IRIS cu ajutorul proiecției Sammon (vezi subcap.5.6.3.7) 4D→2D.....	142
141.Figura 5.48 Reprezentarea valorilor bazei de date Fisher IRIS cu ajutorul proiecției Sammon (vezi subcap.5.6.3.7) 4D→2D.....	142
142.Figura 5.49 Schema bloc a RNP implementate pentru clasificarea setului de date Fisher IRIS.....	143
143.Figura 5.50 Funcțiile triunghiulare asociate neuronilor de intrare și decalajele temporale.....	143
144.Figura 5.51 Domeniile receptive ale neuronilor de intrare.....	143
145.Figura 5.52 Măsurători în timpul antrenării RNP (albastru-variația PM celor trei some, roșu valoarea axonală a acestora).....	144
146.Figura 5.53 Variația ponderilor în timpul antrenării (primele 15).....	144
147.Figura 5.54 Diagrama bloc a sistemului testbench realizat.....	145
148.Figura 5.55 Sistemul de dezvoltare Digilent Nexys2 (XC3S1200E).....	145

149.Figura 5.56 Structura logică a procesorului Xilinx MicroBlaze.....	146
150.Figura 5.57 Harta memoriei procesorului MicroBlaze.....	146
151.Figura 5.58 Forma caracterelor învățate.....	147
152.Figura 5.59 Diagrama bloc a sistemului încorporat care implementează RNP pentru recunoașterea de caractere.....	147
153.Figura 5.60 Structura RNP pentru recunoaștere de caractere, implementată cu procesor MicroBlaze încorporat în circuit FPGA.....	148
154.Figura 5.61 Modelul neuronal pulsativ utilizat în aplicația de recunoaștere caractere.....	148
155.Figura 5.62 Principiul de funcționare a modelului neuronal utilizat în aplicația de recunoaștere caractere... ..	149
156.Figura 5.63 Structura setului de antrenare a stratului ascuns (a se nota, că primele 6 valori sunt pentru neuronii de coloană iar următoarele 8 pentru cele de linie).....	150
157.Figura 5.64 Structura setului de antrenare a stratului de ieșire.....	150
158.Figura 5.65 Diagrama de componente a modului neuronului pulsativ implementat software.....	152
159.Figura 5.66 Diagrama cazurilor de utilizare (use-case) a RNP implementate pe procesorul MicroBlaze.....	152
160.Figura 5.67 Diagrama de secvență a RNP implementate pe procesorul MicroBlaze.....	152
161.Figura 5.68 Variația erorii în timpul antrenării (pentru caracterele 2,3,4,5,6,7).....	153
162.Figura 5.69 Variația potențialelor de membrană a 6 neuroni în timpul primului ciclu de antrenare (al caracterului 0).....	153
163.Figura 5.70 Variația potențialelor de membrană în timpul ciclului 32 de antrenare (al caracterului 0).....	154
164.Figura 5.71 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 0) , care este învățat de neuronul 1.....	154
165.Figura 5.72 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 1), care este învățat de neuronul 2.....	155
166.Figura 5.73 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 2), care este învățat de neuronul 3.....	155
167.Figura 5.74 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 3), care este învățat de neuronul 4.....	155
168.Figura 5.75 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 4), care este învățat de neuronul 5.....	156
169.Figura 5.76 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 5), care este învățat de neuronul 6.....	156
170.Figura 5.77 testarea RNP cu caracterul 0 perturbat de zgomot.....	156
171.Figura 5.78 Figura 5.79 testarea RNP cu caracterul 0 perturbat de zgomot accentuat.....	157
172.Figura 5.80 Structura hardware a sistemului implementat.....	157
173.Figura 5.81 Generarea impulsului decalat din valoarea pe 4 biți citită din BRAM.....	158
174.Figura 5.82 Simbolul bloc al perifericului de codificare a intrărilor.....	158
175.Figura 5.83 Codificarea valorii de intrare în trei impulsuri decalate temporal în decursul unui pas temporal de 15 cicluri de tact.....	158
176.Figura 5.84 Funcțiile triunghiulare ale domeniilor receptive utilizate în cazul RNP fără strat ascuns.....	159
177.Figura 5.85 Funcțiile triunghiulare ale domeniilor receptive utilizate în cazul RNP cu strat ascuns.....	159
178.Figura 5.86 Prima versiune a RNP implementate pentru aplicația benchmark de clasificare a setului de date WDBC.....	160
179.Figura 5.87 Versiunea cu strat ascuns a RNP implementate pentru aplicația benchmark de clasificare a setului de date WDBC.....	160
180.Figura 5.88 Reprezentarea grafică prin proiecție Sammon 9D→2D a clasificării prescise a bazei de date WBCD cu 699 mostre.....	161
181.Figura 5.89 Reprezentarea grafică prin proiecție Sammon 9D → 2D a clasificării prescise a bazei de date WBCD cu 683 mostre.....	162
182.Figura 5.90 Rezultatul clasificării bazei de date WBCD cu RNP fără strat ascuns, implementată hardware, reprezentat grafic prin proiecție Sammon 9D → 2D.....	162
183.Figura 5.91 Stadiul clasificării într-un ciclu de antrenare intermediar (Varianta 2 a RNP-FPGA).....	163
184.Figura 5.92 Stadiul clasificării în ciclul de antrenare 250 (Varianta 2 a RNP-FPGA).....	163
185.Figura 5.93 Pașii proiectării sistemelor încorporate în mediul Xilinx EDK.....	165
186.Figura 5.94 Fazele dezvoltării comune hardware-software.....	165
187.Figura 5.95 Comparatie a rezultatelor clasificării setului de date Fisher IRIS cu diferite metode similare ...	167
188.Figura 5.96 Comparatie a rezultatelor clasificării setului de date WBCD cu diferite metode similare.....	167



# 1. INTRODUCERE

## 1.1. Inteligență artificială și tehnici de machine learning

Cercetătorii au urmărit de multă vreme dezvoltarea conceptului de mașini instruibile și inteligente. La început s-a considerat că, deoarece orice funcție se poate calcula cu un dispozitiv arbitrar compatibil cu o mașină Turing, și inteligența trebuie să fie calculabilă cu aceleași dispozitive. În urma acestui fapt a luat naștere în anii 1960 domeniul de cercetare al Inteligenței Artificiale (IA). Scopul IA a fost de a produce programe inteligente, fără a se afla cum reușește acest lucru creierul biologic, considerându-se, că numai comportamentul unui sistem este important și nu procesul în sine care generează acest comportament. În prima fază rezultatele au fost promițătoare, deși performanțele au fost slabe. S-a afirmat însă, că această performanță slabă se va îmbunătăți drastic odată cu creșterea exponențială a puterii de calcul prevăzută de legea lui Moore. În anii 1980 a devenit evident, că această creștere de performanță foarte așteptată nu se realiza, chiar după 20 de ani de creștere exponențială a puterii de calcul. În multe cazuri, la probleme ce par simplu de rezolvat unui subiect uman, soluțiile cu IA întâmpinau greutăți multiple. Exemple semnificative erau recunoașterea vorbirii, robotică, procesarea imaginilor, prelucrarea semnalelor biologice, controlul sistemelor nestaționare și/sau neliniare. Pentru a putea rezolva o astfel de problemă, în general impunem unui sistem de procesare două cerințe importante, comportament neliniar și abilitatea de învăța din exemple.

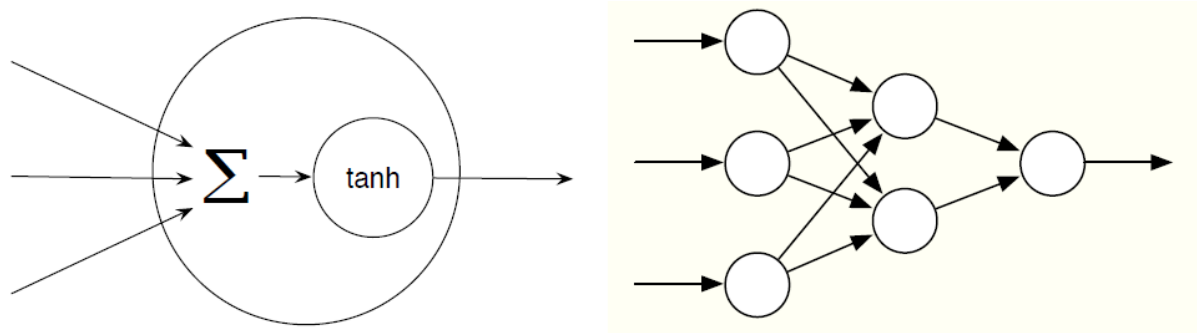


Figura 1.1 Primele modele de neuroni și rețele neuronale  
(stânga – neuron analog cu sumator a intrărilor și o funcție neliniară,  
dreapta – rețea neuronală feed-forward cu trei straturi, putând reprezenta numai funcții independente de timp)

Cercetătorii au început să se orienteze după paradigme și arhitecturi de calcul alternative pentru a crea sisteme inteligente, în mare măsură inspirate de unicul sistem cunoscut a fi cu adevărat inteligent: creierul biologic. Astfel a luat ființă domeniul de cercetare al Rețelelor Neuronale Artificiale (RNA), care de fapt a fost abordat și în anii 1960, dar a pierdut rapid interes datorită rezultatelor negative în non-separabilitate utilizând neuroni individuali. Acestea sunt modele foarte simplificate ale structurilor prezente în creierul natural, în cele mai multe cazuri compunându-se din numai câteva straturi de noduri simple de calcul interconectate (astfel fiind foarte departe de realitatea biologică). Utilizând aceste structuri (vezi Figura 1.1), însă, a devenit posibil să se *încea* un anumit comportament în loc de a fi programat. Acest lucru a fost foarte important, deoarece multe probleme complexe nu sunt definite clar de un model, cum ar fi ecuații analitice sau chiar un algoritm, dar pot fi descrise de un set de măsurători zgomotoase ale sistemului împreună cu răspunsurile predefinite.

Modelele de Rețele Neuronale Artificiale sunt ecuații neliniare simple, ce au fost parțial inspirate de modul în care cercetătorii anilor 1940 presupuneau că funcționează neuronii biologici. există multe variante, dar ecuația de bază este:

$$a_j = \tanh\left(\sum_i w_{ij} a_i\right) \quad \text{Ec. 1}$$

unde ieșirea  $a_j$  a neuronului  $j$  este calculată de o funcție neliniară (cum este funcția tangenta hiperbolică) a sumei ponderate a activărilor a neuronilor sau intrărilor la care este conectat. A fost arătat de

(Cybenko, 1989) că orice rețea cu două straturi creată din astfel elemente cu funcții neliniare simple, sunt capabile a aproxima orice funcție continuă cu o precizie arbitrară.

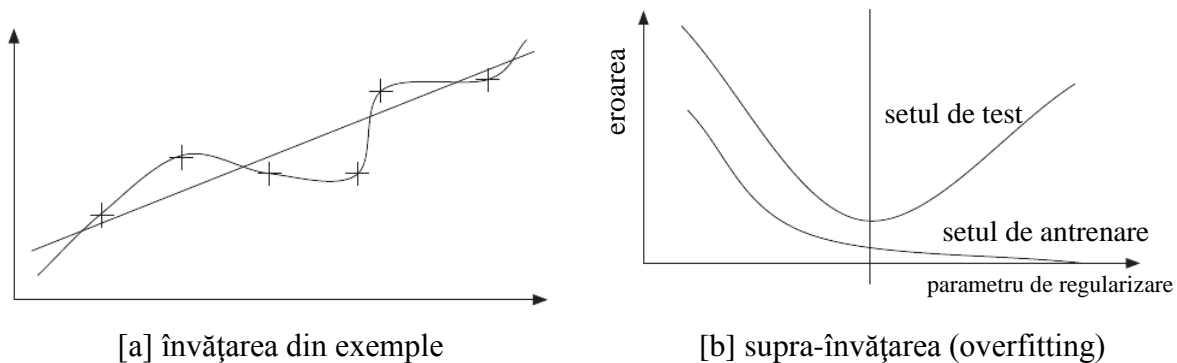


Figura 1.2 Reprezentarea schematică a conceptului de învățare în rețele neuronale artificiale

Ideea învățării este bazată pe faptul, că este posibil să se antreneze un sistem să răspundă corespunzător pentru un set de date de antrenare iar apoi demonstreze același lucru (cu o eroare neglijabilă) și pentru un set de date de test, neîntâlnite apriori. Această proprietate se numește generalizare și este un aspect crucial al realizării sistemelor inteligente. Ca urmare a introducerii acestei proprietăți, putem învăța funcții complexe fără a cunoaște explicit funcția respectivă. Învățarea supervizată este realizată prin specificarea ieșirii impuse pentru fiecare set de date de intrare, ajustând parametrii funcției (proces numit și antrenare) în așa fel, încât o metrică de eroare calculată pe setul de date de antrenare să se minimizeze. Un lucru important de remarcat este capacitatea de control: în cazul în care antrenăm o funcție cu mulți parametri utilizând un set de date de antrenare limitat este posibil să învățăm și zgomotul prezent în acest set de date. Această problemă, numită și supra-învățare (overfitting) apare, când numărul gradelor de libertate (a parametrilor) funcției ce este învățată este prea mare pentru complexitatea funcției care a generat setul de antrenare. Rezolvarea cea mai simplă constă în a limita numărul (efectiv) parametrilor funcției până la o valoare pentru care rezultă o generalizare optimă (acest lucru se poate realiza prin stoparea antrenării mai devreme, modificare arhitecturii, etc.) pe un set de date separat, numit de obicei setul de validare.

Mai recent, cercetarea ce a pornit ca rețele neuronale a fost fertilizată de domenii conexe, cum ar fi statistică și control, astfel obținându-se metode de învățare avansate, bazate pe principii matematice puternice. În anii 1990 s-a pus în evidență faptul, că deși RNA sunt bine înțelese ele nu pot spune multe despre modul de funcționare al creierului natural. Această realizare a dus la o intensificare a interesului cercetătorilor în a găsi modele mai plauzibile din punct de vedere biologic, care să ia în considerare natura pulsativă a neuronilor corticali. Astfel, s-a studiat capacitatea de învățare a acestor modele mai detaliate, numite Rețele Neuronale Pulsative (RNP) (spiking neural networks), unde informația este transmisă prin șiruri de impulsuri. Primele astfel de rezultate (Maas, 1997) (Maas & Natschläger, 1997) au arătat, că unele dintre aceste modele neuronale complexe sunt capabile de performanțe de calcul mai mari decât modelele neuronale clasice.

S-a demonstrat (Delorme & Thorpe, 2001), de asemenea, că se pot implementa cu RNP aplicații ingineresti impresionante, cu un necesar de putere de calcul foarte redus. Aici însă, s-au utilizat modele foarte mult simplificate pentru a reprezenta proprietățile specifice ale RNP pe arhitecturi clasice bazate pe unități centrale de prelucrare (CPU). Pe de altă parte, există și o direcție de cercetare mai apropiată de neuro-biologie, unde există o nevoie acută de tehnici inteligente pentru a procesa multitudinea de semnale măsurate de la neuroni reali, biologici cu metode in-vitro și in-vivo. În acest caz, scopul este de a decoda informațiile stocate în aceste trenuri de impulsuri și de a înțelege ce procesează acei neuroni. Cercetarea metodelor de codare/decodare a informațiilor în impulsuri este importantă din acest punct de vedere.

Unul dintre scopurile generale al acestei teze, este de a arăta, că aceste modele neuronale mai complexe pot fi utilizate în aplicații ingineresti, obținându-se performanțe de vârf și că aceste modele sunt foarte pliabile unor implementări pe suporturi hardware digitale, depășind chiar modelele clasice de RNA. În continuarea acestei introduceri se va face o descriere mai detaliată a RNP, cu o privire generală asupra neajunsurilor acestora din punctul de vedere ingineresc, respectiv se va deschide tema implementărilor hardware.

## 1.2. Rețele neuronale artificiale

Structura creierului uman se compune dintr-o rețea complexă de neuroni care se interconectează între ei în mod paralel prin intermediul unor axoni, respectiv dendrite. Schimbul de informație între diferiții neuroni se realizează prin intermediul așa-numitelor sinapse. Funcționarea creierului uman se caracterizează prin fenomene electro-chimice relativ simple. S-a pus problema explicării modului în care aceste rețele neuronale complexe, dar constituite din elemente de procesare relativ simple, pot să prezinte performanțe notabile. Aceste idei au condus la crearea diferitelor modele matematice ale neuronilor “artificiali”.

Prin studierea rețelelor neuronale artificiale putem încerca înțelegerea operațiilor esențiale ce au loc în rețeaua de neuroni dens interconectată a sistemului nervos central natural. Primul model neuronal propus de (McCulloch & Pitts, 1943) a fost bazat pe neuroni binari simplificați, care implementează o funcție de prag simplă, conform căreia un neuron este sau activ sau inactiv, stări calculate pe baza sumei ponderate a stărilor neuronilor la care acesta este conectat. În acest scop conexiunile dintre neuroni sunt direcționate (de la neuronul  $i$  la neuronul  $j$ ) și au atașate câte o pondere. Legăturile sunt deci niște conexiuni care transportă semnale excitatorii sau inhibitorii de la un element de procesare la altul. Dacă suma ponderată a neuronilor  $i$  conectați la neuronul  $j$  depășește o valoare de prag, atunci starea neuronului  $j$  este activ altfel inactiv. Remarcabil, rețele construite din astfel de elemente de calcul simple, pot implementa o gamă largă de funcții matematice, corelând stările de intrare cu stările de ieșire și variind ponderile conexiunilor pe baza unor algoritmi specifici, astfel rețelele de neuroni artificiali învățând funcțiile respective.

## 1.3. Fundamente biologice

### 1.3.1. Neuronul pulsativ ideal

Un neuron natural tipic se poate împărți în trei componente funcțional diferite: dendrite, soma (corpul celulei) și axonul. Generalizând, se poate afirma, că dendritele sunt intrările, care achiziționează semnale de ieșire de la alte celule pe care le transmite către soma neuronului (Figura 1.3). Soma este unitatea centrală de prelucrare, care realizează un pas important de prelucrare neliniară: dacă ieșirea totală calculată de acesta depășește o valoare de prag, se generează un semnal de ieșire – un impuls de potențial (Gerstner & Kistler, 2002). Acest semnal este transmis spre receptorii altor celule de către axon, dispozitivul de ieșire a neuronului.

Conexiunea dintre doi neuroni este denumită sinapsă. Să presupunem, că un neuron transmite un semnal printr-o sinapsă. Ne vom referi la acest neuron ca neuron pre-sinaptic, iar la cel care va recepționa semnalul ca neuron post-sinaptic. În cazul unei rețele neuronale naturale, ca cea a unui creier uman, un singur neuron pre-sinaptic poate fi conectat până la  $10^4$  neuroni post-sinaptici. Deseori o sinapsă se poate afla în imediata vecinătate a axonului dar se poate afla chiar și la câțiva centimetri depărtare.

### 1.3.2. Impulsuri de activare

Impulsurile de activare sunt impulsuri electrice de scurtă durată, dar care se pot sesiza plasând un electrod sensibil foarte aproape de soma sau axonul neuronului natural. Amplitudinea caracteristică a acestor impulsurile de activare este de aproximativ 100 mV, iar durata lor este de 1-2 ms. Forma impulsului nu se modifică, pe durata propagării acestuia de-a lungul axonului. Un șir de astfel de impulsuri stereotipice, emise

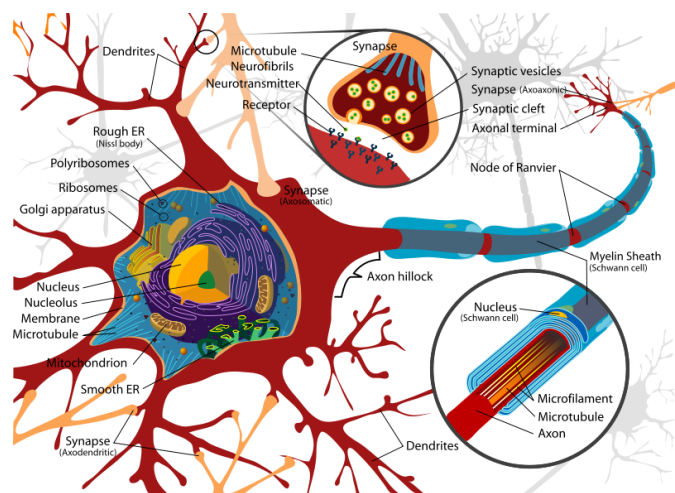


Figura 1.3 Structura unui neuron natural

de un neuron, se pot repeta la intervale de timp regulate sau neregulate. Deoarece toate impulsurile emise de către un neuron sunt asemănătoare, forma acestor impulsuri nu este purtătoare de informație.

Important este numărul și temporizarea acestor impulsuri. Impulsurile sau potențialele de activare sunt elementul de bază în transmisia de date dintre neuroni. Potențialele de activare se pot delimita cu acuratețe în cadrul unui astfel de șir de impulsuri. Distanța minimă dintre două activări (două impulsuri) determină perioada absolută de reprimare a neuronului. Perioada relativă de reprimare este cea care urmează celei anterior menționate, în care este dificil, dar nu imposibil excitarea neuronului până la activare.

### 1.3.3. Sinapsele

Punctul în care axonul unui neuron pre-sinaptic se conectează la dendrita unui neuron post-sinaptic se numește sinapsă. În cazul natural, sinapsele cele mai des întâlnite sunt cele chimice.

La o astfel de sinapsă, terminalul (capătul) axonului se apropie extrem de mult de neuronul post-sinaptic, (la o distanță de

ordinul micrometrilor) spațiul rămas între membranele celulelor pre- și post-sinaptice numindu-se interspațiu sinaptic. (Figura 1.4). Dacă un potențial de activare ajunge la un interspațiu sinaptic, determină startul unui proces complex de prelucrare chimică, care are ca efect eliberarea de către terminalul pre-sinaptic în interspațiu a unei substanțe speciale, numită neuro-transmițător (Trappenberg, 2002). Moleculile acestei substanțe ajungând la senzorii membranei post-sinaptice deschid canale specifice, prin care anumiți ioni, aflați în fluidul intracelular din interspațiul sinaptic, pot pătrunde în interiorul celulei. Ca urmare a acestei migrări de ioni, potențialul membranei post-sinaptice se modifică, astfel, semnalul chimic se transformă din nou în semnal electric.

### 1.3.4. Dinamismul neuronal

Efectul unei activări pre-sinaptice asupra neuronului post-sinaptic se poate măsura cu ajutorul unui electrod intracelular, care măsoară diferența de potențial  $u(t)$  dintre interiorul celulei și mediul înconjurător.

Această diferență de potențial este numită potențialul membranei. În lipsa unui impuls de intrare, neuronul este în stare latentă, având un potențial de membrană specific acestei stări. După sosirea unui impuls de activare la intrare potențialul membranei se modifică ca apoi să revină la potențialul de repaus (vezi Figura 1.5).

Figura 1.5 Neuronul post-sinaptic nr.  $i$  recepționează impulsuri de activare de la neuronii pre-sinaptici  $j=1,2$ . Ambii neuroni emit câte un impuls care modifică potențialul membranei la neuronul post-sinaptic, modificarea se poate măsura cu un electrod ca diferența de potențial  $u_i(t) - u_{rest}$ . A. Variația potențialului de membrană post-sinaptic ( $t - t_1(f)$ ) ca răspuns la impulsul de activare excitator recepționat de la neuronul pre-sinaptic  $j = 1$ . B. Sosirea unui al doilea impuls de activare de la neuronul pre-sinaptic  $j = 2$ , imediat după cel de la neuronul  $j = 1$ , determină un al doilea salt în potențialul membranei post-sinaptice, care se adună cu primul. C. Dacă  $u_i(t)$  atinge pragul, neuronul post-sinaptic emite, la rândul lui, un impuls de activare.

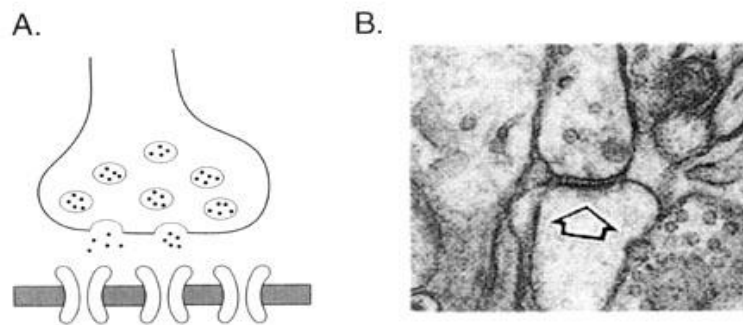
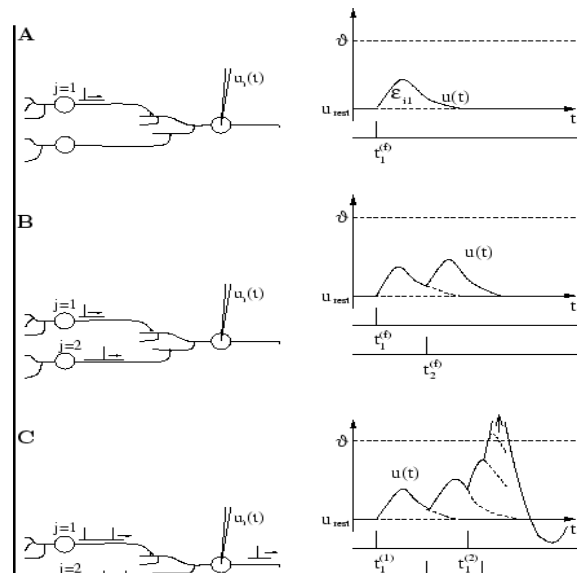


Figura 1.4 A. Schema unei sinapse chimice; B. Fotografia electro-microscopică a unei sinapse (Trappenberg, 2002)



Dacă această variație are gradientul pozitiv, atunci este vorba de o sinapsă excitatoare, iar în cazul negativ de una inhibitorie. În stare latentă, sau de repaus, membrana celulei este deja polarizată negativ, având o valoare de aproximativ  $-65$  mV. Efectul unei intrări sosite pe o sinapsă excitatoare, scăderea stării de polarizare negativă a membranei se numește depolarizare, iar efectul invers, care duce la accentuarea polarizării negative se numește hiperpolarizare.

### 1.3.5. Pragul de activare sau aprindere și potențialul de activare

Să presupunem că avem doi neuroni pre-sinaptici  $j = 1, 2$ , fiecare transmițând câte un impuls către neuronul post-sinaptic  $i$ . Neuronul  $j = 1$  se aprinde la momentele de timp  $t_1^{(1)}, t_1^{(2)}, \dots$ , iar neuronul  $j = 2$  la  $t_2^{(1)}, t_2^{(2)}, \dots$ . Fiecare impuls de intrare generează separat câte un potențial post-sinaptic  $\varepsilon_{i1}$  și  $\varepsilon_{i2}$ . Dacă avem numai câteva impulsuri de intrare, variația potențialului post-sinaptic este aproape liniară, ea putând fi echivalată cu suma salturilor de potențial post-sinaptic (Figura 1.5B):

$$u_i(t) = \sum_j \sum_f \varepsilon_{ij}(t - t_j^{(f)}) + u_{\text{repaus}} \quad \text{Ec. 2}$$

Altfel, dacă avem multe impulsuri de intrare sosite într-un interval de timp relativ scurt, atunci variația menționată își pierde liniaritatea. De îndată ce potențialul post-sinaptic atinge o valoare critică, forma acestui semnal va lua cu o totul altă formă decât cea a sumării impulsurilor intrate, și ea devenind un impuls abrupt, cu o amplitudine de circa  $100$  mV. Acest potențial de aprindere se va propaga de-a lungul axonului neuronului  $i$  către sinapsele altor neuroni din rețea. După evenimentul aprinderii, potențialul membranei nu se ca întoarce imediat la potențialul de repaus, ci va trece printr-o fază de hiperpolarizare. Amplitudinea impulsurilor de intrare este în domeniul milivolților. Nivelul de prag al potențialului membranei, care trebuie atins pentru a declanșa starea de aprindere, este situat cu cca.  $20 \div 30$  mV deasupra potențialului de repaus. După cum se poate observa în Figura 1.5C, patru impulsuri de intrare nu sunt suficiente, pentru a atinge potențialul de prag și implicit pentru a aprinde neuronul. Acest lucru necesită sosirea a cel puțin  $20-50$  de impulsuri pre-sinaptice într-un interval de timp scurt.

## 1.4. Sisteme de control neuronal

În prezent sistemelor de reglare li se cer să posede performanțe dinamice ridicate împreună cu proprietăți de robustețe și utilizare în sisteme din ce în ce mai complexe, cu caracteristici dinamice neliniare. Capacitatea de modelare cu metode matematice pierde teren față de o modelare și reglare inteligentă. În ultimul timp s-a manifestat un larg interes în folosirea modelelor bazate pe structuri biologice și algoritmi de învățare pentru modelare și reglare adaptivă. Aplicațiile de reglare inteligentă au nevoie de algoritmi care sunt capabili (Eldredge & Hutchings, 1994) să:

- opereze într-un domeniu insuficient de bine definit și variabil în timp
- se adapteze la schimbări în dinamica sistemului ca și în efectele mediului
- învețe informații semnificative în mod stabil
- considere câteva restricții pe dinamica sistemului
- să opereze autonom în mediu haotic cu intervenție minimă.

Căutarea unui algoritm care să ne furnizeze o metodă universală pentru toate problemele de reglare inteligentă nu este un țel realist. În ingineria reglării, algoritmi de învățare au importanță mare. În general, algoritmi de adaptare sunt bazați pe sisteme liniare și modele liniare, numărul parametrilor care trebuie aleși determinând flexibilitatea schemelor de adaptare. Rețele neuronale furnizează și metode cu care acești algoritmi pot fi adaptați pentru sisteme neliniare. Majoritatea algoritmilor de învățare supervizată sunt bazați pe metoda gradientului dar recent au apărut strategii de adaptare bazate pe concepte de stabilitate (James-Roxby & Blodget, 2000) (Zhu & Gunther, 1999) (Pérez-Urbe & Sanchez, 1996). În general arhitecturile de modelare și reglare sunt independente de rețea.

În ultima perioadă se manifestă un mare interes în utilizarea rețelelor neuronale în proiectarea sistemelor de control. În proiectele legate de controlul proceselor industriale foarte des se utilizează rețele neuronale atât pentru identificarea sistemului, cât și comanda optimală și adaptivă a lor. În aceste procese de control se utilizează diferite topologii de rețele neuronale cât și diferite procedee de antrenare a acestor rețele.

Sistemele de control neuronale se pot clasifica așa cum urmează:

- control neuronal *indirect*
- control neuronal *direct*

Față de sistemele de control neuronale directe cele cu control indirect nu sunt în legătură directă cu procesul controlat. În general o rețea neuronală este utilizată să modeleze parametrii procesului sau a regulatorului. Dacă parametrii regulatorului sunt furnizați de o rețea neuronală atunci avem de a face cu un regulator auto-tuner care adaptează parametrii regulatorului în conformitate cu comportamentul procesului.

Astfel modelele de control neuronal *indirecte* se pot clasifica:

- control pe bază de model proces utilizând rețele neuronale artificiale
- control bazat pe model invers folosind rețele neuronale artificiale
- control bazat pe autoacordare a parametrilor realizat de rețea neuronală

În schemele de control *directe* cu rețele neuronale, rețeaua este folosită ca regulator, și trimite semnalul de comandă direct către proces. În funcție de conceptul de proiectare a acestor sisteme amintim:

- modelare regulator tip neuronal
- structură de control fără model
- structură de control bazată pe de model
- structură de control robust bazată pe model

În toate clasele descrise mai sus rolul rețelei neuronale este rezolvarea unei probleme de optimizare bazată pe o funcție de criteriu  $J(w)$  care depinde de valorile ponderilor rețelei.

O astfel de funcție de criteriu se poate modela astfel ca:

$$NN : \min_w \{J(w)\} \quad \text{Ec. 3}$$

unde  $NN$ : reprezintă faptul că problema de optimizare implică o rețea neuronală.

O rețea neuronală poate fi privită ca o cutie neagră pentru modelul procesului cum se poate vedea Figura 1.6.

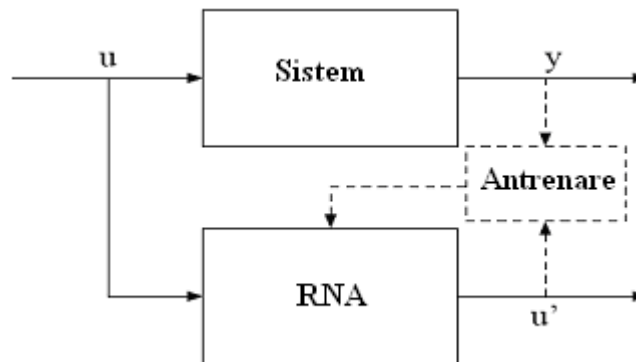


Figura 1.6 Rețea neuronală privită ca o cutie neagră

### 1.4.1. Control neuronal indirect

*Control bazat de modelul procesului, utilizând RNA*

Cea mai cunoscută aplicație a sistemelor de control cu rețele neuronale este de a utiliza rețeaua neuronală ca un model de proces intrare-ieșire. Rețeaua neuronală este antrenată printr-un algoritm de învățare supervizată. În mod obișnuit se utilizează ca și model de sistem modelul NARMAX sau NARX (**N**-neural **AR**-autoregressive **MA**-moving average **X**-eroarea estimată) care exprimă funcții de transfer discrete neliniare.

O alternativă ar fi identificarea sistemului controlat cu o rețea neuronală dinamică Indiferent de structura modelului și a strategiei de control, modelul cu control neuronal în acest caz poate fi exprimat conform relației din Ec. 4.

$$NN : \min_w \{F(y_p - y_n(w, \dots))\} \quad E \quad c. 4$$

unde  $y_p$  reprezintă ieșirea procesului,  $y_n$  ieșirea rețelei neuronale,  $w$  ponderile rețelei neuronale.  $F\{\}$  este o funcțională care măsoară performanța procesului de optimizare. De obicei  $F$  este o integrală din suma erorii de predicție dintre ieșirea procesului  $y_p$  respectiv ieșirea rețelei neuronale  $y_n$ . Sunt măsurate, pe o perioadă finită, intrările respectiv ieșirile procesului  $\{u_p, y_p\}$ , și acest set de date măsurate sunt utilizate pentru antrenarea rețelei neuronale.

Deseori, se utilizează pentru funcția de criteriu:

$$NN : \min_w \left\{ \sum_t |y_p(t) - y_n(t)|^2 \right\}; \quad y_n(t) = N(w, \dots) \quad Ec. 5$$

Odată ce modelul este specificat, se poate utiliza în implementarea regulatorului pe bază de model.

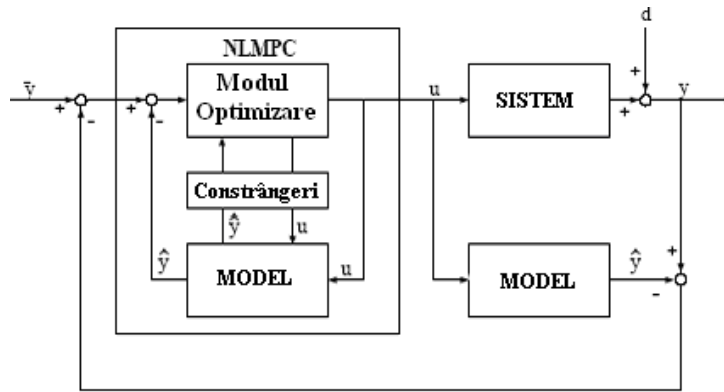


Figura 1.7 O rețea neuronală poate fi încorporată într-o schemă de control predictiv bazat pe model

În stadiul de implementare, de la început modelul de rețea neuronală, nu poate fi utilizat ci trebuie încorporat într-o schemă de control pe bază de model. În procesele industriale chimice rețeaua neuronală este utilizată într-o schemă de control neliniar cu model predictiv, obținând un regulator NL MPC. Blocul **Model** din Figura 1.7 reprezintă rețeaua neuronală care este utilizată pentru identificarea procesului controlat și paralel se desfășoară și predicția stărilor respectiv elaborarea semnalului de comandă pentru eșantionul curent.

De fapt controlul cu model predictiv este tot o problemă de optimizare și poate fi exprimată prin relația:

$$NN : \min_u \{F\{y^* - y_n(u, \dots)\}\} \quad Ec. 6$$

- $y^*$ - ieșirea așteptată a procesului cu buclă închisă
- $y_n$ -ieșirea estimată (cu modelul de rețea neuronală)
- $F$  funcția de criteriu pentru evaluarea performanțelor sistemului cu buclă închisă

Problema de optimizare poate fi descrisă prin relația din Ec. 7.

$$\min_u \left\{ \sum_t |y^*(t) - y_n(t) - d(t)|^2 \right\}; \quad y_n(t) = N(u, \dots) \quad c. 7$$

unde:  $y^*$  – reprezintă valorile pentru traiectoria impusă,  
 $d(t)$ - valoarea perturbației estimate.

### 1.4.2. Model de rețea neuronală parțială sau parametrică

În multe cazuri este posibil să avem informații concrete despre model cum ar fi structura modelului, sau fenomene fizice particulare bine cunoscute care descriu funcționarea unei părți a modelului. În acest caz nu este necesar să realizăm în totalitate modelul ca și o cutie neagră. Dacă



structura cu modelul procesului ne stă la dispoziție, valorile pentru parametrii asociați modelului pot fi determinate cu o rețea neuronală. Acești parametri pot fi constante de timp, factori de amplificare, factori de întârziere sau parametri fizici cum ar fi rata de difuziune, coeficienți de material etc. Când structura de model nu este cunoscută apriori, rețeaua neuronală poate fi antrenată pentru a selecta elemente din structura de model dintr-un set de modele predefinite. Aceste elemente pot fi compuse într-o structură utilizabilă. În alte cazuri unde structura modelului este parțial cunoscută, rețeaua neuronală se poate folosi împreună cu un model parțial și astfel procesul poate fi mai bine modelat (Figura 1.8).

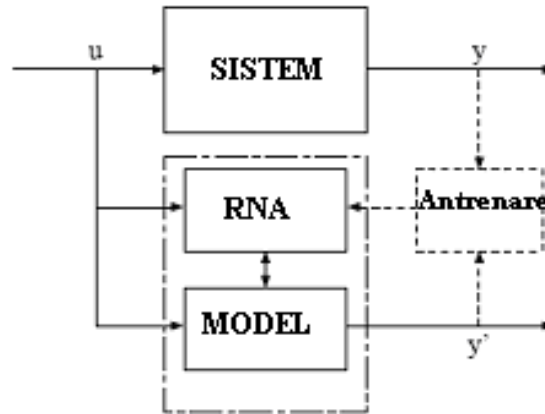


Figura 1.8 Rețeaua neuronală ca estimator de parametri

Problema de modelare cu rețea neuronală parțială sau parametrică se poate formula astfel:

$$NN : \min_w \{F\{y_p - y_m(\theta, \dots)\}\}; \quad \theta = N(w, \dots) \quad \text{Ec. 8}$$

- unde  $y_m$  este o ieșire estimată de model

$\theta$  - parametrii procesului, informații structurale despre model, sau alte elemente pentru completarea modelului.

Din punct de vedere al controlului bazat pe model, această aproximare în esență este identică cu modelul de rețea neuronală full black box cu excepția faptului că rețeaua neuronală nu descrie direct comportamentul procesului.

#### Modelul sistem invers cu RNA

În acest caz rețeaua neuronală este antrenată pentru a obține un model invers al procesului Figura 1.9. Ieșirea procesului este intrarea rețelei neuronale, iar intrarea procesului corespunde cu ieșirea rețelei neuronale. În general problema de optimizare se poate formula:

$$NN : \min_w \{F(u_{p-1} - u_n(w, \dots))\} \quad \text{Ec. 9}$$

Unde  $u_{p-1}$  - semnalul de intrare a procesului

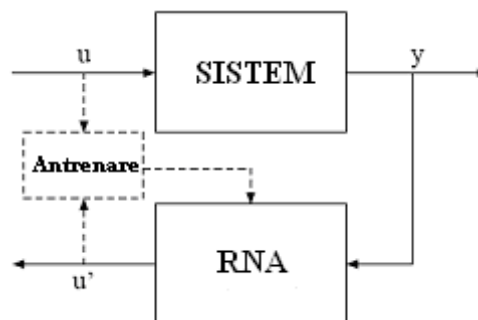


Figura 1.9 Control cu rețea neuronală bazată pe modelul invers al procesului

Pentru o valoare impusă ieșirii sistemului  $y$  și semnalul de comandă se obține din ieșirea rețelei neuronale: ( $u=u'$ )



$$u^* = N(y^*, \dots) \tag{Ec. 10}$$

Modelul invers există numai când procesul se comportă ca un sistem minimal, altfel această structură este inaplicabilă.

În principiu, o rețea neuronală poate învăța cu anumite restricții dinamica inversă a unui sistem. Modelul invers este utilizat într-un mod similar unui controller MPC. În practică, și în special la un model dinamic discret, modelul invers nu este capabil să învețe dinamica inversă pe tot spațiul de comandă inversă.

Modelul invers al sistemelor poate să nu existe dacă sistemul este cu timp mort. Se poate obține un model invers a procesului utilizabil cu restricții impuse pe intrare respectiv parametri de model invers.

*Model cu rețea neuronală cu autoacordare.*

În acest caz rețeaua neuronală poate fi utilizată pentru estimarea parametrilor de acordare a regulatorului a cărui structură este cunoscută apriori. Deseori estimatorul parametrilor de acordare a regulatorului este denumit autoacordare.

Problema de optimizare în acest caz poate fi formulat astfel:

$$NN : \min_w \{F\{\eta^* - \eta_n(w, \dots)\}\} \tag{Ec. 11}$$

$\eta^*$  - parametrii regulatorului

$\eta_n$  - parametrii estimați de rețeaua neuronală

Parametri  $\eta$  nu se pot unic determina din caracteristica procesului. Acești parametri depind de caracteristicile sistemului închis dorit.

De obicei, parametrii regulatorului se obțin din optimizarea funcției de criteriu corespunzătoare sistemului în buclă închisă

$$NN : \min_{\eta} \{F\{y^* - y_{p/m}(u, \dots)\}\}; \quad u = C(\eta, \dots) \tag{Ec. 12}$$

unde C este un regulator cu structură cunoscută  $y_{p/m}$  reprezintă un model sau un proces ce poate fi inclus în sistemul cu buclă închisă pentru obținerea sistemului de reglare cu un C dorit. Avantajul acestei metode este că rețeaua neuronală din buclă se poate antrena în condiții de simulare. Pentru autoacordare în buclă deschisă este suficientă o simulare în buclă deschisă, altfel necesită o simulare în buclă închisă, Figura 1.10.

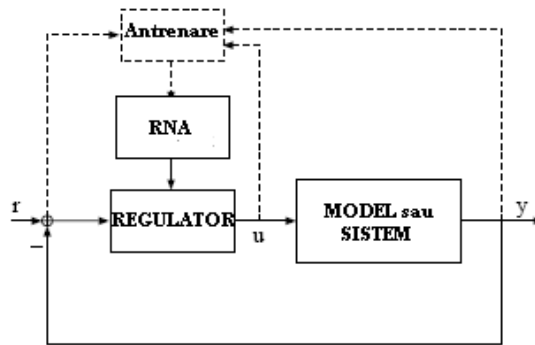


Figura 1.10 Regulator în buclă închisă cu autoacordare cu RNA

Antrenarea trebuie condusă astfel încât să acopere tot spațiul modelului de proces. În astfel de sisteme regulatoarele PID sunt utilizate pe scară largă.

### 1.4.3. Control neuronal Direct

#### Modelarea regulatorului

Dintre cele patru modele de control neuronal direct, cel mai simplu din punct de vedere al proiectării regulatorului este utilizarea unei rețele neuronale pentru modelarea unui regulator existent, Figura 1.11.

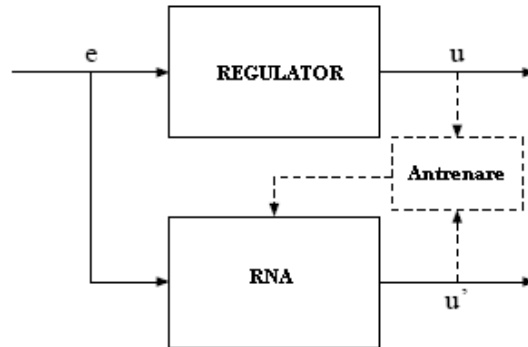


Figura 1.11 Modelare regulator

Setul de antrenare pentru regulatorul neuronal ales este format din perechi de intrare ieșire. Intrarea reprezintă intrarea regulatorului iar ieșirea reprezintă valorile impuse pentru intrările corespunzătoare sistemului.

În mod similar, această structură de control neuronal se poate formula astfel (Ec. 13)

$$NN : \min_w \{F\{u_c - u_n(w, \dots)\}\} \quad E \quad c. 13$$

unde  $u_c$  este ieșirea regulatorului (curent) existent  $C^*$ . De obicei regulatorul curent existent  $C^*$  poate fi și un operator uman, sau se poate obține din:

$$NN : \min_C \{F\{y^* - y_{p/m}(u, \dots)\}\}; \quad u = C(\dots) \quad Ec. 14$$

Ca și modelul procesului, regulatorul este un sistem dinamic și deseori conține elemente integratoare și/sau elemente de diferențiere.

Dacă pentru modelarea regulatorului utilizat se folosește o rețea feed-forward, la intrarea rețelei atât pentru modelare cât și în timpul exploatarei trebuie specificate și intrări dinamice (modul de variație a semnalelor de intrare).

O altă aproximare este similară cu dezvoltarea unui model de proces tip ARX (auto regresiv cu intrări exogene) cu excepția că ieșirea respectiv intrarea sistemului sunt schimbate cu eroarea feed-back și ieșirea regulatorului.

#### 1.4.3.1. Neuro-control fără modelul sistemului

În anumite procese în absența unui regulator, cercetători s-au inspirat din modul cum un operator uman învață să controleze un proces, fără să aibă cunoștințe amănunțite despre dinamica procesului.

Astfel s-au încercat proiectarea unor regulatoare implementate cu rețele neuronale care prin adaptare și învățare pot rezolva probleme dificile de reglare și control în absența unui model al procesului.

În general, acest tip de *control neuronal fără model* se poate exprima astfel:

$$NN : \min_w \{F\{y^* - y_p(u, \dots)\}\}; \quad u = N(w, \dots) \quad Ec. 15$$

- unde  $y_p$  este ieșirea sistemului.

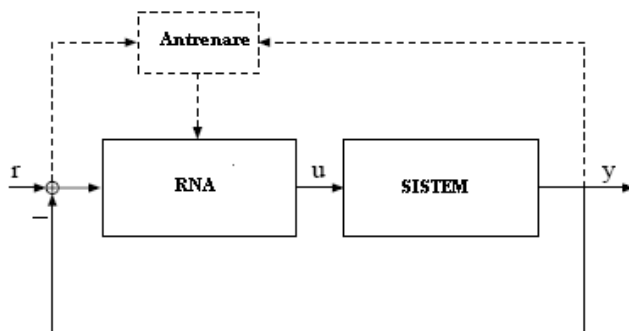


Figura 1.12 Neuro-control fără modelul sistemului

Ideea de mai sus se poate descrie ca o metodă de control adaptivă directă când nici modelul nu se cunoaște dinainte și nici în timpul proiectării sistemului de control.

De multe ori la astfel de modelări de sisteme de control foarte des se utilizează metoda de învățare prin reinformare și din această cauză metoda se numește la fel.

### 1.4.3.2. Control neuronal bazat pe model

Metoda se bazează pe antrenarea rețelei neuronale pe baza utilizării unui model al sistemului în cauză. Astfel putem evita defectarea sistemului în timpul proiectării, dar metoda este utilizabilă numai în cazul în care putem avea un model foarte exact al procesului de condus.

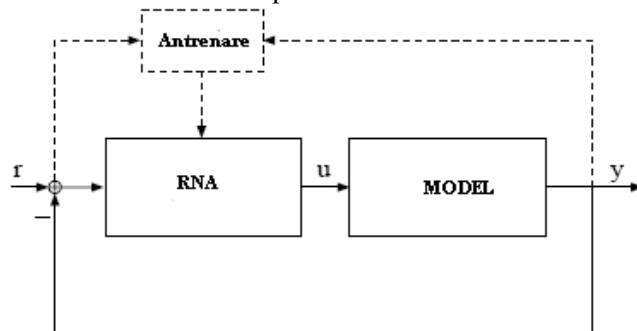


Figura 1.13 Control neuronal bazat pe model

Regula de control în general se poate formula astfel:

$$NN : \min_w \{F\{y^* - y_m(u, \dots)\}\}; \quad u = N(w, \dots) \quad \text{Ec. 16}$$

Dacă nu există model pentru proces, atunci putem utiliza un model al procesului anterior identificat cu o metodă corespunzătoare.

După obținerea modelului pentru sistem, modelul poate fi utilizat în modelarea structurii de control. După o antrenare și testare riguroasă în mod de simulare a rețelei neuronale, controlerul neuronal poate fi instalat în sistemul de control real.

În realitate, această structură de regulator neuronal bazat pe modelul sistemului nu numai că s-a demonstrat în multe studii dar de asemenea utilizarea lui a condus la utilizări eficiente. Aceste abordări pot fi folosite atât în control on-line cât și în control off-line.

### 1.4.3.3. Regulator neuronal robust pe bază de model

Structurile de control cu rețele neuronale discutate în paragrafele anterioare au câteva deficiențe:

- rețeaua neuronală trebuie antrenată pentru fiecare aplicație nouă
- Re-antrenarea rețelei neuronale este necesară și la mici modificări ale criteriilor de reglare sau dacă regulatorul se aplică la un alt proces similar cu procesul anterior reglat.

Pentru a preveni aceste deficiențe în mod natural se introduce conceptul de metodă robustă de proiectare a reguletoarelor neuronale.

În această nouă viziune fiecare proces și modul de reglare a lui nu este considerat ca un proces nominal ci un proces făcând parte dintr-o familie de procese. În majoritatea cazurilor aceste familii sunt

specificate pe baza unor modele de zgomot și de variabilitate a parametrilor proceselor din familie.

Controlul neuronal robust pe bază de model se poate formula astfel:

$$NN : \min_w \{F\{y^* - y_{m_i}(u, \dots)\}\}, \quad u = N(w, \dots) \quad \text{Ec. 17}$$

unde  $m_i$  reprezintă membrul  $i$  din familia de model  $M_i$

În caz ideal robustețea se referă atât la familia de modele cât și la o anumită variabilitate în modele reale. Stabilitatea sistemelor robuste se referă la stabilitatea tuturor sistemelor care fac parte din aceeași familie.

Totodată trebuie să precizăm că metodele robuste de reglare adeseori nu ajung la performanța reglatoarelor prezentate anterior.

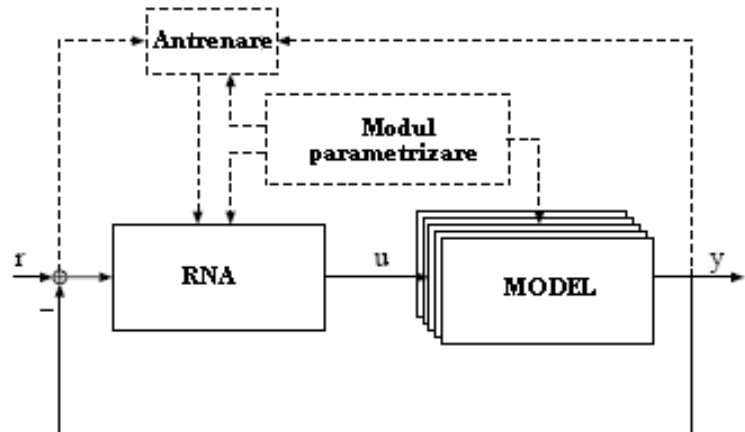


Figura 1.14 Regulator neuronal robust pe bază de model

## 1.5. Organizarea tezei

În acest subcapitol se trece în revistă câte o introducere succintă asupra conținutului celor șase capitole, prezentând astfel structura acestei teze de doctorat.

Capitolul 1 constituie o descriere a fundamentelor biologice ale inteligenței artificiale și a realizării rețelelor neuronale artificiale bazate pe modele neuromorfe. De asemenea se prezintă stadiul actual al cercetărilor în domeniul aplicării rețelelor neuronale în sisteme de control.

În Capitolul 2 este expusă teoria modelării și simulării rețelelor neuronale neuromorfe. În cadrul acestei sinteze sunt detaliate metodele de învățare specifice acestui tip de rețea neuronală și se introduce și un model propriu, adaptat pentru implementare în hardware digital. În partea a doua a capitolului se găsește o prezentare a strategiilor de simulare a rețelelor neuronale pulsative precum și a celor mai importante sisteme de simulare realizate de cercetătorii domeniului și prezente în literatura de specialitate. În finalul capitolului se prezintă rezultatele proprii în simularea acestor sisteme neuronale.

Capitolul 3 cuprinde modalitățile de realizare a rețelelor neuronale artificiale pe sisteme digitale reconfigurabile FPGA și sunt prezentate tehnicile de implementare a logicii de reconfigurare compilată în timp și reconfigurare în timpul utilizării. În acest capitolului se discută maparea diferitelor algoritmi de rețele neuronale pe sisteme FPGA și sunt prezentate mai multe sisteme de rețele neuronale hardware la baza cărora stau sisteme FPGA.

Capitolul 4 prezintă motivații pro și contra în alegerea suportului electronic pentru implementarea total paralelă a rețelelor neuronale artificiale. Sunt descrise o serie de astfel de implementări hardware, citând lucrările științifice care stau la baza acestora. Partea de contribuții personale ale acestui capitol descrie aplicațiile realizate în domeniul implementării totale paralele a rețelelor neuronale neuromorfe, detaliind funcționarea subsamblelor acestora, dar evidențiind și rezultatele experimentale obținute în comparație cu realizări software.

În Capitolul 5 sunt prezentate realizări proprii de rețele neuronale pulsative implementate în hardware cu ajutorul incorporării unor procesoare de tip soft-core în circuitele FPGA utilizate. Sunt prezentate patru astfel de aplicații, câte una aplicativă și una de test (benchmark) utilizând procesoare încorporate pe 8 respectiv pe 32 de biți. Printre acestea se numără o aplicație de detectare automată a componentelor de frecvență a unui semnal analogic, un sistem de recunoaștere de caractere, respectiv testele benchmark de clasificare a seturilor de date Fischer IRIS și a bazei de date Wisconsin Breast Cancer.

Capitolul 6 conține concluziile generale și lista contribuțiilor originale realizate de autor în cadrul acestei lucrări.

## 2. MODELAREA ȘI SIMULAREA REȚELELOR NEURONALE NEUROMORFE

### 2.1. Alegerea modelului optim pentru neuronii pulsativi corticali

#### 2.1.1. Introducere

În ultimii câțiva ani, în comunitatea care se ocupă de rețelele neuronale artificiale, s-a observat o trecere accentuată spre rețele neuronale pulsatorii. Motivate de descoperiri biologice, mai multe studii consideră rețelele neuronale legate de pulsuri cu temporizarea pulsațiilor ca o componentă esențială a studierii și înțelegerii procesării informației din creier.

Rețelele neuronale pulsative dispun de mai multe proprietăți interesante și cu potențial bun pentru a fi utilizate la aplicații ingineresti.

**Comportament temporal intrinsec:** În RNP timpul este utilizat pentru a codifica informația, iar calculele sunt bazate de asemenea pe temporizări. Acest fapt face ca astfel de rețele neuronale să fie în mod intrinsec capabile de a fluxuri de valori de intrare variabile în timp. În rețelele neuronale clasice, acest lucru este realizat prin introducerea de întârzieri sau prin adăugarea conexiunilor recurente. În RNP este posibil să se efectueze procesarea temporală chiar și cu arhitecturi de tip feed-forward.

**Comunicație utilizând impulsuri:** În RNA neuronii comunică prin valori continue, reale. Din punct de vedere al implementărilor, acest lucru poate deveni o problemă importantă, deoarece aceste valori reale trebuie multiplicat cu valori de ponderare, transmise unor alte elemente aflate la distanțe arbitrare (în termeni de implementare hardware) și necesită o capacitate semnificativă de memorie. Pe de altă parte, neuronii din RNP comunică prin impulsuri, adică prin evenimente „rare” reprezentabile pe un singur bit. Astfel, lățimea de bandă necesară pentru comunicația dintre neuroni este drastic redusă.

**Capacitate de calcul mai ridicată decât în cazul RNA:** S-a demonstrat în (Maas & Natschläger, 1997) că neuronii bazați pe modele cu impulsuri sunt mai eficienți din punct de vedere al puterii de calcul decât neuronii bazați pe modele analogice. O RNP perturbată de zgomot poate aproxima orice RNA cu o precizie arbitrară cu un timp de execuție influențat numai de numărul de straturi din RNA și cu un număr de neuroni ce este numai dublul numărul neuronilor din RNA. Pe de altă parte, este necesară o rețea neuronală analogică de dimensiuni mult mai mari (mii de neuroni) pentru a simula unele funcții ce pot fi implementate utilizând un singur neuron pulsativ. Această asimetrie duce la concluzia, că RNP sunt strict mai eficiente ca putere de calcul decât rețelele neuronale analogice, bazate pe modele clasice.

**Inspirat de biologie, neuromorfic:** Rețelele neuronale pulsative sunt mai plauzibile din punct de vedere biologic decât alte rețele neuronale. Utilizarea sistemelor de calcul cu arhitectură inspirată de mecanisme biologice este însă un concept foarte dezbătut. Este adevărat, că pentru multe domenii ingineresti a propune ca unic scop al unei lucrări plauzibilitatea biologică nu este o abordare validă. Un exemplu frecvent este cel al unei aeronave Boeing 747. Păsările naturii sunt foarte diferite ca design, dar construind numai o pasăre mare nu obținem un Boeing 747. Modul de operare al unei astfel de mașini este drastic diferit de cea a unei păsări, deci este nevoie de idei de construcție complet noi. În cazul sistemelor inteligente însă, o mulțime de tehnici de calcul, teorii matematice și statistice, respectiv 50 de ani de cercetări în aceste domenii nu au reușit să se apropie de inteligența ce o putem percepe în interacțiunea dintre ființele umane. Despre nici un singur sistem artificial curent nu se poate spune că este inteligent, comparat cu un sistem biologic. Mai mult decât atât, nici măcar noțiunea de sistem inteligent nu este clar definită. Turing a spus, că un sistem este inteligent dacă nu-l putem distinge de o ființă umană. De aceea, putem afirma cu încredere, că abordarea neuromorfică în domeniul inteligenței artificiale este cu siguranță întemeiată. Acest lucru nu înseamnă, că fiecare detaliu complicat trebuie înțeles și simulat, dar cel puțin arhitecturile de procesare importante pot

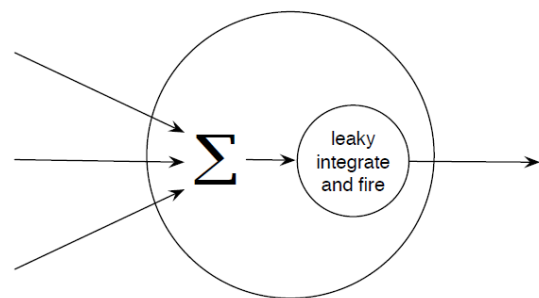


Figura 2.1 Structura de bază a unui neuron pulsativ

Figura 2.1 Structura de bază a unui neuron pulsativ

fi surse de inspirație folositoare. Studiarea arhitecturii simplificate ce se află la baza modului de operare al creierului natural poate ajuta de asemenea, la înțelegerea funcționării acestuia. Aceste concluzii nu ar putea fi extrapolate pentru sisteme nervoase biologice dacă am utiliza modele ce sunt foarte departe de funcționarea propriu-zisă a creierului, cum ar fi de exemplu mașinile Turing.

De curând, au apărut câteva aplicații importante ce utilizează RNP, cum ar fi recunoașterea de fețe prezentată în (Delorme & Thorpe, 2001) care prezintă o viteză de execuție mare, propagând un singur val de impulsuri printr-o RNP sau o aplicație de procesare de imagini ca o metodă de citire de pe buze în (Booij, 2004), precum și un sistem detector de epilepsie bazat pe semnale EEG în (Ghosh-Dastidar & Adeli, 2007) sau sistemul iterativ de găsire a rădăcinilor cu o RNP prezentat în (Iannella & Kindermann, 2005).

### 2.1.2. Dispute curente în domeniul RNP

Rețelele neuronale pulsative sunt modele neuronale noi utilizabile ca instrumente ingineresti foarte promițătoare și atrag o atenție tot mai mare în literatura de specialitate. Totuși, mai există multe teme deschise care necesită a fi adresate pentru ca aceste RNP să ajungă să fie instrumente bine dezvoltate de la stadiul de studiu neuro-biologic în care sunt de fapt la momentul de față.

Primul obstacol de trecut în utilizarea acestor rețele neuronale speciale este faptul că fluxul de informații spre și dinspre acestea se realizează prin impulsuri. Întrucât cele mai multe cantități din lumea reală sunt de natură analogică și variază continuu, reprezentarea acestora cu impulsuri fără amplitudine nu este un lucru simplu. Studiul neuronilor biologici a rezultat în apariția mai multor ipoteze pentru codarea valorilor reale în trenuri de impulsuri. Acestea pot apărea ca o codare a ratei impulsurilor, care ia în considerare numai valoarea medie a timpului de apariție a unui impuls sau ca codarea exactă a momentelor temporale de apariție a acestor impulsuri. Coduri intermediare care utilizează filtre temporale sunt de asemenea utilizate. Totuși, nu este încă clar care sunt proprietățile acestor scheme de codare și dacă este posibilă o codare fără pierdere de informație cu aceste metode. Toate aceste întrebări sunt cruciale și trebuie clarificate înainte de a începe implementarea unor aplicații ingineresti cu neuroni pulsativi.

Când vorbim despre RNP există de fapt o supraabundență de diferite modele, de la cele simple de tipul integrează și activează prin neuroni cu diferite modele sinaptice până la modele comportamentale complete descrise prin ecuații diferențiale. Există câteva studii teoretice (Maass, 1996) care arată că unele dintre aceste modele sunt mai eficiente decât celelalte, dar aceste studii se limitează doar la un set redus de modele și scheme de codare. O explicație generală a capacității de calcul a diferitelor modele neuronale pulsative și a schemelor de codare conexe reprezintă și astăzi o temă deschisă cercetărilor actuale.

Pe de altă parte, schema de codare utilizată pentru a introduce datele într-o RNP și pentru a evalua ieșirile acesteia nu este întotdeauna aceeași cu cea utilizată pentru comunicarea dintre neuronii rețelei. Există chiar și posibilitatea de a utiliza concomitent mai multe scheme de codare pentru comunicațiile interne. Acest lucru face ca analiza funcționării interne a acestor rețele să fie dificilă, precum și traducerea acesteia în termeni ușor de asimilat de către om, cum este cazul modelelor neuronale, unde există legături puternice cu sisteme bazate pe reguli fuzzy (Mantas, Puche, & Mantas, 2006).

Antrenarea acestor rețele este o altă temă în mod notoriu recunoscută a fi dificilă, deoarece literatura de specialitate ne arată doar o singură metodă supervizată, limitată și aceasta la o singură schemă de codare, cu frecvente probleme de convergență. Deoarece sistemele naturale nu cunosc antrenarea supervizată, majoritatea metodelor derivate din studiile neuro-informatică sunt de natură nesupervizată. Cea mai importantă este cea numită Spike Time Dependent Plasticity (STDP) sau plasticitate dependentă de timpii de activare. În cazul acestei metode, eficacitatea unei conexiuni dintre doi neuroni este mărită, dacă activarea acestuia este corelată cu activarea predecesorului său altfel eficacitatea scade. Această învățare de tip Hebb a fost observată în multe zone ale creierului biologic. Un alt exemplu pentru astfel de metode de antrenare este plasticitatea homeostatică sau intrinsecă (Intrinsic Plasticity - IP). Această metodă nu este văzută ca o regulă de antrenare, ea vizând de fapt menținerea la o valoare constantă a ratei medii de activare a unui neuron, funcționând deci ca și un mecanism de optimizare energetică. S-a arătat, însă recent (Triesch, 2007), că această metodă simplă poate fi utilizată pentru antrenare și are un impact drastic asupra proprietăților de învățare a metodei STDP.

### 2.1.3. Caracteristici neuro-computaționale

În orice studiu legat de dinamica unei rețele neuronale pulsative, sunt două probleme specifice: 1. ce model descrie dinamica pulsației fiecărui neuron și 2. cum sunt conectați neuronii. O alegere greșită a modelului sau o conexiune greșită, poate duce la rezultate care nu au nimic în comun cu procesarea ce are loc în creier.

În continuare se elaborează prima problemă, comparând neuronii pulsativi. Se vor prezenta diferite modele de neuroni pulsativi precum și o clasificare după următoarele proprietăți: 1. numărul de caracteristici neuro-computaționale pe care le pot reproduce, și 2. eficiența lor de implementabilitate, de ex. numărul de operații în virgulă flotantă (adunare, înmulțire, etc.) necesare pentru simularea modelului pe o durată de 1ms. Rezultatele comparației sunt rezumate în Figura 2.2.

În Figura 2.2 se prezintă 20 dintre cele mai proeminente caracteristici biologice ale neuronilor pulsativi (Izhikevich, 2004). Țelul urmărit este ilustrarea bogăției și complexității comportamentului pulsativ al neuronilor individuali ca răspuns la impulsuri simple de curent continuu. Ce se întâmplă dacă doar zeci (dar dacă milioane) de astfel de neuroni sunt cuplați - depășește tot ceea ce noi putem înțelege. Utilizând câteva dintre modelele amintite în continuare, nu este dificil să simulăm mii de neuroni corticali în timp real cu o rezoluție de 1ms.

#### ❖ Pulsație tonică (Tonic Spiking)

Majoritatea neuronilor pot fi excitați, ceea ce înseamnă, că sunt pasivi (statici), dar pot să emită impulsuri când sunt stimulați. Pentru a testa această caracteristică, cercetătorii de neurofiziologie injectează pulsuri de curent continuu printr-un electrod legat la un neuron și înregistrează potențialul membranei. Curentul la intrare și răspunsul neuronal sunt înregistrați unul sub celălalt Figura 2.2(A).

Cât timp impulsul de intrare este activ, neuronul continuă să emită un tren de impulsuri. Acest comportament, numit pulsație tonică, poate fi observat la trei tipuri de neuroni corticali: neuronii excitatori cu pulsație regulată (RS), cu pulsații cu prag scăzut sau cu valoare de prag scăzută (LTS), și neuronii inhibitorii (Gerstner & Kistler, 2002) (Gibson, Belerlein, & Connors, 1999), cu pulsație rapidă (FS). O continuă emisie a acestor neuroni indică faptul că sunt supuși unor impulsuri de intrare persistente.

#### ❖ Pulsație fazică (Phasic Spiking)

Neuronul poate să emită un singur impuls la pornirea impulsului de intrare, ca în Figura 2.2(B), și rămâne pasiv după aceea. Acest tip de răspuns se numește pulsație fazică, și poate fi utilizată pentru detectarea începutului stimulării.

#### ❖ Fluxuri de impulsuri tonice (Tonic Bursting)

Anumiți neuroni, ca și neuronii de tip chattering din neocortexul pisicii, emit un rând de impulsuri când sunt stimulați, ca în Figura 2.2(C). Frecvența apariției acestor trenuri poate fi mare, de până la 50 Hz, și se crede că astfel de neuroni contribuie la oscilațiile de frecvență gamma a creierului.

#### ❖ Fluxuri de impulsuri fazice (Phasic Bursting)

Asemănător cu pulsațiile fazice, anumiți neuroni emit impulsuri fazice, ca în Figura 2.2(D). Astfel de neuroni raportează începutul unei stimulări prin emiterea unui tren de impulsuri. Există trei ipoteze importante despre rolul unor trenuri de impulsuri al creierului care ar fi: **1.** trenurile de impulsuri sunt necesare pentru a învinge defecțiunea transmisiei și reducerea zgomotului neuronal, **2.** trenurile de impulsuri pot transmite proeminențele semnalului de intrare, deoarece efectul unui tren de impulsuri este mai mare decât al unui singur impuls asupra neuronului post-sinaptic și **3.** fluxurile de impulsuri pot fi folosite ca un mod de comunicare selectivă între neuroni (Izhikevich, Desai, Walcott, & Hoppensteadt, 2003), unde frecvența dintre impulsuri să fie codificarea canalului de comunicare. Un model bun al rețelei neuronale corticale nu poate să neglijeze aceste „izbucniri” neuronale.

#### ❖ Model combinat (Trenuri de impulsuri urmate de pulsații - Bursting Then Spiking)

Neuronii excitatori intrinseci emit trenuri de impulsuri în neocortexul mamar, pot emite o activitate de pulsație mixtă, descrisă în Figura 2.2(E). Ei emit trenuri de impulsuri fazice în momentul apariției stimulului și trec la modul de pulsație tonică. Nu este clar ce fel de mod de calcul un astfel de neuron poate crea în plus față de detectarea momentului de început al stimulării și raportând lungimea perioadei stimulării.

### ❖ Adaptarea frecvenței pulsației

Cel mai răspândit tip de neuron excitator în neocortexul mamiferelor, numit RS - celulă cu pulsație regulată -, emite impulsuri tonice cu frecvență în descreștere, ca în Figura 2.2(F). Înseamnă că frecvența este relativ mare în momentul stimulării, și ulterior se adaptează. Neuronii inhibitori cu pulsațiile cu valoarea de prag scăzută (LTS), au de asemenea această caracteristică. Frecvența dintre impulsuri la aceste celule este o codificare pentru timpii de la apariția impulsului de stimulare.

### ❖ Excitabilitatea de clasa 1

Frecvența pulsațiilor tonice a neuronilor excitatori neo-corticali RS depinde de puterea impulsului de intrare, și poate să varieze între 2 Hz și 200 Hz, sau chiar o valoare mai mare. Abilitatea de a emite impulsuri de frecvență mică, când impulsul de intrare este slab (dar deasupra valorii de prag) se numește excitabilitate de clasa 1. Neuronii de acest tip pot să codifice puterea impulsului de intrare în rata lor de emisie, Figura 2.2(G).

### ❖ Excitabilitate de clasa 2

Anumiți neuroni nu pot să emită trenuri de impulsuri de frecvență mică. Ei sunt ori pasivi sau emit un tren de impulsuri la o valoare destul de mare de frecvență, de ex de 40 Hz, ca în Figura 2.2(H). Acești neuroni se numesc de excitabilitate clasa 2, rata lor de emisie este un aproximativ slab al puterii stimulării.

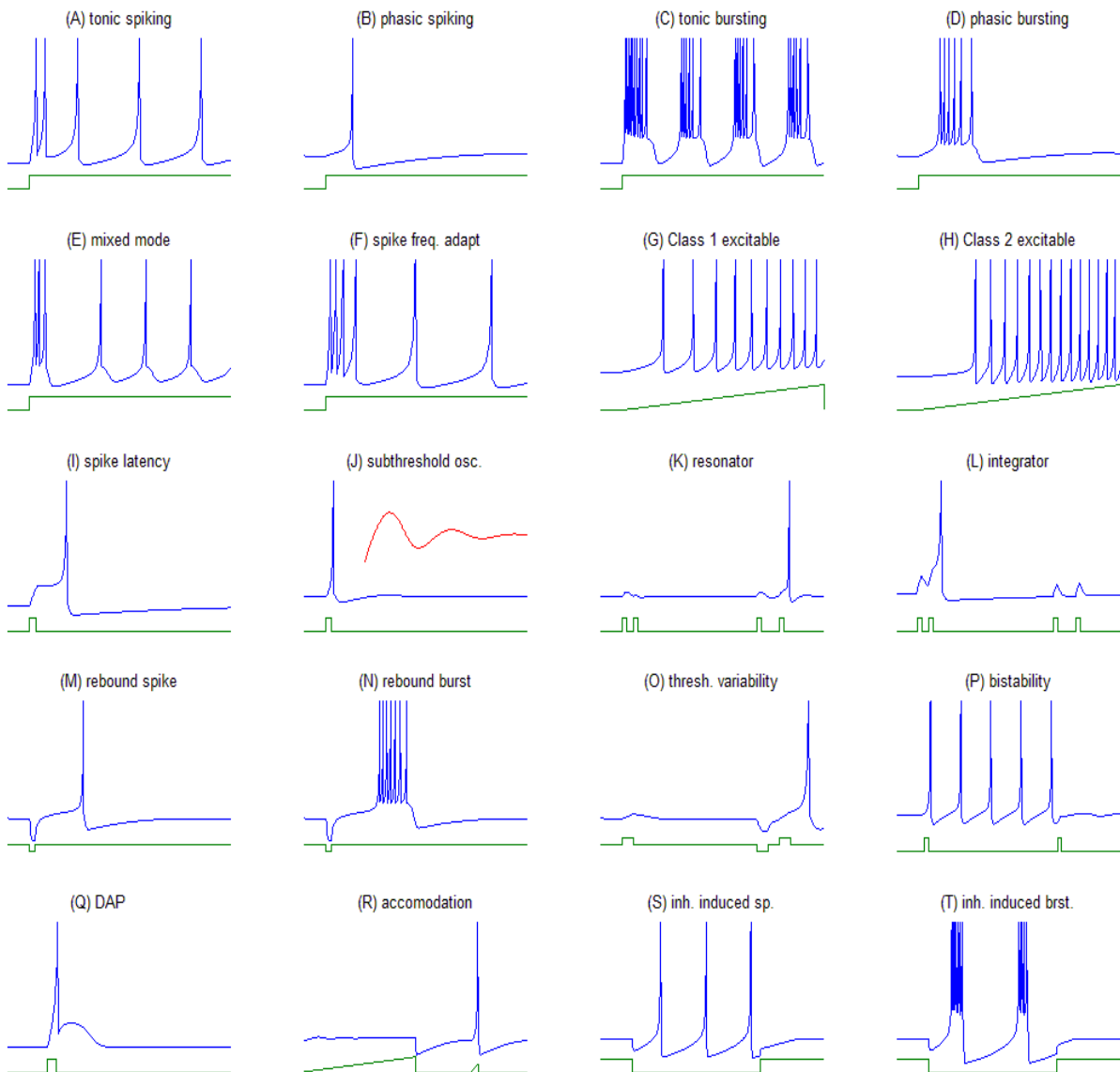


Figura 2.2 Tipuri de neuroni pulsativi



❖ **Latența impulsului (Spike Latency)**

Majoritatea neuronilor corticali emit impulsuri cu o întârziere care depinde de puterea semnalului de intrare. Pentru un impuls relativ slab, dar deasupra valorii de prag, întârzierea, numită și latența impulsului, poate fi destul de mare, după cum se vede în Figura 2.2(I). Celulele RS din neocortexul mamiferelor pot avea o latență de zeci de milisecunde. Astfel de latențe (întârzieri) pot asigura un mecanism de temporizare pentru a coda puterea impulsului de intrare.

❖ **Oscilații sub valoarea de prag**

Practic orice structură a creierului are neuroni capabili să arate comportament oscilatoriu, Figura 2.2(J). Frecvența unor astfel de oscilații joacă un rol important și neuronii de acest tip acționează ca niște filtre trece bandă, după cum va fi discutat mai târziu.

❖ **Preferință și rezonanță a frecvenței**

Datorită fenomenului de rezonanță, neuronii având oscilatori, pot răspunde selectiv la valorile de intrare, având conținut de frecvență similară (asemănătoare) cu frecvența oscilațiilor sub valoarea de prag. Astfel de neuroni pot implementa interacțiuni cu frecvențe modulate (FM) și multiplica semnalele. În Figura 2.2(K), se prezintă un astfel de neuron stimulat cu o pereche de impulsuri, având frecvențe diferite ale impulsurilor. Neuronul răspunde numai la una dintre perechi, la impulsul a cărui frecvență rezonază cu frecvența oscilațiilor sub nivelul de prag. Acești neuroni sunt numiți rezonatori.

❖ **Integrare și detectarea coincidenței**

Neuronii fără potențial oscilator acționează ca integratori. Ei preferă impulsuri de intrare de frecvență înaltă; cu cât este mai mare frecvența cu atât este mai posibilă activarea lor, vezi Figura 2.2(L). Aceștia pot fi folositori pentru detectarea impulsurilor coincidente sau aproape coincidente.

❖ **Impulsuri de revenire (Rebound Spike)**

Când un neuron primește și după aceea este eliberat de impulsul inhibitor, poate să emită un impuls post inhibitor sau de ricoșare, ca în Figura 2.2(M). Acest fenomen este în legătură cu excitația anodică de rupere în membrana excitabilă. Mulți neuroni pulsativi pot emite ca răspuns la intrări scurte, inhibitorii, astfel estompând diferența dintre excitație și inhibiție.

❖ **Trenuri de impulsuri de revenire**

Anumiți neuroni, inclusiv cei din celulele thalamo-corticale, pot emite trenuri de impulsuri post – inhibitorii, ca în Figura 2.2(N). Se crede că astfel de trenuri de impulsuri contribuie la oscilațiile somnului în sistemul thalamo-cortical.

❖ **Variabilitatea valorii de prag**

O concepție greșită în comunitatea cercetărilor rețelelor artificiale neuronale este că neuronii pulsatorii au un voltaj de valoare fixă a valorii de prag. Este binecunoscut, că neuronii biologici au o valoare de prag care depinde de activitatea apriorică a neuronului. În Figura 2.2(O), prima dată stimulăm neuronul cu un puls excitator mic de curent care produce 10 mV de-polarizație. Neuronul nu emite, prin urmare valoarea de intrare este mai mică decât valoarea de prag. După aceea, aplicăm un puls inhibitor scurt, și după aceea același puls de curent sub valoarea de prag. Neuronul emite de această dată, deoarece valoarea lui de prag a devenit mai mică din cauza impulsul inhibitor anterior. Prin urmare, aceeași depolarizare de 10mV poate fi sub valoarea de prag sau deasupra, depinzând de activitatea anterioară. În mod interesant, un următor puls excitator poate să ridice valoarea de prag și să facă neuronul mai puțin excitabil.

❖ **Bi-stabilitatea stării de pulsație și pauză**

Anumiți neuroni pot să arate două modele stabile de operare: pauză și pulsație tonică (sau chiar trenuri de impulsuri). Un puls excitator sau inhibitor, poate să schimbe de la un mod la altul, Figura 2.2(P), astfel creând o posibilitate interesantă pentru bi-stabilitate și memorie de scurtă durată. Trebuie să notăm că, pentru ca schimbarea de la pulsație tonică la modul de pauză să aibă loc, impulsurile de intrare trebuie să sosească într-o fază potrivită a oscilației, astfel punând accent pe importanța temporizării impulsurilor în astfel de procesare a informației.

❖ **Potențiale ulterioare depolarizante**

După emiterea unui impuls, potențialul membranei neuronului poate să arate o prelungită hiperpolarizare ulterioară ca în exemplul din Figura 2.2(B), I sau M, sau un prelungit potențial ulterior (DAP), ca în Figura 2.2(Q). Astfel de DAP-uri pot apărea din cauza influenței dendritice, din cauza unor curenți de prag mari, interiori, activați în timpul impulsului, sau din cauza unei interacțiuni dintre curenții sub valoarea de prag. În orice caz, un astfel de neuron are o perioadă refractară mai scurtă și devine super-excitabil.

❖ **Obișnuința (acomodarea)**

Neuronii sunt foarte sensibili la valorile de intrare scurte și coincidente, dar s-ar putea să nu emită semnal ca răspuns la o valoare de intrare puternică și încet crescătoare, după cum se vede în Figura 2.2(R). Curentul în creștere ușoară nu provoacă un impuls, pe când un curent mai mic dar cu creștere rapidă provoacă un impuls. În timpul unei rampe de creștere mică, curenții interni au timp suficient să se inactiveze iar curenții externi au timp suficient să se activeze. Deci neuronul se obișnuiește, devine mai puțin excitabil și nu poate genera un impuls.

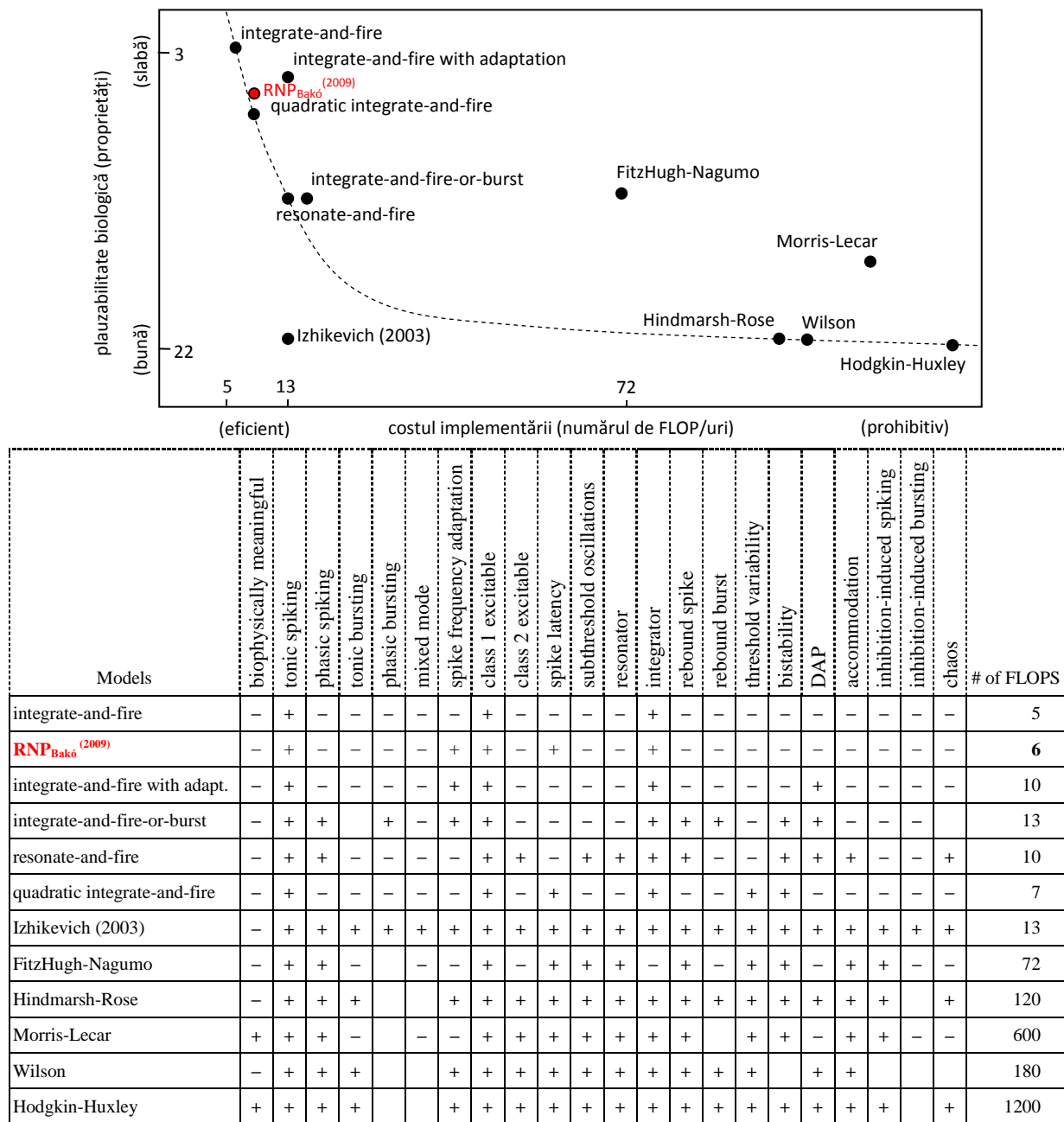


Figura 2.3 Comparație a diferitelor modele de neuroni pulsativi

### ❖ Impuls indus de inhibiție

O caracteristică interesantă a neuronilor thalamo-corticali este că ei sunt pasivi, latenți, când nu primesc impulsuri, dar emit dacă sunt hiper-polarizați de către un impuls inhibitor sau din cauza unui curent injectat, după cum se vede în Figura 2.2(S). Aceasta se întâmplă deoarece curentul injectat activează curentul  $h$  și dezactivează curentul de calciu  $T$ , ducând la pulsație tonică.

### ❖ Trenuri de impulsuri induse prin inhibiție

În loc de pulsație, un neuron thalamo-cortical poate să emită un tren de impulsuri tonice ca răspuns la o hiperpolarizare prelungită, ca în Figura 2.2(T). Se crede că astfel de trenuri de impulsuri au loc în timpul unor oscilații de undă în sistemul thalamo-cortical și joacă un rol important în ritmurile somnului.

Nici un model nu poate manifesta toate aceste 20 de caracteristici neuro-computaționale simultan, deoarece câteva sunt caracteristici mutual exclusive. De exemplu, un neuron nu poate fi un integrator și un rezonator în același timp. Totuși, există modele care cu ușurință pot fi reglate, să arate fiecare dintre aceste caracteristici. De exemplu, toate răspunsurile neurale din Figura 2.2 au fost obținute printr-un simplu model pulsativ având patru parametri ușor de reglat (Izhikevich E. , 2003).

## 2.2. Realizarea învățării Hebbiene competitive

Învățarea Hebbiană, dezvoltarea circuitelor neurale bazate pe activitatea corelată, are două mecanisme de bază. Cea mai bine cunoscută este modificarea sinaptică, bazată pe activitatea pe lângă ideile propuse de Hebb. La fel de important este mecanismul care forțează diferite sinapse să intre în competiție una cu cealaltă în așa fel, încât dacă anumite sinapse ale unui neuron post-sinaptic devin mai puternice, altele devin mai slabe.

De exemplu, regulile bazate pe corelația modificărilor sinaptice pot asigura un număr mare al aspectelor de dezvoltare în cortexul vizual, dar doar atunci când sunt combinate cu introducerea constrângerilor, ca să se asigure competitivitatea. Cu toate că modificările sinaptice Hebbiene au fost susținute de experimente cu algoritmi de excitare și inhibare pe termen lung, se cunoaște foarte puțin despre mecanismele care generează competiția dintre sinapse. Evidențele experimentale de la diferitele preparate, sugerează că atât semnalul cât și gradul de modificare sinaptică, care apare din repetatele împerecheri ale potențialelor acțiunilor pre- și post-sinaptice, depind de temporizarea lor relativă. Acțiunile potențiale pre-sinaptice care urmează impulsuri post-sinaptice, produc o slăbire îndelungată a sinapselor. Cele mai mari modificări în eficacitatea sinaptică apar, când diferența de timp dintre acțiunile potențiale pre- și post-sinaptice este foarte mică, și există o tranziție bruscă din valori în creștere spre valori în scădere, când această valoare tinde spre 0. Această modificare sinaptică se numește plasticitate dependentă de temporizare a impulsului (STDP). Modificarea sinaptică prin reguli STDP a fost studiată prin modele de recunoaștere temporală a modelului (temporal pattern recognition), învățare de secvență temporală de detectarea coincidenței (Gerstner, Kempter, van Hemmen, & Wagner, 1996). Distribuțiile conductanței sinaptice produsă prin STDP forțează neuronul post-sinaptic într-un regim de activare echilibrată, cu emisii neregulate de impulsuri, care este sensibil la temporizarea acțiunilor pre-sinaptice pe care le primește. O astfel de sensibilitate duce la competiție între valorile de intrare pentru controlul temporizării impulsurilor post-sinaptice. Aceasta permite STDP să întărească selectiv sinapsele cu întârzieri relativ mai mici sau cu corelări mutuale mai puternice, în timp ce sinapsele rămase devin mai slabe.

### 2.2.1. Excitație echilibrată

Pentru ca un neuron să funcționeze normal, trebuie să stabilească și să mențină un nivel potrivit de excitație în așa fel, încât să poată răspunde la impulsurile primite prin acțiuni potențiale în proporții rezonabile. Variabilitatea răspunsurilor asigură o constrângere pe impulsurile sinaptice de intrare asupra unui neuron. Răspunsurile primite de la un model care primește multe impulsuri independente pre-sinaptice pot fi mai puțin variabile decât răspunsurile observate în realitate (in-vivo). Corelările temporizării impulsurilor de intrare, ca de exemplu sincronizarea, pot contribui la creșterea variabilității. Totuși, mulți autori au observat că o variabilitate mare poate apărea și dacă impulsurile excitatorii de intrare asupra unui neuron sunt relativ echilibrate față de curenții sinaptici inhibitori și ai membranei. Condiția critică este că impulsul asupra neuronului ar trebui să fie suficient ca să ridice potențialul

membranei la un punct puțin sub sau deasupra valorii de prag pentru generarea potențialului de acțiune, astfel încât timpii impulsului să fie determinați prin fluctuație pozitivă la nivelul total de intrare. După cum se va arăta, STDP asigură un mecanism prin care acest echilibru poate fi stabilit și întreținut asupra unui șir de valori de intrare. Din aceasta rezultă o stare în care impulsurile pre-sinaptice pot controla temporizarea impulsurilor post-sinaptice și o competiție între sinapse poate să apară.

### 2.2.2. Discuții

Cu toate că plasticitatea sinaptică Hebbiană este un concept puternic, ea suferă din cauza apariției mai multor probleme. În primul rând, sinapsele sunt modificate de câte ori activitățile pre- sau post-sinaptice sunt corelate. Astfel de activități corelate pot să apară doar prin întâmplare, reflectând ineficient o relație causală care poate fi studiată. Pentru corectarea acestei probleme, modelele rețelilor neurale adesea utilizează o covarianță față de regulile de modificare sinaptice bazate pe corelare. Totuși, o astfel de regulă nu poate, în general, obține o competiție între sinapse. Aceasta creează o nouă problemă a modificărilor pur Hebbiene: nu este competitivă, deci trebuie să se aplice constrângeri asupra lor, pentru a obține rezultate interesante. STDP pare să rezolve ambele probleme menționate anterior. Coincidențele accidentale, non-cauzale slăbesc sinapsele, după cum am presupus, dacă integrala funcției modificării sinaptice este negativă. Competiția apare într-un mod nou, nu din cauza unui echilibru în mod artificial impus al factorilor ne-specifici de creștere sau descreștere a eficienței sinapselor, ci mai ales prin competiția pentru controlul temporizării impulsurilor post-sinaptice. Impulsurile de intrare, care în mod consistent sunt cele mai previzibile pentru răspunsuri post-sinaptice, devin cele mai puternice impulsuri pentru neuron. Cauzalitatea este elementul cheie în STDP. Așa cum a sugerat Hebb, sinapsele devin mai puternice doar dacă impulsurilor lor pre-sinaptice preced și au contribuit la activarea neuronului post-sinaptic. STDP duce în mod automat la o stare de echilibru neregulat de activare, la care timpii impulsurilor pre- și post-sinaptici sunt corelați.

Acest rezultat depinde de neliniaritatea procesului de generare a impulsului. Într-un model în care probabilitatea impulsurilor depinde de liniaritatea tensiunii de pe membrană, corelația dintre încărcarea pre- și post-sinaptică nu-și schimbă forma în eficacitatea sinaptică, așa cum se întâmplă în Figura 2.4.

Efectele neliniare, care fac corelațiile cazuale a valorilor de intrare-ieșire să crească relativ față de o relație causală ca o descreștere a întregii eficacități sinaptice, sunt cruciale pentru producerea stabilității și ale efectelor competitive al STDP. STDP reglează atât rata cât și coeficientul de variație a încărcării post-sinaptice peste o gamă largă de valori de intrare. Aceasta reprezintă o funcție regulatoare homeostatică a STDP, ceea ce în mod surprinzător, ca și la regula lui Hebb, este destabilizatoare pentru sinapsele individuale. STDP poate în mod diferențial să mărească valorile de intrare de scurtă întârziere date de stimuli.

Există evidențe experimentale, că reducerea întârzierii în răspunsul post-sinaptic apare și în in-vivo. Un fenomen analog cu reducția latenței, discutată aici, prezice că, dacă un șobolan cobai se mișcă într-o anumită parte a spațiului, acele celule care sunt active pentru spațiu, s-ar putea să emită mai repede, după ce șobolanul a trecut de câteva ori prin acel spațiu. Acest efect a fost observat experimental în (Blum & Abbott, 1996). O presupunere a modelului nostru este că slăbirea de către STDP este dominantă față de creșterea sinaptică. Aceasta este critică, pentru stabilitate. Dacă această presupunere nu este adevărată, rezultatele pe care le-am raportat, nu ar apărea niciodată, dintr-o combinație a STDP și depresia homo-sinaptică, de lungă durată (slăbirea impulsurilor pre-sinaptice, care emit în lipsa unui impuls post-sinaptic). Cât timp STDP duce la creșterea impulsurilor de intrare efective în timp ce STDP și/sau alte forme de plasticitate de lungă durată slăbește mai tare impulsurile ineficiente, rezultatele de bază întâlnite aici, ar trebui să fie aplicate.

Cea mai eficientă metodă de a mări puterea sinapsei este de a elibera emisia înainte de un impuls post-sinaptic și după aceea să se oprească din emisie, în așa fel, încât să nu piardă din putere prin emisii subsecvente care apar în urma activității post-sinaptice. O rată mare a depresiei post-sinaptice, care este

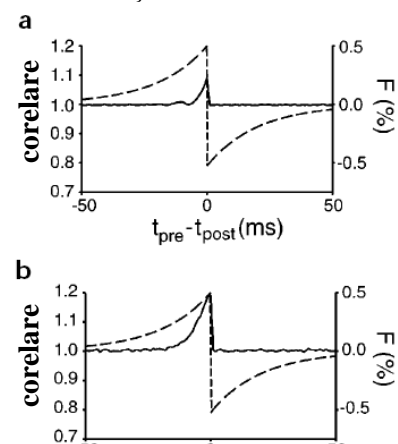


Figura 2.4 Corelația dintre încărcarea pre- și post-sinaptică

o caracteristică a sinapselor puternice din schema de re-distribuție, asigură aceasta. STDP care acționează pentru a modifica probabilitatea de emisie și schimbă rata depresiei sinaptice este astfel un mecanism foarte eficient și competitiv pentru a conduce sinapsele individuale spre limite maxime și minime. STDP, cu toate că contribuie la competiție, probabil nu este singura sursă a plasticității în situațiile de analiză Hebbiene. Ca orice regulă de modificare Hebbiană, STDP nu poate mări eficiența sinapsei fără o emisie post-sinaptică. Dacă sinapsa excitatoare din anumit motiv este prea slabă ca să emită, STDP nu poate să o salveze. Un mecanism non-hebbian, ca și scalarea sinaptică 7-9, poate să înlocuiască această funcție. În modelul STDP utilizăm două seturi de impulsuri care emit la un interval mai mare de timp de 100 ms, generează separat STDP și totuși nu sunt în competiție. Experimentele spun că competiția nu va avea loc în astfel de situații. Astfel de rezultate pot apărea dacă în fereastra temporară STDP, pentru o scădere a valorii sinaptice apare o urmă destul de lungă sau dacă STDP este suplimentat cu o depresie hetero-sinaptică suficient de puternică, sau este indusă competiția printr-o scalare sinaptică 7-9. Dimensiunea ferestrei temporare prin care mărirea și scăderea valorii sinaptice are loc, este critică în determinarea efectelor STDP. Ar fi în avantajul dimensiunii ferestrelor să fie diferite în anumite regiuni ale creierului, să se modifice în timpul dezvoltării, și să fie dinamic ajustabil în timp scurt. Acesta ar permite STDP să rămână compatibilă cu valori de intrare relevante de corelare.

### **2.3. Metode de propagare înapoi a erorilor (back-propagation) în rețelele neuronale pulsative codificate temporal**

Pentru o rețea de neuroni pulsanti care codifică informația în cronometrarea timpilor de impulsuri individuale, se poate formula o regulă de antrenare supravegheată, asemănătoare cu metoda error back-propagation tradițională. Cu acest algoritm, se poate demonstra cum pot executa rețelele de neuroni pulsanti cu potențiale de acțiune biologic rezonabile clasificări neliniare complexe în codificare temporală rapidă la fel de bine ca și rețelele cu codificarea ratei impulsurilor. Există studii, care arată, aplicabilitatea pentru problema clasică XOR, ridicând problema în termenii unui cadru temporal, și totodată pentru alte seturi de date reper (benchmark). Comparând numărul neuronilor pulsanti necesari pentru codificarea problemei interpolate XOR, rețelele antrenate demonstrează că codificarea temporală este o codificare fiabilă pentru o procesare rapidă a informației, și astfel necesită mai puțini neuroni decât în cazul rețelelor cu codificarea ratei de impulsuri. Totodată, putem observa că un calcul temporal sigur se poate obține doar prin folosirea funcțiilor-răspuns pulsative cu o constantă a timpului mai mare decât intervalul de codificare, cum se preconiza prin considerațiile teoretice.

Datorită succesului său în rețelele neuronale artificiale clasice, neuronul cu funcție de activare sigmoid este cunoscut ca fiind un model de succes a comportamentului neuronal biologic. Prin modelarea gradului la care un singur neuron biologic descarcă potențialii de acțiune ca și o funcție monoton crescătoare a intrării, s-au construit mai multe aplicații utile ale rețelelor neuronale artificiale, și s-au obținut mai multe puncte de vedere teoretice cu privire la comportamentul structurilor de conectare. Totodată, natura pulsativă a neuronilor naturali a dus la cercetări asupra puterii computaționale asociată cu codificarea temporală a informației în impulsuri individuale. În (Maas, 1997) s-a dovedit faptul că rețelele de neuroni pulsanti pot simula rețele neuronale sigmoid feed-forward întâmplător, și pot aproxima astfel orice funcție continuă, iar neuronii care propagă informația prin timpii impulsurilor individuale, sunt mai puternici din punct de vedere computațional decât neuronii cu funcții de activare sigmoid.

Cum impulsurile se pot descrie prin a.n. coordonate de „evenimente” (spațiu, timp) și numărul neuronilor activi (pulsanti) este tipic redus, rețelele neuronale de neuroni pulsanti artificiale permit implementarea unor rețele neuronale de dimensiuni mari (Mattia & Del Giudice, 2000). Prelucrarea timpilor mai multor impulsurilor individuale s-a propus de asemenea ca o nouă paradigmă pentru implementările de rețele neuronale VLSI care ar putea oferi o accelerare semnificativă. Arhitecturile de rețele bazate pe neuroni pulsanti care codifică informația în timpii individuali ai impulsurilor au oferit, printre altele, o hartă de auto-organizare asemănător SOM Kohonen și rețelelor de clustering nesupervizat (Natschläger & Ruf, 1998) (Bohte, Kok, & La Poutre, 2000).

Principiul codificării intensității de intrare prin timpii relativi de activare a fost de asemenea aplicat cu succes în cazul rețelelor de recunoaștere a caracterelor (Buonomano & Merzenich, 1999). Pentru aplicările rețelelor neuronale pulsative codificate temporal, totuși, nu s-a dezvoltat un algoritm supervizat practic până în prezent. Chiar și în (Buonomano & Merzenich, 1999), autorii recurg la

metode back-propagation de tip sigmoidal, tradiționale pentru a învăța rețeaua să distingă histogramele diferitelor răspunsuri de timpi ai impulsurilor. Pentru a face posibilă antrenarea supravegheată utilă cu modelul codificării temporale, s-a dezvoltat (Bohte, Kok, & La Poutr, 2002) un algoritm de antrenare pentru impulsuri individuale, care păstrează avantajele neuronilor pulsați, și, în același timp, permite o antrenare la fel de performantă ca și în cazul rețelelor neuronale sigmoid.

Se poate obține astfel un algoritm de antrenare supravegheată bazată pe propagarea înapoi a erorii (error back-propagation), pentru rețelele neuronale pulsative care transferă informația în temporizarea unui singur impuls. Pentru a depăși problema discontinuității neuronilor pulsați, se va aproxima funcția limită. Algoritmul poate să învețe să execute sarcini neliniare complexe, în rețele neuronale pulsative, cu aceeași precizie ca și în cazul rețelelor neuronale sigmoid tradiționale. Acest lucru se poate demonstra prin experimente pentru problema clasică de clasificare XOR, și pentru câteva seturi de date reale.

Aceste rezultate pot prezenta interes pentru comunitatea de specialitate, pentru că posibilitatea codificării informației în timpii impulsurilor se bucură de o atenție din ce în ce mai mare. Se poate demonstra empiric, faptul că rețelele de neuroni pulsați cu un raționament biologic pot să execute clasificarea neliniară complexă printr-o codificare temporală rapidă, la fel ca și rețelele cu codarea ratei de impulsuri (CRI). Cu toate că acest capitol descrie o regulă de antrenare aplicabilă la o clasă a rețelelor neuronale artificiale, neuronii pulsați sunt mai aproape de modelul neuronului biologic, decât cei sigmoid.

Pentru a realiza calcule cu CRI într-un timp foarte scurt, se consideră că în sistemele neuronale biologice se adună răspunsurile unui număr mare de neuroni pulsați, pentru a obține a măsură instantanee a ratei medii de activare. Chiar dacă necesită un număr mare de neuroni, acest răspuns totalizator este robust din cauza numărului mare de neuroni participanți. Dacă comparăm un astfel de cod de rate instant cu codificarea temporală prin impulsuri individuale, se știe că sunt necesari mult mai puțini neuroni, chiar dacă și cu prețul scăderii robusteții.

Conform unor studii din literatura de specialitate, s-au obținut rezultate care susțin anticiparea faptului că lungimea segmentului în creștere al potențialului post-sinaptic trebuie să fie mai mare decât structura temporală relevantă pentru a permite o computație temporală de încredere. Pentru neuronii pulsați, potențialul post-sinaptic descrie dinamica unui impuls care influențează un neuron, și se poate modela ca și o diferență a două funcții descrescătoare exponențial (Gerstner W. , Spiking neurons, 2001).

### 2.3.1. Propagarea înapoi a erorii

Am obținut propagarea înapoi a erorii în mod analog cu rezultatele prezentate de (Rumelhart, Hinton, & Williams, 1986). Ecuațiile au fost obținute pentru o rețea neuronală pulsativă feed-forward total conectată cu straturi notate cu  $H$  (strat de intrare),  $I$  (strat ascuns) și  $J$  (strat de ieșire), unde algoritmul rezultat se aplică la fel de bine pentru rețele cu mai multe straturi ascunse.

Scopul algoritmului este să învețe un set de timpi de activare prescriși, notați  $\{t_j^d\}$ , la neuronii de ieșire  $j \in J$  pentru un set dat de modele de intrare  $\{P[t_1..t_h]\}$ , unde  $P[t_1..t_h]$  ce definește un model de intrare individual descris de timpi de impuls individuali pentru fiecare neuron  $h \in H$ . Am ales ca funcție de eroare media celor mai mici pătrate, dar alte variante ca entropie sunt de asemenea posibile. Având timpii de activare prescriși  $\{t_j^d\}$  și timpii de activare măsurăți  $\{t_j^a\}$ , această funcție de eroare este definită de:

$$E = \frac{1}{2} \sum_{j \in J} (t_j^a - t_j^d)^2 \quad \text{Ec. 18}$$

Pentru propagarea înapoi a erorii vom trata fiecare terminal sinaptic ca o conexiune separată  $k$  cu ponderea  $w_{ij}^k$ . Astfel pentru regula propagării înapoi a erorii, trebuie să calculăm:

$$\Delta w_{ij}^k = -\eta \frac{\partial E}{\partial w_{ij}^k} \quad \text{Ec. 19}$$

unde  $\eta$  este rata de învățare și  $w_{ij}^k$  ponderea conexiunii  $k$  de la neuronul  $i$  la neuronul  $j$ .

Cum  $t_j$  este funcție de  $x_j$  ce depinde de ponderile  $w_{ij}^k$ , derivata din partea dreapta a Ec. 19 poate fi extinsă în:

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{\partial E}{\partial t_j} (t_j^a) \frac{\partial t_j}{\partial w_{ij}^k} (t_j^a) = \frac{\partial E}{\partial t_j} (t_j^a) \frac{\partial t_j}{\partial x_j(t)} (t_j^a) \frac{\partial x_j(t)}{\partial w_{ij}^k} (t_j^a) \quad \text{Ec. 20}$$

În ultimii doi factori din dreapta, exprimăm  $t_j$  ca o funcție mărginită a intrării post-sinaptice  $x_j(t)$  în jurul lui  $t = t_j^a$ . Presupunem că pentru un interval suficient de mic în jurul lui  $t = t_j^a$ , funcția  $x_j$  poate fi aproximată ca o funcție liniară de  $t$ , cum se vede în Figura 2.5. Pentru un astfel de interval mic, aproximăm funcția de prag  $\partial t_j(x_j) = -\partial x_j(t_j) / \alpha$ , cu  $\frac{\partial t_j}{\partial x_j(t)}$  fiind derivata funcției inverse a  $x_j(t)$ . Valoarea lui  $\alpha$  este derivata locală a

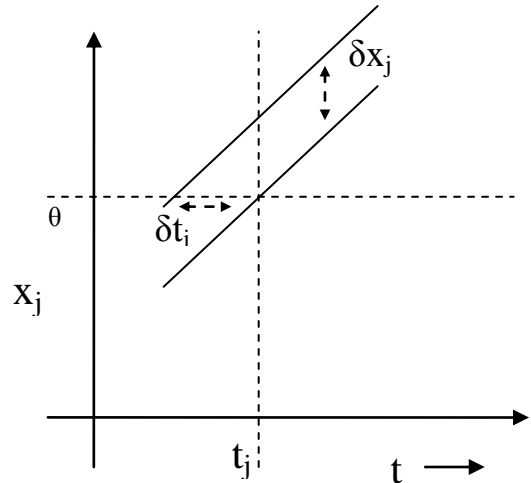


Figura 2.5 Relația dintre  $\partial x_j$  și  $\partial t_j$  pentru un spațiu  $\mathcal{E}$  în jurul  $t_j$ .

$x_j(t)$  în punctul  $t$ , care este  $\alpha = \frac{\partial x_j(t)}{\partial t} (t_j^a)$ .

Al doilea factor din EEc. 20 poate fi scris ca:

$$\frac{\partial t_j}{\partial x_j(t)} (t_j^a) = \frac{\partial t_j(x_j)}{\partial x_j(t)} \Big|_{x_j=\theta} = \frac{-1}{\alpha} = \frac{-1}{\frac{\partial x_j(t)}{\partial t} (t_j^a)} = \frac{-1}{\sum_{i,l} w_{ij}^l \frac{\partial y_i^l(t)}{\partial t} (t_j^a)} \quad \text{Ec. 21}$$

În următoarele calcule vom scrie termeni ca  $\frac{\partial x_j(t_j^a)}{\partial t_j^a}$  pentru  $\frac{\partial x_j(t)}{\partial t} (t_j^a)$ .

Observăm că această aproximare este valabilă numai dacă ponderile spre un neuron nu sunt modificate în așa fel încât potențialul de membrană să nu mai atingă valoarea de prag, și în consecință neuronul să nu se mai activeze. Aceasta este o potențială problemă, dar poate fi combătută prin codificarea intrărilor astfel încât primele impulsuri sunt în mod automat considerate „mai importante” decât următoarele impulsuri. Observăm însă că în cazul în care un neuron nu se activează pentru *oricare* din intrări, nu exista un mecanism care să corecteze ponderile. În experimentele noastre am setat ponderile inițiale în așa fel încât fiecare neuron din rețea a răspuns la cel puțin o parte din modelul de intrare. Cu această prevedere, nu am avut nici o problemă cu neuroni inactivi.

De asemenea observăm că această aproximare ar putea sugera că pentru rate de învățare mari, algoritmul poate fi mai puțin eficient.

Primul factor în EEc. 20, derivata lui  $E$  în punctul  $t_j$ , este:

$$\frac{\partial E(t_j^a)}{\partial t_j^a} = (t_j^a - t_j^d) \quad \text{Ec. 22}$$

Avem:

$$\frac{\partial x_j(t_j^a)}{\partial w_{ij}^k} = \frac{\partial \left\{ \sum_{n \in \Gamma_j} \sum_l w_{nj}^l y_n^l(t_j^a) \right\}}{\partial w_{ij}^k} = y_i^k(t_j^a) \quad \text{Ec. 23}$$

Dacă combinăm aceste rezultate, Ec. 19 devine:

$$\Delta w_{ij}^k(t_j^a) = -\eta \frac{y_i^k(t_j^a) \cdot (t_j^d - t_j^a)}{\sum_{i \in \Gamma_j} \sum_l w_{ij}^l \frac{\partial y_i^l(t_j^a)}{\partial t_j^a}} \quad \text{Ec. 24}$$

Pentru comoditate, definim  $\delta_j$ :

$$\delta_j \equiv \frac{\partial E}{\partial t_j^a} \frac{\partial t_j^a}{\partial x_j(t_j^a)} = \frac{(t_j^d - t_j^a)}{\sum_{i \in \Gamma_j} \sum_l w_{ij}^l \frac{\partial y_i^l(t_j^a)}{\partial t_j^a}} \quad \text{Ec. 25}$$

și EEc. 20 poate fi acum exprimat ca:

$$\frac{\partial E}{\partial w_{ij}^k} = y_i^k(t_j^a) \frac{\partial E}{\partial t_j^a} \frac{\partial t_j^a}{\partial x_j(t_j^a)} = y_i^k(t_j^a) \delta_j \quad \text{Ec. 26}$$

de aici rezultă:

$$\Delta w_{ij}^k = -\eta y_i^k(t_j^a) \delta_j \quad \text{Ec. 27}$$

Ecuțiile Ec. 25 și Ec. 27 oferă funcția simplă de adaptare a ponderilor pentru neuronii din stratul de ieșire.

Continuăm cu straturile ascunse: pentru propagarea înapoi a erorii în alte straturi decât cel de ieșire, eroarea delta generalizată în stratul  $I$  este definită pentru  $i \in I$  cu timpii de activare măsurati  $t_i^a$ :

$$\delta_i \equiv \frac{\partial t_i^a}{\partial x_i(t_i^a)} \frac{\partial E}{\partial t_i^a} = \frac{\partial t_i^a}{\partial x_i(t_i^a)} \sum_{j \in \Gamma^i} \frac{\partial E}{\partial t_j^a} \frac{\partial t_j^a}{\partial x_j(t_j^a)} \frac{\partial x_j(t_j^a)}{\partial t_i^a} = \frac{\partial t_i^a}{\partial x_i(t_i^a)} \sum_{j \in \Gamma^i} \delta_j \frac{\partial x_j(t_j^a)}{\partial t_i^a} \quad \text{Ec. 28}$$

unde  $\Gamma^i$  denotă un set de succesori neuronali direcți în stratul  $J$  conectați la neuronul  $i$ .

Ca și în (Bishop, 1995), în Ec. 28 dezvoltăm eroarea locală  $\frac{\partial E(t_i^a)}{\partial t_i^a}$  conform erorii ponderate corespunzătoare stratului următor  $J$ . Pentru această dezvoltare, aceleași reguli ca și în EEc. 20 sunt folosite cu aceleași restricții, dar pentru  $t = t_i$ .

Termenul  $\frac{\partial t_i^a}{\partial x_i(t_i^a)}$  a fost calculat în Ec. 21, iar pentru  $\frac{\partial x_j(t_j^a)}{\partial t_i^a}$

$$\frac{\partial x_j(t_j^a)}{\partial t_i^a} = \frac{\partial \sum_{l \in I} \sum_k w_{lj}^k y_l^k(t_j^a)}{\partial t_i^a} = \sum_k w_{ij}^k \frac{\partial y_i^k(t_j^a)}{\partial t_i^a} \quad \text{Ec. 29}$$

Deci,

$$\delta_i = \frac{\sum_{j \in \Gamma^i} \delta_j \left\{ \sum_k w_{ij}^k \frac{\partial y_i^k(t_j^a)}{\partial t_i^a} \right\}}{\sum_{h \in \Gamma_i} \sum_l w_{hi}^l \frac{\partial y_h^l(t_i^a)}{\partial t_i^a}} \quad \text{Ec. 30}$$

Astfel pentru un strat ascuns și conform Ec. 27, Ec. 28, Ec. 29 și Ec. 30, regula adaptării ponderilor se compune în:

Analog cu algoritmi tradiționali de propagare înapoi a erorii, regula adaptării ponderilor Ec. 31

$$\Delta w_{hi}^k = -\eta y_h^k(t_i^a) \delta_i = -\eta \frac{\sum_j \left\{ \delta_j \sum_k w_{ij}^k \frac{\partial y_i^k(t_j^a)}{\partial t_i^a} \right\}}{\sum_{n \in \Gamma_i} \sum_l w_{ni}^l \frac{\partial y_n^l(t_i^a)}{\partial t_i^a}} \quad \text{Ec. 31}$$



de mai sus se generalizează pentru o rețea cu mai multe straturi ascunse  $I$  numerotate  $J-1, \dots, 2$  calculând eroarea delta la stratul  $i$  din eroarea delta în stratul  $i+1$ , propagând efectiv eroarea înapoi.

Algoritmul obținut mai sus, numit *SpikeProp*, este rezumat în tabelul următor:

Tabelul 1 Algoritmul *SpikeProp*

<b>1</b>	Se calculează $\delta_j$ pentru toate ieșirile conform Ec. 25
<b>2</b>	Pentru fiecare strat ulterior $I = J-1, \dots, 2$ Se calculează $\delta_j$ pentru toți neuronii din $I$ conform Ec. 30
<b>3</b>	Pentru stratul de ieșire $J$ , ajustăm $w_{ij}^k$ cu $\Delta w_{ij}^k = -\eta y_i^k(t_j) \delta_j$ <span style="float: right;">Ec. 27</span>
<b>4</b>	Pentru fiecare strat ulterior $I = J-1, \dots, 2$ Ajustăm $w_{hi}^k$ cu $\Delta w_{hi}^k = -\eta y_h^k(t_i) \delta_i$ Ec. 31

O simplă modificare a acestei scheme ar fi să includem un termen de inerție:

$$\Delta w_{ij}^k = -\eta \Delta w_{ij}^k(t) + p w_{ij}^k(t-1) \tag{Ec. 32}$$

unde  $p$  este variabila de inerție. Ca urmare a lucrării (Bohte, Kok, & La Poutré, 2000), Xin *et al.* (Xin & Embrechts, 2001) au arătat că asemenea modificări ale algoritmului *SpikeProp* accelerează într-adevăr convergența în mod semnificativ.

### 2.3.2. Problema XOR

În continuare, aplicăm algoritmul *SpikeProp* problemei XOR. Funcția XOR este un exemplu clasic de problemă neliniară care necesită unități ascunse pentru a transforma intrarea în ieșirea dorită.

Pentru a codifica funcția XOR în modele impuls-timp, am asociat valoarea 0 cu o activare „întârziată” și valoarea 1 cu o activare „prematură”. Cu valorile specifice 0 și 6 pentru timpii de intrare respectivi, vom folosi următorul XOR codificat temporal:

Tabelul 2 Codificarea modelelor de intrare și ieșire al problemei XOR rezolvată cu algoritmul *Spikeprop*

Modele de intrare	Modele de ieșire
0	16
6	10
0	10
6	16

Numerele în tabel reprezintă durata impulsurilor, să spunem în milisecunde. Vom folosi un al treilea neuron de intrare (offset) în rețea care se va activa întotdeauna la  $t = 0$  pentru a marca timpul de referință (altfel problema devine banală). Definim diferența dintre timpii echivalenți cu „0” și „1” ca fiind intervalul de codificare  $\Delta T$ , care în acest exemplu corespunde cu 6ms.

Pentru rețea vom folosi o arhitectură de tip feed-forward. Conexiunile au un interval de întârziere de 15 ms; astfel vom avea întârzieri sinaptice între 1 și 16 ms. Potențialul post-sinaptic (PPS) este definit de o funcție  $\alpha$  ca în

$$\varepsilon(t) = \frac{t}{\tau} e^{-\frac{t}{\tau}} \tag{Ec. 33}$$

cu un timp de scădere  $\tau = 7$  ms. Valori mai mari până la cel puțin 15 ms rezultă într-o învățare similară.

Rețeaua a fost compusă din 3 neuroni de intrare (2 neuroni de codificare și 1 neuron de referință), 5 neuroni ascunși (dintre care un neuron inhibitor generând doar PPS-uri negative) și 1

neuron de ieșire. Au fost permise doar ponderi pozitive. Cu această configurație rețeaua a învățat în mod fiabil modelul XOR în timpul a 250 de cicluri cu  $\eta = 0.01$ . Pentru a „învăța” XOR, au trebuit să fie ajustate  $16 \times 3 \times 5 + 16 \times 5 \times 1 = 320$  de ponderi individuale.

În timpul folosirii algoritmului s-a descoperit că este necesară introducerea specifică a neuronilor inhibitori și excitatori, neuronii inhibitori și excitatori fiind definiți ca generând PPS-uri cu valori negative și respectiv pozitive, folosind doar ponderi pozitive. De fapt, algoritmul *SpikeProp* nu va converge spre un rezultat satisfăcător dacă se permite ca conexiunile să conțină un amestec de ponderi pozitive și negative. Cauza acestei probleme ar putea fi că în cazul în care avem valori atât pozitive cât și negative, efectul unei conexiuni individuale către un neuron țintă nu mai este o funcție monoton crescătoare (cum ar fi pentru constante de timp suficient de mari). S-a observat că introducerea de neuroni inhibitori și excitatori nu este o restricție: dacă extindem rețeaua în așa fel încât fiecărui neuron excitator să-i corespundă un neuron inhibitor, vom obține de fapt o conexiune mixtă. În experimente însă, includerea unui sau a doi neuroni inhibitori a fost suficientă pentru a realiza învățarea.

### 2.3.3. Gradientul erorii și rata de învățare

Considerăm influența ratei de învățare  $\eta$  și constanta de timp  $\tau$  pe capacitățile de învățare ale algoritmului *SpikeProp* într-o rețea neuronală pulsativă.

Cum s-a văzut în secțiunea anterioară, aproximarea dependenței de timpul de activare  $t_j^a$  pe intrarea post-sinaptică  $x_j$  este valabilă doar pentru un interval mic în jurul lui  $t_j^a$ . S-a descoperit că pentru rate de învățare mari probabilitatea convergenței scade, însă pentru rate de învățare până la 0.01, creșterea ratei de învățare corespunde cu timpi mai rapizi de învățare. Asta se poate vedea în

Figura 2.6, unde numărul mediu de iterații de învățare necesare pentru funcția XOR este reprezentat pentru un număr de constante de timp  $\tau$ .

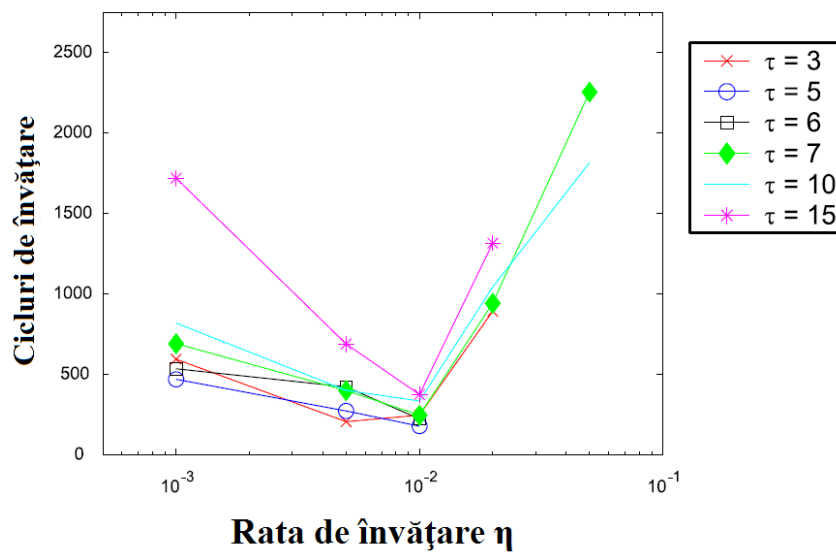


Figura 2.6 Învățare XOR: numărul mediu necesar de iterații de învățare pentru a atinge suma propusă a erorilor pătratice de 1.0. Media a fost calculată pentru cazurile în care s-a converș la un rezultat satisfăcător.

În Figura 2.7 este reprezentată fiabilitatea învățării pentru diferite valori ale constantei de timp  $\tau$ . Reprezentarea arată că cele mai sigure rezultate pentru o convergență optimă au fost obținute pentru valori ale constantei de timp mai mari (cu puțin) decât intervalul de timp  $\Delta T$  în care se codifică problema XOR (aici  $\Delta T = 6$ ).

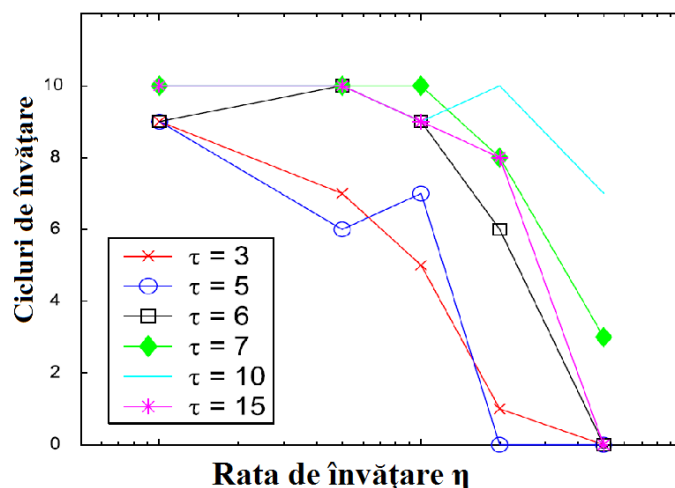
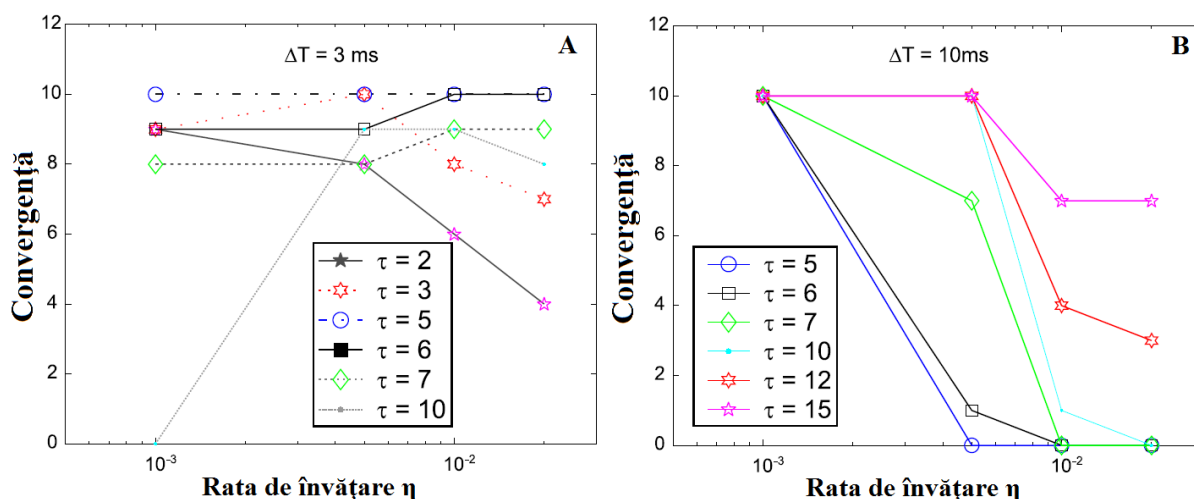


Figura 2.7 Învățare XOR: Numărul de cazuri din 10 care au converș

Graficele de convergență pentru diferite valori ale intervalului de codificare  $\Delta T$  sunt reprezentate în Figura 2.8 și arată același comportament. Aceste rezultate confirmă rezultatele obținute în (Maass, 1996), unde s-a demonstrat că teoretic constanta de timp  $\tau$  trebuie să fie mai mare decât intervalul de codificare relevant  $\Delta T$ .


 Figura 2.8 Numărul de cazuri din 10 care au converș pentru diferite valori ale  $\tau$ .

- A) Pentru un interval de codificare  $t = 0 \dots 3, [0,3]^2 \rightarrow [7,10]$ . Suma propusă a erorilor pătratice de 0.7.  
 B) Pentru un interval de codificare  $t = 0 \dots 10, [0,10]^2 \rightarrow [15,25]$ . Suma propusă a erorilor pătratice de 3.0.

## 2.4. Dezvoltarea modelelor neuronale neuromorfe

În concordanță cu cele prezentate în prealabil, putem afirma că în creierul uman elementul de primordial care stă la baza prelucrării informației este neuronul pulsativ. Acestea sunt unități specializate în efectuarea unor "calculuri" și sunt total diferite de cele utilizate în dezvoltarea rețelelor neuronale artificiale. Un neuron pulsativ emite impulsuri la un anumit moment în timp, adică se aprinde. Astfel, neuronul transmite un semnal electric prin axonul său – spre neuronii la care este conectat – numit potențial de aprindere, sau impuls de aprindere sau de activare (IA). Forma și magnitudinea IA este independentă de semnalul de intrare a neuronului, însă temporizarea acesteia depinde de intrare.

Poate fi vorba de orice fel de calculator artificial, de multe ori ținem ca o anumită ordine în efectuarea pașilor de calcul să fie respectată, ea încadrându-se într-o schemă temporală globală de lucru. De exemplu, în cazul unei rețele neuronale artificiale de tip *feed-forward*, așteptăm de la un neuron al

stratului  $x$ , ca ieșirea lui să se activeze la momentul  $x$  de calcul, independent de semnalele aplicate la intrarea rețelei. În contrast cu acest aspect, în cazul unei rețele neuronale neuromorfe (similare cu cele naturale) momentele de activare ale neuronilor sunt determinate de intrarea rețelei. Rețelele bazate pe astfel de neuroni sunt capabile să exploateze timpul ca o resursă de codare și calcul, mult mai eficient decât majoritatea modelelor de calcul existente.

În consecință, analiza rețelelor neuronale bazate pe astfel de modele ne determină să reevaluăm rolul timpului și spațiului, care sunt dimensiunile de bază ale oricărui calcul. Unele cercetări în acest sens au avut ca efect stimularea dezvoltării de elemente electronice hardware noi, care utilizează impulsuri pentru a efectua calcule. Astfel de elemente sunt de exemplu, circuitele VLSI analogice utilizate în sisteme senzoriale în timp-real.

### 2.4.1. Modelul unei rețele neuronale artificiale formale cu neuroni pulsativi

Să considerăm o rețea neuronală pulsativă (RNP) (Maas W. , 1997) (Maas W. , 1998), formată dintr-o mulțime finită  $V$  de neuroni pulsanți (NP) și o mulțime  $E \subseteq V \times V$  formată de sinapse.

Fiecărui neuron  $i \in V$ , corespunde o funcție de prag  $\mathcal{G}_i : R^+ \rightarrow R^+$  (Figura 2.9) și fiecărei sinapse  $\langle i, j \rangle \in E$  corespunde o funcție de răspuns  $\varepsilon_{i,j} : R^+ \rightarrow R$ , unde  $R^+ = \{x \in R : x \geq 0\}$  (Figura 2.10)

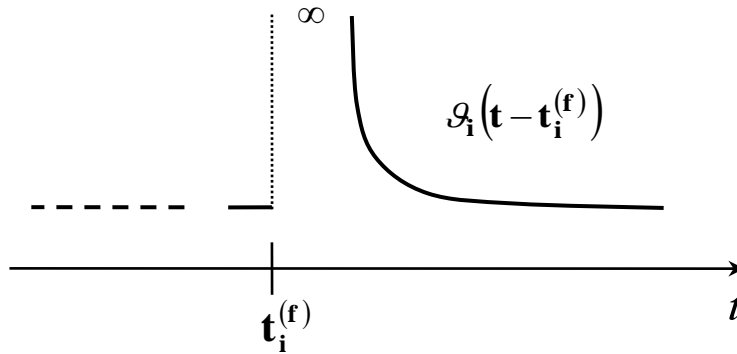


Figura 2.9 Forma funcției de prag al unui neuron natural

Un NP  $j$  emite IA în momentele de timp  $t_j^{(1)}$ ,  $t_j^{(2)}$ , ..., adică emite un impuls stereotipic de-a lungul structurii axonale arborescente. În sinapse, care realizează legătura dintre neuronul  $j$  și alți neuroni, IA-ul emis declanșează producerea unei materii denumite neuro-transmițător. Ca urmare se generează impulsuri pozitive (potențiale post-sinaptice excitatorii, PPSE) sau negative (potențiale post-sinaptice inhibitorii, PPSI) care vor modifica potențialul membranei neuronului post-sinaptic.

Caracteristica de timp a PPSE și PPSI-urilor cauzate de IA al neuronului  $j$  se poate exprima în felul următor  $\omega_{ij} \varepsilon(t - t_j^{(f)})$ , unde valoarea funcției de răspuns  $\varepsilon(t - t_j^{(f)})$  este 0 dacă  $(t - t_j^{(f)})$  este mai mic decât, timpul de propagare  $\Delta_{ij}$ , apoi crește aproape liniar și în final descrește exponențial spre 0. Factorul  $\omega_{ij}$  este o pondere, care are valoare pozitivă în cazul unor PPSE-uri și valoare negativă în cazul PPSI-urilor.

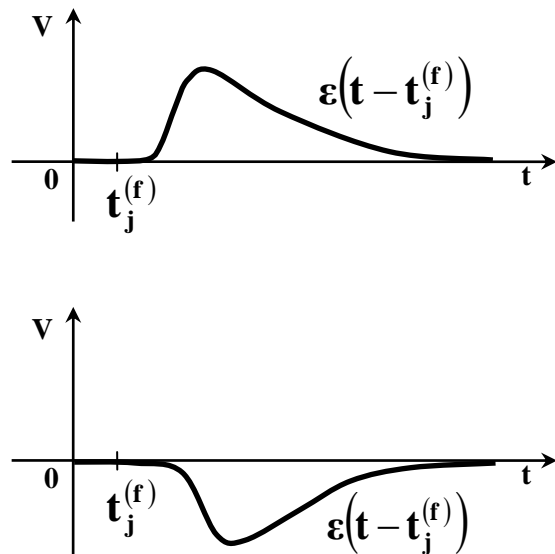


Figura 2.10 Forma funcției de răspuns a NP natural (sus PPSE și PPSI jos)

Există numeroase modele matematice, care încearcă să identifice momentele de timp  $t_i^{(1)}, t_i^{(2)}, \dots$  ale apariției de IA-uri la neuronii biologici, în funcție de momentele  $t_j^{(f)}$  în care un neuron pre-sinaptic  $j$  emite un IA. Datorită relativei simplități matematice care duce și la o mai ușoară implementare software și hardware, *modelul de reacție la impuls* este unul dintre modelele cele mai convenabile pentru aplicarea în analize computerizate.

### 2.4.2. Modelul de reacție la impuls

În cazul modelului de reacție la impuls (MRI) potențialul de membrană  $u_i^{(t)}$  al somei unui neuron  $i$  la un moment dat  $t$ , se poate exprima în felul următor

$$\sum_{j,f} \omega_{ij} \varepsilon(t - t_j^{(f)}) \quad \text{Ec. 34}$$

sumă care alcătuită din PPSE și PPSI-uri emise de neuronul pre-sinaptic  $j$  în momente de timp  $t_j^{(f)}$  premergătoare momentului  $t$ .

Pe de altă parte în MRI definim și un *factor de inhibare sau înăbușire*,  $\eta(t - \hat{t}_i)$ , care este o funcție de  $\hat{t}_i$ , adică momentul premergător lui  $t$  în care neuronul  $i$  a emis ultimul IA. Dacă  $(t - \hat{t}_i)$  este pozitiv, dar de valoare foarte mică (sub 2 ms), atunci factorul  $\eta(t - \hat{t}_i)$  arată o valoare puternic negativă, iar în cazul în care  $(t - \hat{t}_i)$  este mai mare *factorul de inhibare* tinde spre zero.

Deci, în cazul modelului de reacție la impuls neuronul  $i$  se aprinde atunci și numai atunci, când valoarea sumei rezultante atinge un prag  $\mathcal{G}$ , iar derivata de ordinul unu al acestei expresii este pozitivă, cu alte cuvinte ne apropiem de prag în sens crescător.

Modelul se poate utiliza mai eficient în aplicații de calcul, dacă presupunem, că potențialul de repaus al neuronilor este 0, iar valoarea de prag este întotdeauna mai mare decât zero ( $\mathcal{G} > 0$ ). În realitate, potențialul de repaus al unui neuron natural are o valoare aproximativă de  $-70\text{mV}$  iar cel de prag este în jurul valorii de  $-50\text{mV}$ . Efectul PPSE și PPSI-urilor asupra acestor valori este de ordinul câtorva milivolți.

Comparând acest model cu modelul clasic *Hodkin-Huxley* (Trappenberg, 2002) (Maas & Bishop, 1999) putem spune că MRI este mult mai ușor de utilizat în analize de calcul, deoarece în definiția acestuia nu am aplicat ecuații diferențiale. Pe de altă parte, dacă alegem corespunzător valorile pentru funcțiile  $\eta$ ,  $\mathcal{G}$  și  $\varepsilon$ , MRI este capabil de aproxima foarte bine dinamica neuronilor biologici.

## 2.5. Simularea rețelelor neuronale pulsative: o sinteză a instrumentelor și a strategiilor

### 2.5.1. Introducere

Dovezile experimentale au arătat că sincronizarea impulsurilor poate fi importantă pentru a explica modul în care se realizează calculele la nivel de neuron, motivând astfel folosirea modelelor cu neuroni pulsați în locul modelelor tradiționale bazate pe rata impulsurilor. În același timp a apărut un număr din ce în ce mai mare de instrumente ce permit simularea rețelelor neuronale pulsative. Aceste instrumente oferă utilizatorului posibilitatea să obțină simulări precise ale unui sistem computațional dat, precum și rezultate publicabile într-un timp relativ scurt. Problemele de calcul legate de neuronii pulsați sunt foarte variate însă. În unele cazuri este necesară folosirea unei reprezentări biofizice detaliate a neuronilor, de exemplu când se dorește reproducerea măsurătorilor electrofiziologice intracelulare (Destexhe & Sejnowski, 2001). În aceste cazuri se folosesc modele bazate pe conductanță (COBA), cum ar fi tipul de model folosit de Hodgkin și Huxley (Hodgkin & Huxley, 1952). În alte

cazuri nu se dorește să se surprindă exact mecanismele prin care se generează impulsul, și sunt suficiente modele mai simple, cum ar fi modelul „integrează și activează” (integrate-and-fire – IF). Simularea modelelor de tip IF este de asemenea foarte rapidă, făcându-le în special atractive pentru simularea rețelelor neuronale la scară mare.

Există două familii de algoritmi pentru simularea rețelelor neuronale: algoritmi sincroni sau „bazați pe tact”, în care toți neuronii sunt actualizați simultan la fiecare tact, și algoritmi asincroni sau „bazați pe eveniment”, în care neuronii sunt actualizați doar când primesc sau transmit un impuls (există de asemenea strategii hibride). Pentru algoritmi sincroni codul este ușor de scris și se poate aplica oricărui model. Deoarece timpii de impuls sunt legați de o grilă discretă de timp, precizia simulării poate fi o problemă. Algoritmi asincroni au fost dezvoltati în special pentru simulări exacte, ceea ce este posibil pentru modele simple. Pentru rețele foarte mari, durata simulării prin ambele metode crește foarte mult odată cu numărul impulsurilor transmise, dar fiecare strategie are propriile avantaje și dezavantaje.

În cele ce urmează vom trece în revistă diferite strategii de simulare, și vom sublinia în ce măsură precizia în timp a evenimentelor impulsuri au impact asupra unui sau a unei mici rețele de neuroni de tip IF cu sinapse plastice. Apoi vom trece în revistă simulatoarele sau mediile de simulare disponibile în prezent, punând accent numai pe instrumentele publice și non-comerciale pentru simularea rețelelor neuronale pulsative.

## 2.5.2. Strategii de simulare

Ne limităm la algoritmi seriali pentru precizie. Secțiunile specifice NEST și SPLIT oferă informații în plus despre conceptele pentru calculul în paralel.

Există două familii de algoritmi pentru simularea rețelelor neuronale: algoritmi sincroni sau „bazați pe tact”, în care toți neuronii sunt actualizați simultan la fiecare tact, și algoritmi asincroni sau „bazați pe eveniment”, în care neuronii sunt actualizați doar când primesc sau transmit un impuls. Aceste două abordări au unele proprietăți comune pe care le vom descrie întâi exprimând problema simulării rețelelor neuronale în formalismul sistemelor hibride, de ex. ecuații diferențiale cu evenimente discrete (impulsuri). În aceste condiții apar unele strategii comune pentru o reprezentare și simulare eficiente.

Deoarece vom compara algoritmi din punct de vedere al eficienței de calcul, să ne punem întâi următoarea întrebare: cât timp îi poate lua unui algoritm bun să simuleze o rețea mare? Să presupunem că avem  $N$  neuroni cu rata medie de activare  $F$  și numărul mediu de sinapse  $p$ . Dacă toate impulsurile transmise sunt luate în calcul, atunci o simulare care durează 1s (timp biologic) trebuie să proceseze  $N \times p \times F$  transmisii. Scopul unui algoritm eficient este să atingă acest număr minim de operații (desigur până la un factor constant multiplicativ). Dacă simularea nu este restricționată la interacțiuni imediate de impulsuri, de ex. dacă modelul include fante sinaptice sau interacțiuni dendro-dendritice, atunci numărul optim de operații poate fi mult mai mare, dar nu vom lua în seamă aici problema interacțiunilor gradate.

### 2.5.2.1. Formalismul unui sistem hibrid

Matematic, neuronii pot fi descriși ca fiind sisteme hibride: starea unui neuron evoluează continuu conform unor ecuații biofizice, care sunt de obicei ecuații diferențiale (deterministice sau stocastice, ecuații ordinare sau parțial diferențiale), și impulsurile recepționate prin sinapse declanșează schimbări în unele variabile. Astfel dinamica unui neuron poate fi descrisă cum urmează:

$$\frac{dX}{dt} = f(X) \quad \text{Ec. 35}$$

$$X \leftarrow g_i(X) \text{ la un impuls recepționat prin sinapsa } i \quad \text{Ec. 36}$$

unde  $X$  este un vector care descrie starea neuronului. În teorie, luând în calcul morfologia neuronului ar duce la ecuații parțial diferențiale; însă, în practică se poate aproxima arborele dendritic prin compartimente iso-potențiale cuplate, ceea ce va duce de asemenea la un sistem diferențial cu evenimente discrete. Impulsurile sunt emise când o anumită condiție de activare (threshold) este satisfăcută, de exemplu  $V_m \geq \theta$  pentru modelele de tip IF (unde  $V_m$  este potențialul membranei, și ar fi prima componentă a vectorului  $X$ ), și / sau  $dV_m/dt \geq \theta$  pentru modelele de tip Hodgkin-Huxley (HH). Putem rezuma spunând că un impuls este emis de fiecare dată când o anumită condiție  $X \in A$  este

îndeplinită. Pentru modelele de tip IF, potențialul membranei, care ar fi prima componentă a  $X$ , se resetează când se produce un impuls. Resetarea poate fi integrată în formalismul sistemului hibrid dacă considerăm de exemplu că impulsurile de ieșire acționează asupra  $X$  printr-o sinapsă adițională (virtuală):  $X \leftarrow g_0(X)$ .

Cu acest formalism devine evident că *timpii impulsurilor nu trebuiesc stocați* (exceptând desigur cazurile în care sunt incluse întârzierile de transmisie), chiar dacă formulările fenomenologice ar sugera contrariul. Spre exemplu, considerăm următorul model de tip IF descris în (Gütig & Sompolinsky, 2006):

$$V(t) = \sum_i \omega_i \sum_{t_i} K(t-t_i) + V_{rest} \quad \text{Ec. 37}$$

unde  $V(t)$  este potențialul membranei,  $V_{rest}$  este potențialul de repaus,  $\omega_i$  este ponderea sinaptică a sinapsei  $i$ ,  $t_i$  sunt timpii impulsurilor venite de la sinapsa  $i$ , și  $K(t-t_i) = \exp(-(t-t_i)/\tau) - \exp(-(t-t_i)/\tau_s)$  este potențialul post-sinaptic (PSP) la care contribuie fiecare impuls primit. Modelul poate fi descris ca un sistem diferențial de două variabile cu evenimente discrete cum urmează:

$$\begin{aligned} \tau \frac{dV}{dt} &= V_{rest} - V + J \\ \tau_s \frac{dJ}{dt} &= -J \end{aligned} \quad \text{Ec. 38}$$

$$J \leftarrow J + \frac{\tau - \tau_s}{\tau} w_i \text{ la un impuls recepționat prin sinapsa } i$$

Practic putem exprima în acest mod fiecare PSP sau curenți descriși în literatură (de ex. funcții  $\alpha$ , funcții bi-exponențiale). Mai mulți autori au descris transformarea de la expresii fenomenologice la formalismul unui sistem hibrid pentru conductanțe sinaptice și curenți ((Destexhe, Mainen, & Sejnowski, 1994), (Destexhe, Mainen, & Sejnowski, 1994); (Rotter & Diesmann, 1999); (Giugliano, 2000)), depresiuni sinaptice pe termen scurt (Giugliano, Bove, & Grattarola, 1999), și plasticitatea dependenței de timpii impulsurilor (Song, Miller, & Abbott, 2000). În multe cazuri, modelul răspunsului la impuls (Gerstner & Kistler, 2002) este de asemenea expresia integrală a unui sistem hibrid. Pentru a deduce formularea diferențială a unui curent sau conductanță post-sinaptică dată, o metodă ar fi de a considera ultima ca răspunsul la impuls al unui sistem liniar invariabil în timp (care poate fi considerat ulterior ca un filtru (Jahnke, Roth, & Schoenauer, 1998)) și să folosim instrumentele de transformare din teoria procesării semnalelor cum ar fi transformata  $Z$  (Kohn & Wörgötter, 1998) (a se vedea de asemenea (Sanchez-Montanez, 2001)) sau transformata Laplace (transformata  $Z$  este echivalenta transformatei Laplace în domeniul timpului digital, de ex. pentru algoritmi sincroni).

### 2.5.2.2. Folosirea liniarităților pentru simulări sinaptice rapide

În general, numărul de variabile de stare a unui neuron (lungimea vectorului  $X$ ) crește odată cu numărul sinapselor, deoarece fiecare sinapsă are propria sa dinamică. Acest fapt constituie o problemă majoră pentru simularea eficientă a rețelelor neuronale, atât în privința memoriei consumate cât și a timpului de calcul. Cu toate acestea, mai mulți autori au observat că toate variabilele sinaptice caracterizate de aceeași dinamică liniară pot fi reduse la o singură variabilă (Wilson & Bower, 1989); (Bernard, Ge, Stockley, Willis, & Wheal, 1994); (Lytton, 1996); (Song, Miller, & Abbott, 2000). De exemplu, următorul set de ecuații diferențiale, descriind un model de tip IF cu  $n$  sinapse având conductanțe exponențiale:

$$\begin{aligned} C \frac{dV}{dt} &= V_0 - V + \sum_i g_i(t)(V - E_s) \\ \tau_s \frac{dg_1}{dt} &= -g_1 \quad \dots \quad \tau_s \frac{dg_n}{dt} = -g_n \end{aligned} \quad \text{Ec. 39}$$

$$g_i \leftarrow g_i + w_i \text{ la un impuls recepționat prin sinapsa } i$$

este matematic echivalent cu următorul set de două ecuații diferențiale:

$$C \frac{dV}{dt} = V_0 - V + g(t)(V - E_s) \quad \text{Ec. 40}$$

$$\tau_s \frac{dg}{dt} = -g$$

$$g \leftarrow g + w_i \quad \text{la un impuls recepționat prin sinapsa } i$$

unde  $g$  este conductanța sinaptică totală. Aceeași reducere se aplică sinapselor cu dinamică dimensională mai mare, atât timp cât este liniară și schimbările determinate de impulsuri ( $g_i \leftarrow g_i + w_i$ ) sunt aditive și nu depind de starea sinapsei (ex. regula  $g_i \leftarrow g_i + w_i * f(g_i)$  ar cauza o problemă). Unele modele de plasticitate dependentă de timpii impulsurilor (cu interacțiuni liniare între perechile de impulsuri) pot fi de asemenea simulate în acest mod (Abbott & Nelson, 2000). Însă, unele modele biofizice importante nu sunt liniare și astfel nu pot beneficia de această optimizare, în particular interacțiunile mediate NMDA și sinapsele saturate.

### 2.5.2.3. Algoritmi sincroni sau bazați pe tact

Într-un algoritm sincron sau bazat pe tact (vezi pseudocodul din Figura 2.11), starea variabilelor fiecărui neuron (și posibil a sinapselor) este actualizată la fiecare tact:  $X(t) \rightarrow X(t + dt)$ .

```

t=0
while t<duration
  for every neuron
    process incoming spikes
    advance neuron dynamics by dt
  end
  for every neuron
    if vm>threshold
      reset neuron
      for every connection
        send spike
      end
    end
  end
end
t=t+dt
end

```

*Actualizări stare*

*Propagare impulsuri*

Figura 2.11 Un algoritm simplu bazat pe tact

Având ecuații diferențiale neliniare, se poate folosi o metodă de integrare cum ar fi Euler sau Runge-Kutta (Press, Flannery, Teukolsky, & Vetterling, 1993) sau, pentru modele HH, metode implicite (Hines, 1984). Neuronii cu morfologii complexe sunt de obicei separați spațial și modelați ca compartimente interactive: sunt de asemenea descriși matematic de perechi de ecuații diferențiale, pentru care s-au dezvoltat metode de integrare speciale (pentru detalii vezi de ex. secțiunea specifică Neuron în cele ce urmează). Dacă ecuațiile diferențiale sunt liniare, atunci operațiunea de actualizare  $X(t) \rightarrow X(t + dt)$  este de asemenea liniară, ceea ce înseamnă că actualizarea variabilelor de stare se realizează pur și simplu înmulțind  $X$  cu o matrice:  $X(t + dt) = AX(t)$  (Hirsch & Smale, 1974) (vezi de asemenea (Rotter & Diesmann, 1999), pentru o aplicație în rețelele neuronale), ceea ce este foarte convenabil în medii științifice bazate pe vectori cum ar fi Matlab sau Scilab. Apoi, după ce s-au actualizat toate variabilele, condiția de activare este verificată pentru fiecare neuron. Fiecare neuron ce satisface această condiție produce un impuls care se va transmite către neuronii țintă, actualizând variabilele corespunzătoare ( $X \leftarrow g_i(X)$ ). Pentru modele IF, potențialul membranei fiecărui neuron care pulsează va fi actualizat.



**Complexitatea calculului.** Timpul de simulare al unui astfel de algoritm constă în două etape: (1) actualizări de stare și (2) propagare de impulsuri. Presupunând că numărul de variabile de stare pentru întreaga rețea crește odată cu numărul de neuroni  $N$  din rețea (ceea ce se întâmplă când se aplică simplificarea descrisă în secțiunea 2.5.2.2), costul fazei de actualizare este de ordinul  $N$  pentru fiecare pas, deci este  $O(N/dt)$  pe secundă de timp biologic ( $dt$  este durata pasului temporal). Această componentă crește odată cu complexitatea modelelor neuronale și cu precizia simulării. În fiecare secundă (de timp biologic) o medie de  $F \times N$  impulsuri sunt produse de către neuroni ( $F$  este media ratei de activare) și fiecare impuls trebuie propagat către  $p$  neuroni țintă. Astfel faza de propagare constă în  $F \times N \times p$  impulsuri propagate pe secundă. Acestea sunt în esență adăugări de ponderi  $w_i$  variabilelor de stare, în consecință operații simple al căror cost nu crește odată cu complexitatea modelelor. Însușind, costul total al calculului efectuate pe secundă de timp biologic este de ordinul

$$\begin{aligned} \text{Actualizare} &+ \text{Propagare} \\ c_U \times \frac{N}{dt} &+ c_p \times F \times N \times p \end{aligned} \quad \text{Ec. 41}$$

unde  $c_U$  este costul unei actualizări și  $c_p$  este costul unei propagări de impuls; de obicei  $c_U$  este mult mai mare decât  $c_p$  dar asta depinde de implementare. Astfel, pentru rețele foarte dense, totalul este dominat de faza de propagare și este liniar în numărul de sinapse, ceea ce e optim. În practică însă, prima fază este neglijabilă numai când se îndeplinește următoarea condiție:

$$\frac{c_p}{c_U} \times F \times p \times dt \gg 1 \quad \text{Ec. 42}$$

De exemplu, rata medie de activare în cortex poate fi redusă ca  $F = 1\text{ Hz}$  (Olshausen & Field, 2005), și presupunem că  $p = 10000$  sinapse pe neuron și  $dt = 0.1\text{ ms}$ , vom avea  $F \times p \times dt = 1$ . În acest caz, considerând că fiecare operație în faza de actualizare este mai grea decât în faza de propagare (în special pentru modele complexe), de ex.,  $c_p < c_U$ , este probabil ca faza de actualizare să fie partea dominantă în costul total al calculului. Astfel se pare că chiar și în rețele cu conectivitate realistă, creșterea preciziei ( $dt$  mai mic) poate fi în detrimentul eficienței simulării.

#### 2.5.2.4. Algoritmi asincroni sau bazați pe eveniment

Algoritmii asincroni sau bazați pe eveniment nu sunt la fel de folosiți ca și cei bazați pe tact deoarece implementarea lor este mai complicată (vezi pseudocodul din Figura 3) și mai puțin generală. Avantajele cheie ale algoritmilor asincroni constă într-un posibil câștig în viteză prin faptul că nu vom mai calcula mulți pași mici de actualizare pentru un neuron în care nu se produce nici un eveniment și pentru că timpii impulsurilor sunt calculați în mod exact (dar se poate vedea mai jos pentru algoritmi bazați pe eveniment aproximați); în particular timpii impulsurilor nu mai sunt aliniați unei grile (ceea ce este o posibilă sursă de erori).

Problema simulării sistemelor dinamice cu evenimente discrete este un subiect de cercetare bine stabilit în domeniu (Ferscha, 1996); (Sloot, Kaandorp, Hoekstra, & Overeinder, 1999); (Fujimoto, 2000); (Zeigler, Praehofer, & Kim, 2000) (vezi de asemenea (Roche & Martinez, 2003); (Mayrhofer, Affenzeller, Praehofer, Hfer, & Fried, 2002)), cu structuri de date adecvate și algoritmi deja disponibili pentru cercetătorii din domeniu. Vom începe prin a descrie cazul simplu în care interacțiunile sinaptice sunt instantanee, de ex. când impulsurile se pot produce doar în momentul recepționării unui impuls (fără întârzieri); apoi vom aborda cazul cel mai general.

**Interacțiuni sinaptice instantanee.** Într-un algoritm asincron sau bazat pe eveniment simularea avansează de la un eveniment la următorul eveniment. Evenimentele pot fi impulsuri emise de neuronii din rețea sau impulsuri externe (de obicei impulsuri întâmplătoare descrise de un proces Poisson). Pentru modele în care impulsurile pot fi produse de către un neuron doar în momentul recepționării unui alt impuls, simularea bazată pe eveniment este relativ ușoară (vezi pseudocodul din Figura 2.12).

Evenimentele sunt stocate într-o coadă de așteptare (care este un fel de listă ordonată). O iterație constă în:

1. Extragerea următorului eveniment
2. Actualizarea stării neuronului corespunzător (de ex. calcularea conform ecuației diferențiale și adăugarea ponderii sinaptice)
3. Se verifică dacă neuronul satisface condiția de activare, caz în care se introduc în coada de așteptare evenimente pentru fiecare neuron post-sinaptic

```

while queue not empty and t<duration
  extract event with lowest timing           Procesare eveniment
  (= timing t, target i, weight w)
  compute state of neuron i at time t
  update state of neuron i (+w)
  if vm>threshold
    for each connection i -> j           Propagare impuls
      insert event in the queue
    end
    reset neuron i
  end
end
end

```

Figura 2.12 Un algoritm simplu bazat pe eveniment cu interacțiuni sinaptice

În cazul simplu al întârzierilor de transmisie identice, structura de date pentru coada de așteptare poate fi pur și simplu o listă FIFO (First In, First Out – primul intrat, primul ieșit), care se poate implementa rapid (Cormen, Leiserson, Rivest, & Stein, 2001). Când întârzierile iau valori dintr-o mulțime discretă de dimensiuni mici, cel mai ușor este să folosim o listă FIFO pentru fiecare valoare de întârziere, așa cum e descris în (Mattia & Del Giudice, 2000). Este de asemenea mai eficient să folosim o listă FIFO separată pentru a gestiona evenimentele externe diverse.

În cazul întârzierilor arbitrare, este nevoie de o structură de date mai complexă. În informatică, structurile de date care mențin în mod eficient o listă ordonată de evenimente marcate în timp sunt grupate sub numele de liste de prioritate (priority queues) (Cormen, Leiserson, Rivest, & Stein, 2001). Subiectul listelor de prioritate este vast și bine documentat; exemple sunt stive binare, stive Fibonacci (Cormen, Leiserson, Rivest, & Stein, 2001), liste calendar (Brown, 1988); (Claverol, Brown, & Chad, 2002) sau arbori Emde Boas (van Emde Boas, Kaas, & Zijlstra, 1976) (vezi de asemenea (Connolly, Marian, & Reilly, 2003), în care sunt comparate diverse liste de prioritate). Folosirea unei liste de prioritate eficiente este crucială pentru buna funcționare a unui algoritm bazat pe eveniment. Devine și mai crucială când interacțiunile sinaptice nu sunt instantanee.

**Interacțiuni sinaptice non-instantanee.** Pentru modele în care momentul emiterii unui impuls nu coincide neapărat cu momentul recepționării unui impuls, simularea bazată pe eveniment este mai complexă. Descriem întâi algoritmul simplu fără întârzieri și fără evenimente externe (vezi pseudocodul din Figura 2.13). O iterație constă în:

1. Identificarea neuronului care va emite următorul impuls
2. Actualizarea acestui neuron
3. Propagarea impulsului, adică actualizarea neuronilor către care se transmite impulsul

În general acest lucru se realizează păstrând o listă ordonată a timpilor viitoarelor impulsuri pentru toți neuronii. Acești timpi ai impulsurilor sunt doar provizorii deoarece orice impuls din rețea poate modifica toți timpii impulsurilor viitoare. Cu toate acestea, impulsul cu cel mai mic timp din listă este garantat ca fiind următorul. Astfel, următorul algoritm pentru o iterație garantează corectitudinea simulării (vezi Figura 2.13):

1. Se extrage impulsul cu cel mai mic timp din listă
2. Se actualizează starea neuronului corespunzător și se recalculează timpul impulsului viitor
3. Se actualizează starea neuronilor țintă
4. Se recalculează timpii impulsurilor viitoare pentru neuronii țintă

```

for every neuron i
    compute timing of next spike           Inițializare
    insert event in priority queue
end
while queue not empty and t < duration
    extract event with lowest timing       Procesare impuls
    (event = timing t, neuron i)
    compute state of neuron i at time t
    reset membrane potential
    compute timing of next spike
    insert event in queue

    for every connection i -> j          Comunicare impuls
        compute state of neuron j at time t
        change state with weight w(i,j)
        compute timing of next spike
        insert / change / supress event in queue
    end
end
end
    
```

Figura 2.13 Un algoritm simplu bazat pe eveniment cu interacțiuni sinaptice non-instantanee

Pentru simplificare, vom ignora întârzierile de transmisie pentru descrierea de mai sus. Includerea lor într-un algoritm bazat pe evenimente nu este la fel de simplă ca în cazul unui algoritm bazat pe tact, dar este o complicație minoră. Când un neuron produce un impuls, evenimentele sinaptice viitoare sunt stocate într-o altă listă de priorități, în care timpii evenimentelor nu pot fi modificați. Prima fază a algoritmului (extragerea impulsului cu cel mai mic timp) este înlocuită cu extragerea următorului eveniment, care poate fi un eveniment sinaptic sau emiterea unui impuls. Se pot folosi două liste separate sau una singură. Evenimentele externe pot fi tratate în același mod. Deși întârzierile aduc complicații în codul algoritmilor bazați pe eveniment, acestea pot de fapt să simplifice administrarea listei de priorități pentru impulsurile emise. Într-adevăr, principala dificultate în simularea rețelelor cu interacțiuni sinaptice non-instantanee constă în faptul că impulsurile programate spre a fi emise pot fi anulate, amânate sau avansate de către impulsurile ce urmează a fi primite. Dacă întârzierile de transmisie sunt mai mari decât o valoare pozitivă  $\tau_{\min}$ , atunci toate impulsurile programate spre a fi emise în  $[t, t + \tau_{\min}]$  ( $t$  fiind timpul prezent) sunt garantate. Astfel algoritmi pot exploata natura acestor întârzieri pentru a crește viteza simulării (Lee & Farhat, 2001).

**Complexitatea calculului.** Lăsând la o parte costul tratării evenimentelor externe (care este minor), putem împărți costul computațional al tratării unui impuls emis după cum urmează (presupunem că  $p$  este numărul mediu de sinapse per neuron):

- Se extrage evenimentul (în cazul interacțiunilor sinaptice non-instantanee)
- Se actualizează neuronul și țintele acestuia:  $p + 1$  actualizări
- Se inserează  $p$  evenimente sinaptice în listă (în caz că avem întârzieri)
- Se actualizează timpii de impuls pentru  $p + 1$  neuroni (în cazul interacțiunilor sinaptice non-instantanee)
- Se inserează sau se reprogreamază  $p + 1$  evenimente în listă (impulsuri viitoare în cazul interacțiunilor sinaptice non-instantanee)

Deoarece avem  $F \times N$  impulsuri pe secundă de timp biologic, numărul operațiilor este aproximativ proporțional cu  $F \times N \times p$ . Costul computațional total pe secundă de timp biologic poate fi scris în mod concis după cum urmează:

*Actualizare + Impuls + Listă*

$$(c_U + c_S + c_Q) \times F \times N \times p$$

*Ec. 43*

unde  $c_U$  este costul unei actualizări a variabilelor de stare,  $c_S$  este costul calculării timpului următorului impuls (interacțiuni sinaptice non-instantanee) și  $c_Q$  este costul mediu al inserărilor în și al extragerilor din lista (sau listele) de prioritate. Astfel timpul simulării este liniar cu numărul sinapselor, ceea ce este optim. Cu toate acestea observăm ca operațiile implicate sunt mai grele decât cele din faza propagării pentru algoritmi bazați pe tact (vezi secțiunea anterioară), astfel factorul multiplicativ este foarte probabil să fie mai mare. De asemenea am presupus că  $c_Q$  este  $O(1)$ , adică că operațiile de intrare și ieșire din coadă pot fi efectuate într-un timp mediu constant cu structura de date aleasă pentru lista de prioritate. În cazul simplu al interacțiunilor sinaptice instantanee și având întârzieri omogene, se poate folosi o simplă listă FIFO (First In, First Out), caz în care inserările și extragerile sunt foarte rapide și au durată de timp constantă. Pentru cazul general, există structuri de date pentru care operațiile de intrare și ieșire din coadă pot fi efectuate într-un timp mediu constant ( $O(1)$ ), de ex. listele calendar (Brown, 1988); (Claverol, Brown, & Chad, 2002), dar sunt destul de complexe, și deci  $c_Q$  este o constantă mare. În implementări mai simple de liste de prioritate cum ar fi stivele binare, operațiile de intrare și ieșire din coadă se pot efectua în  $O(\log m)$  timp, unde  $m$  este numărul de evenimente din coadă. În concluzie, se pare că administrarea cozii este componenta crucială pentru algoritmi bazați pe eveniment în general.

### 2.5.3. Prezentare generală a mediilor de simulare

#### 2.5.3.1. NEURON

**Domeniul de utilitate al NEURON.** NEURON este un mediu de simulare folosit pentru crearea și folosirea modelelor empirice a neuronilor biologici și a circuitelor neuronale. Inițial și-a câștigat reputația de a fi bine adaptat pentru modele COBA de celule cu anatomie ramificată complexă, inclusiv potențial extracelular în jurul membranei, și proprietăți biofizice cum ar fi tipuri de canale multiple, distribuție neomogenă a canalelor, acumulare și difuzie ionică, și neurotransmitter-i secundari. La începutul anilor 90, NEURON era deja folosit în unele laboratoare pentru modele de rețele cu multe mii de celule, iar în ultima decadă a avut parte de multe îmbunătățiri care au permis ca construcția și simularea rețelelor la scară mare să fie mai ușoară și mai eficientă.

Până la această dată au apărut mai mult de 600 articole și cărți care descriu modele NEURON, de la bucați de membrană până la rețele la scară mare cu zeci de mii de celule pulsative artificiale sau COBA<sup>1</sup>. În 2005 au fost publicate peste 50 de articole cu subiecte cum ar fi mecanismele care stau la baza transmisiei și plasticității sinaptice (Banitt, Martin, & Segev, 2005), modularea integrării sinaptice prin intermediul curenților activi sub pragul de activare (Prescott & De Koninck, 2005), excitabilitate dendritică (Day, Carr, Ulrich, Ilijic, Tkatch, & Surmeier, 2005), rolul fantelor sinaptice în rețele (Migliore, Hines, & Shepherd, 2005), efectele plasticității sinaptice în dezvoltarea și modul de operare al rețelelor biologice (Saghatlyan, et al., 2005), consecințele zgomotului de canal și sinaptic asupra procesării informației în neuroni și rețele (Badoual, Rudolph, Piwkowska, Destexhe, & Bal, 2005), mecanisme celulare și ale rețelei pentru codificare și recunoaștere temporală (Kanold & Manis, 2005), stări și oscilații ale rețelei (Wolf, et al., 2005), efectele îmbătrânirii asupra funcțiilor neuronale (Markaki, Orphanoudakis, & Poirazi, 2005), imprimarea corticală (Moffitt & McIntyre, 2005), stimularea profundă a creierului (Grill, et al., 2005), și epilepsia rezultată în urma mutațiilor de canal (Vitko, Chen, Arias, Shen, Wu, & Perez-Reyes, 2005), și traumatisme cerebrale (Houweling, Bazhenov, Timofeev, Steriade, & Sejnowski, 2005).

<sup>1</sup> <http://www.neuron.yale.edu/neuron/bib/usednrn.html>

**Prin ce diferă NEURON de alte neuro-simulatoare.** Principalul motiv pentru care se folosesc simulatoare specializate în locul celor de uz general constă în perspectiva unui control conceptual îmbunătățit, și posibilitatea de a exploata structura ecuațiilor care definesc modelul pentru a obține robustețe, acuratețe și eficiență computațională. Câteva diferențe cheie între NEURON și alte neuro-simulatoare constă în modul în care se încearcă atingerea acestor scopuri.

**Control conceptual.** Ciclul formulării ipotezei, testării, și revizuirii, care stă la baza tuturor cercetărilor științifice, presupune că se pot deduce consecințele unei ipoteze. Principalul motiv pentru care se folosește modelarea computațională este utilitatea acesteia în abordarea acelor ipoteze a căror consecințe nu pot fi nu pot fi determinate intuitiv sau analitic. Valoarea oricărui model ca metodă de evaluare a unei anumite ipoteze depinde în mod critic de potrivirea cât mai exactă dintre model și ipoteză. Fără o asemenea potrivire rezultatele simulării nu pot fi un bun test al ipotezei respective. Din punctul de vedere al utilizatorului, primul obstacol din calea modelării computaționale constă în dificultatea cu care se obține controlul conceptual, adică să ne asigurăm ca modelul computațional reflectă cu fidelitate ipoteza noastră.

NEURON are mai multe proprietăți care facilitează controlul conceptual, și continuă să obțină noi astfel de proprietăți în timp ce evoluează pentru a satisface nevoile în continuă schimbare ale cercetătorilor din domeniu. Multe din aceste proprietăți intră în categoria generală a specificațiilor de „sintaxă nativă” a proprietăților modelului: astfel atributele cheie ale neuronilor biologici și ale rețelelor au echivalente directe în NEURON. De exemplu utilizatorii NEURON pot specifica proprietățile de obturare ale canalelor ionice prin scheme kinetice sau familii de ecuații diferențiale de tip HH. Un alt exemplu ar fi că modelele pot include circuite electronice construite cu LinearCircuitBuilder, un instrument GUI a cărui paletă include rezistențe, condensatori, surse de tensiune și curent, și amplificatori operaționali. Cel mai elocvent exemplu de aplicare a conceptului de sintaxă nativă a proprietăților NEURON ar fi modul în care tratează proprietățile de înlănțuire ale neuronilor, care este diferit față de orice alt neuro-simulator. Utilizatorii NEURON nu trebuie să trateze în mod direct compartimentele. Astfel, celulele sunt reprezentate prin neurite separate, numite secțiuni, care pot fi asamblate în arhitecturi de tip arbore (topologia unei celule model). Fiecare secțiune are propriile proprietăți anatomice și biofizice, plus un parametru de discretizare care specifică rezoluția locală a grilei spațiale. Proprietățile unei secțiuni pot varia încontinuu odată cu lungimea sa, și variabilele neomogene spațial sunt accesate în termeni de distanță normalizată de-a lungul fiecărei secțiuni (Hines & Carnevale, 1997) (Capitolul 5 din (Carnevale & Hines, 2006)). Odată ce utilizatorul a specificat topologia celulei, geometria, proprietățile biofizice, și parametrul de discretizare pentru fiecare secțiune, NEURON construiește automat structurile de date interne care corespund unei familii de ODE-uri pentru ecuația discretizată de înlănțuire a modelului.

### 2.5.3.2. GENESIS

**Filozofia de design și capacitățile GENESIS.** GENESIS (the General Neural Simulation System) a primit acest nume deoarece a fost proiectat de la bun început ca un sistem general de simulare extensibil, folosit la modelarea realistă a sistemelor neuronale și biologice (Bower & Beeman, 1998). Cu GENESIS s-au executat simulări de la componente subcelulare și reacții biochimice (Bhalla, 2004) până la modele complexe de neuroni singulari (De Schutter & Bower, 1994), simulări de rețele mari (Nenadic, Ghosh, & Ulinski, 2003), și modele la nivel de sistem (Stricanne & Bower, 1998). Aici, „modelele realiste” sunt definite ca fiind acele modele care sunt bazate pe organizarea anatomică și fiziologică cunoscută a neuronilor, circuitelor și rețelelor (Bower, 1995). De exemplu, modelele de celulă realiste includ în mod tipic morfologie dendritică și o mare varietate de conductanțe ionice, iar modelele de rețea realiste încearcă să copieze tiparele cunoscute de proiecție axonală.

*Parallel GENESIS* (PGENESIS) este o extensie a software-ului GENESIS care poate rula pe aproape orice cluster paralel, SMP, super-computer, sau rețea de calculatoare care suportă MPI și / sau PVM, și pe care poate rula și software-ul GENESIS serial. Este folosit de obicei la simularea rețelelor mari ce implică zeci de mii de modele de celulă realiste (Hereld, Stevens, Teller, & van Drongelen, 2005).

Utilizatorii GENESIS au la dispoziție un proces bine documentat prin care îi pot extinde capacitățile prin adăugarea de noi tipuri de obiecte (clase) GENESIS definite de utilizator, sau comenzi de limbaj de scriptare, fără să fie nevoie să înțeleagă sau să modifice codul de simulare GENESIS.

GENESIS vine deja echipat cu mecanismele necesare pentru a crea ușor modele de rețele la scară mare realizate din modele de neuroni singulari care au fost deja implementate cu GENESIS.

Deși utilizatorii au adăugat, de exemplu, modelul simplificat de neuron pulsativ (Izhikevich E. M., 2003) (acum încorporat în GENESIS), și ar putea de asemenea adăuga modele IF sau alte forme abstracte de modele de neuroni, aceste forme de neuroni nu sunt suficient de realiste pentru a satisface interesul majorității celor ce modelează cu GENESIS. Din acest motiv, GENESIS nu include în mod normal modele IF de neuroni, și nu se oferă implementări GENESIS pentru modele IF. Neuronii tipici GENESIS sunt modele multi-compartimentale cu o varietate de conductanțe de tip HH dependente de voltaj și / sau calciu.

**Modelarea cu GENESIS.** GENESIS este un sistem de simulare orientat-obiect, în care o simulare se construiește din blocuri de construcție de bază (elemente GENESIS). Aceste elemente comunică între ele prin transmiterea de mesaje și fiecare conține informația propriilor variabile (câmpuri) și metode (acțiuni) folosite la efectuarea calculului și a altor sarcini din timpul simulării. Elementele GENESIS sunt create ca instanțieri ale unui tip particular de obiecte precompilate care se comportă ca un șablon. Neuronii model sunt construiți din aceste componente de bază, cum ar fi compartimente neuronale și canale ionice cu conductanță variabilă, legate între ele prin mesaje. Neuronii pot fi legați între ei prin conexiuni sinaptice pentru a forma circuite neuronale și rețele. Această abordare orientată-obiect asigură generalitatea și flexibilitatea sistemului, deoarece permite utilizatorilor să facă schimb sau să refolosească modelele sau componentele modelelor. Mulți utilizatori GENESIS își bazează scripturile pe exemplele care sunt puse la dispoziție cu GENESIS sau în pachetul GENESIS Neural Modeling Tutorials (Beeman, 2005).

GENESIS folosește un interpretor și un limbaj de simulare de nivel înalt pentru a construi neuronii și rețelele acestora. Folosirea unui astfel de interpretor cu tipuri de obiecte pre-compilate, în locul unui pas în plus pentru a compila scripturile în cod mașină binar, are avantajul de a permite utilizatorului să interacționeze cu și să modifice simularea în timpul derulării acesteia, fără să sacrifice viteza simulării. Comenzile pot fi date fie interactiv prin linia de comandă, fie prin folosirea scripturilor de simulare, sau prin intermediul interfeței grafice. Cele 268 de comenzi ale limbajului de scriptare și cele 125 de tipuri de obiecte puse la dispoziție cu GENESIS sunt suficient de puternice astfel încât este nevoie de doar câteva rânduri de script pentru a specifica o simulare sofisticată. De exemplu, „citorul de celule” al GENESIS permite construirea modelelor neuronale complexe prin citirea specificațiilor dintr-un fișier de date.

GENESIS pune la dispoziție o varietate de mecanisme pentru a modela difuzia calciului și conductanțe dependente de calciu, precum și plasticitatea sinaptică. Are de asemenea un număr de „obiecte dispozitiv” care pot fi folosite ca interfață în simulare pentru a asigura simulării diverse tipuri de intrări (generatoare de puls și impuls, circuite de limitare a potențialului, etc.) sau măsurători (histograme de peristimulus și interval inter-impuls, măsurători de frecvență a impulsurilor, histograme de auto-corelare și inter-corelare, etc.). Se oferă de asemenea tipuri de obiecte pentru modelarea căilor biochimice (Bhalla & Iyengar, 1999). O listă și descriere a tipurilor de obiecte GENESIS, cu legături la documentația completă, poate fi găsită în secțiunea „Obiecte” a documentului GENESIS Reference Manual, care poate fi descărcat sau vizionat pe site-ul web al GENESIS.

### 2.5.3.3. NEST

**Inițiativa NEST.** Problema simulării rețelelor neuronale de mărime și complexitate realiste a fost multă vreme subestimată. Acest lucru este reflectat de numărul limitat de articole având ca subiect structuri de date și algoritmi corespunzători apărute în revistele cu impact mare. Lipsa conștientizării de către cercetători și agențiile de finanțare a nevoii de progres în tehnologia simulării și durabilitatea investițiilor poate fi cauzată parțial de faptul că un simulator corect din punct de vedere matematic pentru un anumit model de rețea neuronală poate fi implementat de un individ în câteva zile. Însă, datorită acestei rutine s-a ajuns la un ciclu de cod program nescalabil fiind continuu rescris într-un mod greu de întreținut de către novici, obținându-se un progres minim al fundamentelor teoretice.

Datorită disponibilității în creștere a resurselor computaționale, studiile din domeniul simulării devin din ce în ce mai ambițioase și mai populare. Într-adevăr, răspunsurile la multe întrebări din domeniul neuro-științei sunt în prezent accesibile doar prin intermediul simulării. O consecință negativă

a acestei tendințe este că reproducerea și verificarea rezultatelor acestor studii devine din ce în ce mai grea. Instrumentele de simulare ad-hoc ale trecutului nu ne pot oferi un grad de cunoaștere potrivit. În locul acestora avem nevoie de simulatoare de rețele neuronale atent dezvoltate, validate, documentate și expresive având o bază mare de utilizatori. Mai mult decât atât, actualul progres spre modele mai realiste impune necesitatea unor simulări corespunzătoare mai eficiente. Acest lucru este valabil mai ales pentru noul domeniu de studiu asupra modelelor de rețele la scară mare care încorporează plasticitatea. Aceste cercetări nu sunt deloc fiabile fără simulatoare paralele cu proprietăți scalabile excelente, ceea ce soluțiile ad-hoc nu pot oferi. De asemenea, pentru a fi considerate folositoare de către un număr cât mai mare de cercetători pe o perioadă lungă de timp, aceste simulatoare trebuie să fie ușor de întreținut și de extins.

Pe baza acestor considerații, inițiativa NEST a fost fondată ca un proiect de colaborare pe termen lung care să încurajeze dezvoltarea tehnologiei pentru simulatoarele sistemelor neuronale (Diesmann & Gewaltig, 2002). Instrumentul de simulare NEST este implementarea de referință a acestei inițiative. Programul este oferit comunității științifice sub o licență open-source prin site-ul web al inițiativei NEST<sup>2</sup>. Licența cere cercetătorilor să facă trimitere către inițiativa NEST în munca lor derivată de la codul original și, mai important, în rezultatele științifice obținute cu ajutorul acestui program. Site-ul web oferă de asemenea referințe la material relevant pentru simularea rețelelor neuronale în general, și vrea să devină o sursă științifică importantă de informații pentru simulări de rețele. Ajutorul este oferit prin intermediul site-ului NEST și o listă de discuții prin email. În prezent NEST este folosit în predare la școli de vară internaționale și în cursurile regulate ale Universității din Freiburg.

#### 2.5.3.4. SPLIT

**Simulatoare paralele.** Dezvoltarea simulatoarelor paralele în știința neuronală computațională a fost relativ înceată. În prezent există puține simulatoare paralele disponibile publicului larg, dar sunt departe de a fi la fel de generale, flexibile și documentate ca și simulatoarele seriale folosite în general cum ar fi NEURON (Hines & Carnevale, 1997) și GENESIS (Bower & Beeman, 1998). Pentru GENESIS există PGENESIS, iar pentru NEURON a început dezvoltarea unei versiuni paralele. De asemenea există simulatoare ca NCS<sup>3</sup> (Ermentrout, 2004), NEST (Morrison, Mehring, Geisel, Aertsen, & Diesmann, 2005), și simulatorul de paralelizare SPLIT (Hammarlund & Ekeberg, 1998). Însă acestea sunt încă în stadiu experimental sau de dezvoltare din multe puncte de vedere.

#### 2.5.3.5. Mvaspike

**Modelarea cu evenimente.** S-a discutat de multe ori că potențialele de acțiune așa cum sunt produse de multe tipuri de neuroni pot fi considerate ca fiind *evenimente*: ele constă în impulsuri stereotipice care apar suprapuse peste dinamica voltajului intern al neuronilor. Ca rezultat, multe modele de neuroni oferă moduri de definire a timpilor evenimentelor asociați cu fiecare potențial de acțiune emis, deseori prin definirea unui prag de activare<sup>4</sup>. Instrumentele de simulare neuronală au profitat de acest fapt de multă vreme, prin folosirea unui *algorithm bazat pe eveniment* (vezi secțiunea 2). Într-adevăr când se vorbește despre *evenimente* în contextul simulării rețelelor neuronale, ne gândim în primul rând la algoritmi *bazați pe eveniment*, și avem impresia că folosirea evenimentelor upstream, în timpul etapei de modelare, este deseori subestimată.

Mvaspike a fost proiectat ca un cadru de modelare și simulare bazată pe eveniment. Este bazat pe o abordare a modelării teoriei mulțimilor bine stabilită – specificație de sistem bazată pe evenimente discrete (discrete event system specification – DEVS) (Zeigler & Vahie, 1993); (Zeigler, Praehofer, & Kim, 2000). Modelele țintă sunt sisteme bazate pe evenimente discrete: dinamica lor poate fi descrisă

<sup>2</sup> <http://www.nest-initiative.org>

<sup>3</sup> <http://brain.cse.unr.edu/ncsdocs>

<sup>4</sup> Pragul de activare aici trebuie înțeles într-un sens foarte larg, de la un simplu prag de detectare impuls într-un model continuu (ex. HH) până la un prag activ care este folosit în expresia matematică a dinamicii (modelul IF).



prin schimbările variabilelor de stare la momente alese arbitrar în timp<sup>5</sup>. Mvaspike realizează apropierea dintre expresia mai familiară a dinamicii continue, folosită în general în neuroștiință, și folosirea de modele centrate pe eveniment în simulator. Lucru convenabil pentru modele simple care reprezintă majoritatea modelelor în Mvaspike (mai ales modele IF sau de fază, și SRM-uri). Watts (Watts, 1994) a notat deja că multe proprietăți neuronale pot fi reprezentate în mod explicit și ușor în sisteme bazate pe evenimente discrete. Dacă ne gândim la *perioade* refractare absolute, *timi* de creștere ai PSP-urilor, *întârzieri* de propagare axonală, acestea sunt noțiuni legate în mod direct de intervale de timp (și în concluzie, de evenimente) care sunt folositoare pentru a descrie multe aspecte ale dinamicii neuronale. Astfel fiind în mod evident destul de departe de modelele bine stabilite bazate pe conductanță, mai corecte din punct de vedere electro-fiziologic, un alt scop al Mvaspike este să ia în considerare cât mai mult posibil aceste modele mai complexe, prin susținerea explicită a evenimentelor discrete în timp, și, posibil, discretizarea spațiului de stare pentru integrarea dinamicilor continue sau hibride.

Formalismul DEVS face posibilă modelarea unor sisteme ierarhice sau modulare mari (de ex. rețele cu populații cuplate de neuroni, coloane corticale, etc.), printr-un sistem de cuplare și compoziție bine definit. Asta ajută la modelarea rețelelor mari și complexe, și favorizează re folosirea codului sursă, crearea de prototipuri, și folosirea nivelelor diferite de modelare. Au fost implementate instrumente adiționale în Mvaspike pentru a lua în considerare de exemplu întârzierile de propagare sinaptică sau axonală, precum și arhitectura rețelelor structurate sau conectate aleator într-un mod eficient pentru a descrie conectivitatea (Rochel & Martinez, 2003).

## 2.6. Implementări software proprii, validarea modelelor dezvoltate

### 2.6.1. Idei de bază

Programul de simulare dezvoltat are la bază modelul de reacție la impuls prezentat anterior, aplicând codarea cu un impuls pe neuron.

Fiind vorba de o simulare într-un mediu digital (realizată pe un calculator personal), a fost necesară dezvoltarea unui model digital adecvat, care să aproximeze cu o precizie suficient de bună modelul original și în același timp, să exploateze la maxim capacitatea sistemului de calcul utilizat. În consecință, trebuie găsită o soluție pentru atingerea preciziei dorite a parametrilor și variabilelor utilizate, precum și pentru construirea unei simulări bazate pe divizarea timpului în unități de timp corespunzător alese. Această divizare temporală necesită o acordare sensibilă, pentru a asigura precizia temporizărilor impulsurilor de activare.

### 2.6.2. Algoritmul simulării

Scopul principal al procedurilor ce compun mediul de simulare este calcularea acelor parametri ai rețelei care sunt necesari pentru a genera fiecare IA într-un moment dat al procesării continue a rețelei. Impulsurile de activare sunt calculate cu ajutorul potențialelor de membrană ale neuronilor, respectiv utilizând funcțiile de prag aferente. Este necesar, deci, să prelucrăm în cursul simulării impulsurile pre- și post-sinaptice, calculând cu ajutorul lor potențiale de membrană valide pentru fiecare neuron, determinând care dintre aceștia va emite un IA în pasul temporal următor. Totodată trebuie procesate ponderile și decalările temporale valide pentru sinapse. Toți acești parametri sunt stocați într-o structură definită în memoria calculatorului, structură ce definește starea rețelei neuronale.

Datorită complexității modelelor de NP, algoritmi de calcul tradiționali (ai rețelelor bazate pe modelul perceptronului) care utilizează matrici și vectori de date nu sunt utilizabili. În acest caz neuronii (somele acestora) și sinapsele dintre ei necesită calcule separate. Deoarece într-un calculator personal nu avem un număr de unități de calcul separate necesare pentru a efectua aceste calcule în paralel, suntem nevoiți să recurgem – cel puțin parțial – la prelucrare secvențială de date.

Este necesar, însă, ca aceste prelucrări secvențiale, să garanteze rezultate temporale consecutive. Pentru a realiza acest lucru, am utilizat procedura de divizare a fiecărui pas temporal în două faze: în prima fază se calculează IA al tuturor neuronilor din populația rețelei, IA ce vor fi utilizate în faza a doua pentru a

---

<sup>5</sup> Spre deosebire de sistemele discrete în timp, în care schimbările de stare se produc în mod periodic, și sistemele continue unde schimbările de stare se produc încontinuu.



efectua operațiile de transmitere a acestora conform structurii rețelei.

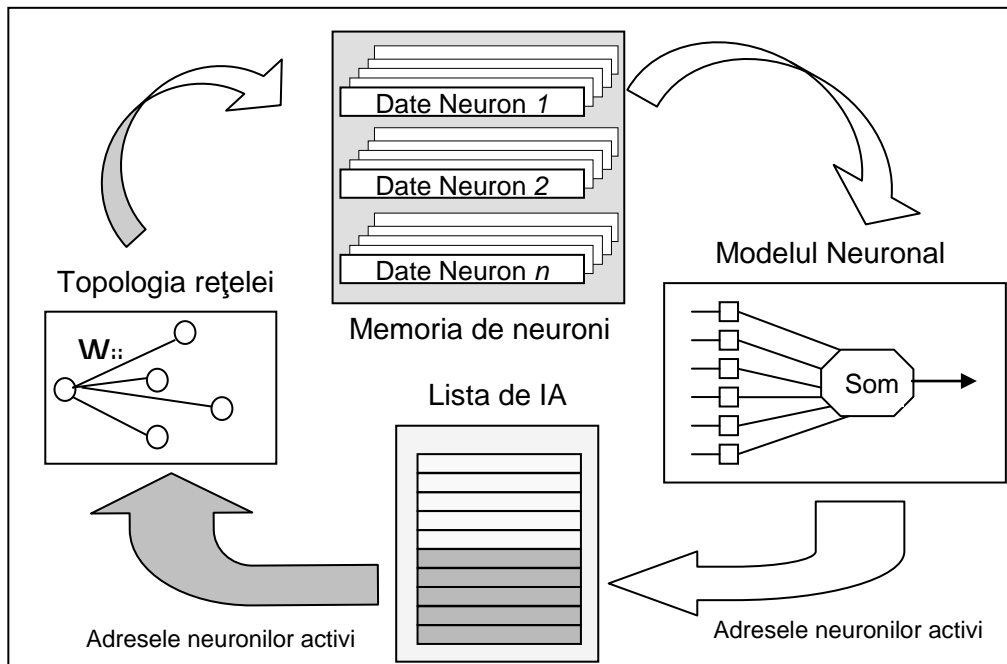


Figura 2.14 Schema algoritmului de simulare

Această divizare duce la realizarea unui algoritm trivial Figura 2.14:

- În primul pas calculez potențialele de membrană și de prag pe baza valorilor citite din memoria de neuroni. Dacă o valoare de potențial de membrană se situează peste valoarea de prag aferentă, atunci se produce un IA care se salvează în lista de IA, împreună cu codul de identificare al neuronului emitent. Toate acestea se efectuează urmând cu strictețe cerințele modelului neuronal utilizat.
- În pasul al doilea, valorile  $w_{ij}$  sunt utilizate pentru a pondera IA citite din lista în care au fost salvate, urmând a fi transmise spre neuronii post-sinaptici, ținând cont de topologia rețelei construite. În acest moment se poate aplica algoritmul de învățare dorit, modificând ponderile și întârzierile sinapselor.

Acum starea rețelei este complet determinată, poate începe următorul pas temporal.

Viteza de calcul a parametrilor neuronali și sinaptici se poate ridica – se poate ajunge chiar și la calcul în timp real – dacă utilizăm prelucrare paralelă sau aplicăm tehnici pipe-line (Hopfield, 1995). De exemplu, într-un sistem multiprocesor, rolul unei unități centrale de prelucrare (CPU) ar fi prelucrarea parametrilor neuronali, iar un al doilea procesor ar rula procedurile de calcul a parametrilor sinaptici. Dacă nu avem la îndemână un astfel de calculator, atunci putem utiliza o unitate specială de prelucrare, foarte rapidă, care utilizează tehnica pipe-line (alimentare continuă cu date) și reușește să efectueze aceste calcule în doar câteva cicluri de tact.

### 2.6.3. Încorporarea algoritmilor de învățare în procesul de simulare

Fiind vorba de rețele neuronale artificiale, ne referim la algoritmi de învățare ca la mecanisme de modificare a intensității conexiunilor dintre neuroni în vederea optimizării funcționării rețelei. În cazul de față însă, lucrul cu rețele neuronale bazate pe NP ne determină să ne concentrăm asupra procedurilor de învățare observate în rețelele neuronale naturale. Cercetări neurobiologice (Trappenberg, 2002) au scos la iveală mai multe astfel de mecanisme, care modifică eficacitatea sinapselor de la o perioadă de timp foarte scurtă (de ordinul milisecundelor) până la modificări care determină schimbări structurale durabile.

Proiectând un neuro-calculator bazat pe codare prin impulsuri, avem nevoie de un model de acordare a sinapselor care realizează – cel puțin parțial – fenomenele naturale înglobând în același timp caracteristici considerate importante de către autor. Scopul principal este de a găsi o cale de mijloc, care va conduce la un sistem aplicabil în diferite probleme de control dar care să fie ușor implementabil în

hardware, cu ajutorul circuitelor reconfigurabile FPGA. Aceste considerații sunt foarte importante, deoarece în cazul unor probleme complexe de control automat, dimensiunile rețelei neuronale pot crește semnificativ, rezultând o sarcină de calcul dificil de rezolvat.

În această ordine de idei am recurs la utilizarea *algoritmului de învățare de tip Hebb*. În rețele neuronale neuromorfe, algoritmul Hebb se poate considera ca o colecție de reguli care ajustează eficacitatea sinapselor în funcție de activitatea neuronilor pre- și post-sinaptici. La modul general, învățarea de tip Hebb se poate exprima cu următoarele patru componente:

1. O conexiune neuronală este întărită (creșterea eficienței sinaptice), dacă neuronii pre- și post-sinaptici sunt simultan activi (emit IA).
2. O conexiune neuronală este slăbită, (scăderea eficienței sinaptice), dacă numai neuronul pre-sinaptic este activ.
3. O conexiune neuronală este slăbită, (scăderea eficienței sinaptice), dacă numai neuronul post-sinaptic este activ.
4. O conexiune neuronală rămâne neschimbată, dacă nici unul dintre neuronii între care se află sinapsa nu este activ.

Modificările eficiențelor sinaptice se realizează exclusiv pe baza parametrilor locali, parametrii globali, ca activitatea medie a rețelei, nu influențează acest proces. Nu există funcționare *corectă* sau *incorectă* a rețelei, deci nu există un *expert*, care să judece acest lucru. În consecință, algoritmul Hebb implementează o învățare nesupervizată.

Ca rezultat al cercetărilor din acest domeniu au apărut mai multe variante metodei de învățare Hebb.

$$\Delta w_{ij}(t) = (x_i(t) - \bar{x}_i)(x_j(t) - \bar{x}_j) \quad \text{Ec. 44}$$

Acordarea sinapselor este dată de corelația dintre activitățile de ieșire ale neuronilor  $x_i$  și  $x_j$ , expresie ce se poate interpreta și ca *factorul mediei de activare* a acestor neuroni.

Având ca scop realizarea RNP cu ajutorul echipamentelor hardware, o soluție mai potrivită în ceea ce privește alegerea metodei de învățare ar fi cea bazată pe observații neurofizice, și anume *excitarea pe termen lung* (ETL) respectiv *inhibare pe termen lung* (ITL). Ca scurtă descriere a metodei ETL-ITL se poate spune că eficiența unei sinapse active crește, dacă potențialul membranei neuronului post-sinaptic se depolarizează, respectiv scade în caz contrar. Primul caz este cel de ETL iar cel din urmă ITL. Ambele cazuri se pot compara cu punctele 1. respectiv 2. din descrierea metodei Hebb. Punctul 3. din descrierea metodei Hebb lipsește însă la analiza metodei ETL-ITL, deoarece aici luăm în considerare doar sinapsele active.

Pentru a putea dezvolta un mediu de simulare flexibil, adăugarea acestui punct 3. ar fi totuși, necesară în pofida faptului, că acest lucru nu se conformează cu rezultatele experimentale din literatură (Földiák, 1990), (Natschläger & Ruf, 1998), (Bohte, Kok, & La Poutré, 2002). Autorii acestor lucrări au arătat, că utilizarea algoritmului Hebb în rețele neuronale cu codare al factorului mediu de activare, poate duce la apariția unor efecte neuromorfe de genul auto-organizării. Este demonstrat, de asemenea, că în mediul natural, sinapsele care nu influențează pe termen lung activitatea post-sinaptică vor înceta să existe.

Studiind în amănunțime rezultatele experimentale din (Natschläger & Ruf, 1998), (Bohte, Kok, & La Poutré, 2002) în ceea ce privește ETL-ITL, putem trage două concluzii importante:

1. O anumită valoare de prag a depolarizării post-sinaptice este aceea care determină activarea ETL sau a ITL.
2. “Retro-propagarea” IA post-sinaptice creează un comportament de învățare, în funcție de decalajul temporal dintre IA pre- și post-sinaptice.

Pe baza acestor rezultate și considerente am dezvoltat o strategie de învățare, care va fi prezentată în cele ce urmează.

### 2.6.4. Învățare bazată pe valori de prag

Modificarea eficacității sinapselor este funcția măsurii de depolarizare a potențialului membranei neuronului post-sinaptic. (Între valorile de prag  $\mathcal{G}_{ITL}$  și  $\mathcal{G}_{ETL}$  eficacitatea sinaptică scade.) Dacă potențialul membranei neuronului post-sinaptic se află în acest interval, atunci are loc ITL [14]. Dacă acest potențial depășește  $\mathcal{G}_{ETL}$ , atunci se va activa ETL. Comparând acest model cu modelul de învățare Hebb, putem spune, că aici activitatea pre-sinaptică este dată de emiterea unui IA pre-sinaptic, iar activitatea post-sinaptică este dată de valoarea potențialului de membrană (Figura 2.15).

Deoarece această metodă descrie numai felul modificării sinapselor active, un impuls pre-sinaptic se poate percepe ca un semnal de start al procesului de învățare. Acest fapt îndeplinește cerința unei simulări eficiente, deoarece frecvența IA nu este prea mare. Deci, într-un anumit pas al simulării nu este necesar a **calcula**, numai parametrii unei fracțiuni din populația de neuroni ai rețelei.

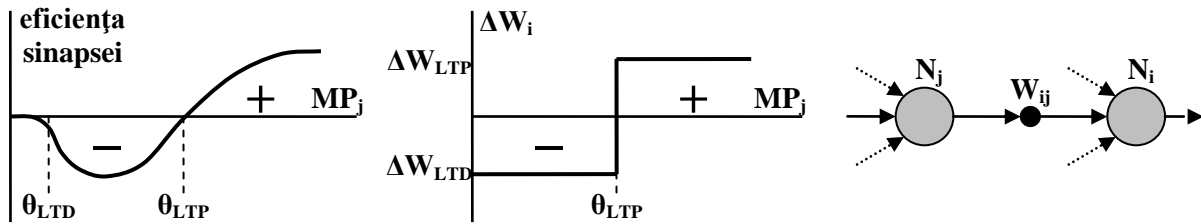


Figura 2.15 Regula de învățare bazată pe valori de prag

( $\Delta w_{ij}$  variația ponderii sinaptice,  $MP_i$  potențialul membranei postsinaptice)

A –modelul biologic, B – modelul adaptat

Având în vedere necesitatea implementării hardware, este posibilă o simplificare adițională, setând valoarea de prag  $\mathcal{G}_{ITL}$  la zero. Procesul de învățare rezultat se poate exprima în pseudo-cod, cu algoritmul următor:

```

if neuronul pre-sinaptic  $j$  emite IA
    if  $MP_i > \mathcal{G}_{ETL}$ 
        incrementează  $w_{ij}$ 
    else
        decrementează  $w_{ij}$ ,
    
```

unde  $w_{ij}$  este ponderea conexiunii dintre neuronii  $j$  și  $i$ , iar  $MP_i$  este potențialul membranei post-sinaptice.

Acest algoritm, însă nu descrie numai punctele 1. și 2. din regula de acordare sinaptică Hebb. Punctul 3. care se referă la cazul în care numai neuronul post-sinaptic este activ lipsește, pentru moment. Pentru a completa algoritmul, va trebui să adăugăm următoarele rânduri:

```

if  $MP_i > \mathcal{G}_{ETL}$ 
    if neuronul pre-sinaptic  $j$  nu emite IA
        decrementează  $w_{ij}$ .
    
```

Deci, acest proces de învățare se activează numai în cazul în care potențialul de membrană post-sinaptic depășește  $\mathcal{G}_{ETL}$ , eveniment care este mult mai frecvent, decât apariția unui IA.

## 2.7. Rezultate experimentale a simulărilor software

Datorită complexității problemei, am acordat o atenție deosebită alegerii mediului de dezvoltare a programului de simulare. După cum s-a văzut în capitolele anterioare, în rețelele neuronale naturale precum și în cele artificiale neuromorfe, componenta temporală joacă un rol extrem de important. Era nevoie, deci de o soluție care să asigure temporizări precise și previzibile concomitent cu un mediu de programare de nivel înalt, corespunzător sarcinii de rezolvat. În această ordine de idei, și după câteva experimente efectuate, am renunțat la calea utilizării sistemului de operare Windows și a unui mediu interpretor (de ex. Matlab) sau compilator (de ex. Visual C++). Am rămas la limbajul de programare C++, acesta oferind mijloacele cele mai dinamice, însă am trecut la utilizarea sistemului de operare RedHat Linux (respectiv Fedora Core 3). Pe lângă suportul substanțial mai bun al temporizărilor în comparație cu Windows, în Linux am beneficiat și de alte avantaje

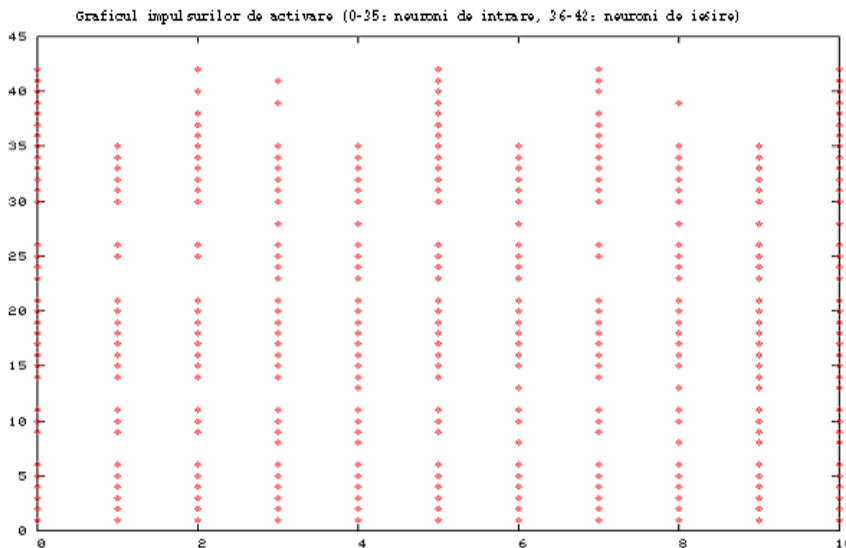


Figura 2.18 IA în timpul funcționării rețelei neuronale prezentate în Figura 2.16.

componentele sistemului de simulare dezvoltat ele ar fi: modelul neuronal implementat, modelul sinaptic implementat, modulele de construcție a straturilor respectiv cele ce realizează conexiunile dintre straturi. Orice configurație reprezentată în 2D sau 3D este realizabilă, construind rețeaua neuron cu neuron sau strat cu strat. Am implementat și posibilitatea de realizare de reacții, astfel putându-se construi rețele compuse din câteva straturi de NP, al căror ieșiri sunt conectate la intrările celui de-al doilea strat de neuroni. Primul strat este format din neuroni speciali, fără sinapse care posedă numai câte o singură intrare. Putem ajunge astfel la structuri care vor emite impulsuri de ieșire la intervale regulate, ele fiind adecvate pentru utilizare ca generator de tact sau pentru sincronizarea altor grupuri de neuroni.

Programul utilizează, deci, trei tipuri de neuroni, și anume:

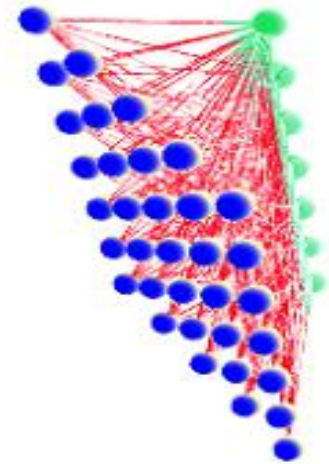


Figura 2.16 O rețea neuronală realizată în mediul de simulare dezvoltat. Aplicația: recunoaștere de caractere dintr-o matrice de 7x5 puncte (albastru – neuroni de intrare, verde – neuroni de ieșire, roșu – sinapse)

precum reprezentarea grafică 3D mai ușoară a rețelelor neuronale realizate (pachete grafice Povray, Gnuplot) sau salvări de structură a rețelelor în format modern XML.

Dezvoltarea de software a urmat regulile impuse de *modelul de reacție la impuls*, descris la subpunctul 2.4.2, respectiv învățarea bazată pe valori de prag. Trecând în revistă pe scurt

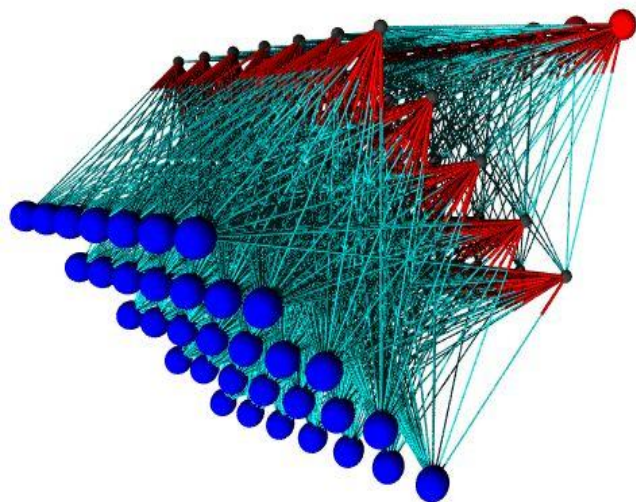


Figura 2.17 Aceași problemă, rezolvată cu o rețea cu un strat ascuns

- neuroni de intrare (pot ocupa loc numai în stratul de intrare)
- neuroni ascunși (neuronii straturilor interioare)
- neuroni de ieșire (neuronii stratului de ieșire)

Semnalele de intrare și de ieșire ale rețelelor de NP construite sunt considerate evenimente ale căror temporizări se pot regla, ele fiind de fapt valori binare - 1 și 0 logic.

În Figura 2.17 este prezentată o RNP cu un strat ascuns, care este capabilă să categorizeze caractere prezentate chiar și sub formă deteriorată.

Cei șapte neuroni de ieșire corespund segmentelor unui afișor cu 7 segmente. Ieșirile date de rețea la intrările zgomotoase date la intrare se pot urmări în Figura 4.6, luând în considerare, că pe axa X s-a reprezentat numărul pasului temporal în care s-a efectuat măsurătoarea respectiv, pe axa Y identificatoarele neuronilor din rețea. Identificatorii neuronilor sunt: (ID's): 0\_34 strat de intrare, 34\_69 strat ascuns, 70\_76 strat de ieșire.

În fiecare pas temporal la intrarea rețelei este prezentat un nou vector de intrare, după care sunt efectuate calculele sinaptice aferente, care au ca rezultat acordarea ponderilor sinapselor rețelei.

Convergența procesului de învățare este confirmată de faptul că, după un număr suficient de mare de pași temporali efectuați, pentru un model anume de intrare se va activa numai un anumit grup de neuroni de ieșire. În Figura 2.19 au fost reprezentate primii treizeci de pași temporali din cei o mie ai simulării în care s-a efectuat măsurătoarea.

După ~800 de cicluri de învățare, caracterele de la 0 la 9 au fost identificate corect în proporție de 95%.

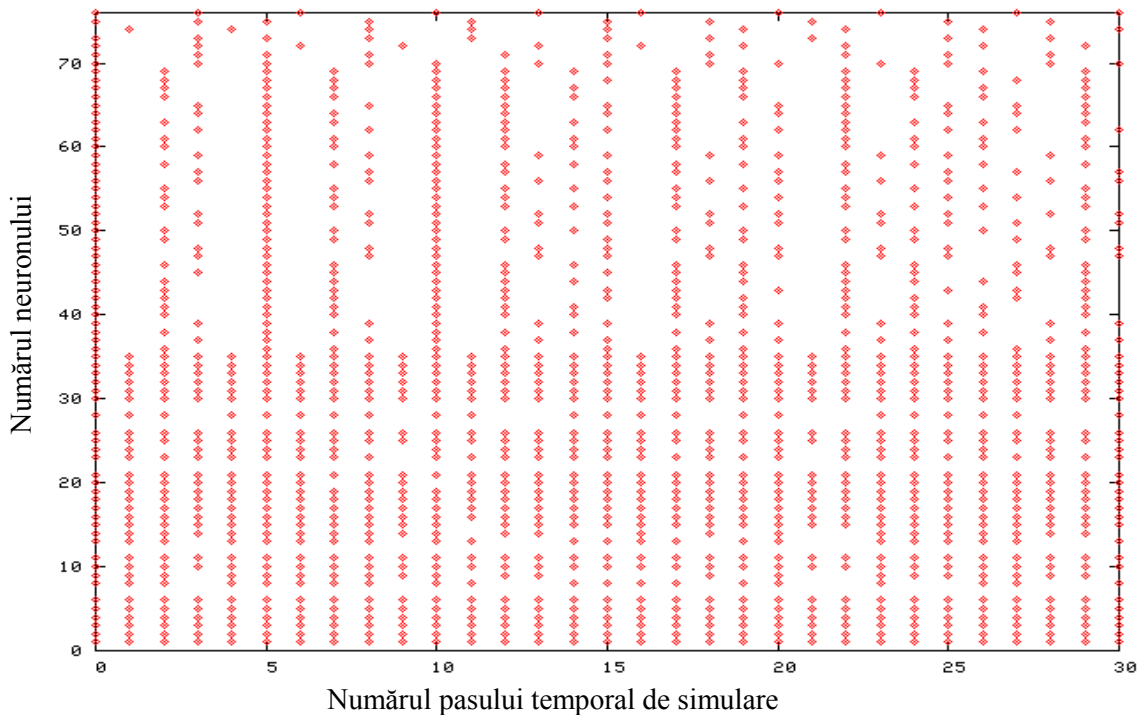


Figura 2.19 Rezultatul măsurătorii efectuate pe ieșirea RNP din Figura 2.17. Primii treizeci de pași temporali

## 2.8. Concluzii

În prima parte a acestui capitol s-a discutat modelarea rețelelor neuronale neuromorfe, expunând comportamentul intrinsec temporal al acestora, utilizarea impulsurilor în comunicarea datelor în cadrul acestora și evidențiind capacitatea de calcul mai ridicată decât în cazul rețelelor neuronale bazate pe modele clasice.

Studiul neuronilor biologici a rezultat în apariția mai multor ipoteze pentru codarea valorilor reale în trenuri de impulsuri. Acestea pot apărea ca o codare a ratei impulsurilor, care ia în considerare

numai valoarea medie a timpului de apariție a unui impuls sau ca codarea exactă a momentelor temporale de apariție a acestor impulsuri.

Când vorbim despre RNP există de fapt o supraabundență de diferite modele, de la cele simple de tipul integrează și activează prin neuroni cu diferite modele sinaptice până la modele comportamentale complete descrise prin ecuații diferențiale. Există câteva studii teoretice (Maass, 1996) care arată că unele dintre aceste modele sunt mai eficiente decât celelalte, dar aceste studii se limitează doar la un set redus de modele și scheme de codare. O explicație generală a capacității de calcul a diferitelor modele neuronale pulsative și a schemelor de codare conexe reprezintă și astăzi o temă deschisă cercetărilor actuale.

În orice studiu legat de dinamica unei rețele neuronale pulsative, sunt două probleme specifice, și anume ce model descrie dinamica pulsației fiecărui neuron respectiv cum sunt conectați neuronii. O alegere greșită a modelului sau o conexiune greșită, poate duce la rezultate complet eronate.

În cadrul acestui capitol s-a elaborat prima problemă, comparând neuronii pulsatori. Se vor prezenta diferite modele de neuroni pulsatori precum și o clasificare a acestora.

În continuare am prezentat cele mai des utilizate metode de antrenare (ne-supervizate și supervizate) utilizate în învățarea rețelelor neuronale neuromorfe. Astfel, în această parte a capitolului se regăsesc discuții asupra realizării învățării Hebbiene competitive precum și asupra modalităților specifice de aplicare a regulii back-propagation în rețele neuronale pulsative codificate temporal (algoritmul SpikeProp).

După acest studiu amănunțit al acestui domeniu de inteligență artificială prezintă modelul neuronal teoretic ales pentru implementările hardware ce urmează a fi prezentate în capitolele următoare.

A doua contribuție importantă a acestui capitol este prezentare generală a strategiilor de simulare software a RNP împreună cu descrierea celor mai importante astfel de medii existente la ora actuală, dar și a simulatorului neuronal implementare proprie, cu care am testat funcționarea modelului dezvoltat. Se detaliază încorporarea algoritmilor de învățare în procesul de simulare software și se expun rezultatele experimentale obținute pe această cale, care au stat la baza primelor publicații ale autorului (Bakó, Székely, & Brassai, Development of Advanced Neural Models. Software And Hardware Implementation, 2004) (Bakó, Analiza și simularea sistemelor cu inteligență artificială neuromorfă, 2005) (Bakó, Székely, Dávid, & Brassai, 2004), (Bakó & Brassai, 2004) în acest domeniu de cercetare.

### 3. DISPOZITIVE FPGA ȘI IMPLEMENTAREA REȚELELOR NEURONALE ARTIFICIALE

În acest capitol sunt prezentate sistemele FPGA și modalitățile de implementare ale rețelelor neuronale pe aceste sisteme. Urmează să fie examinată evoluția sistemelor hardware reconfigurabile, în special a sistemelor FPGA.

În continuare vom descrie sistemul de dezvoltare FPGA Xilinx XC3S1000 utilizat la implementarea rețelelor neuronale. Sunt prezentate elementele de bază ale sistemului FPGA, și accentul se va pune în special asupra modulelor de memorii Block RAM, module care stau de fapt la baza structurilor de rețele neuronale neuromorfe care vor fi prezentate în Capitolul 4.

Sunt prezentate moduri de abordare ale cuplării și configurării pentru sisteme FPGA care să funcționeze în combinație cu procesorul gazdă. De asemenea sunt introduse cele două moduri de implementare a logicii reconfigurabile: reconfigurarea compilată în timp și reconfigurarea în timpul utilizării.

În următoarea parte a capitolului se discută despre maparea diferiților algoritmi de rețele neuronale pe sisteme FPGA și sunt prezentate diferite sisteme de rețele neuronale hardware la baza cărora se află sisteme FPGA. Sunt trecute în revistă câteva sisteme, care se pot considera de bază, cum sunt RRANN, GANGLION, sistemele REMAP, RAPTOR2000 și de asemenea sunt referiri la câteva teze de doctorat și lucrări științifice publicate în ultimii ani.

Modalitatea de implementare a funcțiilor de activare respectiv alegerea unei aritmetici corespunzătoare joacă un rol foarte important în implementarea rețelelor neuronale pe sisteme hardware. Ultima parte a Capitolului prezintă aspecte cu privire la modalitățile de implementare a funcțiilor de activare și aritmetica neuronală utilizată pentru reprezentarea datelor.

Sistemele FPGA sunt circuite integrate construite din configurații de blocuri logice (CLB, configuration logic blocks) interconectate printr-o rețea de conexiuni. Atât blocurile logice cât și interconexiunile sunt configurabile (Bondalpati, 2001). Aceste dispozitive pot fi configurate și reconfigurate de proiectant printr-un software CAD. Evoluția sistemelor FPGA seamănă mai mult cu evoluția structurilor memoriilor ROM. Prima dată au apărut structurile FPGA bazate pe elemente fuzibile analogice structurilor PROM. Alte structuri bazate pe elemente fuzibile au fost Programmable Logic Array (PLA) și Programmable Array Logic (PAL).

La modul general, toate aceste structuri bazate pe elemente fuzibile sunt clasificate sub denumirea de dispozitive logice programabile (PLD Programmable Logic Device). Limitările acestor structuri PLD au dus la apariția structurilor FPGA bazate pe memorii RAM. Structurile FPGA disponibile în prezent sunt de acest fel, care permit realizarea unei configurații prin programarea memoriei RAM a blocurilor CLB. Structurile FPGA au început să fie disponibile la sfârșitul anilor 80 și la începutul anilor 90.

Din acel moment cercetătorii și proiectanții au ajuns la concluzia că este posibilă realizarea sistemelor de calcul reconfigurabile prin sisteme FPGA. Arhitectura reconfigurabilă a început să-și arate fezabilitatea însemnată de aplicare numai recent, dar ideile de bază s-au conturat aproape cu ideea calculatorului programabil cu scop general. În anul 1966, János Neumann, recunoscut ca părintele modelului serial al calculatoarelor programabile, a conceput de asemenea un automat de calcul, compus dintr-o rețea simplă de celule, care pot fi configurate pentru executarea unor operații computaționale (Dehon, 2000).

Minnick în 1971 a descris o matrice celulară programabilă utilizând bistabile pentru a memora un context de configurație.

Posibilitatea de utilizare eficientă a sistemelor de hardware reconfigurabil realizat într-un FPGA a apărut când era deja posibilă introducerea sutelor de elemente programabile pe un singur chip.

După cum considera Moore “numărul transistoarelor într-un circuit integrat se va dubla aproximativ în 18 luni”. Astăzi este posibil să utilizăm sisteme FPGA cu peste 10 milioane de porți într-un singur chip.

În general, pentru a implementa un circuit dorit, sistemele FPGA sunt reconfigurate de un controller extern. Reconfigurarea se poate realiza cu o viteză corespunzătoare, ce se apropie de o funcționare în timp real, depinzând de timpul necesar pentru reconfigurare respectiv timpul necesar pentru încărcarea datelor de configurare.

Ca și la toate sistemele VLSI, densitatea și viteza în creștere a circuitelor pe bază de siliciu este limitată la modulele de comunicație și de la această regulă nici sistemele FPGA nu constituie o excepție. Hardware-ul reconfigurabil poate fi clasificat în una din următoarele categorii, presupunând că un astfel de dispozitiv are câteva nivele proprii de configurabilitate:

*Hardware-ul configurabil* poate fi configurat o dată sau de două ori.

*Hardware-ul reconfigurabil* – are avantajul că poate fi reconfigurat de mai multe ori conform funcțiilor pe care trebuie să le îndeplinească. Aceste sisteme conțin elemente ca Dispozitive Logice Programabile (PLA), care se comportă ca un tabel de adevăr, funcții Booleene, automat de stare, sau memorii PROM, care pot fi utilizate ca memorii nevolatice pentru microprocesoare și microcontrolere. Câteva variante ale sistemelor reconfigurabile sunt *in-system programmable* (ISP) și pot fi reprogramate în timp ce sunt încă rezidente pe placa de circuit.

*Hardware-ul dinamic reconfigurabil* a apărut odată cu apariția sistemelor FPGA bazate pe memorii SRAM. Celulele SRAM permit configurarea logicii unui sistem FPGA prin simpla încărcare a unui șir de biți în modulul SRAM. Deci acestea sunt sisteme care poate fi rapid reprogramate pentru realizarea unei sarcini. De exemplu la pornire, dispozitivul FPGA poate fi configurat pentru a realiza o diagnosticare a propriilor circuite înainte de reconfigurarea dinamică pentru realizarea sarcinii principale pentru care a fost proiectat.

*Hardware-ul parțial reconfigurabil*. La majoritatea sistemelor hardware dinamic reconfigurabile problema care apare este că, în timpul reconfigurării, operațiile sale trebuie întrerupte, și conținutul memoriilor de configurare SRAM reîncărcate, rezultând o pierdere irecuperabilă de date din registrele sistemului FPGA. A devenit astfel necesară dezvoltarea unui nou tip de dispozitiv FPGA, care permite reconfigurarea dinamică a unei porțiuni a logicii interne fără deconectarea conexiunilor I/O sau a semnalelor de tact de sistem, și care asigură funcționarea continuă a operațiilor dispozitivului în porțiunile care nu sunt supuse reconfigurării. O caracteristică care oferă un avantaj important este conținutul registrelor interne care nu se pierd în timpul reconfigurării. Astfel rezultatul unei instanțieri poate fi preluat de instanțierea dinamic reconfigurată.

### 3.1. Evoluția dispozitivelor reconfigurabile

Există mai multe metode aflate la dispoziția proiectanților de sisteme reconfigurabile pentru integrarea diferitelor componente discrete într-un singur chip VLSI ASIC (Hauck, 1995), fiecare metodă având o capacitate specifică, performanță, flexibilitate și cost unitar atât din punctul de vedere al timpului, cât și financiar. Acest subcapitol evidențiază funcțiile diferitelor tehnologii, subliniind noile oportunități oferite de folosirea FPGA-urilor.

Un *sistem full-custom* este acel sistem la care întregul circuit este construit după cerințele proiectantului. Rezultatul este un chip cu performanțe foarte bune, cu o utilizare optimă a substratului de siliciu și este câteodată singura metodă de implementare luată în calcul în cazul unor aplicații speciale, cum ar fi microprocesoarele realizate cu tehnologii de vârf. Costurile de fabricație și timpul inițial de punere în funcțiune sunt de asemenea mari, deoarece flexibilitatea oferită proiectantului rezultă într-un design al chipului foarte complex și care necesită mult timp.

O proiectare pe bază *de celule* sacrifică flexibilitatea și performanța unui sistem full-custom pentru a accelera procesul de proiectare. Aceasta se realizează prin folosirea unui set de restricții de proiectare și celulele de format standard îngăduind folosirea instrumentelor soft la automatizarea procesului de proiectare.

Problema cu ambele metode, full custom și pe bază de celule, este faptul că ele trebuie începute de la zero. În cazul MPGA-rilor (Mask Programmable Gate Array), marea parte a procesului de fabricație se realizează mai înainte. De exemplu, un MPGA poate consta din seturi de transistoare în locații specifice care sunt parțial interconectate pentru a forma baza mapării unor tipuri specifice de matrici de porți. Astfel, tot ce trebuie să facă un proiectant este interconectarea elementelor necesare pentru a atinge funcționalitatea dorită. Această metodă reduce atât costul, cât și timpul de producție, deoarece o astfel de bază poate stoca cantități mari de proiect standard, chipuri parțial realizate. Proiectantul are o flexibilitate destul de mare și este limitat doar de rutările aflate la dispoziție pe chip. Din păcate, este inevitabil să nu fie anumite ineficiențe, deoarece unele proiecte necesită resurse specifice sau lățime diferită a benzii de rutare, decât cele care există pe chip.



Unul dintre primele chipuri de uz general prefabricate care au ajuns să fie folosite la scară largă a fost memoria read only programabilă (PROM). Ea constă dintr-un șir de celule preprogramabile în care prima dată se scriu datele. Informația de pe un PROM se stochează folosind una dintre următoarele metode:

- *Fuses/antifuses*- se poate folosi pentru a stoca biți, condiția cărora va determina care valoare de biți se va citi. Ele se activează prin aplicarea unei tensiuni de programare.
- *Erasable PROM (EPROM)*- ele sunt programate prin folosirea unei tensiuni mai mari decât tensiunea de alimentare, dar se deosebesc de prima categorie prin faptul că ele se pot șterge prin expunerea la lumină ultravioletă (UV). Acest lucru permite ca EPROM să fie reprogramabil.
- *Electrically erasable PROMs (EEPROMs)*(PROM care se poate șterge pe cale electrică)-sunt programate prin folosirea tensiunii înalte, dar în acest caz chipul se poate șterge prin metode pur electrice. O altă tehnologie care merită amintită, dar care nu este clasificată strict ca și PROM, este memoria cu acces aleator (*RAM*). Caracteristicile sale sunt identice cu cele ale PROM-ului, cu deosebirea că acesta se poate citi și scrie, și astfel, reprograma în timp ce sistemul este în funcțiune.
- *Dispozitivele logice programabile (PLD)* au fost proiectate pentru a implementa circuite logice și constau dintr-un șir de porți AND urmat de un șir de porți OR într-o sumă-de-produse.
- *Programmable array logic (PAL)* este cel mai larg răspândit tip de PLD cu un plan AND programabil urmat de un plan OR fixat.
- *Programmable logic arrays (PLA)* constituie un tip mai flexibil al PLD-urilor, care permit atât planului AND, cât și OR să fie programabile.
- *Generic array logic (GAL)* constituie un tip mai avansat de PLD care se poate configura pentru implementarea diferitelor tipuri de PAL cu inversia ieșirii opțională.
- *Complex programmable logic devices (CPLDs)* sunt sisteme complexe, care conectează mai multe elemente PLD folosind o matrice de comutare.
- PLD-urile se pot implementa prin folosirea tehnicilor de programare fuse/antifuse, EPROM sau EEPROM.
- *Field programmable gate arrays (FPGA)* sunt dispozitive prefabricate construite pentru implementarea unor circuite multistrat în loc de termeni sumă-a-produsului simple PLD, singura limitare fiind resursele aflate la dispoziție în circuitele chipului. Din păcate, efectul acestei tehnologii este apariția unei întârzieri a propagării care nu se poate estima. Trebuie făcută mențiunea că în unele cazuri se referă la FPGA ca fiind dispozitive CPLD, dar cu toate acestea noi ne vom referii la ele ca FPGA-uri, pentru a evita ambiguitățile. FPGA-urile nu sunt disponibile doar prin tehnologiile fuse/antifuse, EPROM sau EEPROM; dar și în SRAM sau RAM dinamic (DRAM) (Motomura, Aimoto, Shibayama, Yabe, & Yamashina, 1998). Versiunile SRAM și DRAM permit reconfigurarea rapidă “in-circuit” a chipului, dar sunt dezavantajate din punctul de vedere al substratului de siliciu necesar pentru tehnologia RAM.

Reconfigurarea poate fi parțială sau completă, ceea ce depinde de dispozitiv, și nu trebuie să rezulte neapărat în distrugerea blocurilor.

FPGA-ul (Brown, Francis, Rose, & Vranesic, 1993) (Oldfield & Dorf, 1995) s-a introdus pentru prima dată în 1985 de Xilinx și de atunci multe alte companii au lansat produse similare.

### 3.1.1. Evaluarea performanței aplicațiilor bazate pe circuite FPGA

Inovațiile remarcabile din domeniul tehnologiei computaționale (Asanovic, 2006) îndeplinesc previziunile pe termen lung ale NASA (Sobieski & Storaasli, 2004) pentru calcule științifice și ingineresti mai rapide. O inovație de prim plan este valorificarea circuitelor FPGA (Field Programmable Gate Arrays) pentru accelerarea aplicațiilor de calcul de înaltă performanță (High-Performance Computing - HPC) cu unul sau mai multe ordine de mărime comparativ cu procesoarele tradiționale. FPGA, o formă complet nouă de logică programabilă, a fost inventată în 1984 de către Ross Freeman, co-fondator al Xilinx Corporation. Arhitecturile FPGA sunt extrem de flexibile și interconectează mii de module încorporate cum ar fi sumatoare, multiplicatoare, module de memorie și structuri de elemente logice utilizabile de exemplu în prelucrarea digitală a semnalelor și comunicația de mare viteză (Hypertransport, PCI-EXPRESS). Spre deosebire de microprocesoarele programabile convenționale

„fixe”, circuitele FPGA sunt reconfigurabile. Unele circuite FPGA pot fi parțial reconfigurate chiar în timp ce alte porțiuni ale aceluiași circuit funcționează.

Circuitele FPGA Virtex-4 pot avea unul sau mai multe procesoare PowerPC încorporate. Piața circuitelor FPGA de 2 miliarde de dolari, cu o creștere rapidă (15-20% / an) este dominată de Xilinx și Altera. Folosirea circuitelor FPGA se extinde rapid și în aplicații mai mici, aeronautice și de calcul încorporat de înaltă performanță (High-Performance Embedded Computing - HPEC).

### 3.1.1.1. Caracteristici FPGA

Structura circuitelor FPGA este strict reglementată în comparație cu cea a microprocesoarelor, simplificând astfel procesul de fabricație, și a permis circuitelor FPGA să fie printre primele care și-au redus dimensiunile (90nm => 65nm => 45nm). La fiecare ciclu de tact, algoritmi FPGA (când au fost codificați cu scopul de a mări la maximum numărul de operații executate în paralel) își folosesc aproape 100% din resursele de siliciu, comparativ cu microprocesoarele mai puțin eficiente care își folosesc mai puțin de 2% din resursele de siliciu (în timp ce consumă de 10 ori mai mult decât FPGA pentru a executa doar una sau două operații).

Figura 3.1 și Figura 3.2 arată câteva caracteristici cheie ale circuitelor FPGA. Spre deosebire de microprocesoare, circuitele FPGA continuă să avanseze la viteza legii lui Moore, și mai au mult până să atingă pragul limită de unități logice și viteză de tact (Figura 3.1). Viteza de tact a circuitelor FPGA (frecvent în jur de 100-200 MHz pentru multe aplicații) mai are mult de crescut până să se lovească de problemele de supraîncălzire ce au determinat microprocesoarele să adopte o structură multi-core cu chipuri rulând la viteză de tact redusă. Când aplicațiile FPGA sunt programate cu un grad ridicat de paralelism, viteza lor de calcul depășește cu mult pe cea a microprocesoarelor (Figura 3.2). Cum era de așteptat, pentru sisteme de comunicație de mare viteză, lățimea de bandă pentru memorie și I/O depășește de asemenea cu mult pe cea a microprocesoarelor (Figura 3.2).

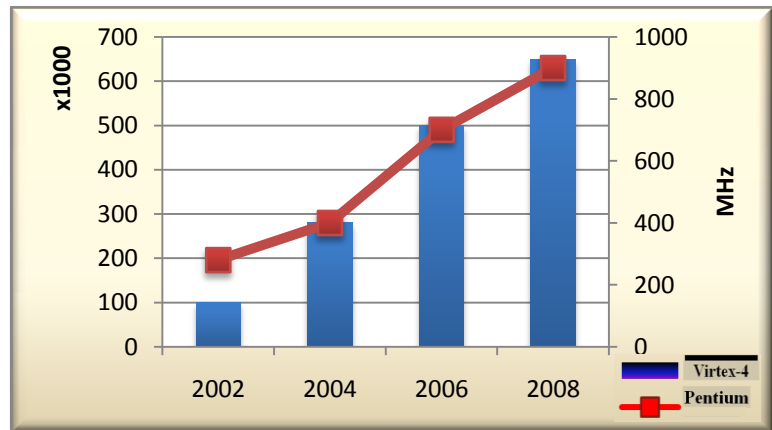


Figura 3.1 Caracteristici FPGA: evoluția creșterii unităților logice și a vitezei de tact

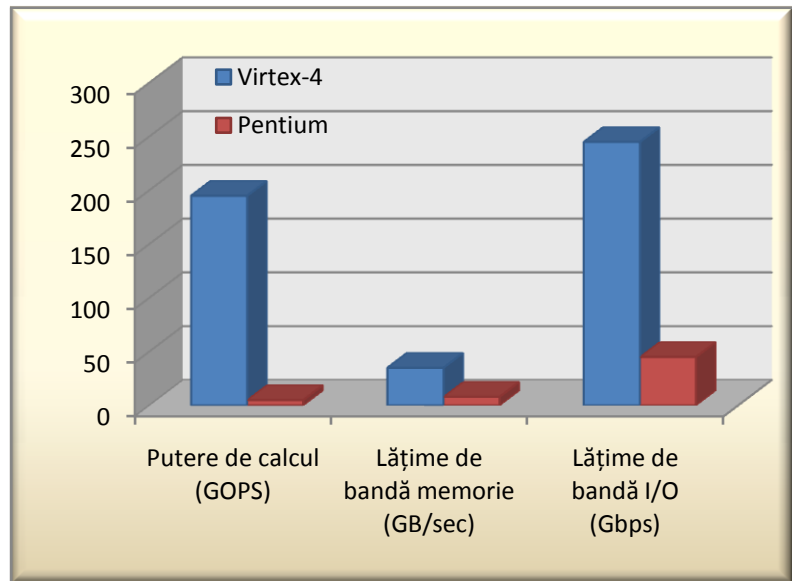


Figura 3.2 Caracteristici FPGA: evoluția creșterii vitezei de calcul și a lățimii de bandă

### 3.1.1.2. Programarea circuitelor FPGA

Cum circuitele FPGA au fost dezvoltate de proiectanți de circuite logice, se programează în mod tradițional folosind limbaje de programare de descriere hardware cum ar fi VHDL și Verilog. Aceste limbaje necesită cunoștințele și experiența unui proiectant de circuite logice, învățarea lor necesitând luni de zile. Chiar și cu experiența și cunoștințele necesare, programarea în VHDL sau Verilog este

extrem de greoaie, dezvoltarea unor prototipuri timpurii necesită luni de programare, iar perfecționarea și optimizarea chiar ani. Programarea FPGA, spre deosebire de compilatoarele HPC, este mult încetinită de pașii mari în plus necesari pentru a sintetiza un proiect, a plasa componentele și a concepe conexiunile dintre acestea.

Odată efortul depus pentru programarea de aplicații specifice în VHDL, performanța lor pe FPGA este greu de atins. În particular, aplicații care folosesc doar operații simple cu numere întregi sau operații logice (comparație, adunare, înmulțire) cum ar fi comparații de secvențe ADN, criptografie, logica jocului de șah, rulează extrem de rapid pe FPGA. Cum aplicațiile care folosesc operații de dublă precizie în virgulă mobilă au epuizat rapid numărul de structuri de elemente logice disponibile pe primele circuite FPGA, acestea au fost evitate adesea pentru calcule de înaltă precizie. Această situație s-a schimbat însă pentru ultimele modele de FPGA de la Xilinx (Virtex-4 și Virtex-5) care au suficiente unități logice pentru a acomoda în jur de 80 module de multiplicare pe 64 de biți (Strenski, 2007).

Capabilitățile primelor circuite FPGA erau foarte potrivite pentru calculul încorporat de înaltă performanță (HPEC) cu destinație specială. Folosirea lor în HPC cu destinație generală a fost inițial restrânsă la supercomputere reconfigurabile de primă generație. (de ex. Starbridge Systems, SRC, Cray XD1). Această primă generație a fost caracterizată de lipsa unui sistem de I/O de mare viteză și a unei infrastructuri (compilatoare, librării) care să faciliteze aplicații supercomputer cu destinație generală, incluzând cod deja existent. Această situație se schimbă rapid odată cu ultima generație de supercomputere reconfigurabile și circuitele FPGA pe care le folosesc. DRC Computer, Xtreme Data și Xilinx în colaborare cu Intel oferă acum pe piață module conținând cele mai noi circuite FPGA care se potrivesc în același socket folosit și de microprocesoare și folosesc aceeași linie de comunicație de mare viteză. Cray a ales circuitul DRC pentru a accelera performanțele liniei de supercomputere XT.

Alternative concurente la FPGA în domeniul HPC: Creșterea vitezei aplicațiilor HPC este atât de necesară, astfel au apărut pe piață mai multe alternative. Deși multe coduri bazate pe fizica tradițională au fost scrise în limbajul secvențial Fortran acum mai bine de 30 de ani, au supraviețuit remarcabil de bine mai multor generații de HPC: vector (via compilatoare), paralel (via MPI, OpenMP) și acum primele generații de microprocesoare multi-core. Despre unele se presupune că ar putea suferi degradări masive ale performanței, sau chiar va fi nevoie să fie rescrise pentru a folosi la maximum capacitatea a 8 sau mai multe „core”-uri / chip. Marii producători de chip-uri (Intel și AMD) au depus eforturi substanțiale pentru a satisface cerințele pieței cu acceleratoare, concentrându-se în principal pe circuite FPGA ca un mod de a recâștiga performanță. Cum microprocesoarele multi-core înfruntă probleme cu consumul de curent, răcirea, dimensiunile și numărul porturilor de I/O, circuitele FPGA par mai atractive.

Alte trei tipuri de acceleratoare sunt disponibile arhitecților HPC: Cell (IBM), Array (ClearSpeed) și procesoare grafice (Graphical Processors - GPUs). Ca și circuitele FPGA, procesoarele grafice și Cell au o piață vastă de desfacere (grafică și jocuri video), fapt ce duce la scăderea costurilor, promovarea concurenței și încurajarea dezvoltării în domeniu, făcându-le din ce în ce mai atractive pentru HPC. În schimb procesoarele Array sunt dispozitive specializate care trebuie amortizate pe un număr relativ redus de utilizatori. Procesoarele GPU necesită multă putere/răcire și au probleme complexe de programare și precizie a datelor de rezolvat înainte de a putea intra pe piața HPC. Programarea procesoarelor Cell 8+1 e probabil să fie considerată mai greoaie decât programarea circuitelor FPGA în VHDL sau Verilog, care are deja o bază mare de utilizatori. Odată cu dezvoltarea hardware-ului circuitelor FPGA, software-ul și instrumentele FPGA devin mai ușor de folosit pentru HPC, incluzând spre exemplu MitrionC, Viva, DSPlogic, ImpulseC, Celoxica, Aldec, și CHiMPS de la Xilinx, toate având ca țintă specifică piața HPC.

### 3.1.1.3. Lățimea de bandă și localizarea datelor

Performanța multor algoritmi este limitată mai degrabă de lățimea de bandă decât de puterea de calcul. Astfel de limitări ale lățimii de bandă pot apărea între microprocesor și dispozitivul reconfigurabil, sau între dispozitivul de procesare și propria sa memorie. Lățimea de bandă dintre microprocesor și dispozitivul reconfigurabil tinde să fie fixă și probabil mai eficientă când se transferă un volum mare de date. Lățimea de bandă a memoriei tinde să crească cu cât memoria este mai aproape de microprocesor. De exemplu, memoria cache internă a microprocesoarelor moderne este substanțial mai rapidă decât memoria cache externă, care de asemenea este mai rapidă decât memoria SDRAM externă. La fel se întâmplă și în cazul memoriei dispozitivelor reconfigurabile. Algoritmii care sunt în

prezent limitați de lățimea de bandă SDRAM a unui microprocesor pot fi limitați în mod similar pe un dispozitiv reconfigurabil.

Mai departe, pentru a putea profita de paralelism, se determină dacă există suficientă lățime de bandă între microprocesor și dispozitivul reconfigurabil, și de asemenea între dispozitivul reconfigurabil și resursele de memorie ale acestuia. De exemplu, pentru a calcula punctajul maxim de aliniament (maximum alignment score) între secvențe de ADN, să considerăm o arhitectură în care microprocesorul transmite către circuitul FPGA secvența de interogare, secvența de date și mai mulți parametri de punctaj. Numărul de punctaje pe care circuitul FPGA trebuie să le calculeze pentru a afla maximum este egal cu lungimea secvenței de interogare înmulțită cu lungimea secvenței de date. Pentru fiecare caracter transmis în secvența de date, circuitul FPGA trebuie să calculeze un rând întreg de punctaje. Determinarea fiecărui punctaj în parte necesită multe calcule pentru fiecare caracter transmis către circuitul FPGA și astfel „blocajele” în lățimea de bandă pe care se transmit secvențele către circuitul FPGA devin puțin probabile. De asemenea, deoarece singura dată retransmisă de către circuitul FPGA este punctajul maxim, nici lățimea de bandă de la circuitul FPGA către microprocesor nu va suferi limitări. Astfel, singura limitare a algoritmului este viteza cu care circuitul FPGA poate calcula punctajele.

#### 3.1.1.4. Memorie suficientă?

Analiza exemplului de mai sus este valabilă doar dacă circuitul FPGA poate calcula și stoca întregul tabel de punctaje într-o singură trecere. Fapt puțin probabil, deoarece ar necesita ca punctajele dintr-un rând întreg al tabelului să fie calculate și stocate în paralel. Asta ar limita foarte mult dimensiunea secvenței de interogare. Pentru ca circuitul FPGA să poată calcula tabela de punctaje pe segmente, trebuie să stocheze datele intermediare în memoria locală. Pentru a împărți tabela de punctaje în segmente verticale, trebuie stocată în memorie ultima coloană a unui segment pentru a o folosi la calculul primei coloane din următorul segment. Deoarece numai o coloană trebuie stocată, e puțin probabil ca lățimea de bandă a memoriei necesare va fi factorul limitant. Însă, dimensiunea acestei memorii va limita lungimea maximă a secvențelor de interogare și de date.

### 3.1.2. Clase de FPGA

Există patru mari clase de modele FPGA, așa cum apare în Figura 3.3. Toate aceste patru arhitecturi sunt disponibile în comerț, fiecare descriind poziționarea liniilor de rutare și a blocurilor logice. Se folosesc instrumente CAD sofisticate la proiectarea circuitelor logice care „traduc” desenul, arhitectura aflată la baza lor devenind transparentă. Alegerea arhitecturii depinde de aplicație, unele fiind mai potrivite în cazul anumitor sarcini decât altele.

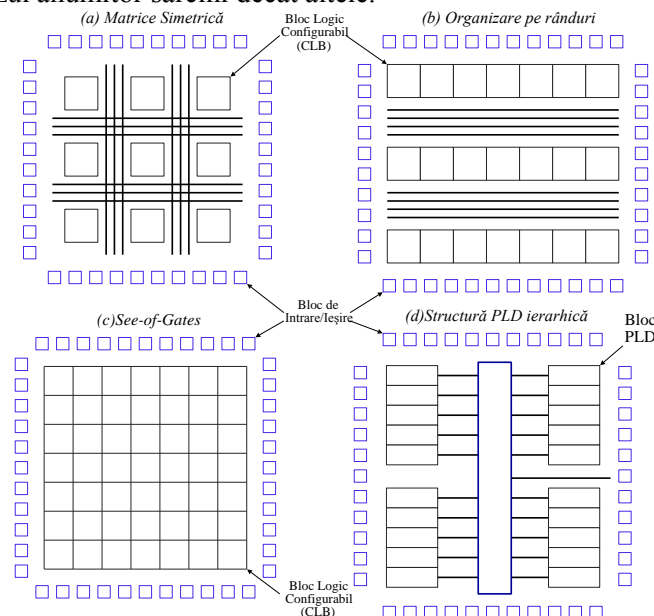


Figura 3.3 Principalele clase de FPGA

*Blocurile logice configurabile* (CLB) asigură elementele funcționale pentru construcția proiectului în FPGA. Cel mai important merit al CLB-ului este funcționalitatea sa. Funcționalitatea sa mărită permite implementarea unor funcții logice mult mai complexe pe un singur CLB, dar totuși, cum crește funcționalitatea unui CLB, tot așa crește și dimensiunea sa, rezultatul fiind mai puține CLB-uri pe fiecare FPGA. Un CLB conține un număr de blocuri de construcție standard, aranjate conform metodei specifice a fabricantului, pentru a permite implementarea diferitelor funcții logice. Blocurile de construcție sunt:

- *Tabele de căutare* (LUTs) –acceptă M intrări binare și dau N ieșiri, unde N este de obicei 1. Se poate nota cu M-LUT-N un LUT pentru a evidenția caracteristicile sale: ex. 3-LUT-1 este un tabel de căutare cu 3 intrări și o ieșire. Un astfel de tabel de căutare poate fi folosit pentru implementarea unui port AND cu trei intrări, ilustrat în Tabelul 3. Numărul biților de ieșire care trebuie stocați pentru orice tabel de căutare sunt o funcție directă a intrărilor:  $2^M$ . Astfel, un 4-LUT-1 necesită 16 biți stocați, iar un 3-LUT-1 are nevoie de doar 8. S-a arătat faptul că, pentru cele mai multe aplicații, soluția optimă este un 4-LUT-1.
- *Matricile logice programabile* (PLA) au fost introduse datorită faptului că tabelele de căutare se dovedesc a fi ineficiente în cazul  $M > 5$ . S-a examinat aria de eficiență pentru PLA, cu M intrări, N ieșiri și K termeni-produs, și s-a dovedit că valorile optime sunt  $K=10$ ,  $M=3$  și  $N=12$ . Acest dispozitiv oferă cu 4% mai multă funcționalitate decât o implementare 4-LUT-1.
- *Bistabile* – sunt extrem de utile, deși nu necesare, atunci când trebuie făcută o operație logică secvențială. S-a arătat prin experimente că fără utilizarea unui bistabil, numărul CLB-urilor necesare implementării unei funcții se dublează.

Tabelul 3 **3-LUT-1**

<b>A</b>	0	1	0	1	0	1	0	1
<b>B</b>	0	0	1	1	0	0	1	1
<b>C</b>	0	0	0	0	1	1	1	1
<b>X</b>	0	0	0	0	0	0	0	1

Exemplu pentru un tabel de căutare configurat pentru a funcționa cu trei intrări și un port de ieșire. CLB-urile pot să conțină alte componente, cum ar fi porți logici, registre și multiplexere, module sumatoare, totuși, neexistând o implementare standard.

Problema pe care o întâlnește orice producător la proiectarea unui FPGA este găsirea echilibrului optim între granularitatea CLB-ului și funcționalitatea sa. O posibilă rezolvare a acestei probleme este crearea unui șir neomogen de blocuri logice, dar un astfel de șir va avea întotdeauna elemente care nu se potrivesc unei anumite aplicații. Indiferent de configurația CLB-ului, eficiența FPGA-ului depinde în mare măsură de instrumentele CAD. Dacă instrumentele CAD nu mapează circuitul logic în mod optim pe CLB, avantajul arhitectural va fi irosit.

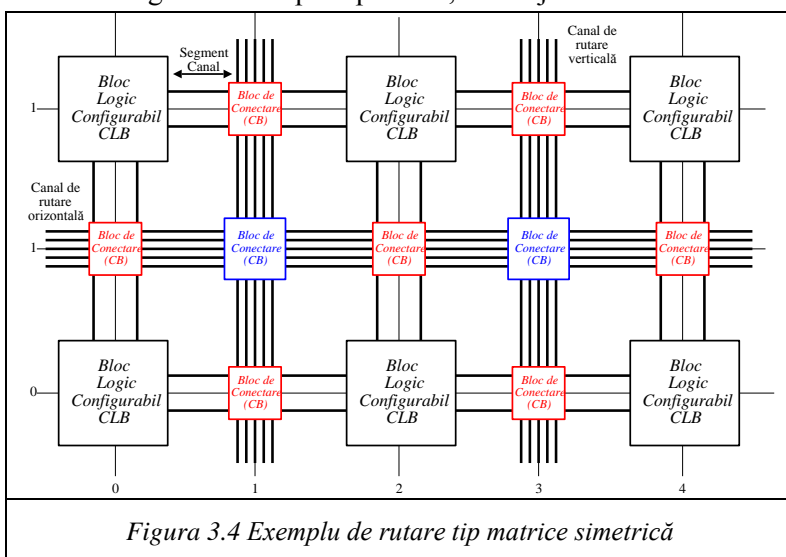


Figura 3.4 Exemplu de rutare tip matrice simetrică

Rutare în FPGA. Rutarea în FPGA se ocupă în mod esențial de conectarea CLB-urilor cu alte CLB-uri și cu blocurile de intrare/ieșire (IOB), oferind astfel funcționalitatea necesară. Echilibrul dintre spațiul oferit rutării și cel dedicat CLB-urilor este critic: dacă vor fi prea multe CLB-uri, conectarea sistemelor complexe va fi imposibilă, iar dacă vor fi prea puține, rutarea va rămâne nevalorificată. Pe un FPGA standard, rutarea va ocupa apr. 70-90% din spațiul total disponibil. Rutarea necesită mult spațiu și timp.

Există două tipuri de blocuri pe un FPGA, definite ca și bloc de conectare (CB) și bloc comutator (SB). În Figura 3.4 este reprezentată construcția unei arhitecturi CLB simetrice, care folosește

două blocuri adiționale și esențiale. Conexiunile interne posibile atât în CB și SB depind de producător. Un producător trebuie să ia în considerare rutarea necesară arhitecturii respective, și să optimizeze conectivitatea CB și SB pentru a permite o flexibilitate maximă a interconexiunilor, în același timp păstrând o redundanță minimă.

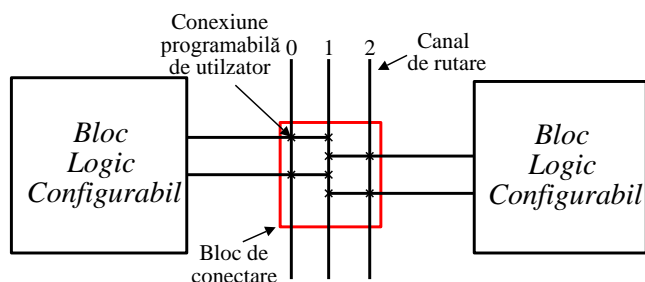


Figura 3.5 Exemplu Bloc de conectare

comutatori sau conexiuni programabile ar putea însemna imposibilitatea interconectării a două linii. În Figura 3.6 vom prezenta un exemplu de bloc comutator.

Flexibilitatea unui bloc comutator este definită ca și numărul segmentelor de cablaj la care orice linie de intrare poate fi conectată, în exemplul arătat aceasta fiind de 6. La proiectarea unei topologii de bloc comutator, producătorul trebuie să aleagă o topologie care să nu împiedice conectarea a două linii CLB.

Resursele de rutare în FPGA se clasifică în general în trei categorii, prin definirea utilizării lor primare. Aceste trei categorii sunt:

- *Interconectare cu scop general.* Se folosește pentru interconectări care cuprind una sau mai multe CLB-uri locale. Se implementează prin folosirea canalelor de rutare, blocurilor de conectare și a blocurilor comutatoare, așa cum s-a descris mai înainte. Un factor negativ îl constituie faptul că fiecare bloc comutator sau de conectare prin care trebuie să treacă un semnal, are o întârziere RC ca și o rezultantă a problemelor de propagare a semnalelor la frecvențe mai înalte ale dispozitivului.
- *Interconectare directă.* Asigură o conexiune directă cu unul sau cu toate CLB-urile învecinate la dreapta, stânga, în sus sau în jos.
- *Liniile lungi* se folosesc pentru rutarea conexiunilor care necesită interconectarea mai multor CLB-uri cu un defazaj minim, care asigură soluții parțiale pentru semnalele care altfel ar traversa mai mulți comutatori de rutare, rezultând întârzieri RC cumulative.

Un ultim și foarte important aspect al rutării în FPGA sunt IOB-urile, fără de care un FPGA ar fi inutil. IOB-urile trebuie alese astfel încât să susțină mai multe moduri de ieșire, cum ar fi TTL (transistor-transistor logic) sau CMOS (complementary metal-oxid semiconductor), și să asigure curent necesar, dar nu excesiv, pentru a permite funcționarea corectă a oricărui dispozitiv atașat. Din păcate, I/O este veriga slabă într-un FPGA și se dovedește a fi disproporțională față de densitatea chipului.

*FPGA-urile reconfigurabile în mod dinamic (DRFPGA)* combină, în principiu, viteza unei soluții hardware dedicate cu flexibilitatea softului. Astfel de FPGA-uri se pot reconfigura de un controller extern pentru implementarea oricărui set de funcții Booleene dorit. Acesta este conceptul din spatele computării reconfigurabile dinamice: combinarea flexibilității softului cu viteza hardware-ului.

Există, în general, două motive pentru adoptarea unei reconfigurări dinamice:

- *Creșterea vitezei* prin folosirea unei arhitecturi specializate pentru o anumită problemă, rezultând o performanță mult mai bună
- *Toleranța față de eroare* îmbunătățește procesul de fabricare, prin faptul că permite schimbări ale designului sau depășirea ușoară a unor probleme.
- Dimensiunile și complexitatea a celui mai mic bloc dintr-un dispozitiv reconfigurabil se pot clasifica în funcție de granularitatea sa:
  - *De granularitate fină* - aceste arhitecturi constau din blocuri logice mici și simple care sunt configurate pentru a realiza operații mai complexe.

Blocurile de conexiune (CB) și comutatoare (SB) se pot folosi pentru a ruta I/O de pe elementele CLB în canale de rutare. Un exemplu de bloc de conexiune apare în Figura 3.5 în care este prezentată un punct de comutare de rutare configurabilă folosind o intersecție, CB-ul având la rândul lui o topologie de interconectare predefinită.

Flexibilitatea este o problemă importantă aici, deoarece prea puțini

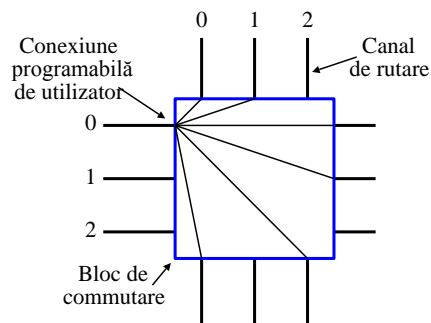


Figura 3.6 Exemplu de bloc comutator



- *De granularitate medie* - aceste arhitecturi constau din blocuri logice complexe, din care fiecare realizează o anumită parte semnificativă a calculului. Calculele nestandardizate se pot face prin reconfigurarea rețelei de interconectare.
- *De granularitate dură* - constau dintr-un număr de unități de execuție dintre care fiecare are propriul set de instrucțiuni. Fiecare unitate e mai simplă decât un microprocesor, dar integrată, astfel încât permite o viteză mare de comunicare. (Nagami, Oguri, Shiozawa, Ito, & Konishi, 1998).

*Integrarea dispozitivului* specifică cât de strâns este cuplat sistemul reconfigurabil la gazda sa:

- *Sistemele dinamice* sunt sisteme de inspirație biologică care nu sunt controlate de un dispozitiv extern. Aceste sisteme auto-evoluază.
- *Sisteme statice, cuplate strâns* - leagă elementele reconfigurabile strâns ca și elemente de execuție pe calea de date a unui procesor-gazdă.
- *Sisteme statice, cuplate larg* - se situează pe o placă separată de gazdă. Acest lucru se desfășoară în general în detrimentul vitezei, deoarece datele trebuie transferate spre și de la sistemele fiică.

Reconfigurabilitatea unei rețele de interconectare externă între unități reconfigurabile se poate de asemenea clasifica în două categorii:

- *Rețea externă reconfigurabilă* - extinde conceptul reconfigurării asupra mai multor circuite reconfigurabile în mod eficient, având ca și rezultat o unitate reconfigurabilă de dimensiune mare. Totuși, penalitățile depășirii chipului sunt mari.
- *Rețele externe fixe* - folosesc conexiuni statice între circuitele reconfigurabile, limitând flexibilitatea în vederea creșterii vitezei și minimizării costului.

### 3.1.3. Elementele unui circuit FPGA din familia Xilinx Spartan3

Arhitectura circuitului FPGA din familia Spartan-3 conține următoarele elemente funcționale programabile:

- **Blocuri logice configurabile (CLBs)** (Configurable Logic Blocks) –conține tabele de căutare LUT cu scopul implementării de funcții logice, și elemente de stocare care se pot utiliza ca bistabile sau registre.
- **Blocuri de intrare ieșire IOB** (Input/Output Blocks) –acestea există pentru controlul fluxului de date între pini de intrare/ieșire și logica internă a circuitului FPGA. Fiecare bloc IOB suportă flux bidirecțional de date și de asemenea permite operații cu trei stări. Recunoaște 65 standarde de semnale diferite. Include registre Double Data-Rate.
- **Bloc de control digital al impedanței** (DCI Digitally Controlled Impedance) asigură terminații on-chip, simplificând astfel complexitatea proiectului.
- **Blocuri RAM** permit stocarea de date în blocuri de memorie cu poartă duală (dual-port blocks) sau cu poartă simplă (single port) de 18-Kbit.
- **Blocuri multiplicator** –acceptă două numere binare pe 18 biți ca intrare, și calculează produsul.
- **Manager digital de semnal de tact** - asigură soluții digitale pentru distribuirea, întârzierea, multiplicarea, divizarea și modificarea fazei semnalului de tact.

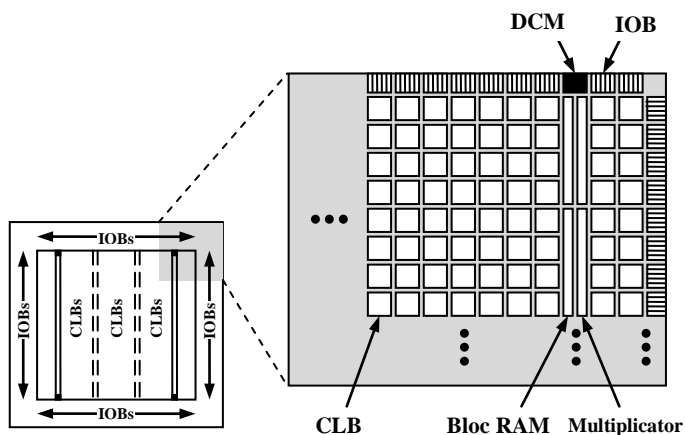


Figura 3.7 Organizarea elementelor în circuitul FPGA

Organizarea acestor elemente în circuitul FPGA este prezentată în Figura 3.7.

Toate aceste elemente constructive sunt conectate printr-o magistrală de control extrem flexibil programabilă, care este realizată pentru a asigura un număr foarte ridicat de procese de reconfigurare. Șirul de biți care determină configurarea acestor conexiuni între elementele de bază ale FPGA se poate descărca prin porturi seriale, paralele sau dintr-o memorie PROM externă. Un aspect important este faptul că aceste circuite sunt capabile să funcționeze la frecvențe de până la 200 MHz, ceea ce permite realizarea unor rețele neuronale artificiale foarte rapide.

Utilizarea memoriilor are un rol foarte important în implementarea rețelei neuronale de tip CMAC și RBF. Circuitul FPGA descris înglobează două tipuri de memorii:

- blocuri de memorie RAM- de 18Kbit.
- RAM distribuit în blocurile logice, de 64 de biți cu poartă simplă sau de 32 de biți cu poartă duală (Xilinx, DS099, 2006).

Modulele de memorii Block RAM s-au utilizat cu succes la implementarea unei rețele neuronale de tip CMAC și RBF, reprezentând elementele de bază pentru stocarea ponderilor rețelei și a funcțiilor de activare. Deoarece atât memoriile Block RAM cât și memoriile distribuite au un rol esențial în realizarea rețelelor neuronale, în următorul subcapitol sunt descrise în amănunt aceste două elemente.

### 3.1.3.1. Elementele Bloc RAM

Pentru aplicațiile care necesită memorii extinse, on-chip, FPGA-urile de generația Spartan™-3 asigură o multitudine de plăci Select RAM™ foarte eficiente. Utilizând diferite opțiuni de configurare, plăcile SelectRAM creează RAM, ROM, FIFO, tabele extinse de căutare, buffer-e circulare, registre, fiecare susținând o varietate de date în ceea ce privește extinderea și adâncimea lor. Prin utilizarea sistemului Xilinx Core Generator™ VHDL, sau Verilog Instantiation se pot specifica diferite caracteristici de proiectare.

Fiecare placă FPGA de generație Spartan™-3, inclusiv Spartan-3, Spartan-3L și Spartan-3E se conturează ca și memorii RAM de plăci multiple, organizate în coloane. Cantitatea memoriei RAM depinde de mărimea FPGA-ului de generație Spartan-3. Fiecare placă RAM conține 18432 biți de RAM static rapid, din care 16 Kbiți au fost alocați pentru stocare de date, și, în câteva configurații de memorie, 2 Kbiți adiționali au fost alocați ca biți de paritate sau biți adiționali de date.

Sistemul Xilinx CORE Generator susține module variate care conțin Block RAM pentru mecanisme Spartan-3, incluzând:

- A. module RAM înglobate cu poartă simplă sau duală
- B. module ROM
- C. module FIFO sincron sau asincron
- D. module de memorii adresabile prin conținut (CAM)

Memoria bloc RAM se poate încorpora în orice proiect prin folosirea unui modul „RAMB16” adecvat din colecția de proiecte Xilinx.

**Locația memoriei Block RAM și mediul înconjurător.** Memoria Block RAM este organizată în coloane. XC3S50-ul are doar o singură coloană de memorie bloc RAM, așezată la două coloane CLB de marginea din stânga a dispozitivului. Dispozitivele Spartan-3 mai mari decât XC3S50 au două coloane de memorie Block RAM, adiacente marginilor din stânga respectiv dreapta ale dispozitivului, la două coloane CLB de I/O (intrare/ieșire) de la margine. Pe lângă aceste coloane de Block RAM de la margine, XC3S4000 și XC3S5000 au câte două coloane în plus- patru coloane în total- distribuite în mod egal de la coloanele de la margini.

Adiacent la fiecare memorie Block RAM este un multiplicator hardware încorporat pe 18\*18 biți. Așezarea memoriei bloc RAM și a multiplicatorului încorporat unul lângă celălalt îmbunătățește performanța unor funcții de procesare/prelucrare a semnalelor digitale.

Conectarea specială din jurul memoriei Block RAM asigură o distribuție eficientă a semnalelor referitoare la liniile de adresă și de date. Pe lângă aceasta, unele dotări speciale asigură ca memoriile Block RAM multiple să poată fi așezate în cascadă pentru a crea memorii mai largi sau mai adânci.

#### Fluxurile de date

Memoria Block RAM Spartan-3 este construită din memorie dual-port și sprijină simultan următoarele operații:

- Porturile A și B funcționează separat ca și un RAM cu port simplu independent, susținând operații simultane de citire și scriere folosind câte un singur set de linii de adresă și de date.



- Portul A este un port de scriere cu o adresă de scriere separată și portul B este portul de citire cu o adresă de citire separată. Lărgimea semnalelor de date poate să difere la cele două porturi.
- Portul B este un port de scriere cu o adresă de scriere separată și portul A este portul de citire cu o adresă de citire separată. Lărgimea semnalelor de date poate să difere la cele două porturi.

Semnalele conectate la un primitiv al memoriei Block RAM se împart în patru categorii:

- Ințrările și ieșirile de date.
- Ințrările și ieșirile de ponderi, accesibile când portul de date este de o mărime de un byte sau mai mare.
- Ințrări de adrese pentru selectarea unei locații specifice a memoriei.
- Semnale de control pentru a administra diferite operații de scriere, citire sau de setare/resetare.

**Ințrări și ieșiri de date.** Lățimea totală de date a unui port include magistrala de date și magistrala de paritate.

**Magistrala de ințrări de date -DI[#:0] (DIA[#:0], DIB[#:0])**

Magistrala de ințrări de date este sursa datelor care vor fi înscrise (introduse) în RAM.

Datele de la magistrala de intrare DI sunt înscrise în celulele RAM specificate de magistrala de adresă de intrare, ADDR, prin transmiterea unui semnal cu front crescător la semnalul CLK, când ințrările de validare semnal de tact EN și validare scriere WE sunt la nivel logic 1.

**Magistrala de ieșire de date -DO[#:0] (DOA[#:0], DOB[#:0])**

Magistrala de ieșire de date, DO, reprezintă conținutul celulelor de memorie selectate de magistrala de adresă, ADDR, la un front crescător al semnalului de tact în cursul unei operații de citire. În cursul unei operații de scriere simultană, comportamentul registrului de date de ieșire este controlat de atributul WRITE\_MODE.

**Ințrările și ieșirile de paritate**

Cu toate că ne referim la ele aici ca și biți de „paritate”, ințrările și ieșirile de paritate nu au o funcționalitate specială, și pot fi folosite ca și biți adiționali de date.

**Ințrări de adrese.** Fiind un RAM dual-port, ambele porți operează în mod independent, accesând același set de celule de memorie de 18 K-biți.

**Magistrala de adrese - ADDR[#:0] (ADDRA[#:0], ADDR B[#:0])** Magistrala de adrese selectează celulele de memorie pentru operațiile de citire sau scriere.

**Semnale de control**

**Ceas - CLK (CLKA, CLKB)** Fiecare port este sincronizat în totalitate cu semnalele de ceas independente. Toate semnalele de intrare ale portului au fixați timpii de validare relativ față de semnalul de tact CLK.

**Semnalul de validare EN**

Intrarea semnalului de validare, EN, controlează operațiile de citire, scriere, setare/resetare. Când EN este pe nivel logic 0, nu sunt înscrise date, și ieșirile DO și DOP păstrează ultima valoare. Polaritatea EN este configurabilă și este activă pe nivel logic 1.

**Semnalul de validare de scriere — WE (WEA, WEB)**

Semnalul de validare de scriere, WE, controlează înscrierea datelor în RAM. Când ambele semnale de validare, EN și WE sunt active la frontul pozitiv al semnalului de ceas, semnalele din magistralele de intrare a datelor și de paritate sunt înscrise în zona de memorie selectată de magistrala de adresă.

Registrul de ieșire a datelor este încărcat sau nu în funcție de atributul WRITE\_MODE.

**Setare/resetare sincron - SSR (SSRA, SSRB)** Intrarea de setare/resetare sincron, SSR, forțează registrele de ieșire de date să ia valoarea specificată de atributul SRVAL. Când SSR și semnalul de validare sunt la nivel logic 1, registrele de ieșire pentru ieșirile DO și DOP sunt setate la „0” sau „1” în funcție de parametrul SRVAL.

O operație de setare/resetare sincron nu influențează celulele de memorie RAM și nu perturbă operațiile de scriere la celălalt port.

**Semnale de control de inversare/schimbare.** Pentru fiecare port, toate cele patru semnale de control-CLK, EN, WE și SSR- au fiecare câte o opțiune individuală de inversare. Fiecare semnal de control

poate fi configurat să fie activ la nivel logic ‘1’ sau ‘0’, iar ceasul poate fi activ la un front pozitiv sau negativ fără a fi necesar utilizarea altor resurse logice.

**Organizarea memoriei.** Organizarea datelor sau aspectul proporțional al unui bloc RAM sunt configurabile, așa cum apare în Tabelul 4. Dacă calea de date este de o lățime de 1 byte sau mai largă, blocul RAM asigură biți adiționali pentru a susține paritatea pentru fiecare byte. În consecință, o organizare a memoriei de 1Kx18 este de o lățime de 18 biți cu 16 biți (2 byte) alocați datelor plus 2 biți de paritate, câte unul pentru fiecare byte. De asemenea, cantitatea fizică de memorie accesibilă de la un port depinde de organizarea memoriei. Pentru memorii de o lățime de 1 byte sau mai largi, există 18 Kbiți de memorie accesibilă. Pentru memorii de o lățime mai mică, numai 16 biți de memorie sunt accesibile din cauza lipsei biților de paritate în această organizare. De regulă, 16 Kbiți sunt alocați semnalelor de date și 2 Kbiți semnalelor de paritate în cazul blocului RAM de 18 Kbiți.

Tabelul 4 Organizarea memoriei Bloc RAM

Organizare	Nr biți adrese	Nr biți date	Nr biți paritate	DI-DO	DIP/DO P	ADDR	Primitivă port simplu	Total RAM KBit
512x36	512	32	4	(31 : 0)	(3:0)	(8:0)	RAMB16_S36	18K
1Kx18	1024	16	2	(15 : 0)	(1:0)	(9:0)	RAMB16_S18	18K
2Kx9	2048	8	1	(7 : 0)	(0:0)	(10:0)	RAMB16_S9	18K
4Kx4	4096	4	-	(3:0)	-	(11:0)	RAMB16_S4	18K
8Kx2	8192	2	-	(1:0)	-	(12:0)	RAMB16_S2	18K
16Kx1	16384	1	-	(0:0)	-	(13:0)	RAMB16_S1	18K

## Modurile de scriere/citire

### Comportamentul de citire în timpul scrierii simultane - WRITE\_MODE

Pentru a maximiza performanțele datelor și ale utilizării memoriei dual-port, la fiecare front al ceasului, memoria bloc RAM susține unul dintre cele trei moduri de scriere pentru fiecare port de memorie. Aceste moduri determină datele care sunt disponibile la registrele de ieșire după un front de ceas valid la scriere la același port. Modul inițial, WRITE\_FIRST, asigură compatibilitatea în retrospectivă cu arhitecturile Virtex™/Virtex-E și Spartan-IIE FPGA, și este de asemenea starea inițială pentru echipamentele Virtex-II/Virtex-II Pro™. Totuși, modul READ\_FIRST este cel mai util deoarece crește eficiența memoriei Block RAM la fiecare ciclu al ceasului, permițând echipamentelor să folosească lățimea de bandă la maxim. În modul READ\_FIRST, un port de memorie susține operații simultane de citire și scriere la aceeași adresă la același front al ceasului, independent de orice complicații de temporizare. Tabelul 5 subliniază modul în care aranjamentul WRITE\_MODE influențează registrele de ieșire la același port și cum influențează registrele de ieșire la portul opus în timpul accesului simultan la aceeași adresă.

Tabelul 5 Influența modului WRITE\_MODE

Modul de scriere	Efectul asupra aceleiași port	Efectul asupra portului opus (doar mod dual-port, aceeași adresă)
WRITE_FIRST Citește după scriere (din oficiu)	Datele de la intrările DI, DIP sunt scrise într-o zonă specificată a memoriei RAM și apar simultan la ieșirile DO, DOP.	Invalidează datele la ieșirile DO, DOP.
READ_FIRST Citește înainte de scriere (recomandat)	Datele din zona specificată RAM apar la ieșirile DO, DOP. Datele de la intrările DI, DIP sunt scrise într-o zonă specificată.	Datele din Zona specificată a RAM apar la ieșirile DO, DOP.
NO_CHANGE Fără citire asupra scrierii	Datele de la ieșirile DO, DOP rămân neschimbate. Datele de la intrările DI, DIP sunt scrise în zona specificată.	Invalidează datele la ieșirile DO, DOP.

Selectarea modului este setat prin configurație. Unul dintre aceste trei moduri este setat în mod individual pentru fiecare port printr-un atribut.

**Modul WRITE\_FIRST** este modul de operare din oficiu din considerente de compatibilitate în retrospectivă. Pentru proiecte noi se recomandă modul READ\_FIRST.

În acest mod, datele de intrare sunt scrise în zona de memorie RAM adresată și sunt stocate simultan în registrele de ieșire de date, rezultând în operații transparente de scriere, așa cum apare în Figura 3.9. În Figura 3.8 este prezentată ciclul de semnale pentru o operație de scriere WRITE\_FIRST.

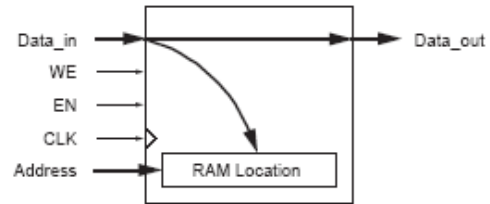
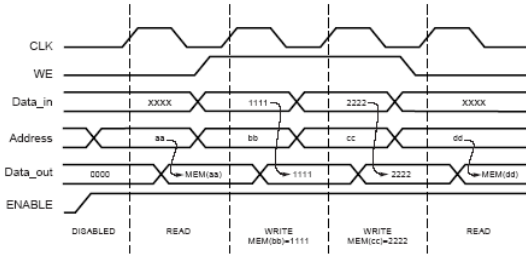


Figura 3.8 Ciclul de scriere în modul WRITE\_FIRST      Figura 3.9 Fluxul de date în modul WRITE\_FIRST

În modul READ\_FIRST, datele stocate anterior la adresa de scriere apar la registrele de ieșire, în timp ce datele de intrare sunt stocate în memorie, rezultând o operație de citire-înainte-de-scriere așa cum apare în Figura 3.10. Datele RAM mai vechi apar la ieșirea de date în timp ce datele RAM noi sunt stocate în zona de memorie RAM specificată.

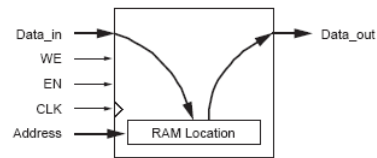
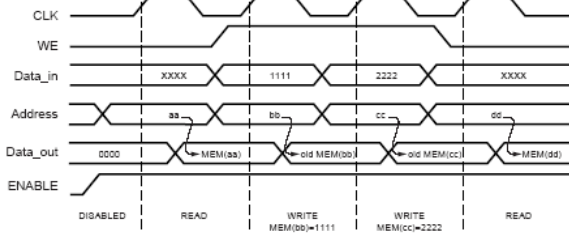


Figura 3.10 Ciclul de scriere în modul READ\_FIRST      Figura 3.11 Fluxul de date în modul READ\_FIRST

### Modul NO\_CHANGE

În modul NO\_CHANGE, registrele de ieșire sunt blocate și rămân neschimbate în timpul unei operații simultane de scriere, așa cum apare în Figura 3.12. Acest comportament imită cel al unei memorii sincron simple când zona de memorie este citită sau scrisă în timpul unui ciclu de ceas.

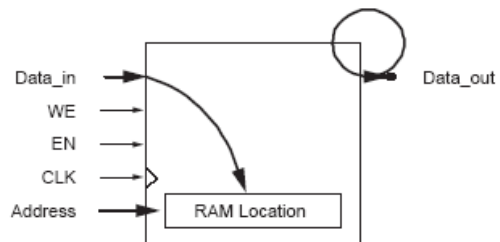
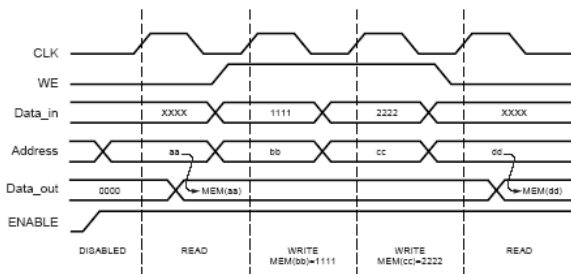


Figura 3.12 Ciclul de scriere în modul NO\_CHANGE      Figura 3.13 Fluxul de date în modul NO\_CHANGE

Modul NO\_CHANGE este util în cazul mai multor aplicații, incluzând acelea la care blocul RAM conține forme de unde, tabele de funcții, coeficienți, etc. Memoria poate fi actualizată fără afectarea ieșirilor memoriei.

Fluxul de date în modul NO\_CHANGE arată că ieșirea de date reține ultima dată citită dacă există o operație simultană de scriere la același port (Xilinx, XAPP463, 2005).

#### 3.1.3.2. RAM distribuit în CLB

Adițional față de Block RAM de 18Kbit FPGA-ul Spartan-3 conține blocuri de memorie distribuită în fiecare bloc logic configurabil (CLB). Fiecare generator de funcție SLICEM sau LUT din resursele CLB poate să implementeze opțional un 16x1bit RAM sincron. Tabelele de căutare din SLICEL nu conțin RAM distribuit.

Scrierea datelor în memoriile RAM distribuite se realizează în mod sincron, iar citirea asincron. Fiecare 16x1 bit RAM se poate conecta în cascadă pentru a obține memorii cu o magistrală mai lată sau îngustă de adrese și/sau date prin resursele logice aferente cu o penalizare minimă de timp. Blocurile logice configurabile Spartan-3 suportă o varietate largă de primitive până la 64 biți de adresare x 1 bit lățime linie de date.

### Fluxul de date Single-Port și Dual-Port RAM distribuite

Memoriile RAM distribuite susțin următoarele tipuri de configurare:

- RAM cu port simplu cu înscriere sincron și citire asincron. Și citirea sincron este permisă prin utilizarea bistabilelor asociate memoriei RAM distribuite.
- Dual-port RAM cu o scriere sincron și două operații de citire a porturilor în mod asincron. De asemenea este posibilă citirea sincron.

Cum este ilustrat în Figura 3.14, memoria RAM distribuită conține un port pentru scriere și citire și un port independent pentru citire. Orice operație de scriere prin intrarea D sau de citire prin ieșirea SPO poate să apară simultan și independent de o operație de citire din portul secund DPO.

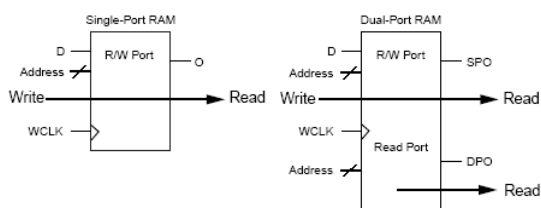


Figura 3.14 Fluxul de date Single-Port și Dual-Port RAM distribuit

### Operații de citire/scriere

**Operația de scriere** se realizează pe frontul semnalului de tact controlat prin intrarea validare scriere WE. Din oficiu WE este activ pe '1', dar acesta se poate inversa. Când WE este 1, pe frontul semnalului de tact este validată adresa și datele din intrarea D sunt înscrise în locația de memorie selectată. Dacă WE='0' datele nu sunt înscrise în memorie.

**Operația de citire** este pur combinațional. Portul de

adresare, atât pentru memoria cu port simplu sau port dublu, este asincron cu un timp de acces echivalent cu întârzierile logicii tablei de căutare LUT.

**Citire în timpul scrierii.** Când sunt înscrise date noi, ieșirea reflectă datele care au fost înscrise în celula de memorie adresată. În Figura 3.15 este prezentată diagrama de timp, care ilustrează o operație de scriere cu datele anterior citite prin portul de ieșire, urmate de datele noi după apariția frontului crescător pe semnalul de tact.

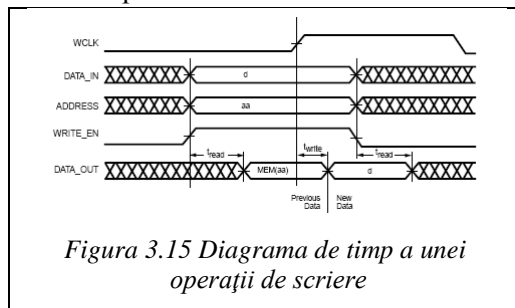


Figura 3.15 Diagrama de timp a unei operații de scriere

### Caracteristici

- operația de scriere necesită numai un singur front de semnal de tact
- durata operației de citire depinde numai de timpul de accesare a elementelor logice
- ieșirile sunt asincron și depind numai de întârzierea logicii tablelor de căutare
- intrările de adrese și date sunt activate pe frontul pozitiv al semnalului de tact. Nu necesită timp de

menținere.

- pentru RAM cu port dual, portul A[#:0] este adresă pentru scriere citire, și DPRA[#:0] este adresă independentă pentru citire

Spartan-3 suportă următoarele primitive prezentate în tabelul Tabelul 6.

Tabelul 6 Primitive pentru RAM distribuit cu Port Simplu și Port Dual

Primitive	Dimensiune memorie	Tip	Semnale de adrese
RAM16x1S	16x1	Port Simplu	A3,A2,A1,A0
RAM32x1S	32x1	Port Simplu	A4,A3,A2,A1,A0
RAM64x1S	64x1	Port Simplu	A5,A4,A3,A2,A1,A0
RAM16x1D	16x1	Port Dual	A3,A2,A1,A0

Datele de intrare ieșire sunt de un bit. Mai multe memorii RAM distribuite conectate în paralel permit implementarea unor elemente cu funcție de memorie cu un număr mai mare de biți de date.

**Semnalele de port.** Figura 3.16 reprezintă primitivile pentru memoriile RAM distribuite cu port simplu sau port dual.

**Semnalele pentru controlul memoriilor distribuite.** Ambele porturi de RAM distribuit operează independent unul față de celălalt în timp ce citește același set de celule de memorie.

**Clock — WCLK** Semnalul de tact WCLK este utilizat pentru scriere sincron. Liniile de adrese și de date au timpii de setare specificați față de semnalul WCLK.

**Validare scriere — WE.** Semnalul WE validează funcționalitatea de scriere a portului. Un semnal inactiv de WE previne orice operație de scriere a datelor în celulele memoriei.

Un semnal activ WE='1' permite ca apariția unui front pozitiv pe semnalul de tact să înscrie semnalele de date în locația de memorie specificată de semnalele de adrese.

**Address — A0, A1, A2, A3 (A4, A5).** Liniile de adrese selectează celulele de memorie pentru scriere sau citire. Lățimea portului determină numărul liniilor de adrese.

**Date de intrare — D** intrarea de date furnizează valorile noilor date care să fie înscrise în RAM

**Date de ieșire — O, SPO, și DPO.** Ieșirea de date O în cazul memoriei RAM distribuite cu port simplu, sau SPO și DPO în cazul memoriei RAM cu port dual reflectă conținutul celulelor de memorie selectate prin liniile de adrese. După un semnal de tact activ de scriere, liniile de date de ieșire (O sau SPO) reflectă datele noi care au fost înscrise (Xilinx, XAPP464, 2005).

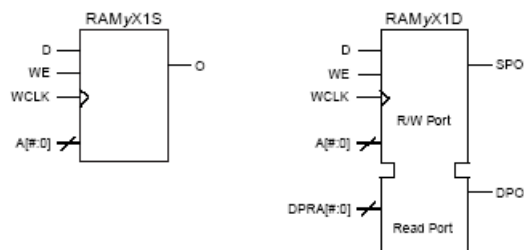


Figura 3.16 Semnalele de control la RAM distribuit

## 3.2. Implementări RNA pe FPGA

Progresul în curs al microelectronicii este forța motrice a dezvoltării continue a produselor tehnice noi. Sistemele FPGA sunt inovații de acest fel în sistemele microelectronice pentru aplicațiile computerizate.

Creșterea continuă în densitate a sistemelor FPGA a făcut posibilă realizarea unor proiecte de sisteme on-chip cu o complexitate mai mare de un milion de porți și RAM intern. De aceea, sistemele FPGA s-au dovedit în prezent a fi o platformă hardware atractivă pentru algoritmi RNA care necesită mult spațiu.

Un alt avantaj al acestui dispozitiv este capacitatea de a combina programabilitatea cu viteza crescută a operațiilor asociate cu soluțiile hardware paralele (Lysaght, Stockwood, Law, & Girma, 1994).

### 3.2.1. Moduri de cuplare și configurări pentru FPGA

Sistemele reconfigurabile sunt de obicei formate din combinarea sistemelor logice reconfigurabile și a unui microprocesor de uz general. Procesorul efectuează operațiile care nu pot fi realizate eficient în logica reconfigurabilă, ca de exemplu comenzi în bucle sau în derivații, accesări ale memoriei.

Sunt trei moduri diferite de abordare pentru ca un FPGA să funcționeze în combinație cu procesorul gazdă. În prima abordare, o unitate reconfigurabilă poate fi utilizată ca un coprocesor. În acest caz coprocesorul execută calculele fără a apela la interogarea procesorului gazdă. Procesorul gazdă inițializează hardware-ul reconfigurabil și trimite datele necesare pentru calcule și coprocesorul prelucrează independent aceste date (Compton & Hauck, 2000). O astfel de integrare permite, în majoritatea cazurilor, atât pentru procesorul gazdă, cât și pentru hardware-ul reconfigurabil, să funcționeze simultan datorită nesupunerii comunicării.

În a doua abordare, o unitate reconfigurabilă poate fi utilizată ca o unitate de procesare atașată. Aceasta este o configurație mai puțin strânsă în comparație cu prima abordare. Memoria cache a procesorului gazdă nu este vizibilă de către unitatea reconfigurabilă atașată. De aceea există o întârziere mai mare între comunicarea dintre procesorul gazdă și hardware-ul reconfigurabil, între datele de intrare și rezultate (Compton & Hauck, 2000).

Totuși această configurație permite o mai mare independență pentru calcule, deoarece se pot transmite mai multe operații CPU către hardware-ul atașat.

A treia abordare, caracterizată de conexiunea cea mai slabă, este un hardware extern, reconfigurabil, de sine stătător. Există o comunicare rară între hardware și procesorul gazdă.

Fiecare mod de abordare are avantajele și dezavantajele lui. Cu cât este mai strânsă integrarea hardware-ului reconfigurabil cu procesorul gazdă, cu atât poate fi mai mult utilizat în comunicări mai puțin suprapuse. Cu toate acestea nu poate opera în mod independent o perioadă îndelungată de timp fără întrerupere de către procesorul gazdă. Și la acest tip de abordare suprafața de logică reconfigurabilă este destul de mică. Cu cât este mai liberă cuplarea, cu atât permite mai multe paralelisme în executarea programelor, suferă de suprapunerea comunicării (Compton & Hauck, 2000). Această formă de conexiune se potrivește pentru aplicațiile la care o parte a calculului poate fi efectuată cu un hardware reconfigurabil o perioadă de timp destul de îndelungată, fără prea multă comunicare cu procesorul gazdă. Într-un sistem în care există o cuplare strânsă între gazdă și hardware-ul reconfigurabil, este extrem de necesar ca sarcinile de efectuat să fie împărțite între software și hardware. Implementarea acestei metode cade în sarcina unui calcul proiectat atât pentru hardware, cât și pentru software.

În chipurile FPGA există două moduri de implementare a logicii reconfigurabile. Sunt binecunoscute CTR (reconfigurarea compilată în timp) și RTR (reconfigurarea în timpul utilizării). CTR este o strategie de implementare statică, la care fiecare aplicare constă într-o singură configurare. RTR este o strategie de implementare dinamică unde fiecare aplicație constă în multiple configurații cooperative (Dehon, 2000).

În CTR, o dată ce o operație începe, configurația nu se mai schimbă în timpul procesului. Această abordare este asemănătoare cu implementarea unei aplicații în ASIC, în timp ce RTR utilizează o schemă de alocare dinamică care realocă partea de hardware în timpul executării aplicației (Hutchings & Writhlin, 1995). Fiecare aplicație constă în mai multe configurații pentru FPGA.

Aceste configurații sunt încărcate succesiv în timpul derulării unei aplicații pe chipul FPGA. Acest lucru înseamnă realizarea unui nou hardware pentru fiecare configurație pentru porțiuni particulare ale aplicației. Există două modele RTR: globală și locală. RTR global alocă întreaga resursă a întregului chip FPGA pentru fiecare pas al configurației. Aceasta duce la o nouă realizare pe FPGA în timpul fiecărei configurații în timpul utilizării. Proiectantul trebuie să implementeze o aplicație în RTR global și să împartă aplicația în părți aproximativ egale pentru a utiliza în mod eficient resursele FPGA (Hutchings & Writhlin, 1995). Aceasta se numește partiționare temporală a unei aplicații. În RTR local, o aplicație reconfigurează local niște subseturi ale logicii în timpul executării aplicației. Poate configura orice procent din resursele reconfigurabile în orice moment. Organizarea aplicațiilor în RTR local se bazează mai mult pe o divizare funcțională a muncii decât pe partiționarea utilizată de aplicațiile RTR global (Dehon, 2000). Odată ce partiționarea temporală manuală este decisă, poate fi foarte ușor implementată în varianta RTR global. Implementarea RTR local nu este posibilă în instrumente CAD.

### **3.2.2. Maparea algoritmilor de rețele neuronale artificiale pe sisteme FPGA**

Eldredge (Eldredge, 1994) a implementat cu succes algoritmul backpropagation folosind o placă proprie construită din FPGA-uri Xilinx XC3090, denumit Run-Time Reconfiguration Artificial Neural Network (RRANN).

Influențat de arhitectura RRANN a lui Eldredge, Beuchat et al. (Beuchat, Haenni, & Sanchez, 1998) au elaborat o placă FPGA, denumită RENCO (Reconfigurable Network Computer). Sistemul RENCO conține patru FPGA-uri Altera FLEX 10K130, care se pot reconfigura și monitoriza printr-o conexiune LAN (ex. Internet sau altul) printr-o interfață 10Base-T. RENCO a fost destinat recunoașterii caracterelor scrise de mână.

Ferucci și Martin (Ferrucci, 1994) (Marcelo & Martin, 1994) au construit o placă proprie, denumită „Adaptive Connectionist Model Emulator” (ACME) care constă din FPGA-uri Xilinx XC4010 multiple. ACME a fost validat cu succes prin implementarea unei rețele cu 3 intrări, 3 unități ascunse și 1 ieșire pentru rezolvarea problemei XOR cu 2 intrări. Skrbek a folosit această problemă pentru a demonstra că propria sa placă FPGA bazată pe backpropagation funcționează (Skrbek, 1999). Placa FPGA a lui Skrbek, denumită cardul ECX, poate implementa de asemenea rețele neuronale cu funcții de bază radiale, și a fost validată prin folosirea unor aplicații de recunoaștere a modelelor cum ar fi problema parității, recunoașterea cifrelor, recunoașterea semnalelor sonare.

GANGLION-ul este o aplicație dezvoltată de IBM Research Division, care se bazează pe circuite FPGA programabile de serie Xilinx XC3000. Sistemul se compune din cinci asemenea elemente programabile, care au fost așezate pe o placă VME. Rețeaua obținută este un MLP cu 12 intrări, 14 neuroni ascunși, 4 ieșiri. Aplicația este integral paralelă, operează pe date de intrare de 8 biți și ponderi de 8 biți. Astfel puterea atinsă este de 4,48 GCPS.

În ceea ce privește funcțiunile sale, un RNA se poate evidenția prin asociere, clasificare și aproximare. Autorii (Poormann, Witkowski, Kalte, & Ruckert, 2002) au implementat pe sisteme FPGA un algoritm pentru fiecare funcție. Pentru acest scop a fost folosit un accelerator hardware reconfigurabil în mod dinamic, RAPTOR2000 care este compus dintr-o placă de bază și până la șase module specifice aplicației. Placa de bază asigură infrastructura de comunicare pentru modulele găzduite prin magistrala PCI. Cele șase module sunt conectate în inel. Fiecare modul se compune dintr-un chip Xilinx Virtex XCV1000 FPGA și 256 Mbyte SDRAM. Toate modulele sunt conectate la o magistrală locală comună pentru a comunica cu alte dispozitive sau module și pentru comunicarea cu gazda printr-o punte PCI. Puntea de magistrală PCI poate să funcționeze în mod slave serial în care sistemul gazdă configurează inițial modulele prin transferarea șirurilor de biți în FPGA. Se poate folosi o magistrală broadcast în plus pentru comunicarea simultană dintre module.

Pentru accesarea rapidă a memoriei sistemului gazdă, se folosește un SRAM cu poartă dublă, care poate fi accesată de fiecare modul. Reconfigurarea modului pe RAPTOR2000 este inițiată de calculatorul gazdă. Se susține că RAPTOR2000 are un timp de reconfigurare de 20 ms pentru fiecare modul. S-au ales rețelele SOM (hartă Kohonen) pentru sarcina de clasificare, NAM (memorie neurală asociativă) pentru problema de asociere și RBF pentru aproximarea funcției la implementarea sistemului RAPTOR2000. Pentru implementarea SOM se folosesc cinci module. Patru module sunt utilizate la implementarea matricei PE, în timp ce cel de-al cincilea modul este folosit ca și controler de matrice. SRAM cu port dublu se utilizează pentru stocarea vectorilor de intrare. Folosind dispozitive XCV1000 se poate ajunge la 64 elemente de procesare pe un singur modul.

În acest caz (Poormann, Witkowski, Kalte, & Ruckert, 2002) a atins o performanță de 11,3 GCPS ceea ce este mult mai bun decât cele 80 MCPS atinse cu un calculator personal (AMD athlon, 800 MHz). Poormann (Poormann, Witkowski, Kalte, & Ruckert, 2002) a implementat de asemenea algoritmul SOM pe sistemul RAPTOR2000, folosind diferite dispozitive Xilinx Virtex. Operațiile aritmetice se realizează cu reprezentare în virgulă fixă cu o precizie de 16 biți. Pentru implementarea NAM, sistemul RAPTOR2000 este echipat cu șase module. În acest caz, cu fiecare modul se pot realiza 512 neuroni.

Unitățile neurale se compun dintr-o unitate de control, o unitate de prelucrare neurală, codificator și decodor de priorități și 128 Mbyte SRAM. Sinteza rezultată a arătat o folosire a spațiului de 3200 CLB pentru 512 neuroni. Neuronii din această implementare pot să funcționeze la 50 MHz și necesită 5,4  $\mu$ s pentru o singură asociere.

A existat o implementare FPGA inovativă bazată pe modelul rețelei neurale care își schimbă dimensiunile în mod dinamic, numită FAST (topologie de mărime flexibil adaptabilă) (Pérez-Uribe & Sanchez, 1996).

Cu cele mai multe modele de rețele neurale principalele probleme care se ivesc, sunt: determinarea numărului straturilor, a numărului neuronilor dintr-un strat și a faptului cum ele se interconectează (Pérez-Uribe & Sanchez, 1996).

Rețelele neurale ontogenice au fost create cu scopul de a înlătura aceste probleme prin oferirea posibilității de a schimba în mod dinamic topologia. ART (teoria rezonanței adaptabile) și GAR (creștere și reprezentare) sunt rețele neurale de acest tip, iar FAST derivă din aceste concepte.

Astfel, în lucrarea (Pérez-Uribe & Sanchez, FPGA Implementation of an Adaptable Size Neural Network, 1996) s-a implementat o rețea numită rețea neuronală ontogenică pe o placă FPGA proprie, denumită Flexible Adaptable-Size Topology (FAST). FAST a fost folosit pentru a implementa trei tipuri diferite de rețele neuronale nesupravegheate, ontogenice- adaptive resonance theory (ART), adaptive heuristic critic (AHC) și Dyna-SARSA.

Primele implementări FAST au folosit rețele neuronale bazate pe ART. Când se folosește la o problemă de segmentare a culorilor de imagini, patru neuroni FAST au segmentat cu succes o imagine 294x353, de 61 culori înfățișând Floarea soarelui de Van Gogh în patru clase de culori.

A doua implementare FAST a folosit o rețea neuronală bazată pe AHC (Pérez-Uribe & Sanchez, 1997). În această implementare, denumită FAST-AHC, s-au folosit opt neuroni pentru a controla un

pendul invers. Problema pendulului invers este exemplul clasic pentru un sistem instabil, folosit la testarea unor abordări ale controlului antrenării (Pérez-Urbe, 1999). FAST-AHC nu a putut să generalizeze la fel de bine ca și algoritmul de propagare înapoi a erorii, dar antrenarea este mai rapidă și mai eficientă. Acest lucru se datorează faptului că tehnica de antrenare AHC se poate generaliza ca și o formă a antrenării locale, unde doar nodurile active ale rețelei neuronale sunt actualizate, contrar algoritmului backpropagation, unde antrenarea este globală.

A treia implementare FAST a folosit a rețea neuronală Dyna-SARSA (Pérez-Urbe, 1999). Dyna-SARSA constituie un alt tip de antrenare de întărire, și a fost mult mai puțin intensivă din punct de vedere computațional, în comparație cu AHC, și foarte potrivită pentru implementarea digitală. Placa FAST Dyna-SARSA a fost integrată pe un robot mobil de sine stătător, și folosită ca și un neurocontroller pentru a demonstra o sarcină de antrenare-navigare. Neurocontrollerul FAST Dyna-SARSA a fost o încercare reușită, ajutând robotul mobil să evite obstacolele, care s-au adaptat la schimbări ușoare ale poziției obstacolelor.

Arhitectura FAST a fost prima din categoria sa, care folosește rețele neuronale cu antrenare nesupravegheate, ontogenice, dar se poate spune că arhitectura FAST are limitările sale, deoarece poate să rezolve doar probleme simple care necesită categorizare dinamică sau clasificare on-line.

Rețeaua FAST se compune din două straturi feed-forward conectate total și antrenarea este nesupervizată. Această rețea este potrivită pentru gruparea sau clasificarea datelor de intrare.

Se adaugă un neuron în stratul de ieșire dacă se găsesc vectori de intrare suficient de diferiți. Există o fază de curățire în care, în funcție de suprapunerea regiunilor sensibile ale unor neuroni învecinați, un neuron din stratul de ieșire este șters.

Rețeaua neurală FAST se compune dintr-un microcontroler 68331 și patru chipuri FPGA Xilinx XC4013. Pe aceste chipuri se implementează neuroni FAST, registre de mapare, generator de secvențe, și I/O.

Neuronul FAST se compune din trei blocuri diferite care execută câte un algoritm: calculul distanței, antrenare și curățire. Fiecare neuron include nouă sumatori de 8 biți și un singur multiplicator shift-add de 8 biți. Conține de asemenea un generator de numere aleatoare pentru procesul de eliminare. Generatorul de secvențe este un automat cu număr finit de stări care controlează executarea algoritmilor din cele trei stadii diferite.

Modelul de rețea FAST este aplicat la o problemă de învățare și recunoaștere a culorilor la imagini digitale. Aceasta necesită gruparea pixelilor de imagine după proprietățile de asemănare cromatică. În acest experiment, coordonatele de pixeli ale unei imagini sunt prezentate rețelei în mod aleator, din care rezultă segregarea imaginii în patru categorii, câte unul pentru fiecare neuron de ieșire. Se poate observa că fiecare vector de intrare se poate clasifica în 8  $\mu$ s. Rezultatul obținut cu implementarea hardware seamănă foarte mult cu cel obținut prin simularea software.

Ca o aplicație, problema cu N regine cu rețea Hopfield pe FPGA apare în lucrarea (Abramson, Smith, Logothetis, & Duke, 1998), în care este menționat faptul că implementarea rețelei Hopfield pe FPGA pentru rezolvarea acestei probleme diferă de alte implementări în mai multe aspecte.

În primul rând, ponderile sunt mici și se pot reprezenta folosind numere întregi. Astfel se reduce semnificativ dimensiunea unităților aritmetice. În al doilea rând, valorile ponderilor neurale se reduc la 0 sau 1. Acest lucru înlătură necesitatea unităților de multiplicare pentru că produsul vectorial devine simplă operație logică ȘI.

Arhitectura acestei implementări constă din 16 XC4010 conectate de 4 comutatori programabili (FPIC) pe o placă reconfigurabilă Aptix AP4. Folosind această configurație, este posibilă încărcarea unui proiect de până la 160000 porți. Neuronii și interconexiunile lor sunt specificați în VHDL, care este sintetizat în Exempler Logic's Galileo System Tools.

Performanța acestui sistem se compară cu alte două simulări software în „C” pentru problema celor 4, 6 și 8 regine. Rezultatul indică faptul că este posibilă atingerea unei creșteri în viteză a soluționării problemei de 2 sau 3 ori prin realizările hardware.

A existat de asemenea o încercare de a implementa o rețea neurală probabilistică pe sisteme FPGA pentru probleme de clasificare. Clasificarea automată a imaginilor multispectru în spațiu necesită un proces îndelungat din punctul de vedere al calculelor. Imaginile multispectru se obțin prin sateliți sofisticăți ale Sistemului de Observare a Pământului (Earth Observing System (EOS)) lansați de NASA pentru studierea mediului înconjurător. Unul din scopurile clasificării acestor imagini este divizarea solului pământesc în diferite categorii cum ar fi: urban, agricol, silvic, neproductiv etc.



În articolul (Abramson, Smith, Logothetis, & Duke, 1998) s-a implementat o rețea neurală probabilistă (PNN) pe FPGA pentru clasificarea imaginilor multi-spectru obținute de satelitul LANDSAT-2 EOS. Imaginile multi-spectru sunt formate de scanere și reprezintă seturi de imagini, fiecare corespunzând unei bande spectrale. Algoritmul de clasificare PNN implică calcularea probabilității pentru fiecare pixel pentru a-l încadra în una dintre clase. Ecuația pentru verificarea probabilității necesită calcule complexe de scădere, înmulțire și calcul exponențial a vectorilor de imagine.

Arhitectura clasificatorului PNN constă din două sisteme FPGA XC 4013E (numindu-le XFPGA și YFPGA), fiecare cu câte 13000 porți pe o placă acceleratoare X213. Datorită numărului limitat de porți, pentru efectuarea calculelor se folosește aritmetica în virgulă fixă. Lățimea căii de date în virgulă fixă este determinată prin simularea unor variabile cu operații pe bit în limbajul de programare C. În acest fel, valorile ponderilor de 10 biți sunt satisfăcătoare. Memoria ponderilor este realizată pe SRAM-ul de pe YFPGA.

Modulul YFPGA este compus din unitatea de scădere, unitatea de ridicare la putere și unitatea de acumulare. Unitatea XFPGA constă dintr-o unitate exponențială, o unitate de înmulțire și o unitate de acumulare a claselor.

Datorită lipsei de spațiu pe XFPGA, o unitate de comparație a claselor a fost mutată pe calculatorul gazdă. Este folosit un tabel de căutare pentru calcularea exponențialei negative. Ambele sisteme FPGA pot fi interfațate cu calculatorul gazdă printr-o magistrală PCI. Astfel calculatorul gazdă este capabil să configureze datele inițiale. Performanța acestui sistem este comparată cu două simulații software efectuate pe sisteme DEC și Pentium.

Timpul necesar pentru simularea algoritmului scris în „C” pe DEC rulând la 200 MHz și pe Pentium la 166 MHz s-a dovedit a fi de 22 respectiv 30 minute. Utilizând o placă acceleratoare X213 cu două dispozitive FPGA a redus timpul necesar implementării la 77 secunde, ducând la o accelerare mai mare cu un ordin de mărime.

De Garis et al (de Garis, Gers, & Korkin, 1997) (de Garis & Korkin, 2002) au implementat o rețea neuronală evoluționară bazată pe tehnici evoluționare, și au reușit să obțină o antrenare on-chip. De Garis a proiectat o placă bazată pe FPGA, denumită CAM-Brain Machine (CBM) unde se folosește un algoritm genetic (GA) pentru a antrena o rețea neuronală bazată pe un automat celular (cellular automata, CA). Chiar dacă CBM este privit ca o rețea care permite antrenare on-chip, nu a fost inclus nici un algoritm de antrenare în CA. În locul acesteia, antrenarea localizată se realizează în mod indirect printr-un algoritm genetic (în acest caz acesta inițializează datele de configurare pentru fiecare automat celular, care stabilește cum va crește rețeaua), urmat de „creșterea” topologiei rețelei neuronale, care este o caracteristică funcțională a automatului celular.

CBM susține până la 75 milioane de neuroni, fiind cea mai mare rețea neuronală evolutivă a timpului, unde mii de neuroni evoluează în fiecare secundă. CBM-ul s-a dovedit a fi un succes în aplicații precum aproximarea de funcții. Scopul pe termen lung al lui de Garis este folosirea CBM-ului la crearea unor rețele neuronale modulare foarte rapide, la scară largă, care se pot folosi la aplicații de construcție a creierului. De exemplu, de Garis plănuiește folosirea CBM-ului ca și un neurocontroller în reglarea în timp real al unei pisici-robot de mărime naturală denumit „Robokitty”.

Susținerea rețelelor neuronale modulare prin CBM este oarecum limitată, deoarece interacțiunea în aceste module sau conexiunile intermodulare trebuie definite manual, off-line.

Nordstorm (Nordstrom, 1995) a încercat să conceapă o rețea neuronală modulară pe o placă bazată pe FPGA, denumit REMAP (Real-time, Embedded, Modular, Adaptive, Parallel Processor). Nordstorm a considerat că calculul reconfigurabil se poate folosi ca și o bază potrivită pentru a susține cu ușurință diferite tipuri de module (diferiți algoritmi de rețele neuronale). Acest tip de mediu se poate folosi în crearea unor rețele neuronale modulare eterogene, cum ar fi „ierarhia experților” („hierarchy of experts”) propus de Jordan și Jacobs.

Cu excepția CAM-Brain Machine al lui de Garis, observațiile lui Nordstorm cu privire la acest domeniu rămân valabile chiar până astăzi. Dat fiind faptul că densitatea FPGA-urilor aflate la dispoziția lui Nordstorm la timpul cercetării sale era limitată, el a reușit să implementeze doar aplicații de modul singular pe REMAP. Încăzul ideal, Nordstorm ar fi dorit să susțină folosirea rețelelor neuronale modulare pe REMAP, însă acest scop a rămas neatins.

Diferența dintre REMAP- $\alpha$  și REMAP- $\beta$  este mărimea densității FPGA-ului folosit. La REMAP- $\alpha$  s-a folosit FPGA Xilinx XC3090 pentru prototipizarea diferiților algoritmi neuronali, în timp ce la REMAP- $\beta$  s-a folosit inițial Xilinx XC4005, dar care a fost înlocuit cu Xilinx XC4025. REMAP- $\beta$

s-a folosit ca și emulator hard de rețele neuronale și instrument de antrenare, care a putut să efectueze implementarea a următoarelor tipuri de rețele neuronale: rețele Adaline și Madaline, algoritmul backpropagation, rețele de memorie asociativă bidirecțională, memoria asociativă Hopfield, rețele counterpropagation, hartă autoorganizantă (SOM), teoria rezonanței adaptive (ART).

Pentru a crește viteza de procesare, s-a construit un al treilea și ultim prototip, REMAP- $\gamma$ , prin folosirea unei abordări bazate pe ASIC.

În următoarele paragrafe sunt prezentate două lucrări recente (Pandya, 2005) (Nichols, 2003) cu privire la posibilitățile de implementare a rețelelor neuronale artificiale.

Rețelele neuronale artificiale, și algoritmul backpropagation în special, reprezintă o formă a inteligenței artificiale care este caracterizată de antrenarea lentă și lipsa unei metodologii clare pentru determinarea topologiei rețelei înainte de începerea antrenării. Cercetătorii anteriori au folosit calculul reconfigurabil ca o metodă a accelerării testării rețelelor neuronale artificiale. În lucrarea (Nichols, 2003) este prezentat modul în care au ajutat avansarea în domeniu îmbunătățirile recente în instrumentele și metodologiile folosite la calculul reconfigurabil, și au susținut astfel aplicabilitatea lor la accelerarea RNA-urilor. S-a creat o arhitectură RNA nouă, bazată pe FPGA, denumită RTR-MANN, pentru a demonstra performanțele mai bune obținute prin folosirea instrumentelor și metodologiilor de generație curentă. RTR-MANN are o scalabilitate și o densitate funcțională de un ordin mai mare decât arhitecturile RNA bazate pe FPGA mai vechi. În plus, folosirea unei metodologii de proiectare a sistemelor noi (limbaj de nivel înalt) a condus la o fază de verificare/validare mult mai intuitivă, care a fost cu un ordin de magnitudine mai rapid decât simulatorii HDL tradiționali.

Contribuțiile pentru a atinge obiectivele lucrării au fost următoarele:

Analiza și concluzia faptului că folosirea unei aritmetici în virgulă flotantă de 32 biți nu este atât de fezabilă ca și o aritmetică în virgulă fixă de 16 biți, în cazul proiectelor FPGA de generație curentă. Această concluzie nu reprezintă o noutate în domeniu, dar este o concluzie actualizată în ceea ce privește FPGA-urile de generație curentă.

Aritmetica utilizată s-a proiectat în mod special pentru a fi folosit în RTR-MANN, care a conținut operatori aritmetici în virgulă fixă de 16 biți, care au fost optimizate din punctul de vedere al suprafeței pentru Xilinx XCV2000E Virtex-E FPGA. Această bibliotecă a ajutat la obținerea unei scalabilități și densități funcționale mult mai bune pentru RTR-MANN.

Au apărut alte probleme în cazul simulatorilor HDL de generație curentă. În special, timpii de simulare ai simulatorilor HDL de generație curentă s-au dovedit a fi foarte lungi (de ordinul zilelor sau săptămânilor în cazul proiectelor VLSI), ceea ce a dus la faze de verificare/validare anevoioase.

Algoritmul backpropagation (BP), folosit la construirea rețelelor neuronale artificiale a devenit foarte populară de la apariția sa la sfârșitul anilor 1980. Structura regulată și lărgimea domeniului de aplicabilitate a algoritmului BP a dobândit interesul cercetătorilor în tentativele lor de a obține implementări eficiente din punctul de vedere al utilizării timpului. Procesoarele de utilizare generală (General Purpose Processors, GPP) și circuitele integrate specifice aplicației (ASIC) sunt exemple relevante pentru RNA-uri bazate pe algoritmul BP. Totuși, aceste dispozitive de calcul suferă în mod constant de căutarea echilibrului între flexibilitate și performanță. În ultima decadă s-au obținut progrese semnificative în dezvoltarea unui anumit tip de hardware, o platformă reconfigurabilă bazată pe FPGA. FPGA-urile au o flexibilitate excelentă în termeni de reprogramabilitate ai aceluiasi hardware, și, în același timp, ating performanțe foarte bune în termeni de calcul paralel.

Cercetarea descrisă în lucrarea (Nichols, 2003) propune arhitecturi parțial paralele și o arhitectură total paralelă pentru realizarea algoritmului BP pe un FPGA. Proiectele propuse sunt codificate în Handel-C și verificate din punctul de vedere al funcționalității prin sintetizarea sistemului pe un chip FPGA Virtex2000e. Validarea proiectelor se realizează după diferite considerente. Arhitecturile parțial paralele și varianta total paralelă se dovedesc a fi de 2,5, respectiv 4 ori mai rapide decât implementările soft.

Într-o altă lucrare sunt propuse mai multe arhitecturi ANN pentru implementarea algoritmului BP pe FPGA Virtex2000e. Pe parcursul cercetărilor efectuate, s-au realizat inițial două arhitecturi în serie. Scopul dezvoltării acestor arhitecturi în serie a fost identificarea problemelor și stabilirea bazelor dezvoltării arhitecturilor paralele. După aceea, s-au realizat două arhitecturi parțial paralele. Scopul acestor arhitecturi a fost studiul accelerării în timpul acomodării unor rețele de dimensiuni variabile. Unele din problemele întâlnite la aceste arhitecturi parțial paralele au servit ca și îndrumare pentru dezvoltarea unei arhitecturi total paralele.

Operația de multiplicare-acumulare este operația de bază necesară executării algoritmului BP. Conceptele de mod Branch-In și Branch-Out ale operației de multiplicare-adunare au fost introduse în detaliu. Aceste două moduri formează baza dezvoltării paralele. S-au propus și s-au prezentat următoarele arhitecturi:

- SIPEX și SIPOB: abordare serială, diferă în stocarea parametrilor inițiali și tabele de căutare (LUT) pentru funcția sigmoid.
- PAROI, PARPP și PARIO: abordare parțial paralelă, diferă în atribuirea operațiilor modului Branch-In și Branch-Out la fazele algoritmului BP.
- FPAR: un proiect total paralel.

Generarea aleatoare a numerelor prin LFSR și funcția sigmoid, aproximată liniar prin trei variabile, a fost prezentată în detaliu.

Trei repere XOR, Iris și Cancer, respectiv mic, mediu, mare, au fost alese pentru validarea arhitecturilor. Contribuția majoră adusă este elaborarea proiectelor parțial și total paralele care ating performanțe de 2,5 și 4 ori mai rapide decât varianta cu implementare software. Implementările parțial paralele încercă să obțină un echilibru între spațiul de pe chip și performanță, atingând totuși CUPS de  $12 \times 10^6$ , mai mare decât la oricare dintre arhitecturile prezentate în bibliografie. O altă contribuție care merită menționată aici este implementarea cu succes a reperului Cancer care este o problemă reală pentru toate arhitecturile cu excepția FPAR. Validarea reperului Cancer pentru arhitectura FPAR necesită un număr foarte mare de module aritmetice, ceea ce are ca rezultat un consum de spațiu foarte mare.

S-a arătat prin calcule pentru un sistem multi-FPGA ipotetic, că FPAR-ul poate să atingă CUPS de  $51 \times 10^6$  și o viteză de 10 ori mai mare decât versiunea software pentru setul de date Cancer. O altă contribuție importantă este implementarea cu succes a problemelor XOR și Iris pentru arhitectura FPAR cu necesități hardware reduse de 0,4 respectiv 0,6 milioane porturi pe chipul Virtex XCV2000e FPGA. Pe parcursul cercetărilor prezentate în lucrarea (Pandya, 2005), s-a folosit Handel-C ca și limbaj de descriere hardware. Structura generală Handel-C contribuie la înlesnirea muncii de dezvoltare pentru un inginer hardware începător.

Handel-C permite o simulare rapidă a unor proiecte de dimensiuni mari în comparație cu simulatorul VHDL Modelsim. Timpul de dezvoltare este mult mai scurt decât în VHDL. Operațiile de depanare în Handel-C nu sunt ușor de realizat pentru arhitecturile RNA. S-a folosit metoda reprezentării numerelor în virgulă fixă. Problema cu compilatorul Handel-C este că nu poate să afișeze numerele în virgulă fixă, în schimb afișează valorile decimale pentru partea întreagă și partea fracțională. Deoarece algoritmul BP necesită operații aritmetice elaborate, depanarea parametrilor în numere în virgulă fixă devine un proces foarte complex.

Caracteristica de reconfigurare a unui FPGA în timpul rulării se poate utiliza pentru implementarea unor RNA-uri de dimensiuni mari (Pandya, 2005). Algoritmul BP este foarte potrivit pentru o astfel de reconfigurare deoarece se poate diviza în trei faze distincte în timp. Totuși, timpul necesar reconfigurării diferitelor faze pe un FPGA poate fi un impediment în asemenea proiecte. În același timp, faptul că un anumit FPGA poate să aibă caracteristica de reconfigurabilitate sau nu, este determinat de furnizorul plăcii respective. Arhitecturile descrise execută cele trei faze ale algoritmului BP în mod secvențial.

Chiar dacă paralelismul complet s-a aplicat la nivel de strat, un progres în creșterea vitezei în cazul versiunii soft nu a fost posibilă. Investigațiile se pot extinde asupra celor trei faze ale metodei pipeline: calculul ieșirii rețelei, propagarea înapoi a erorii și adaptarea ponderilor. Posibilitatea combinării parametrilor și a topologiei inițiale și descărcarea sa directă în Block RAM-ul unui FPGA trebuie cercetate. Instrumentul sintetic Xilinx asigură asemenea posibilități pentru platforme hardware specifice. Acest lucru va elimina necesitatea de a accesa off-chip RAM-ul ceea ce poate să rezulte într-o execuție mai rapidă a algoritmului. Arhitecturile prezentate în lucrarea (Pandya, 2005) pot utiliza alte limbaje descriptive hardware cum ar fi VHDL și performanțele lor în termeni de spațiu și viteză se pot compara. Viitoarele lucrări bazate pe arhitecturile FPAR se pot realiza pentru implementarea setului de date Cancer pe un hardware multi-FPGA dedicat.

Sistemele multi-FPGA sunt sisteme utilizate în mod frecvent la executarea algoritmului BP. Hardware-ul multi-FPGA nu doar permite implementarea unor rețele de dimensiuni mari, dar atinge de asemenea viteze mari de procesare și de adaptare. În lucrarea (Pandya, 2005) se propune în locul

reprezentării numerelor în virgulă fixă oferită de compilatorul Handle-C, reprezentarea în virgulă fixă separat și investigarea criteriului spațiu-viteză pentru operațiile aritmetice în virgulă fixă.

Ca și o concluzie, se poate spune că rețeaua neuronală folosită la implementarea bazată pe FPGA este o caracteristică importantă la clasificarea diferitelor arhitecturi. Tipul rețelei neuronale implementate depinde de aplicația care se va folosi la rezolvarea problemei în cauză. Ultimele tendințe din acest domeniu au arătat că au existat puține încercări pentru implementarea rețelelor neuronale bazate pe FPGA. Densitatea și viteza FPGA-urilor a crescut, până la punctul unde este rentabilă susținerea implementării rețelelor neuronale pe astfel de sisteme.

### 3.2.3. Implementarea funcțiilor de activare

Există multe tehnici cu privire la evaluarea funcțiilor elementare sau aproximativ elementare cum sunt aproximările polinomiale, algoritmi CORDIC, aproximări raționale, aproximări pe bază de tabele și multe altele (Omondi, 1994) (Muller, 1997). Pentru implementarea hardware a funcțiilor de activare, precizia, performanța și costul sunt cei mai importanți factori. Prin ultimii doi se înțelege că majoritatea tehnicilor performante care au fost dezvoltate în analiza numerică, și sunt ușor de implementat în software, nu sunt potrivite la implementare hardware.

*Dezvoltare în serie Taylor:*

Acesta se realizează prin aproximarea funcției sigmoid utilizând primele 4, 5 termene dintr-o serie Taylor. Publicația (Murtagh & Tsoi, 1992) propune următoarea formulă de serie Taylor pentru aproximarea funcției sigmoid, care poate fi utilizată în domeniul  $[-1, 1]$ .

$$y(x) = \frac{1}{2} + \frac{3}{4}x - \frac{1}{4}x^3 \quad \text{Ec. 45}$$

În graficele din Figura 3.17 sunt reprezentate funcția sigmoid

$$y(x) = \frac{1}{1 + e^{-ax}} \quad \text{Ec. 46}$$

$a=3.5$  respectiv funcția sigmoid aproximată prin dezvoltare în serie Taylor. Cu linie punctată (1) este reprezentată funcția sigmoid, iar cu linie continuă (2), aproximarea prin dezvoltare în serie Taylor. Dacă pentru parametrul  $a$  se ia valoarea 3.5 eroarea de aproximare este 0.0338.

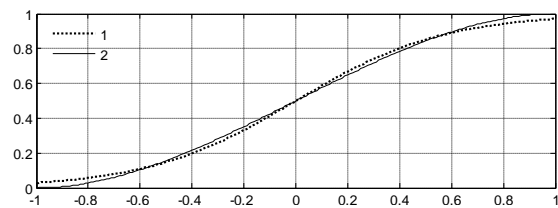


Figura 3.17 Aproximare Taylor a funcției sigmoid

*Tabele de căutare (LUT)*

Tabelele de căutare este un mod de abordare simplă, valorile funcției de activare sunt stocate off-chip sau on-chip în memorii de tip RAM sau Block RAM în FPGA.

Valorile potrivite sunt citite din valorile stocate prin introducerea adresei adecvate, adresă ce poate fi generată de un circuit simplu compus din comparatoare și registre de deplasare.

Avantajul acestei aproximări este costul redus de logică necesară, iar dezavantajul este suprafața mare necesară pentru a realiza LUT-ul în memorie.

*Computer digital de rotație coordonată (CORDIC)*

Schema CORDIC este o metodă iterativă bazată pe operații aritmetice și deplasare la nivel de bit. Necesită suprafață de substrat de siliciu mare. Eroarea de aproximare depinde de numărul de biți utilizați pentru codificarea valorii funcției de activare.

*Aproximarea liniară pe porțiuni (PWL)*

Aceasta se realizează prin aproximarea liniară pe domenii a funcției de activare. Această aproximare s-a dovedit a fi o alegere bună în ceea ce privește cerințele hardware, respectiv precizia de aproximare.

Aproximările PWL pot fi divizate în aproximări liniare de ordin întâi și aproximări de ordin superior. Publicațiile (Alippi & Storti-Gajani, 1991) (Myers & Hutchison, 1989) au fost primele care au prezentat schema pentru aproximarea prin metoda PWL a funcției sigmoid. Autorul (Alippi & Storti-

Gajani, 1991) a lucrat foarte mult la simplificarea funcției sigmoid pentru a obține un model realist, capabil să fie implementat în tehnologie VLSI. Publicația (Alippi & Storti-Gajani, 1991) descrie aproximarea funcției sigmoid cu metoda PWL prin selectarea unui set de valori pentru valorile “x” și calcularea valorilor “y” ca numere în funcție de putere a 2. Astfel luăm în considerare un singur segment pentru fiecare interval întreg, funcția se poate liniariza și descrie după cum urmează:

$$y(x) = 1 + \frac{1}{2^{|\lfloor x \rfloor|}} \left( \hat{x} - \frac{1}{2} \right) \tag{Ec. 47}$$

$\hat{x}$  : partea zecimală a valorii lui x cu semnul său

$\lfloor x \rfloor$  : partea întreagă a valorii lui x

Cu notațiile de mai sus, pentru axa negativă:

$$y(x) = 1 + \frac{1}{2^{|\lfloor x \rfloor|}} \left( \hat{x} + \frac{1}{2} \right) \tag{Ec. 48}$$

În Figura 3.18 este reprezentată funcția sigmoid (linie punctată 1) și aproximarea funcției prin segmente liniare (linie continuă 2) în intervalul [-4, 4]. Eroarea maximă de aproximare este 0.0189. Versiunea simplificată prezentată anterior poate fi simplu implementată în hardware cu un registru de deplasare controlat de un numărător. Acesta elimină necesitatea unui circuit dedicat de multiplicare. Publicația (Alippi & Storti-Gajani, 1991) a derivat formula de mai sus prin descompunerea funcției sigmoid în intervalul [-8, 8] în 15 segmente liniare.

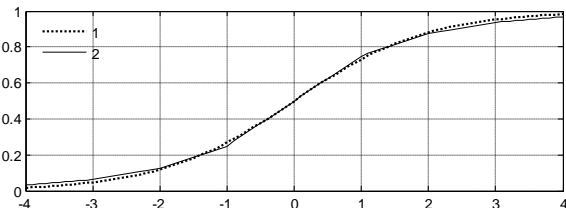


Figura 3.18 Aproximarea liniară pe porțiuni a funcției de activare sigmoid

Aproximarea propusă de lucrarea (Myers & Hutchison, 1989) este o curbă modificată care se bazează pe principiul legii A-law pentru sisteme PCM (Pulse Code Modulation).

Curba modificată compusă din 7 segmente a fost utilizată la aproximarea funcției. Publicațiile (Alippi & Storti-Gajani, 1991), (Myers & Hutchison, 1989) au prezentat de asemenea formula pentru calculul derivatei funcției sigmoid bazată pe segmentarea liniară a curbei.

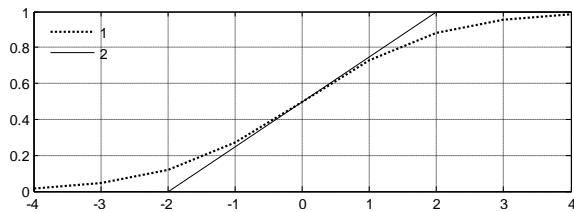


Figura 3.19 Aproximarea funcției sigmoid prin metoda CRI

Lucrarea (Basterretxea & Tarela, 2004) a propus o schemă de Aproximare Liniară Centrată (CRI) pentru o implementare optimă a funcției sigmoid și derivata acestuia. CRI este o schemă computațională secvențială pentru generarea funcției PWL. Autorii susțin că aproximarea PWL bazată pe modelul CRI se poate aplica la aproximarea oricărei funcții neliniare.

Figura 3.19 prezintă structura minimă inițială unde metoda CRI este utilizată la aproximarea funcției sigmoid. Este o aproximare simplă cu trei segmente, unde două segmente reprezintă saturația pentru valori de intrare mici și mari, în timp ce segmentul intermediar este tangenta funcției sigmoid de referință în punctul x=0.

Ecuatiile care descriu aproximarea sunt descrise mai jos.

$$y_1(x) = 0, \quad y_2(x) = \frac{1}{2} \left( 1 + \frac{x}{2} \right), \quad y_3(x) = 1 \tag{Ec. 49}$$

Cum se poate observa și din ecuațiile prezentate și conform Figura 3.19 liniarizarea cu trei segmente este o aproximare dură a funcției, deoarece eroarea de aproximare este mare la sfârșitul curbei înainte de saturație.

Autorii lucrării (Basterretxea & Tarela, 2004) au observat că numărul segmentelor poate fi crescut și poate fi calculat pentru fiecare nivel de interpolare astfel:

$$Nr. \text{deg mente} = 2^{q+1} + 1 \quad \text{Ec. 50}$$

Numărul de segmente se poate mări prin introducerea de perechi de tangente pe ambele părți a axei  $y$  la diferite valori  $|x|$ .

Deci numărul de segmente crește până la 5, 9, 17, 33 pentru  $q=1, 2, 3$  respectiv 4.

În modul de abordare prezentat de lucrarea (Sammur & Jones, 1991) s-a susținut utilizarea circuitelor de adunare și înmulțire rapidă pentru aproximarea funcției sigmoid prin liniarizare. Curba sigmoid este descompusă în patru segmente, pozitiv saturat, pozitiv nesaturat, negativ nesaturat, și negativ saturat. Partea neliniară dintre limitele de saturație este aproximată prin utilizarea unui polinom pătratic de gradul doi.

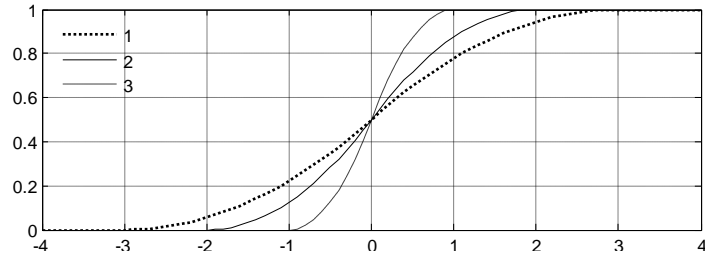


Figura 3.20 Aproximarea funcției de activare sigmoid prin polinom pătratic de gradul doi

$$\begin{aligned} y(x) &= 1 && \text{dacă } x \geq hl \\ y(x) &= 1 - g(hl - x)^2 && \text{dacă } x \geq 0 \text{ și } x < hl \\ y(x) &= 0 + g(hl + x)^2 && \text{dacă } x < 0 \text{ și } x \geq ll \\ y(x) &= 0 && \text{dacă } x < ll \end{aligned} \quad \text{Ec. 51}$$

unde valorile pentru  $ll$  și  $hl$  se pot obține din limitele de saturație a funcției sigmoid și

$$g = \frac{1}{2(hl)^2} \quad \text{Ec. 52}$$

Modelul propus de autorul lucrării (Sammur & Jones, 1991) reprezentat în Figura 3.20 este o aproximare de gradul doi a funcției sigmoid. În reprezentarea grafică, limitele de saturație a curbelor sunt 1, 2 și 3, corespunzând liniei punctate (1), liniei continue (2), și liniei întrerupte (3).

Publicația (Zhang, Vassiliadis, & Fris, 1996) a propus de asemenea un model care îmbunătățește aproximarea funcției în comparație cu varianta sa liniară de gradul I.

În particular, reduce eroarea medie și eroarea maximă a aproximării și necesită numai o singură operație de multiplicare. În acest model, intrările funcției sigmoid sunt descompuse pe segmente și o schemă de aproximare de ordinul II este utilizată pentru calcularea fiecărui segment.

Pentru un segment  $[\alpha, \beta]$ , dacă  $x$  este din intervalul  $x \in [\alpha, \beta]$  și  $y(x)$  este ieșirea aproximată, în general o aproximare de ordinul doi se poate exprima astfel:

$$y(x) = A + B(x + C)^2 \quad \text{unde } |C| = 2^{-n} \quad \text{Ec. 53}$$

unde  $C$  este o putere a lui 2, multiplicarea cu  $C$  se poate calcula prin operații de deplasare.

De asemenea este prezentat în detaliu un algoritm pentru definirea valorilor  $A$ ,  $B$  și  $C$ . Autorul a prezentat o formulă de aproximare a funcției împărțită pe două segmente  $(-4, 0)$  și  $(0, 4)$  care arată astfel:

$$y(x) = \begin{cases} 2^{-1}(1 - |2^{-2}x|)^2 & -4 < x < 0 \\ 1 - 2^{-1}(1 - |2^{-2}x|)^2 & 0 \leq x < 4 \end{cases} \quad \text{Ec. 54}$$

În modelul de mai sus se poate folosi o reprezentare cu bit de semn sau complement față de doi. Pentru reprezentarea numerelor s-a utilizat o reprezentare în virgulă fixă, folosind 4 biți pentru partea întregă, un bit pentru semn, zece biți pentru partea fracționară. În graficele din Figura 3.21 sunt reprezentate funcția sigmoid (linia punctată 1) și aproximarea funcției de activare prin două segmente neliniare (linia continuă 2). Eroarea

maximă de aproximare pentru funcțiile reprezentate este 0.0180. Eroarea medie în aproximarea funcției sigmoid cu această metodă este de ordinul  $10^{-3}$ , mult mai mică comparativ cu metoda de aproximare PWL de gradul întâi, caz în care eroarea medie este de ordinul  $10^{-2}$ .

În publicația (Beiu, Peperstraete, & Vandewalle, 1994) s-a subliniat faptul că, chiar dacă există metode de aproximare pentru simplificarea funcției sigmoid, calculele implicate sunt totuși complexe. Alți autori au introdus o funcție specială, o funcție sigmoid particulară. În această lucrare s-a arătat că funcția sigmoid particulară este echivalentă cu funcția sigmoid clasică dacă amplificarea este modificată cu o constantă. Funcția este descrisă astfel:

$$f^*(n) = \frac{1}{1 + 2^{-n}} \quad \text{Ec. 55}$$

Diferența dintre ecuația de mai sus și funcția sigmoid clasică este un factor constant de scalare, care se poate obține prin relația

$$e^{-z} = 2^{\frac{-z}{\ln(2)}} \quad \text{Ec. 56}$$

La implementarea funcției, întregul este scalat cu acest factor constant și adăugat la valoarea obținută prin formula prezentată mai sus pentru funcția sigmoid particulară. În graficele din Figura 3.22 sunt reprezentate funcția sigmoid clasică (linia punctată 1) și funcția sigmoid particulară pe bază de puterile lui 2 (linia continuă). Eroarea maximă de aproximare a funcției sigmoid clasice prin metoda menționată este 0,0808.

Publicația (Simard & Graf, 1993) a experimentat pentru valori de intrare limitate la intervalul  $[-8,8]$ , reprezentând cu 4 biți partea întregă și cu 8 biți partea fracționară. Eroarea calculată, adică diferența dintre funcția continuă și funcția cuantizată este în intervalul  $[-0.16, 0.16]$ .

Unii autori au propus versiuni simplificate ale funcțiilor cum ar fi tangenta hiperbolică și sigmoida rapidă. Lucrarea (Beiu, Peperstraete, & Vandewalle, 1994) a descris de asemenea un algoritm pentru derivata funcției sigmoid și a extins formula de calcul asupra altor funcții de activare menționate.

În cazul aproximărilor funcției de activare sigmoid clasice prin metodele prezentate apare o eroare de aproximare care este diferită de zero. Din punctul de vedere al rețelelor neuronale artificiale această eroare de aproximare nu prezintă nici o semnificație datorită faptului că rețeaua neuronală este antrenată, și nu este niciun indiciu din care să rezulte că erorile de aproximare ale funcțiilor de activare să se regăsească la ieșirea rețelei neuronale.

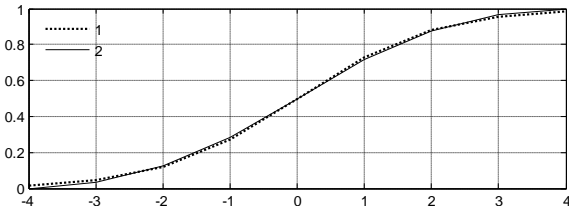


Figura 3.21 Aproximarea funcției de activare prin două segmente neliniare

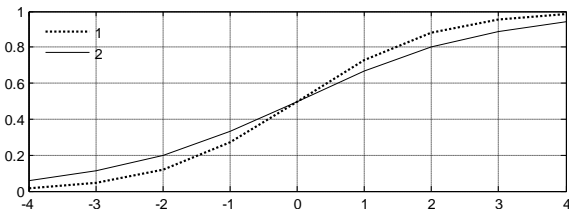


Figura 3.22 Aproximarea funcției sigmoid clasice printr-o funcție particulară pe bază de puterile lui 2

Din descrierile anterioare despre aproximarea funcțiilor de activare, rezultă că soluțiile clasice analizate sau implementate de cercetători sunt următoarele: tabele de căutare în RAM sau ROM, aproximarea PWL, aproximări cu grad superior, dezvoltare în serii Taylor, alte aproximări dedicate cum ar fi CORDIC.

Algoritmul CORDIC este cea mai studiată tehnică în privința implementării hardware a funcțiilor de activare, dar rareori o întâlnim utilizată. Avantajul este că același suport hardware poate fi utilizat pentru implementarea diferitelor tipuri de funcții de activare, dar performanța rezultantă de obicei nu este satisfăcătoare.

Aproximările funcțiilor de activare cu polinoame de grad ridicat duc la erori de aproximare reduse, dar nu sunt adecvate pentru implementare hardware din cauza numărului mare de operații aritmetice care trebuie efectuate pentru fiecare valoare, sau trebuie utilizate multe elemente hardware sau performanța este compromisă. O observație asemănătoare este valabilă și pentru implementările bazate pur pe metode tabelare. În cazul utilizării unui tabel de căutare de dimensiune redusă performanțele sunt nesatisfăcătoare, iar implementarea funcției printr-un tabel de dimensiuni semnificative costă. Așa stau lucrurile atât în cazul implementărilor în ASIC, cât și pe sisteme FPGA. Implementarea funcției prin combinarea funcțiilor polinomiale de ordin redus nu este o tehnică nouă. Întrebarea este cum să se aleagă cel mai bine punctele de interpolare în așa fel încât dimensiunea tabelului de căutare să rămână redusă.

Interpolările cu funcții polinomiale de grad redus au trei avantaje majore: aceeași structură hardware poate fi utilizată pentru implementarea diferitelor funcții în care numai coeficienții polinomului (conținutul tabelii de căutare) trebuie redefinite; pot fi ușor implementate pe sisteme FPGA care conțin module de multiplicare, sumatoare și tabele de căutare, respectiv se pot implementa relativ ușor prin utilizarea unui număr redus de elemente logice (tabele de căutare, sumatoare).

### 3.2.4. Aritmetică neurală/Reprezentarea datelor

Scopul acestei părți este clasificarea aritmeticii neurale folosite la sistemele FPGA și examinarea efectului său asupra resurselor sistemelor FPGA. Deoarece executarea algoritmilor de rețele neuronale (de ex. backpropagation) necesită multe operații aritmetice, alegerea schemei aritmetice și a reprezentării numerelor joacă un rol important în cazul în care avem în vedere proiectarea unui hardware de rețele neuronale.

Performanța RNA pe hardware depinde în foarte mare măsură de domeniul aplicațiilor în care se utilizează și precizia numerelor folosite pentru reprezentarea semnalelor. Aceste semnale includ parametrii inițiali cum ar fi cei de intrare, de ieșire, valorile ponderilor, și valori intermediare cum ar fi funcția de activare a fiecărui neuron, derivata funcției de activare și valorile erorii. Scopul creării unei scheme aritmetice pentru implementările hardware este de a prezenta semnalele într-un număr și cu o precizie suficientă care pot să ducă la convergența rețelei pe lângă păstrarea consumului de spațiu în limitele resurselor FPGA.

Vom descrie mai multe scheme de reprezentare a datelor și scheme aritmetice aplicate la arhitecturile prezentate în capitolele anterioare:

*Aritmetica fluxului de impulsuri (Pulse Stream Arithmetic):*

Modularea frecvenței impulsurilor (Pulse Frequency Modulation, PFM) este o schemă de codificare unde valorile circuitului sunt reprezentate de frecvența unor impulsuri de bandă îngustă constantă (Lysaght J. , Stockwood, Law, & Girma, 1994).

Semnalele astfel codificate se pot multiplica și aduna la fiecare nod folosind porți logice simple. Această tehnică se numește aritmetica fluxului de impulsuri și funcționează la sisteme FPGA fin granulate.

Figura 3.23 reprezintă valoarea fracționară a  $7/16$  prezentat prin flux de impulsuri. În lucrarea (Lysaght P. , Stockwood, Law, & Girma, 1994) s-a elaborat o arhitectură a neuronului bazată pe flux de impulsuri, pornind de la principiul aritmetic menționat. Valoarea de intrare (valoarea de activare) pentru fiecare neuron este un flux constant de impulsuri, pe când ponderile sinaptice sunt construite ca și funcții restrictive prin mascare OR selectivă a unei serii de semnale de tact.



Semnalele de tact sunt obținute de la ceasuri binare sincrone, fără suprapuneri, cu cicluri de  $\frac{1}{2}$ ,  $\frac{1}{4}$  etc.

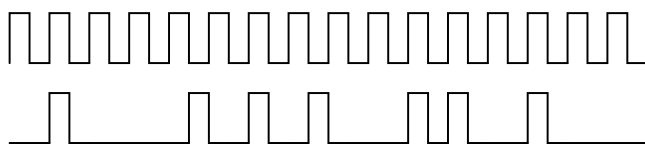


Figura 3.23 Reprezentarea valorii 7-16 prin flux de impulsuri

S-a folosit un generator de tact de 4 biți care poate fi folosit pentru construirea unor ponderi între 0 și 15/16.

Înmulțirea fluxului de impulsuri de intrare cu valorile ponderilor se poate face prin efectuarea simplei operații de AND între funcția de intrare și de mascare.

### Reprezentarea în virgulă flotantă

Această reprezentare a numerelor la implementările RNA este foarte asemănătoare cu reprezentările la un calculator de uz general. Reprezentarea în virgulă flotantă a implementărilor hardware ale RNA oferă un domeniu și o precizie mare la costul spațiului de circuit valoros de pe FPGA. De aceea, în trecut, cei mai mulți dintre cercetători au evitat folosirea unei aritmetici bazate pe reprezentarea punctelor în virgulă flotantă din cauza inaccesibilității sistemelor FPGA de capacitate mare.

Autorul (Simard & Graf, 1993) a folosit reprezentarea în virgulă flotantă pentru mai multe valori la implementarea algoritmului de propagare înapoi a erorii. Autorii au folosit mantisa de un 1 bit și exponent de 4 biți pentru rata de antrenare, mantisa de un 1 bit și exponent de 3 biți pentru valorile de stare activare și mantisa de un 1 bit și exponent de 5 biți pentru valorile gradientului. Chiar dacă această rețea a fost validată pe o platformă software prima dată, autorii susțin în mai multe ipoteze că se poate implementa foarte ușor pe un chip.

Publicația (Cloutier, Pigeon, & Boyer, 1996) a construit un procesor de imagini virtuale (VIP) care folosește sistem FPGA, Altera EPF81500 cu 1,5 MB RAM static pentru a realiza algoritmi neurali, aplicațiile de prelucrare a imaginilor și de recunoaștere a formelor.

Autorii (Sahin, Gloster, & Doss, 2000) a verificat flexibilitatea aritmeticii în virgulă flotantă aplicată la sisteme reconfigurabile. Au arătat în această lucrare că dezvoltarea recentă a tehnologiei FPGA oferă utilizatorului mai multe resurse pe un singur dispozitiv și, astfel, posibilitatea de a construi sisteme reconfigurabile complexe. În aceeași publicație (Sahin, Gloster, & Doss, 2000) s-au construit modulele lor în virgulă fixă, vectorul de adunare, scădere și înmulțire folosind VHDL și realizate pe dispozitiv FPGA Xilinx XC4044XL. Aceste module au viteze de 5 ori mai mari decât modelul software care rulează pe un calculator Pentium II 300MHz.

### Reprezentarea în virgulă fixă

Reprezentarea în virgulă fixă este alegerea cea mai răspândită pentru implementarea algoritmilor RNA pe hardware. Aritmetica în virgulă fixă a fost mai efektivă din punctul de vedere al spațiului utilizat în detrimentul reprezentărilor în virgulă flotantă. Deoarece multe aplicații de procesare a imaginilor și RNA funcționează în mod satisfăcător într-un domeniu restrâns și precizie mică, și reprezentarea în virgulă flotantă necesită o suprafață de siliciu mare, cercetătorii au folosit pe scară largă reprezentarea în virgulă fixă pentru aceste tipuri de aplicații.

Reprezintă un echilibru ideal între suprafața și necesitățile de domeniu-precizie pentru arhitecturile RNA bazate pe FPGA.

În aritmetica în virgulă fixă, există mai multe posibilități prin care calculele de bază pentru algoritmul backpropagation cum ar fi înmulțirea, adunarea și derivarea funcției de activare pot fi făcute.

### Aritmetica bit-serial

În continuare sunt prezentate câteva dintre tehnicile aritmetice pentru realizarea de hardware în virgulă fixă. Acest tip de aritmetică calculează câte un bit, în timp ce aritmetica paralelă calculează toți biții în același timp, simultan. Ocupă suprafață mică pe chip dar este foarte lentă. Această metodă s-a utilizat la multe implementări RNA pe FPGA. Arhitecturile clasice RRANN (Eldredge & Hutchings, 1994) și RENCO, binecunoscute implementări de RNA pe sisteme FPGA, au folosit acest tip de reprezentare.

Acest mod de abordare ajută la încadrarea cât mai multor elemente de procesare pe un FPGA de dimensiuni mici. Liniile de comunicare ale operațiilor aritmetice necesită o lărgime de un singur bit, și sumatorii necesită doar un bloc logic combinațional Xilinx (Xilinx Combinational Logic Block).

#### *Operații de decalare /fără multiplicare*

Această tehnică necesită reprezentarea numerelor în funcție de putere de 2.

Astfel înmulțirea/multiplicarea este implementată prin operații de decalare. A fost introdusă pentru prima dată de autorii publicației (Simard & Graf, 1993). Dezavantajul metodei constă în găsirea unor numere potrivite pentru minimizarea erorii. Publicația (Molz, Molz, Moraes, Torres, & Robert, 2000) a folosit acest mod de abordare pentru a prezenta intrările și ieșirea funcției de activare în format în virgulă fixă de 4 biți.

#### *Aritmetica on-line*

Această metodă folosește o reprezentare redundantă a numerelor ceea ce face posibilă realizarea unor operații aritmetice foarte rapide. Această reprezentare se bazează pe aritmetica în serie. De aceea ea prezintă toate avantajele aritmeticii bit-serial cu o viteză sporită. Ea efectuează o operație aritmetică în mod MSDF, prima dată transmițând cel mai semnificativ bit, spre deosebire de aritmetica în serie obișnuită care transmite prima dată datele cele mai puțin semnificative, printr-un mod de abordare LSDF.

Publicația (Girau & Tisserand, 1996) a explicat în detalii algoritmi pentru adăugare, înmulțire și funcția tangentă hiperbolică folosind transmisii MSDF (aritmetică on-line) cu reprezentarea numerelor redundantă. Prezentarea numerelor redundantă exprimă numerele în format radix-2.

Un astfel de format, definit ca și împrumută-salvează, conține setul de cifre -1,0,1.

În reprezentarea împrumută-salvează fiecare cifră  $a_i$  a unui număr  $a$  este reprezentat de 2 biți,  $a_i^+$  și  $a_i^-$  astfel încât  $a_i = a_i^+ - a_i^-$ .

De exemplu, 1 este reprezentat prin (1,0), pe când pentru 0 există două posibilități de reprezentare: (0,0) și (1,1).

Autorul din publicația (Skrbek, 1999) a propus o aritmetică neurală „deplasare adunare”, care oferă un set complet de funcții approximate potrivite pentru implementarea MLP pe FPGA. Elementele (calculule) de bază ale acestui mod de abordare sunt  $2^x$  și  $\log_2 x$ . Aceste funcții se bazează pe operații de deplasare în combinație cu aproximarea liniară.

Rețelele neuronale artificiale au, din natura lor, o arhitectură paralelă și sunt potrivite pentru implementare pe sisteme FPGA. O importantă problemă de implementare este să se determine formatul numeric utilizat care asigură un optim între precizie și suprafața de implementare. Reprezentarea standard în virgulă flotantă de precizie simplă sau dublă minimizează erorile de cuantizare, dar necesită resurse hardware semnificative. Reprezentarea în virgulă fixă, mai puțin precisă, necesită un număr mai redus de elemente hardware, dar adaugă erori de cuantizare care îngreunează antrenarea rețelei, în special în probleme de regresie.

### **3.3. Concluzii**

Creșterea continuă în densitate (număr porți/suprafață) a sistemelor FPGA a făcut posibilă realizarea unor proiecte de sisteme on-chip cu o complexitate mai mare de un milion de porți și RAM intern. De aceea, sistemele FPGA s-au dovedit în prezent a fi o platformă hardware atractivă pentru algoritmi RNA care necesită mult spațiu. Un alt avantaj al acestui dispozitiv este capacitatea de a combina programabilitatea cu viteza crescută a operațiilor asociate cu soluțiile hardware paralele.

Arhitecturile bazate pe circuite FPGA pot fi exploatate bine pentru realizarea de rețele neuronale artificiale, deoarece utilizând capacitatea lor de reconfigurare rapidă și concurentă se poate ajunge la performanțe bune în adaptabilitatea ponderală și chiar topologică. Cu toate acestea, luând în calcul că prima astfel de aplicație (Cox & Blanz, 1992) a apărut numai cu aproximativ un deceniu în urmă, realizarea de rețele neuronale de dimensiuni mari, cu mulți neuroni în circuite FPGA este o mare provocare chiar și în zilele noastre. Acest lucru se datorează faptului că majoritatea modelelor neuronale clasice necesită un număr mare de operații de multiplicare, ceea ce duce la realizări hardware foarte costisitoare (multe circuite multiplicatoare complexe). În pofida acestei probleme, există tehnologii care

exploatează proprietatea de reconfigurabilitate a acestor circuite și reușesc să implementeze ieftin și eficient RNA în acestea.

Orice realizare de rețea neuronală cu FPGA trebuie să se străduiască să utilizeze într-un anumit fel configurabilitatea acestor circuite, implicit optimizarea lor pentru problema propusă. Această proprietate care se poate aplica într-un timp scurt și de nenumărate ori, este de folos în special în dezvoltarea de *prototipuri și simulări*. Această flexibilitate permite realizarea rapidă de rețele neuronale cu diferite structuri și strategii de învățare, dar și demonstrarea unor principii prin simulări primordiale. Proiectul GANGLION (Cox & Blanz, 1992) este un bun exemplu în acest sens.

Procedurile de *creștere a densității* realizează o mai bună funcționalitate raportată la unitatea de arie FPGA, utilizând reconfigurarea. Acest lucru se poate valorifica dacă dispunem de un circuit FPGA capabil de reconfigurare – cel puțin parțială – în timpul funcționării (run-time/partial reconfigurability). Primul tip de astfel de procedură este *multiplexarea în timp*, ceea ce înseamnă că simularea RNA este descompusă în mai multe faze secvențiale, care corespund fiecare unei configurări diferite ale aceluiași circuit FPGA. Astfel se poate realiza un sistem în care circuitul este întotdeauna optimizat pentru faza care este în execuție. De exemplu, autorii din lucrarea (Eldredge & Hutchings, Density enhancement of a neural network using FPGAs and run-time reconfiguration, 1994) au construit o rețea neuronală feed-forward care utilizează regula de învățare back-propagation. Ei au împărțit algoritmul simulării în faza de reprezentare a intrărilor feed-forward, și faza de învățare back-propagation, care au rulat pe același circuit FPGA reconfigurat secvențial.

Al doilea tip de procedură pentru creșterea densității de implementare se bazează pe execuția de operații cu multiplicatori constanți și se numește *schimbare dinamică de constante*. Această metodă a fost utilizată în lucrarea (James-Roxby & Blodget, 2000), unde s-a arătat că, folosind ponderi de valoare constantă pentru a efectua foarte rapid operațiile necesare, s-a putut efectua ulterior acordarea sinapselor, în timpul reconfigurării circuitului ceea ce a durat mai puțin de 69 μs. Deoarece aceste două proceduri au avut ca scop primordial creșterea densității de implementare, în ceea ce privește o performanță superioară, nu putem avea pretenții față de aceste sisteme, numai în cazul în care timpul de reconfigurare a circuitului este neglijabil în comparație cu timpul utilizat pentru efectuarea calculelor. În circuitele FPGA reconfigurabile dinamice, există posibilitatea de implementare a rețelelor neuronale cu *adaptare topologică*. Altfel spus, putem construi RNA care suferă modificări arhitecturale succesive. În faza de învățare, pe lângă topologia rețelei avem posibilitatea de a acorda și precizia calculelor, conform funcției de criteriu utilizate de algoritmul de învățare.

Câteva grupuri de cercetători și-au propus realizarea RNA utilizând numere întregi pentru valorile ponderilor. Interesul față de această strategie a crescut deoarece complexitatea circuitelor de multiplicare necesare pentru efectuarea calculelor scade vertiginos, dacă utilizăm numere întregi în locul celor în virgulă mobilă. Se cunosc chiar și reguli de învățare, care utilizează puterile lui doi ca valori ale ponderilor. Avantajul mare al acestei idei este că operațiile de multiplicare necesare algoritmilor de învățare a RNA se pot efectua – în acest caz – foarte simplu, doar cu operații de deplasare a biților la stânga sau dreapta. Acele încercări, care au vizat implementarea RNA pe FPGA utilizând reprezentare în virgulă mobilă, au ajuns doar în faza de experimentare.

Neurohardware-ul digital diferă de asemenea în reprezentarea numerică a operațiilor aritmetice. Implementările în virgulă flotantă oferă o precizie mai bună, dar mai complexă și cu necesități de suprafață mai mare. Aritmetica în virgulă fixă este mai puțin complexă, cu necesități de suprafațe mai reduse, dar oferă o precizie scăzută față de implementarea în virgulă flotantă. Aritmetica în virgulă fixă poate fi utilă în cazul în care un algoritm RNA nu necesită precizie ridicată pentru implementarea unor aplicații simple.

Alegerea preciziei ponderilor este un pas foarte important în procesul de construcție a RNA cu ajutorul FPGA. Măsura preciziei ponderilor presupune întotdeauna un compromis între creșterea performanțelor rețelei și costul implementării acesteia. Aplicarea de ponderi cu precizie ridicată duce la erori de cuantizare mai reduse în implementarea finală, iar utilizarea de ponderi cu precizie mai scăzută rezultă în viteze de execuție mai mari, mai puține resurse utilizate și consum de energie mai moderat. În mod tradițional, valoarea optimă se află prin testări și evaluări de erori succesive (metoda empirică). Implementarea directă a unor funcții de activare neuronale, neliniare, în circuite FPGA este prea costisitoare. Există însă soluții mai practice, care garantează o precizie suficientă în cazul implementării acestor funcții în FPGA.



## 4. POSIBILITĂȚI DE IMPLEMENTARE TOTAL PARALELĂ A MODELELOR NEURONALE PULSATIVE CU CIRCUITE FPGA

### 4.1. Considerații inițiale

Există sisteme neuro-adaptive în timp real – clasificatoare, de control – la care învățarea este proces critic, a cărui finalizare trebuie garantată într-un timp finit, bine determinat. Exemple elocvente de astfel de sisteme pot fi roboții de recunoaștere, sisteme senzoriale ale sateliților aflați în medii necunoscute care caută obiecte, semnale interesante. Într-un astfel de mediu străin, sunt imposibil de prevăzut toate modelele posibile de tipuri de obiecte pe care automatul nostru clasificator le va întâlni. Rezultă, că sistemul nostru trebuie să posedă capacitatea de a încorpora cunoștințe noi, de a învăța să clasifice în timp real tipurile noi de obiecte întâlnite de-a lungul misiunii de recunoaștere.

În mod evident, s-a demonstrat de multe ori, că RNA pot fi utilizate excelent ca clasificatoare de precizie rapide. Cu toate acestea, însă RNA bazate pe modele de tip perceptron sau RBF (Radial Basis Function), necesită algoritmi de învățare iterativi. Dacă acestea sunt implementate pe mașini de calcul secvențiale, atunci însăși procesul de învățare va deveni un impediment în calea obținerii unor performanțe edificatoare (Chakrabarti, Roy, & Soundalgekar, 2001). De aici rezultă că procesul de învățare nu va fi convergent sau va dura prea mult pentru a fi de folos în sisteme în timp real. Acesta este motivul, pentru care RNA bazate pe învățare iterativă nu se pot aplica în sisteme în timp real la care procesul de învățare este unul critic. Concluzionând, putem afirma, că RNA oferă o paradigmă promițătoare: rezolvarea de probleme învățând din exemple. Ele realizează prelucrări de date foarte rapide, utilizând elemente masiv paralelizate. În timp ce o implementare software pe o platformă de tip Neuman este foarte utilă în evaluarea performanțelor unui anumit tip de rețea neuronală, în cazul unei implementări hardware este imperativ a ne concentra pe exploatarea paralelismului lor natural, pentru a le folosi în aplicații în timp real.

O rezolvare a problemei puse este algoritmul Kak (Tang & Kak, 2002), care permite învățarea instantanee. De asemenea, un algoritm plauzibil din punct de vedere al implementării hardware este și algoritmul FAST (Flexible Adaptable Size Topology) (Chakrabarti, Roy, & Soundalgekar, 2001) (Tang & Kak, 2002) (Pérez-Uribe, 1999). Rețeaua neuronală FAST este o rețea cu învățare nesupervizată, cu o topologie flexibil adaptabilă. De fapt este o rețea feed-forward, formată din două straturi complet conectate, unul de intrare și unul de ieșire. Constituirea claselor în care împarte rețeaua intrările primite se face pe baza corelației dintre acestea. În cazul unei intrări suficient de diferite de cele de până atunci, dimensiune rețelei se poate mări, adăugând un neuron la stratul de ieșire. Dimensiune rețelei se poate diminua eliminând neuronii cu activitate redusă din stratul de ieșire.

### 4.2. Caracteristici generale ale sistemelor hardware pentru implementarea rețelelor neuronale. Motivații pro și contra

#### 4.2.1. Implementări analogice

##### 4.2.1.1. Avantajele implementărilor analogice:

- viteză de procesare foarte mare
- integrabilitate foarte bună, densitate mare, neuronul analog are o complexitate mult mai redusă
- asigură o funcționare total paralelă, este foarte ușor de implementat din cauza suprafeței de circuit integrat utilizată (chip) reduse

##### 4.2.1.2. Dezavantajele implementărilor analogice:

- **precizie redusă** - (și la cea mai bună tehnologie precizia maximă este de 8 biți) fapt ce rezultă din următoarele cauze:

- curgerile de curenți în cazul capacitiv de stocare a ponderilor, suprafața condensatorului este redusă  $30 \times 30 \mu\text{m}$  de o capacitate  $1 \text{pF}$ , unde un curent de ordinul  $1 \text{pF}$  produce o modificare a tensiunii de  $1 \text{V/s}$
  - tranzistoarele utilizate în modulele de multiplicare analogică la curenți reduși nu sunt liniare, la curenți mai ridicați liniaritatea este mai bună dar apare o nouă problemă rezultată din cauza drift-ului termic.
  - imperfecțiunea procesului de fabricare a circuitelor analogice
  - **sensibilitate**, zgomot, diafonie, asuprirea tensiunii de alimentare, drift-ul de temperatură,
  - **modificarea ponderilor** ridică probleme
  - deficiența adaptivității - din cauza modificării ponderilor și din lipsa algoritmilor continue de antrenare (majoritatea algoritmilor cunoscuți de antrenare discretă în timp)
    - posibilitatea de **legare în cascadă** a chip-urilor este greoaie
    - cost ridicat de fabricare - modificarea parametrilor de fabricație trebuie să fie foarte redus.
- Testul după fabricare este complex, și este posibil numai cu introducerea unui modul de autotestare care ridică suprafața chip-ului și costurile de fabricare și proiectare
- **proiectare greoaie** - necesitând experiență de proiectare ridicată, și o mulțime de simulări din cauza numărului ridicat de parametri a căror efect trebuie luat în considerare.

Precizia obținută cu circuitele analogice rămâne sub nivelul preciziei ce se poate obține cu sistemele digitale. Deoarece în sistemele analogice prezența zgomotului influențează direct valoarea semnalului purtător de informație trebuie avut grijă la reducerea sensibilității la zgomot a rețelei și algoritmului de antrenare. La implementările analogice una dintre problemele majore constă în stocarea ponderilor modificabile. Aceste unități de stocare trebuie să satisfacă simultan următoarele condiții:

- rezoluție și precizie ridicată
- stabilitate ridicată pe toată durata de exploatare, posibilitatea de stocare și în lipsa tensiunii de alimentare
- posibilitate de modificare rapidă și accesibilitate rapidă din exterior

La cele mai multe sisteme ponderile rețelei sunt constante sau pot fi modificate numai din exterior prin intermediul unui calculator.

## 4.2.2. Implementări digitale

Implementările digitale sunt mult mai răspândite față de sistemele analogice. Implementările digitale se pot realiza cu sisteme hardware relativ simple, astfel verificarea sistemului obținut poate fi mai ușor efectuată față de sistemele analogice. Totodată cu implementare hardware se obține o viteză de procesare ridicată, dar totuși sub viteza de procesare obținută prin sistemele analogice sau optice.

O parte din dezvoltările de hardware digitale de rețele neuronale de fapt înseamnă proiectarea unor sisteme de procesare paralelă multiprocesor, a căror arhitectură ia în considerare proprietățile speciale a structurii rețelei neuronale utilizate. O altă posibilitate a rețelelor neuronale este implementarea neuronilor prin circuite digitale ca elemente de circuite logice programabile respectiv circuite VLSI cu un grad ridicat de integrare.

### 4.2.2.1. Avantajele implementărilor digitale

- precizie mare - precizia este în funcție de numărul de biți utilizați care teoretic poate fi foarte mare.
- este ușoară stocarea și modificarea parametrilor rețelei
- în varianta digitală pot fi implementate sisteme mai mari față de sistemele analogice, permite o proiectare flexibilă
- proiectare simplă și rapidă - proiectarea și testarea este simplificată prin utilizarea unor programe evoluate pot fi utilizate elemente logice programabile, circuite „gate array”, microprocesoare
- stabilitate ridicată – rezultată din cauza modului de funcționare sincronă
- insensibilitate la zgomot.

#### 4.2.2.2. Dezavantajele implementărilor digitale

- viteză de funcționare mai redusă față de sistemele analogice
- elementele de procesare (multiplicatoarele) ocupă o suprafață de chip mare, realizarea conexiunilor, și a unui număr ridicat de semnale digitale ocupă o suprafață însemnată din suprafața circuitului integrat
- Proprietățile mai importante a sistemelor optice din punctul de vedere a realizării rețelelor neuronale
- proprietățile de liniaritate a sistemelor optice relativ sunt bune
- prin utilizarea luminii ca purtător de informație se ajunge la viteză ridicată de procesare
- realizarea operațiilor de combinări liniare rapide este simplă
- se poate realiza un număr mare de conexiuni
- implementarea neliniarităților este problematică
- sistemul este încetinit de interfața dintre sistemul optic și mediu
- realizarea ponderilor cu valori modificabile este greoaie

#### 4.2.3. Calificarea implementărilor

Evaluarea implementărilor neuronale prin indicatoare standard (MFLOPS, MIPS) ne dă informație limitată despre sistem. Caracteristicile mai importante a chip-urilor neuronale este suprafața chip-ului, consumul și viteza. Pentru caracterizarea implementării din punctul de vedere a vitezei s-au introdus următoarele indicatoare:

- ❖ Connections Per Second (CPS), - calculul sinapselor într-o secundă
- ❖ Connection Update Per Second (CUPS) - Numărul de ponderi modificate în faza de învățare, este caracteristic vitezei de învățare a rețelei.

Față de cele două indicatoare cel mai mult răspândite se mai obișnuiește utilizarea acestor indicatori normați cu numărul sinapselor, respectiv înmulțit cu numărul de biți din reprezentarea intrărilor și a ponderilor.

#### 4.2.4. Implementări RNA pe FPGA

Progresul în curs al microelectronicii este forța motorică a dezvoltării continue a produselor tehnice noi. Sistemele FPGA sunt inovații de acest fel în sistemele microelectronice pentru aplicațiile computerizate. În ultima decadă a devenit o metodă realizarea unor algoritmi bogați din punctul de vedere al calculului pe FPGA-uri.

Creșterea continuă în densitate a sistemelor FPGA a făcut posibilă realizarea unor proiecte de sisteme on-chip cu o complexitate mai mare de un milion de porți și memorii RAM interne. De aceea, sistemele FPGA s-au dovedit în prezent a fi o platformă hardware atractivă pentru algoritmi RNA care necesită mult spațiu.

Un alt avantaj al acestui dispozitiv este capacitatea de a combina programabilitatea cu viteza crescută a operațiilor asociate cu soluțiile hardware paralele.

În chip-urile FPGA există două moduri de implementare a logicii reconfigurabile. Sunt binecunoscute CTR (reconfigurarea compilată în timp) și RTR (reconfigurarea în timpul utilizării). CTR este o strategie de implementare statică, la care fiecare aplicare constă într-o singură configurare. RTR este o strategie de implementare dinamică unde fiecare aplicație constă în multiple configurații cooperative.

În CTR, odată ce o operație începe, configurația nu se mai schimbă în timpul procesului. Această abordare este asemănătoare cu implementarea unei aplicații în ASIC. Pe când RTR utilizează o schemă de alocare dinamică care re-allocă partea de hardware în timpul executării aplicației. Fiecare aplicație constă în mai multe configurații pentru FPGA.

Aceste configurații sunt încărcate succesiv în timpul derulării unei aplicații pe chip-ul FPGA. Acest lucru înseamnă realizarea unui nou hardware pentru fiecare configurație pentru porțiuni particulare ale aplicației. Există două modele RTR: Globală și locală. RTR global alocă întreaga resursă a întregului chip FPGA pentru fiecare pas al configurației. Aceasta duce la o nouă realizare pe FPGA în timpul fiecărei configurații în timpul utilizării. Proiectantul trebuie să implementeze o aplicație în RTR global și să împartă aplicația în părți aproximativ egale pentru a utiliza în mod eficient resursele

FPGA. Aceasta se numește partiționare temporală a unei aplicații. În RTR local, o aplicație reconfigurează local niște subseturi ale logicii în timpul executării aplicației. Poate configura orice procent din resursele reconfigurabile în orice moment. Organizarea aplicărilor în RTR local se bazează mai mult pe o divizare funcțională a muncii decât pe partiționarea utilizată de aplicațiile RTR global. Odată ce partiționarea temporală manuală este decisă, poate fi foarte ușor implementată în varianta RTR globală.

Implementarea RTR local nu este posibilă în instrumente CAD.

Tabelul 7 de mai jos oferă o recapitulare a celor discutate în această secțiune.

Tabelul 7 Performanțele diferitelor platforme pentru implementări de rețele neuronale hardware

	<i>ASIC analog</i>	<i>ASIC digital</i>	<i>FPGA</i>	<i>Bazat pe procesor</i>	<i>Computer paralel</i>
<i>Viteză</i>	+++	++	+	-	+
<i>Resurse</i>	+++	++	+	-	--
<i>Cost</i>	--	--	++	++	--
<i>Timp de design</i>	--	--	++	+++	+
<i>Fiabilitate</i>	--	+	++	++	++

### 4.3. Sinteză asupra diferitelor implementări hardware de rețele neuronale artificiale

Acest capitol va prezenta câteva exemple de implementări hardware de rețele neuronale artificiale, cu accent pe cele neuromorfe, acestea din urmă reprezentând obiectul primar al tezei de doctorat.

În primii ani ai cercetărilor legate de rețelele neuronale artificiale, s-a presupus, că va fi nevoie de implementarea unor sisteme hardware specializate, pentru a beneficia de funcțiile promițătoare ale acestor modele. De asemenea, s-a prevăzut, că aceste sisteme vor fi probabil circuite analogice, și vor conține multiple elemente de procesare paralele cu un număr mare de conexiuni între ele.

Aceste presupuneri nu s-au adeverit în totalitate, însă, deoarece creșterea imensă a capacității de calcul a mașinilor convenționale von Neuman în domeniul digital a permis implementărilor software să atingă succese răsunătoare în nenumărate aplicații.

Între timp, dezvoltarea echipamentelor hardware de rețele neuronale a fost mult mai lentă, având un succes comercial mult mai modest. Sunt însă câteva contraexemple, dintre care voi prezenta pe scurt, în cele ce urmează, cele mai relevante.

S-ar putea pune întrebarea: de ce să implementăm hardware rețele neuronale, știind că procesoarele convenționale PC (de ex. cele von Neuman, seria Pentium) continuă să arate un ritm de dezvoltare dramatic? Evident, motivația este în continuare VITEZA, care poate fi dusă dincolo de limita procesării în timp real doar de circuitele neurale specializate, ASIC-urile analogice, digitale sau cele hibride.

Aplicațiile de rețele neuronale software au un succes comercial din ce în ce mai semnificativ. De exemplu, orice scanner nou cumpărat se livrează cu un software de recunoaștere a caracterelor (OCR - Optical Character Recognition) care cel mai probabil utilizează algoritmi bazați pe rețele neuronale. Acești algoritmi nu sunt însă realizați exclusiv cu inteligență artificială, procesul de recunoaștere a caracterelor fiind mult mai complex, presupunând mai mulți pași de prelucrare a imaginii scanate (vezi Figura 4.1).

În termeni generali se poate afirma, că sunt foarte puține acele probleme care se pot rezolva exclusiv cu o singură rețea neuronală.



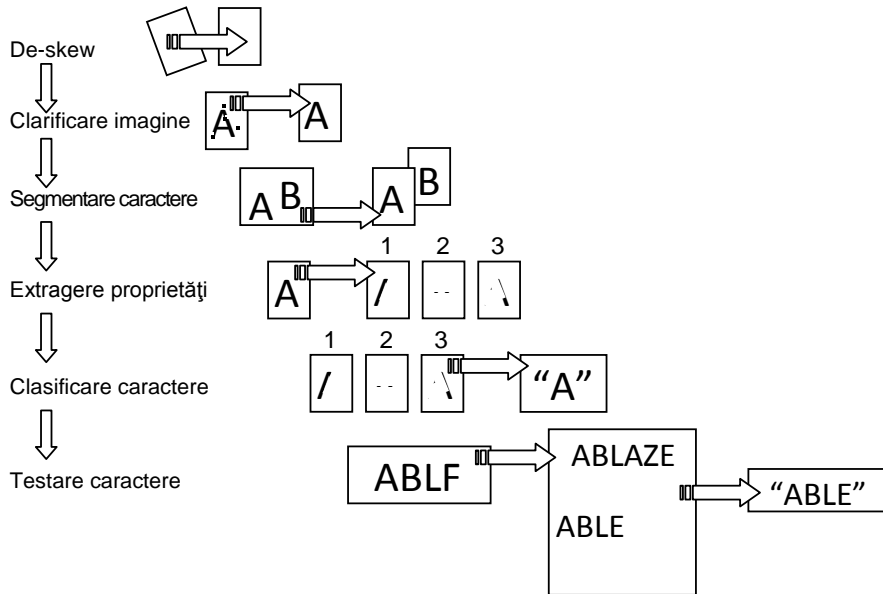


Figura 4.1 Pașii OCR, conform manualului de utilizare al Adaptive Solutions CNAPS

factorul de accelerare a execuției este de numai 2. Numai dacă 90%, sau mai mult din program este executat în paralel, va ajunge factorul de accelerare la valori de peste 10. Aceste studii au sta la baza dezvoltării rapide a sistemelor cu procesare masiv paralelizată (MPP - massively parallel processing).

#### 4.3.1.2. Legea lui Gustafson

Cunoscută de asemenea ca legea Gustafson-Barsis, este o lege din domeniul ingineriei calculatoarelor care afirmă ca orice problemă de dimensiuni suficient de mari poate fi paralelizată în mod eficient. Legea lui Gustafson este strâns legată de legea lui Amdahl, care dă o limită gradului cu care un program poate fi accelerat cu ajutorul paralelizării. Această lege a fost prima dată descrisă de John L. Gustafson în 1988 (Gustafson, 1988).

$$S(P) = P - \alpha \cdot (P - 1)$$

Ec. 57

unde  $P$  este numărul de procesoare,  $S$  este accelerarea, iar  $\alpha$  este partea procesului ce nu poate fi paralelizată.

Legea lui Gustafson abordează deficiențele legii lui Amdahl, care nu poate scala astfel încât să egaleze disponibilitatea puterii computaționale odată cu creșterea în dimensiuni a mașinii. Eliminată problema dimensiunilor fixe sau a sarcinii computaționale fixe pe procesoarele paralele: în locul acestora propune un concept de timp fix care duce la scalarea accelerării pentru probleme de dimensiuni mai mari (adică scalare slabă sau moale).

Legea lui Amdahl (Amdahl, 1967) se bazează pe volum de lucru fix, sau dimensiuni ale problemei fixe (adică scalare puternică sau tare). Ea implică faptul că partea secvențială a unui program nu se schimbă în raport cu dimensiunea mașinii (adică numărul de procesoare). Însă partea paralelă este distribuită în mod egal la  $n$  procesoare.

Impactul legii a fost simțit prin schimbările în cercetare pentru dezvoltarea de compilatoare care paralelizează și reducerile în partea serială a soluției pentru a impulsiona performanțele sistemelor paralele.

#### 4.3.1.1. Legea lui Amdahl

Prezentată în 1967 (Amdahl, 1967) iar apoi versiunea acesteia conform lui Gustafson din 1988 (Gustafson, 1988), au arătat, că paralelizarea soluției unei probleme este eficientă doar dacă partea cea mai semnificativă a acesteia se poate implementa astfel. Spre exemplu, să presupunem, că 50% din aplicația mai sus amintită, programul OCR, poate fi executată de un sistem paralel. În acest caz

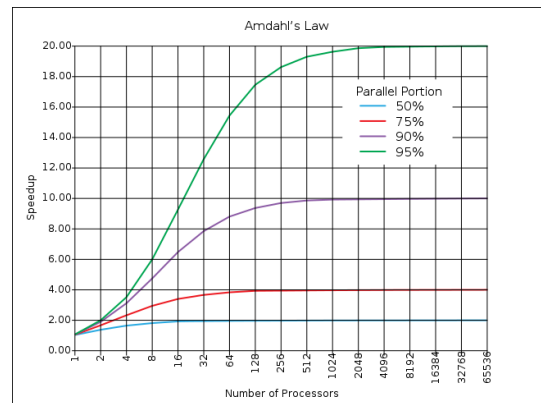


Figura 4.2 Legea lui Amdahl

## Implementarea legii lui Gustafson

Fie  $n$  o măsura a dimensiunilor problemei.

Execuția programului pe un computer paralel se descompune în:

$$a(n) + b(n) = 1 \quad \text{Ec. 58}$$

unde  $a$  este fracțiunea secvențială și  $b$  este fracțiunea paralelă.

Pe un computer secvențial, timpul relativ ar fi  $a(n) + p \cdot b(n)$  unde  $p$  este numărul de procesoare în cazul paralel.

Astfel, accelerarea este:

$(a(n) + p \cdot b(n))$  în cazul paralel, și luând în considerare relația  $a(n) + b(n) = 1$ , obținem:

$$S = a(n) + p \cdot (1 - a(n)) \quad \text{Ec. 59}$$

unde  $a(n)$  este funcția serială.

Presupunând că funcția serială  $a(n)$  se micșorează odată cu dimensiunea problemei  $n$ , atunci accelerarea se apropie de  $p$  când  $n$  se apropie de infinit, exact cum s-a dorit.

Legea lui Gustafson argumentează că chiar și folosirea sistemelor de calculare masiv paralele nu influențează partea serială, și tratează această parte ca o constantă. În comparație, ipoteza legii lui Amdahl rezultă din ideea că influența părții seriale crește odată cu numărul de procese.

O metaforă de conducere auto pentru a explica aceste concepte ar fi următoarea:

*Legea lui Amdahl sugerează aproximativ că:*

„Presupunem că un autoturism circulă între două orașe aflate la o distanță de 60 mile între ele, și a petrecut deja o oră parcurgând jumătate din distanță la viteza de 30 mph. Indiferent de cât de rapid veți conduce pe ultima jumătate de drum, este imposibil să atingeți viteza medie de 90 mph înainte de a ajunge la destinație. Deoarece a trecut deja o oră și distanța este de 60 mile în total; chiar și circulând cu o viteză infinit de mare nu veți putea depăși viteza medie de 60 mph.”

*Legea lui Gustafson sugerează aproximativ că:*

„Presupunem că un autoturism a circulat deja pentru o anumită perioadă de timp cu o viteză mai mică de 90 mph. Având la dispoziție suficient timp și distanță de parcurs, viteza medie a autoturismului poate eventual atinge oricum valoarea de 90 mph, indiferent de cât de mult sau de încet a circulat deja până acum. De exemplu dacă autoturismul a petrecut deja o oră la viteza de 30 mph, ar putea atinge media dorită circulând cu 120 mph pentru încă două ore, sau cu 150 mph pentru încă o oră, și așa mai departe.”

### 4.3.2. Implementări hardware neuronale analogice

Implementarea hardware analogică se remarcă prin abilitatea de a efectua calcule (adunări, înmulțiri, transformări neliniare) extrem de rapid. *Viteza ridicată de procesare* și de eficiența calculului bazate pe principiile clasice ale teoriei circuitelor electrice contrabalansează principalul dezavantaj al metodelor analogice: *imposibilitatea* obținerii unei *precizii* ridicate.

Unul din elementele cele mai importante pentru o rețea îl reprezintă *interconexiunile* dintre neuroni. În mod tipic numărul de neuroni ce pot fi integrați pe un chip este limitat de numărul de conexiuni posibile. Aceste interconexiuni pot să fie *fixe*, implementate prin *rezistori* de o anumită valoare sau *programabile*. Prima alternativă se folosește în aplicațiile în care este cunoscută arhitectura RNA și valorile ponderilor.

Implementările analogice a sistemelor de rețele neuronale sunt stimulate de viteza mare de prelucrare. De obicei în aceste sisteme semnalele sunt funcții continue de timp, iar purtătorul de informație este o mărime analogică - de exemplu intensitatea curentului. VLSI analog este un instrument potrivit pentru implementarea sistemelor neuronale inspirate din modele biologice (de exemplu sisteme neuronale de procesare sunet și imagine). Cu instrumente analogice se poate realiza foarte ușor însumarea ponderată respectiv caracteristica neliniară.

Pentru implementarea analogică a neuronului tipic se utilizează un amplificator operațional, iar ponderarea intrărilor este realizată de raportul rezistențelor legate în serie corespunzătoare intrării. Ponderile sinaptice sunt fixe (realizate prin rezistențe cu valori fixe), potrivite pentru realizarea rețelelor fixe.

Principiul de bază a construirii sistemelor neuronale analogice antrenabile și utilizării lor în sistemele adaptive este implementarea setului de ponderi cu valori modificabile. În literatura de specialitate se regăsesc două tehnici specifice:

- ❖ valorile ponderilor sunt stocate în condensatoare. În cazul utilizării condensatoarelor ca elemente pentru stocarea ponderilor precizia de implementare și timpul de stocare este redus din cauza scurgerilor de curenți, scurgeri care pot descărca condensatorul.
- ❖ utilizarea tranzistoarelor MOS cu poartă grilă flotantă: sunt tranzistore cu mai multe porți grilă, iar una dintre porți este izolată. Sarcina trimisă pe poarta izolată se conservă, iar rezistența canalului MOSFET poate fi continuu modificată în funcție de sarcină. Timpul de stocare este mult mai mare și scurgerile de curenți sunt mai mici.

Această tehnologie este utilizată și la chip-urile ETANN proiectate de INTEL.

Pentru multe aplicații precizia de stocare a ponderilor prin metodele analogice amintite mai sus nu este îndeajuns. Din acest motiv ponderile sunt stocate digital. Ponderile stocate digital sunt trimise pe intrarea multiplicatorului analogic printr-un convertor digital analog, sau se utilizează convertoare digital analogice cu înmulțire MDAC (Multipliyaing D/A converter).

Abordări VLSI analogice ("AVLSI") mai recente marchează o orientare spre mimarea în detaliu a comportamentului neuronului biologic. Implementarea s-a realizat în tehnologie de 10μm nMOS.

În același context ca și abordarea precedentă se regăsește și implementarea în siliciu a dendritelor unui neuron. Conductanța axială este implementată prin intermediul unei capacități comutate iar conductanța trans-membrană este realizată printr-un amplificator transconductanță.

### 4.3.3. Implementări hardware neuronale digitale

Există două principale avantaje ale acestei implementări față de versiunea analogică:

**Ușurința proiectării și realizării.** Folosirea tehnologiei VLSI digitale oferă o precizie ridicată, stocarea facilă a ponderilor și un raport performanță/cost mai bun decât varianta analogică.

**Flexibilitatea.** Implementarea VLSI digitală permite utilizarea unor algoritmi complecși de o mare diversitate, măbind astfel gama aplicațiilor posibile.

Unul dintre *dezavantajele* majore ale tehnologiei VLSI digitale îl constituie **consumul ridicat** de putere și arie de siliciu a blocurilor de multiplicare. Una dintre soluțiile pentru reducerea suprafeței necesare o constituie multiplexarea interconexiunilor.

În cazul implementărilor digitale ajungem la sisteme cu eșantionare, fiecare semnal este discret în timp și amplitudine. Implementările digitale sunt motivate de concepția conecționistă vizavi de realizările analogice motivate de modelele biologice. Se poate privi ca o realizare digitală acel hardware, pe care rulează un program neuronal. Avantajele oferite de rețele neuronale pot fi exploatate cu adevărat numai în cazul unui hardware special proiectat. Din punctul de vedere de al arhitecturii chip-urilor neuronale acestea se pot clasifica astfel:

- **chip neuronal multiprocesor** - chip-urile neuronale multiprocesoare sunt realizat cu ajutorul mai multor procesoare care pot funcționa în diferite moduri de lucru: SIMD, MIMD, matrici sistolici.
- **Arhitectură felie (slice).** Arhitectura de proiectare a sistemelor digitale bazată pe slice-uri este o metodă bine testată și utilizată cu succes. Elementele funcționale sunt descompuse pe slice-uri care ușor pot fi combinate în cascadă. Din slice-uri se pot construi elemente de dimensiuni preferabile. Avantajul principal al acestei abordări este că permite o proiectare modulară, reduce timpul și costurile de proiectare.
- **LSI waffer-scale.** În cazul tehnologiei WSI, analog fabricației circuitelor integrate, slice-urile sunt descompuse pe chipuri dar se realizează și conexiunile dintre chipuri și de alimentare. Aceste resurse la sistemele WSI LSI sunt supra-proiectate și cu redundanță ridicată. La sfârșitul procesului de fabricare chipurile cu defecte sunt neutralizate, și circuitul este utilizat întreg fără dezmembrare în chip-uri.
- **Chip-uri speciale** - Aceste chip-uri au fost dezvoltate pentru rețele și aplicații specifice. Avantajul lor este că numai blocurile și unitățile necesare sunt implementate față de neurocalculatoarele de uz general.
- **Soluții dedicate.** Scopul soluțiilor dedicate este creșterea vitezei de procesare. Posibilitatea de creștere a vitezei de procesare depinde de utilizarea unei proiectări orientată pe metoda țintă prin

care se mărește viteza de prelucrare a unităților individuale totodată este posibilă amplasarea mai multor unități de procesare pe un singur chip

Soluția optimală poate fi remarcată prin proiectarea de către utilizator a unor circuite orientate pe echipament. Fabricarea acestora însă este rentabilă numai în cazul unui număr ridicat de unități. Din această cauză - și în principal la sistemele fabricate în număr de bucăți reduse - a crescut rapid utilizarea sistemelor logice programabile. Deoarece știința rețelelor neuronale în mare parte este în faza de cercetare, dezvoltarea soluțiilor hardware dedicate sunt bazate pe sisteme logice programabile.

În faza de cercetare-dezvoltare dintre sistemele logice programabile în primul rând sunt utilizate acelea care sunt reprogramabile. Din această categorie fac parte chipurile FPGA XILINX conținând blocuri logice reconfigurabile și conectarea acestora prin comutatoare programabile.

Fiecare tip de rețea neuronală conține unități aritmetice, unitatea de comunicare dintre neuroni, unitatea de comandă a întregii rețea, etc. Problema cea mai mare constă în realizarea circuitului de înmulțire, necesitând un număr foarte mare de elemente logice și astfel acesta ocupă o suprafață însemnată din suprafața chip-ului. Scopul este introducerea unui număr ridicat de unități de procesare pe un singur circuit integrat, astfel și acele rețele care nu conțin elemente de înmulțire pot fi implementate eficient.

#### 4.3.4. Aplicații ale implementărilor hardware de rețele neuronale

Deși nu atât de reușite precum rețelele neuronale software, există și rețele neuronale hardware care funcționează la parametri optimi zi de zi.

##### 4.3.4.1. Accurate Automation Corp (AAC)

**Descriere:** Procesorul de rețele neuronale al Accurate Automation (Neural Network Processor - NNP) utilizează o adevărată arhitectură cu instrucțiuni și date multiple (multiple-instruction multiple-data - MIMD), capabilă de a rula în paralel pe mai multe chipuri fără ca performanța sa să aibă de suferit.

Fiecare chip conține un procesor de o viteză înaltă de 16 biți cu stocare de date în circuitul integrat (CI) pentru ponderile sinaptice. Procesorul efectuează doar nouă instrucțiuni lingvistice de asamblare. De exemplu, numai o singură instrucțiune este necesară pentru a transmite ieșirea unui neuron prin funcțiile de transfer bazate pe valori de prag, specifice majorității metodelor de învățare neuronale." Comunicarea interprocesorală dintre multiple NNP-uri este dirijată de către o magistrală interprocesorală proprie. Un toolbox software este inclus care deja a codificat programe în limbaj de asamblare pentru toate metodele cunoscute de învățare neuronale optimizate pentru procesoare paralele multiple. De exemplu, toolbox-ul implementează algoritmi de back-propagation, rețele Hopfield, metode de învățare Adaline și Madaline, învățarea recurentă și cea Kohonen fiind în curs de dezvoltare în cadrul Accurate Automation.

Chipurile sunt disponibile doar de la Accurate Automation fiind deja montate pe circuite imprimate. O versiune ISA pentru PC (\$4,995) vine împreună cu un chip și spațiu pentru a adăuga circuite imprimate de extensie cu până la încă nouă chipuri (\$3,995 bucata). Versiunea VME (\$19,995) a circuitului - cea folosită pentru LoFLZTE - include de asemenea două TMS340 DSP-uri plus spațiu pentru a instala încă nouă chipuri neuronale. DSP-urile permit circuitului să efectueze o preprocesare de semnale pentru rețeaua neuronală, cum ar fi filtrarea digitală, recunoaștere de model, optimizare și procesare de matrice rară.

**Învățare:** Poate fi programat pentru a implementa orice algoritm de antrenare neuronală.

**Performanțe:** 140 MCPS pentru un singur chip. Până la 1.4 GCPS pentru un sistem de 10 procesoare.

##### 4.3.4.2. Adaptive Solutions CNAPS

**Descriere:** sistemul CNAPS este un sistem de dezvoltare a rețelelor neuronale, bazat pe circuitul propriu CNAPS-1064, Procesor Paralel Digital care are 64 de sub-procesoare operând în mod SIMD. Fiecare sub-procesor are proprii 4k biți de memorie locală și o unitate aritmetică în virgulă fixă pentru a efectua operații cu parametrii întregi reprezentați pe 1-bit, 8-biți sau 16-biți. Fiecare sub-procesor poate emula unul sau mai mulți neuroni și mai multe circuite pot fi conectate împreună.

Serverul CNAPS II este mediul de dezvoltare software și hardware. Are o magistrală VME cu 4 slot-uri și un controller Ethernet. Un singur card VME al serverului CNAPS II poate fi setat pentru 64, 128, 256 și 512 de configurații, iar un card are un spațiu de stocare de 16M Bytes.

Procesoarele trebuie să fie programate pentru a executa un anumit algoritm de rețea neuronală. Utilitățile CNAPS includ CNAP-C (un compilator C și un debugger), Quicklib (funcții codate manual apelabile din CNAPS-C), BuildNet (algoritmi pre-codați ai rețelelor neuronale) și CodeNet (debugger pentru limbaj de asamblare).

Cardul CNAPS/PC ISA folosește 1,2 sau 4 circuite integrate paralele ale noului procesor CNAPS-1016 sau două din circuitele integrate 1064 pentru a obține procesoare CNAPS de 16, 32, 64 sau 128. În cazul acestora, ratele de executare a operațiilor de multiplicare și adunare sunt de 320M, 640M, 1.28G și 2.56G pe secundă. Viteza datelor magistralei ISA este de 20MByte/s. Sistemul de dezvoltare CNAPS este similar cu cel pentru stația de lucru: compilator CNAPS-C, modul de asamblare, Debugger, API, QuickLib și programe demonstrative (McCartor, 1991).

Software-ul BrainMaker a fost de asemenea portat pentru a rula card-uri CNAPS PC. **Învățare:** Algoritmii de învățare pot fi programați. Cei clasici ca back-propagation și alți algoritmi de acest fel vin cu pachetul BuildNet.

**Performanțe:** Cu un singur chip, se poate atinge 1.28 miliarde ca rată de multiplicare și adunare pe secundă. Cu 4 chipuri, se ajunge la 10.24 miliarde rată de multiplicare și adunare pe secundă. Algoritmii backpropagation pentru rețele de tip feedforward au performanțe de 1.16 miliarde ca rată de multiplicare și adunare pe secundă și 293 milioane modificări de ponderi pe secundă cu un *chip*.

#### 4.3.4.3. IBM ZISC036

**Descriere:** ZISC este un circuit integrat (CI) cu 64 intrări de 8-biți și 36 neuroni RBF. Mai multe CI pot fi puse laolaltă pentru a crea rețele de mărime arbitrară. Vectorul de intrare  $V$  este comparat cu un vector prototip  $P$  stocat pentru fiecare neuron. Ieșirea neuronului de 14-biți (pentru posibile categorii de 16K) este valoarea de distanță calculată conform a două norme selectabile: (1)  $dist = \sum \text{abs}(V_i - P_i)$ ,  $i=1,64$ ; (2)  $dist = \text{Max}((V_i - P_i))$ ,  $i=1,64$ . Intrarea este serială pe 8-biți. Durează 3,5 microsecunde pentru a încărca cele 64 de elemente și încă 0,5 microsecunde pentru ca semnalul de clasificare să apară.

**Învățare:** Învățare on-chip, stocare a vectorilor de instruire ca prototipuri. Procesarea învățării a unui vector durează încă 2 microsecunde în afara celor 4 microsecunde pentru încărcare și evaluare.

**Performanțe:** La 16MHz, clasificarea unui vector de 64 componente reprezentat pe 8-biți se efectuează în 4 microsecunde.

**Comentarii:** Prima generație de chip-uri utilizează o tehnologie VLSI foarte conservatoare. Generațiile următoare sunt planificate și vor folosi progresiv o tehnologie mai agresivă pentru a mări viteza și densitatea neuronilor.

#### 4.3.4.4. Intel 80170NX Electrically Trainable Analog Neural Network (ETANN)

**Descriere:** rețele neuronale analogice cu 64 de ieșiri (0÷3V), 16 offseturi interne și 64 neuroni cu funcții de transfer sigmoidal. Rețelele de tip feedforward cu două straturi pot fi implementate cu 64 de intrări, 64 de neuroni din stratul ascuns și 64 neuroni de ieșire, folosind două matrici de pondere de 80x64. Ieșirile stratului ascuns vor trece printr-o altă matrice de ponderi pentru a procesa stratul de ieșire. Alternativ, o singură rețea de 64 de straturi și 128 ieșiri poate fi implementată folosind ambele matrici și prezentarea intrărilor în 2 faze de câte de 64 de intrări. Ponderile sunt stocate în sinapse non-volatile de tip *floating gate* synapses (multiplicatoare analogice Gilbert). Ponderile au o precizie de numai 6-biți.

Sistemul de învățare INNTS este un mediu de dezvoltare PC care include software (emularea și interfața chip-ului), o extensie PC a magistralei card-ului și un modul de instruire pentru un singur chip. De asemenea, un circuit imprimat de 8 chip-uri este disponibil.

**Învățare:** Fără învățare on-chip. Emularea este făcută în software și este descărcată pe circuit. Cu toate acestea, pentru a aduce performanța aproape de nivelul emulării este necesară introducerea unei faze speciale de învățare.

**Performanțe:** Un timp de propagare de aproximativ 8 microsecunde pentru o rețea cu 2 straturi. Acesta este echivalent cu numai 2 miliarde de rate de multiplicare/adunare pe secundă.

#### 4.3.4.5. Irvine Sensors 3DANN

**Descriere:** Irvine Sensors au dezvoltat o tehnologie de stratificare FET de mare densitate care permite circuitelor să fie construite atât pe dimensiuni verticale cât și pe dimensiuni orizontale. Folosind aceste circuite stratificate se poate construi o rețea neuronală artificială 3DANN - 3D (NNW) de foarte mare densitate. Compania a câștigat un contract cu Armata SUA pentru a "demonstra fezabilitatea unui sistem

de interconectare cu densitate ridicată pentru a îmbunătăți tehnologia rețelei neuronale artificiale”. Noua schemă de interconectare face parte din planul de îmbunătățire a tehnologiei Irvine Sensors în vederea atingerii unui “Silicon Brain”, un sistem de recunoaștere conceput pentru a emula performanțele sistemului nervos central al omului.

Bazându-se pe o dezvoltare similară al CalTech, care emula cortexul vizual utilizând circuite FPGA, s-a ales o structură stratificată de FPGA-uri pentru proiectul de emulare Silicon Brain Architecture SBIR Phase II. Utilizarea de FPGA este importantă, deoarece se pot implementa multe configurații de circuite integrate fără costuri de design și executare a acestora.

Dificultatea impusă de circuitele FPGA convenționale la implementarea unor funcționalități electronice complexe este faptul, că conexiunile de care dispun nu sunt suficiente. Conceptul 3DFET permite interconexiuni dense între diferite straturi.

În noua implementare a Silicon Brain Architecture SBIR Phase II se va include 3DFET, precum și circuite FPID (Field Programmable Interconnect Devices), care vor permite ca interconexiunile dintre circuitele FPGA stratificate să fie și ele total programabile.

Astfel, se vor putea implementa diferite scheme de interconectare, cum ar fi cea de fluture, crossbar sau mesh. Cu o asemenea flexibilitate a sistemului stratificat, se poate construi orice prototip de procesor digital. Potențialul de piață al unui asemenea produs este, evident, imens.

#### 4.3.4.6. National Semiconductor NeuFuz/COP8 Microcontrollers

**Descriere:** Acest sistem folosește o combinație de rețea neuronală și software cu logică fuzzy pentru a genera cod către microcontrolerele National Semiconductors COP8. Logica fuzzy oferă o metodă foarte flexibilă și puternică pentru aplicațiile de control. Atunci când regulile “if-then” ale sistemului sunt cunoscute, ele pot fi exprimate într-un mod sistematic și direct utilizând funcția de apartenență fuzzy și metodologia regulilor fuzzy. Totuși, pot exista multe situații când nu se cunosc regulile “if-then” dinainte. În asemenea cazuri, o rețea neuronală poate fi folosită pentru a învăța regulile și apoi “cunoștințele” rețelei neuronale sunt înglobate în regulile fuzzy și funcțiile de apartenență. National Semiconductors oferă mai multe pachete:

- NeuFuz Kit de Învățare (NF2-C8A-Kit): software cu rețea neuronală PC/AT (2 intrări, 1 ieșire) și generator de legi și funcții de apartenență fuzzy (maximum 3 funcții), generator de cod COP8, assem/linker COP8.
- NeuFuz4 (NF2-C8A): software cu rețea neuronală PC/AT (4 intrări, 1 ieșire) și generator de legi și funcții de apartenență fuzzy (maximum 3 funcții de apartenență), generator de cod COP8, assem/linker COP8.
- NeuFuz4 Sistem de Dezvoltare (NF2-C8A): software cu rețea neuronală PC/AT (4 intrări, 1 ieșire) și generator de legi și funcții de apartenență fuzzy (maximum 3 funcții de apartenență), generator de cod COP8, assem/linker COP8, emulator în circuit cu programare PROM COP8.

**Învățare:** doar învățare software.

#### 4.3.4.7. Nestor NI1000

**Descriere:** Aceasta este o rețea cu neuroni Radial Basis Function. În timpul învățării, vectorii prototip sunt stocați prin presupunerea faptului că aceștia sunt aleși aleatoriu din distribuția principală. Pot fi stocați până la 1024 de prototip cu 256 de dimensiuni, 5-biți pe dimensiune. Apoi fiecare prototip este desemnat unui anumit neuron din stratul mijlociu. Acest neuron, la rândul lui, este desemnat unui neuron de ieșire care reprezintă clasa specifică pentru acel vector (până la 64 de clase admise). Toți neuronii din stratul de mijloc, care corespund aceleiași clase sunt desemnați aceluiași neuron de ieșire. În modul de redare, un neuron de intrare este comparat cu fiecare prototip, în paralel, și dacă distanța dintre ei este deasupra pragului stabilit, acesta se activează, și activează la rândul lui ieșirea, clasa sau neuronul corespondent.

Algoritmii de învățare pot varia și nu salvează pur și simplu fiecare prototip. De exemplu, algoritmul aparte poate decide dacă un vector de intrare este salvat ca un nou prototip sau e considerat a fi prea asemănător unui prototip deja existent. Setarea pragului poate să depindă și de algoritm.

**Învățare:** Doi algoritmi de învățare on-chip sunt disponibili: Rețeaua Neuronală Probabilistică (RNP) și Restricted Coulomb Energy (RCE). De asemenea, micro-codarea poate fi modificată pentru algoritmi definiți pentru utilizatori.

**Performanțe:** 40k modele a 256 elemente pe secundă (25 micro-secunde/model).

#### 4.3.4.8. Philips L-Neuro chips

**Descriere:** Philips a oferit anterior chipul L-Neuro 1.0 și a anunțat dezvoltarea următoarei generații, chipul L-Neuro 2.3.

**L-Neuro 1.0** conține 16 elemente de procesare (EP) cu registre de 6-biți. Totuși, cei 16-biți pentru un EP pot fi considerați ca 16 neuroni de 1-bit, 4 de 4-biți, 2 de 8-biți sau 1 de 16-biți. Așadar, un singur chip poate implementa, spre exemplu, o rețea neuronală cu 256 neuroni de 1-bit. Există o memorie tampon pentru ponderi de 1kByte pe chip pentru a implementa 1024 de ponderi reprezentate pe 8-biți sau 512 ponderi reprezentate pe 16-biți. Funcțiile de transfer sunt implementate în afara chip-ului pentru ca mai multe chip-uri să poate fi conectate unul de altul. Cele 16 ieșiri EP sunt citite în serie. Chip-urile sunt cel mai ușor interfațate pe procesoare gazdă Transputer. O rată de procesare de 100MCPS pentru modul cu reprezentare pe un 1-bit și una de 26MCPS sunt raportate pentru modul cu reprezentare pe 8-biți pentru un singur chip. Circuitul poate fi programat să învețe și rate de 160MCUPS, respectiv 32 MCUPS este dat pentru modul de 1-bit și 8-biți. Un sistem bazat pe transputer cu până la 112 L-Neuro 1.0 este disponibil și de la Telmat care vinde sisteme Transputer. De asemenea, este disponibil și un circuit imprimat PC cu un singur circuit L-Neuro.

**L-Neuro 2.3** este a doua generație care a beneficiat de experiența lui L-Neuro 1.0. Este conectabil în cascadă ca și L-Neuro 1.0 și constă din 12 procesoare care pot opera fie în mod paralel (SIMD), fie în mod pipelined. Are ponderi pe 16-biți și ieșiri neuronale în modul basic. Fiecare din cele 12 procesoare conține 128 de registre cu reprezentare pe 16-biți pentru a stoca ponderi și stări, un multiplicator de 16 pe 32 biți, o unitate Aritmetică și Logică (ALU) pe 32-biți. Micro-codarea circuitului integrat este disponibilă utilizatorului pentru a putea fi proiectat conform unei anumite aplicații. Cu o viteză de tact de 60MHz, circuitul poate calcula ieșirile, la fiecare 17ns, ponderile fiind reprezentate pe 32 de biți, iar valorile de intrare pe 16. Acesta asigură 720MCPS, adică de 27 ori modul cu reprezentare pe 8-biți al L-Neuro 1.0. În procesul învățării, cele 12 ponderi sunt updatate în paralel în 34ns. Circuitul integrat se poate folosi nu doar ca aplicație pentru rețele neuronale, dar și ca aplicație pentru procesare de semnal, procesare de imagine și logică fuzzy (folosindu-și abilitatea de a extrage valori minime și maxime din vectorii celor 12 elemente).

#### 4.3.4.9. Sensory Circuits RSC-164 Speech Recognition Chip

**Descriere:** RSC-164 este un CI recunoaștere a vorbirii cu cost redus, proiectat pentru uzul produselor electronice. Acesta combină un procesor pe 8-biți cu algoritmi de rețele neuronale pentru o recunoaștere a vorbirii, atât independentă de interlocutor cât și dependentă de acesta, de cea mai bună calitate. Circuitul integrat conține de asemenea funcțiile de înregistrare a vocii, redare, sinteză de muzică, sinteză de vorbire și sistem de control. Dispozitivul CMOS include memorii RAM și ROM înglobate pe chip, 16 linii I/O, convertor D/A și un procesor dedicat de 4-MIPS.

RSC-164 utilizează o rețea neuronală apriori antrenată să execute recunoașterea vorbirii, în timp ce sinteza de înaltă calitate a vorbirii este dobândită folosind o schemă de compresie timp-domeniu care îmbunătățește ADPCM-ul de tip obișnuit. Filtrarea digitală pe chip ameliorează precizia recunoașterii, pre-procesând semnalele primite. Controlul dinamic AGC compensează persoanele care nu stau la o distanță corespunzătoare față de microfon sau pe cele care vorbesc prea încet sau prea tare.

RSC-164 include o interfață externă de memorie care permite conectarea cu dispozitivele de memorie pentru înregistrare audio/playback, recunoașterea dependentă de vorbitor și lungimi extinse de mesaje pentru sinteza de vorbire.

#### 4.3.4.10. Siemens MA-16 chip, SYNAPSE-3 Neurocomputer

**Descriere:** circuitul integrat Siemens MA-16 este un multiplicator rapid de matrici care poate fi combinat pentru a forma *serii sistolice*, adică intrările și ieșirile sunt trimise de la un modul la altul ca într-o linie de asamblare. Un singur modul poate procesa 4 modele, de 16 elemente fiecare (pe 16-biți), cu 16 valori neuronale (pe 16-biți) la o rată de 800 multiplicări-adunări/secundă la 50Mhz. Ponderile sunt încărcate de pe RAM-ul din afara circuitului iar funcțiile de transfer neuronal sunt calculate cu tabele de căutări (look-up table) externe circuitului integrat.

SYNAPSE-1 este un sistem hardware/software complet care folosește 8 circuite integrate MA-16. Acesta se află în propria lui carcasă și comunică via Ethernet cu o stație de lucru gazdă. Cardul magistralei de accelerație SYNAPSE2\*PC conține 1 MA-16, în timp ce card-ul SYNAPSE3-PC PCI conține 2 MA-16.

SYNAPSE3-PC: Până la trei circuite SYNAPSE3\*PC pot fi inserate într-un PC Pentium. Gazda PC și circuitul(ele) pot executa operații simultane sau independente una de cealaltă. Performanțele cele mai înalte ale unui circuit sunt de aproximativ 2560 MOPS ( $1,28 \times 10^9$  MultAcc/sec); 7160 MOPS ( $3,58 \times 10^9$  MultAcc/sec) când se folosesc 3 circuite! O asemenea configurație obține performanțe mai bune decât SYNAPSE1, având de asemenea performanțe I/O mult mai ridicate și un raport optim preț/performance. SYNAPSE3-PC este utilizat ca accelerator pentru rețelele neuronale MediaInterfaces ale SynUse-Base Software.

**Învățare:** Circuitul integrat nu are implementat hardware nici un algoritm specific dar poate fi programat să facă calcule, cum ar fi pentru back-propagation. Cardurile SYNAPSE pot fi programate pentru aproape orice fel de algoritm de rețele neuronale.

**Performanțe:** SYNAPSE-1 cu 8 MA-16 a avut performanțe de vârf de 3,2 miliarde de multiplicări (16-biți x 16-biți) și adunări (48-biți) pe secundă la o frecvență de bază de 25Mhz. SYNAPSE3-PC: 2560 MOPS ( $1,28 \times 10^9$  MultAcc/sec); 7160 MOPS ( $3,58 \times 10^9$  MultAcc/sec) când se folosesc trei circuite.

## 4.4. Implementări software și hardware ale rețelelor neuronale bazate pe impulsuri

S-a demonstrat deja de mulți cercetători, pornind de la studiile sistemelor biologice, că structura temporală a șirurilor de impulsuri din aceste modele neuronale neuromorfe are o putere de procesare demnă de luat în considerare. Aplicațiile acestor modele, spre exemplu în probleme de segmentare de imagini, necesită simularea unor rețele largi într-un timp rezonabil de mic. În cele ce urmează se vor prezenta diferite platforme hardware pe care aceste rețele neuronale au fost deja implementate. După cum se va vedea, majoritatea acestor sisteme sunt foarte complexe și costisitoare, acest lucru făcând necesară dezvoltarea de acceleratoare neuronale neuromorfe implementate hardware mult mai accesibile atât ca preț, cât și ca timp de dezvoltare.

### 4.4.1. Neuronii pulsativi

Să examinăm în primul rând particularitățile de bază ale neuronilor pulsanti (Pérez-Urbe, 1999) (Gerstner W. , 2001) (Trappenberg, 2002):

- Neuronii comunică unul cu celălalt doar prin impulsuri temporizate
- Impulsurile care intră sunt ponderate și induc un potențial post-sinaptic conform unei funcții impuls-răspuns.
- Funcțiile impuls-răspuns sunt construite din unul sau mai multe integratoare permeabile care sunt descrise într-o versiune discretă de  $PI(n+1) = r * PI(n)$ , unde  $PI$  este un modul potențial intern și  $r$  este un factor de dispersie.
- Modulele potențiali interni ( $PI$ ) sunt combinați în funcțiile de ieșire conținând adunare, scădere și/sau multiplicare și cel puțin o funcție de prag.
- Rețelele pot consta din  $10^5$  neuroni și pot fi distribuiți în straturi de neuroni cu proprietăți identice. Fiecare neuron este conectat cu până la  $10^4$  alți neuroni.
- Structura de conectivitate a unei rețele poate fi împărțită în două grupuri: 1.) *conexiuni regulate (CR)* care urmează niște reguli deterministice simple (ex. câmpuri receptive) și 2.) *conexiuni neregulate (CNR)* (conexiuni rare, aleatorii).

### 4.4.2. Probleme de implementare

Când implementăm asemenea rețele pe computere digitale, timpul  $t$  decurge în unități temporale  $T$  proprii discrete (de obicei  $T = 1\text{ms}$ ). Simularea unei etape de timp proprie pentru întreaga rețea ar trebui să fie un *pas temporal* (Time Step - TS). A calcula un pas temporal în mai puțin de o milisecundă va fi numită *simulare în timp real*. Procesarea în fiecare pas temporal implică trei operații:

**Intrare:** impulsurile trebuie să fie distribuite și neuronii-receptori de impulsuri trebuie să incrementeze potențialele lor interne,  $PI$ .

**Ieșire:** fiecare neuron trebuie să-și combine potențialii interni și să decidă dacă să trimită un impuls la ieșire sau nu.



**Dispersie:** potențialii interni trebuie să fie dispersați pentru fiecare neuron. Trebuie menționat faptul că datorită necesității stocării PI-urilor, calcularea acestor rețele este *I/O-limitată* (nr. I/O > nr. de operațiuni). Putem aplica diferite metode pentru a mări viteza simulării:

- *Activitatea rețelei* (numărul mediu al impulsurilor pe TS / numărul de neuroni) este de obicei foarte scăzută. O schemă foarte eficientă de comunicare pentru aceste rețele este *protocolul listă-eveniment* (event-list protocol). Doar adresele neuronilor pulsați trebuie înregistrate într-o listă-eveniment a impulsurilor și trebuie să distribuie impulsuri pentru intrare (Lazarro & Wawrzynek, 1993).
- În mod caracteristic, rețelele neuronale pulsative nu sunt complet conectate. O metodă de a reprezenta această conectivitate risipită (sparse connectivity) este folosirea listelor, una pentru fiecare neuron  $n_i$ . Unitățile de informație din liste sunt seturi de date conținând ponderi și adrese. Adresele  $a_j$  din listă pot indica neuronul  $n_j$  căruia  $n_i$  îi trimite un impuls sau din partea căruia  $n_i$  primește un impuls. Cel dintâi reprezintă o *conectivitate emițător-orientată*, iar cel din urmă o *conectivitate receptor-orientată* (Frank & Hartmann, 1995). Pentru o activitate de rețea scăzută metoda emițător-orientată este avantajoasă: doar neuronii receptori de impulsuri trebuie să calculeze funcția de intrare (Jahnke, Roth, & Klar, 1995).
- Pentru o conectivitate uzuală putem calcula conexiunile unui neuron independent de poziția sa în spațiu. *Calcularea on-line* a conexiunilor reduce drastic cantitatea necesară de stocare pentru listele de conexiune. Presupunând că vectorii de pondere sunt similari pentru toți neuronii, atunci nu trebuie să stocăm decât un vector de pondere pentru întreaga rețea.
- În loc de a utiliza aritmetica în virgulă mobilă, este recomandată folosirea aritmeticii rapide în virgulă fixă. Conform rezultatelor analizei de rezoluție (Roth, Jahnke, & Klar, 1995), acuratețea de la 8 biți (ponderi) până la 18 biți (potențiali interni) este suficientă pentru problemele de vizibilitate redusă.

În ceea ce privește simularea pe un computer paralel, vom face o deosebire între trei metode de mapare a unei rețele care pot fuziona una cu cealaltă:

- **n(neuron)-paralel:** toate sinapsele ale unui neuron sunt mapate aceluiași element de procesare PE. Diferiți neuroni sunt procesați în paralel, dar sinapsele sunt procesate în serie.
- **s(inapsă)-paralelă:** sinapsele unui neuron sunt distribuite către PE-uri diferite. Acum sinapsele unui neuron sunt procesate în paralel și neuronii în serie.
- **m(odel)-paralel:** PE-urile calculează răspunsul pentru un model de intrare în paralel, într-un segment anume al rețelei, în timp ce diferite segmente de rețea vor fi procesate în serie.

#### 4.4.3. Implementare pe calculator paralel

SP2 este un computer MIMD cuplat în paralel, cu până la 256 R6000 PE-uri, arhitectură cu memorie locală și comutatoare 16 pe 16 de viteză înaltă pentru comunicare inter-PE. Planificarea n-paralel și conectivitatea emițător-orientată au fost folosite pentru implementare (Delorme & Thorpe, 2003). Pe de-o parte, un singur PE manifestă o mult mai slabă performanță decât un P90 care poate fi rezultatul unor setări nesatisfăcătoare ale compilatorului și necesită investigații suplimentare. Pe de altă parte, observăm o accelerare de numai 1.12/1.23/1.25 folosind PE-uri de 2/4/8 peste implementarea unui PE individual. Așadar, performanța este îndeosebi limitată de către comunicare inter-PE (calculul și comunicarea nu pot fi realizate în paralel pe SP2) care rezultă în comportamentul slab la scalarea impulsurilor PE observate.

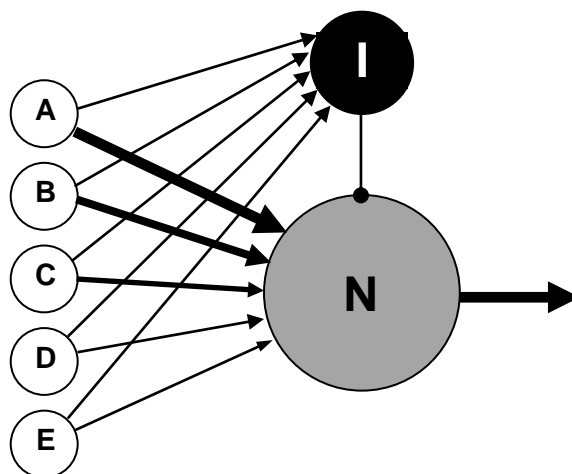


Figura 4.3 Modelul neuronal conceput de Thorpe

#### 4.4.4. SpikeNET

Simon Thorpe și echipa sa (Arnaud Delorme, Rufin VanRullen, Marie Fabre-Thorpe) în Toulouse, Franța, sunt unii dintre acei cercetători care au ajuns să fie la cel mai înalt nivel – în zilele noastre – în ceea ce privește studiul rețelelor neuronale neuromorfe (spiking). Alături de multe alte realizări, ei au inspirat apariția primului produs comercial bazat pe aceste cercetări, și anume **SpikeNET**, un sistem de recunoaștere vizuală.

În lucrarea lor, se argumentează, că un sistem bazat pe codarea ordinii rangurilor, *rank order coding*, - unde primele impulsuri recepționate de neuroni sunt luate în considerare, iar cele care sosesc mai târziu sunt inhibate – poate livra o cantitate extraordinară de informații (Delorme & Thorpe, 2003). Pe lângă acest fapt, simplitatea modelului de dinamism neuronal, care poate realiza ordonarea este foarte atractivă.

În acest model (vezi Figura 4.3), neuronul  $I$  inhibă neuronul  $N$ . Astfel, o acumulare de impulsuri ce activează  $N$ , va stimula spre activare și neuronul  $I$ , a cărui factor de inhibare, va împiedica  $N$  să mai emită impulsuri, chiar dacă va mai recepționa intrări ulterioare.

Alte cercetări foarte importante în acest domeniu, se desfășoară în Austria, la Graz, sub conducerea lui Wolfgang Maass, împreună cu Berthold Ruf, Heinrich Schmitt și alții.

Se poate vedea, că această tematică este foarte tânără, dar cu mare potențial, conturându-se deja, mai multe direcții promițătoare de cercetare.

##### 4.4.4.1. Sistemul de simulare rețele neuronale neuromorfe SpikeNET

În prezent sunt disponibile o multitudine de sisteme, care oferă posibilități variate de simulare a rețelelor neuronale de orice fel. Unele dintre acestea au fost dezvoltate pentru a aproxima în medii artificiale, detaliile biofizice ale funcționării neuronilor biologici. Tipic pentru unitățile ce s-au născut astfel este o structură cu o ieșire simplă (de multe ori 0 sau 1). Aplicare acestora este destul de răspândită, în domenii ingineresti și financiare. La celălalt capăt al spectrului se află programe sofisticate cum ar fi GENESIS (Bower & Beeman, 1998) sau NEURON (Hines & Carnevale, 1997). Acestea au o excelentă capacitate de a modela și simula procese biofizice, spre exemplu structura dendritică sau cinematica canalelor ionice ale membranelor celulelor. Acest nivel ridicat al detalierii funcțiilor reprezintă însă, un impediment în calea implementării eficiente a unor rețele mari conform acestor modele.

Unul dintre avantajele SpikeNET (Figura 4.4) este viteza. Utilizând un calculator G3 Macintosh (PowerPC 750, cu procesor la 266 Mhz), SpikeNET poate update aproximativ 20 de milioane de conexiuni pe secundă, chiar și în cazul utilizării unui parametru de sensibilitate pentru a modula efectul fiecărei intrări sinaptice. Astfel se poate realiza simularea în timp real a unei rețele formate din 400.000 neuroni, cu un pas temporal de 1ms (presupunând 49 de conexiuni pe neuron și o rată medie de activare de 1 pe secundă, acestea fiind o aproximare rezonabilă a ratei de activare a neuronilor corticali). De notat este faptul, că utilizând o simulare mai convențională de rețele neuronale, ar fi nevoie de procesarea fiecărui element la fiecare pas temporal. Acest lucru ar însemna, că utilizând aceeași putere de calcul, numai 20.000 de conexiuni ar putea fi procesate pe milisecundă, ceea ce ar limita mărimea rețelei simulate în timp real la circa 400 de neuroni (49 de conexiuni/neuron).

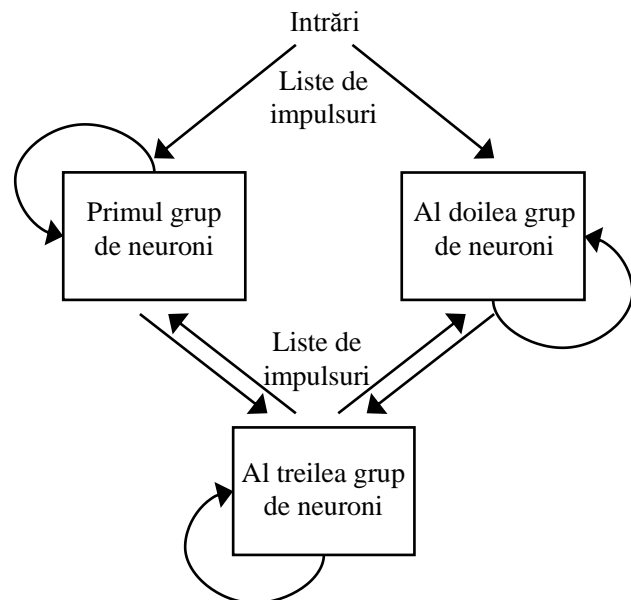


Figura 4.4 Structura de bază a sistemului SpikeNET  
Se redirecționează liste de impulsuri între diferite grupe de neuroni, organizați într-o matrice bidimensională.

Deoarece numai un număr mic de neuroni emit impulsuri în fiecare pas temporal, traficul comunicației este minimalizat.

#### 4.4.5. Rețele neuronale pulsative, dezvoltări hardware

Cel mai bun model neuronal, care asigură fidelitatea cea mai ridicată față de fenomenele biofizice din celulele nervoase biologice, este cel propus de Hodgkin și Huxley (Hodgkin & Huxley, 1952). Acest model se bazează pe patru ecuații diferențiale neliniare cuplate, care implică timpi de simulare software mari, iar implementările hardware ar fi total ineficiente din punct de vedere al raportului preț/performanță (Gerstner & Kistler, 2002). Implementările hardware necesită deci, un model mai simplu, cum ar fi modelul numit Integrare-și-Activare, I&A, (Integrate-and-Fire) studiat în (Koch, 1999) (Dayan & Abbott, 2001) și care a stat și la baza implementărilor realizate și prezentate în primul referat al acestui program de doctorat (Bakó, Analiza și simularea sistemelor cu inteligență artificială neuromorfă, 2005). Modelul se bazează pe o ecuație diferențială de ordinul unu, unde rata de schimbare a stării unui neuron  $v$ , este corelată cu curenții de membrană. Modelul I&A utilizat este următorul:

$$c_m \frac{dv(t)}{dt} = g_l (E_l - v(t)) + \sum_j \frac{w^j g_s^j(t)}{A} (E_s^j - v(t)) \quad \text{Ec. 60}$$

Pentru o sinapsă anume  $j$ , apariția unui potențial pre-sinaptic la un moment  $t_{ap}$  declanșează activarea acesteia la momentul  $t_{ap} + t_{delay}^j$ , prin efectuarea unei schimbări discontinue în conductanța sinaptică de forma:

$$g_s^j(t_{ap} + t_{delay}^j + dt) = g_s^j(t_{ap} + t_{delay}^j) + q_s^j \quad \text{Ec. 61}$$

altfel,  $g_s^j(t)$  este controlat de:

$$\frac{dg_s^j(t)}{dt} = \frac{-1}{\tau_s^j} g_s^j(t) \quad \text{Ec. 62}$$

Folosind schema de integrare directă Euler cu un pas temporal de  $dt=0.125ms$  pentru a rezolva ecuațiile modelului I&A, rezultă:

$$v(t + dt) = v(t) + \frac{1}{c} ((g_l (E_l - v(t)) + \sum_j \frac{w^j g_s^j(t)}{A} (E_s^j - v(t)))) dt \quad \text{Ec. 63}$$

$$g_s^j(t + dt) = g_s^j(t) + \left(\frac{-1}{\tau_i} g_s^j(t)\right) dt \quad \text{Ec. 64}$$

Acest model a fost utilizat de către un **grup de cercetători din Irlanda de Nord, sub conducerea lui B. Glackin**, pentru a implementa rețele neuronale pulsative pe circuite FPGA. Tabelul 8, de mai jos prezintă explicații, respective valori utilizate în această lucrare (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005), pentru ecuațiile modelului I&A.

Tabelul 8. Descrieri și valori ale simbolurilor utilizate în ecuațiile modelului I&A

Parametru	Valori uzuale	Valori în celule piramidale	Valori în celule inhibitorii	Descriere
$c_m$	8nF/mm <sup>2</sup>			Capacitatea membranei
$E_l$	-70mV			Potențialul de inversare a membranei
$v_{th}$	-59mV			Tensiunea de prag
$v_{reset}$	-70mV			Tensiunea de reset
$\tau_{ref}$	6.5ms			Perioada de hiperpolarizare
$t_{delay}^j$	0.5ms			Decalajul temporal de propagare
$g_l$		1μs	1μs	Conductanța de scurgere a membranei
$A$		0.03125mm <sup>2</sup>	0.015625mm <sup>2</sup>	Suprafața membranei
$E_s^j$		0mV	-75mV	Potențialul de inversare a sinapsei
$\tau_s^j$		1ms	4ms	Timpul de descreștere al sinapsei

Implementarea hardware prezentată utilizează o reprezentare a datelor în virgulă fixă. Au fost utilizați 18 biți pentru reprezentarea potențialului de membrană, din care 10 biți sunt partea fracționară. Parametrii de multiplicatori și divizori au fost aleși ca valori multiple ale puterilor lui 2, pentru a putea implementa aceste operații cu registre de deplasare.

Algoritmul de antrenare implementat în hardware în această lucrare este unul cu fidelitate biologică ridicată, și anume algoritmul de Plasticitate dependentă de temporizările impulsurilor (Spike Timing Dependant Plasticity - STDP), bazat pe studiul realizat de Song și Abbott în (Song, Miller, & Abbott, 2000) (Song & Abbott, 2001). Astfel, utilizând o rată de conectare de 1:1 între neuroni și sinapsele STDP, s-a reușit implementarea a unui total de 168 de neuroni și 168 de sinapse pe un circuit FPGA Xilinx, din familia Virtex, XC2V8000, al cărei structură se poate vedea în Figura 4.5. Din capacitatea acestui circuit această rețea neuronală a ocupat aproximativ 60% din resursele reconfigurabil. Dacă s-a aplicat o rată de conectare mai realistă de 1:10, atunci rețeaua a avut 41 de neuroni și 410 sinapse, iar dacă rata s-a ridicat la un nivel de 1:100, atunci valorile erau de 4 neuroni și 400 de sinapse.

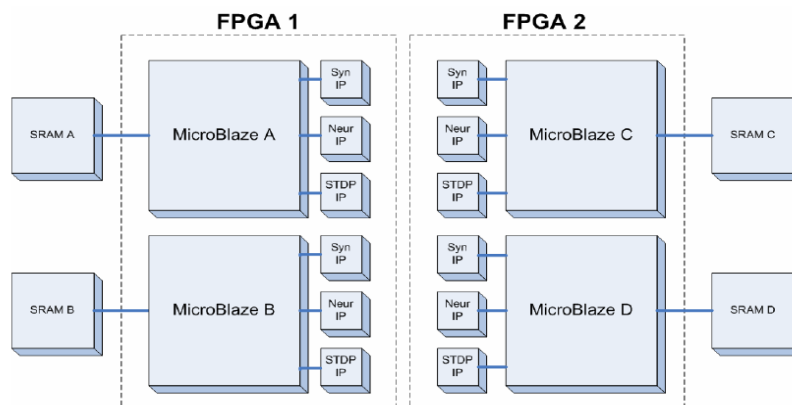


Figura 4.5.- Structura sistemului din (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005) pe un circuit FPAG Xilinx Virtex XC2V8000

Utilizând această arhitectură complet paralelizată la o viteză de tact de 100 MHz și un pas de integrare Euler de 0.125ms pe ciclu de tact, autorii acestei lucrări (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005) au reușit să atingă executarea unei perioade de funcționare în timp real de o secundă în numai 0.8ms. Acest lucru înseamnă o rată de accelerare de 12500 față viteza de timp real. O alternativă pentru mărirea mărimii rețelei implementate este partajarea sistemului între o implementare hardware și una software, însă această formulă implică un compromis între utilizarea mai eficientă a resurselor reconfigurabile și viteza de execuție a rețelei. Astfel, s-ar putea implementa pe una sau mai multe dintre controlerile MicroBlaze înglobate în circuitul Virtex, un algoritm de multiplexare în timp a unui singur neuron sau a unei sinapse. O diagramă a pașilor de execuție a acestui sistem este dată în Figura 4.6.

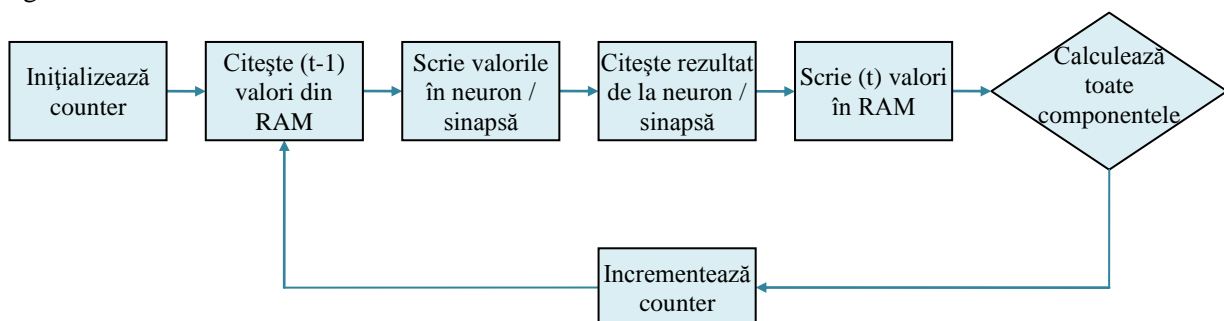


Figura 4.6 Pașii de execuție a sistemului cu elemente de rețea multiplexate prezentată în (Glackin, McGinnity, Maguire, Wu, & Belatreche, 2005)

Versiunea cu mai multe controlere utilizate din FPGA, atinge o îmbunătățire de performanță de 107.25 față de o simulare în Matlab implementată pe un PC Pentium 4 la 2GHz (vezi Tabelul 9).

Tabelul 9 Timpii de simulare pentru o secundă de timp real

Mărimea rețelei	PC (Matlab)	1 MicroBlaze	4 MicroBlaze
100*100*100	2917.2 s	115.7s	27.2s
1400*1400*1400	454418s	N/A	4237s

Un sistem care reprezintă o categorie diferită de implementări de rețele neuronale de tip spiking a fost dezvoltată de un grup de cercetători din Spania, (Agis, Díaz, Ros, Carrillo, & Ortigosa, 2006).

Acest motor de simulare bazată pe evenimente, implementat pe un circuit FPGA Xilinx Virtex-II, poate procesa 2.5 milioane de impulsuri pe secundă, el conținând și cinci memorii SRAM externe. Acest lucru face ca acest sistem să poată fi folosit cu succes în experimente de simulare care necesită motoare de simulare de tip embedded (de exemplu experimente robotice cu agenți autonomi). Această schemă este implementată de obicei în software (Delorme & Thorpe, SpikeNET:an event-driven simulation package for modelling large networks of spiking neurons, 2003) (Ros, Carrillo, Ortigosa, Barbour, & Agis, 2005) (Mattia & Del Giudice, 2000) (Reutimann, Guigliano, & Fusi, 2003), dar a fost abordată și hardware (Makino, 2003) (Schoenauer, Atasoy, Mehrtash, & Klar, 2002) (Mehrtash, Jung, Hellmich, Schoenauer, Lu, & Klar, 2003). Astfel, motorul de simulare poate face salturi de la un impuls de activare la celălalt. Ca urmare, toată activitatea rețelei se organizează în liste cronologice de evenimente (impulsuri) care se procesează în mod secvențial. Această arhitectură se pretează foarte bine calculatoarelor convenționale, datorită naturii sale secvențiale.

În ceea ce privește performanțele simulării și utilizarea de resurse hardware, se poate afirma, că procesarea totală a căii de date consumă 74 cicluri de tact ( în cazul în care lista de evenimente are mai puțin de 1024 de elemente). Utilizând o structură pipe-line, pentru a procesa aceste liste, se reduce sarcina la un impuls de procesat pe fiecare 13 cicluri de tact. Astfel, la o viteză de tact de 33 MHz se obțin cele aproximativ 2.5 milioane de procesări de impulsuri pe secundă.

Comparația acestor implementări cu altele din domeniu este dificilă, deoarece fiecare utilizează alt model neuronal.

Un alt centru important de cercetare al rețelelor neuronale neuromorfe este laboratorul de Sisteme Logice al *Swiss Federal Institute of Technology*, Lausanne, Elveția. Autorii lucrării (Upegui, Peñ a-Reyes, & Sánchez, 2004), pun accent pe proprietățile de reconfigurare dinamică parțială (dynamic partial reconfiguration - DPR) a celor mai recente circuite FPGA, care permite re folosirea aceluiași resurse logice pentru diferite configurații. Se poate reduce astfel necesarul de resurse hardware pentru un anumit proiect și se poate optimiza numărul neuronilor precum și conexiunile dintre aceștia. Există și încercări de a realiza adaptare topologică utilizând DPR, cum ar fi cea descrisă în (Upegui, Peña-Reyes, & Sánchez, 2003)[66]. O topologie modulară, care utilizează DPR, permite mai multe configurații pentru fiecare modul, făcând posibilă și transmiterea bidirecțională a impulsurilor în cazul rețelelor recurente.

Modulele conțin straturi de neuroni cu conexiuni predefinite, iar un algoritm genetic ar putea căuta combinația optimă de straturi.

În (Upegui, Peñ a-Reyes, & Sánchez, 2004) se prezintă un model neuronal funcțional, bazat pe învățare Hebbiană, implementat în hardware pe baza căruia s-a implementat și o rețea neuronală, cu scopul studierii necesității de resurse logice pe un anumit suport. Aplicația de test a fost discriminare între două semnale de frecvențe diferite.

Autorii lucrării mai sus menționate, au utilizat un circuit FPGA Xilinx Spartan II XC2S200 pe care au implementat un neuron cu 30 de intrări, cu și fără învățare ale căror necesități de resurse sunt date de Tabelul 10.

Tabelul 10 Resurse utilizate dintr-un circuit XC2S200 de către un neuron implementat de (Upegui, Peñ a-Reyes, & Sánchez, 2004)

Neuroni fără învățare			Neuroni cu învățare	
Neuron cu 14 intrări	Neuron cu 30 intrări	Neuron cu 62 intrări	Neuron cu 30 intrări	Întreaga rețea (30 de neuroni cu 30 intrări)
17 slice-uri 0.72%	23 slice-uri 0.98%	46 slice-uri 1.95%	41 slice-uri 1.74%	1350 slice-uri 57.4%

Utilizând modelul neuronal cu 30 de intrări ei au implementat o rețea cu trei straturi, (Figura 4.8) fiecare strat conținând 10 neuroni, complet conectați. Fiecare neuron are câte zece intrări, din stratul propriu, din stratul precedent și din stratul următor. Această rețea a ocupat 57.4% din resursele circuitului FPGA, un strat necesitând 19.13%.

În continuare se prezintă un procesor de rețea neuronală cu fidelitate biologică și funcționare în timp real, bazat pe circuit FPGA.

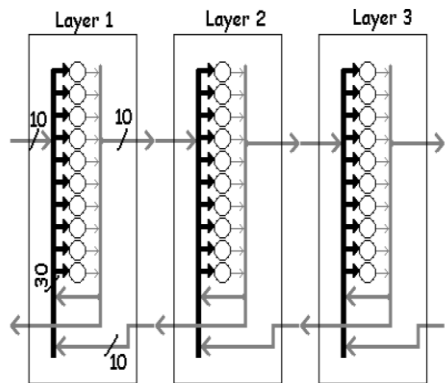


Figura 4.8 Structura rețelei din (Upegui, Peñ a-Reyes, & Sánchez, 2004)

biologice. Pentru a modela rețele spațial dispersate, s-au implementat decalaje inter-neuronale specifice fiecărei sinapse. Atât perioada de dispersie absolută, cât și cea relativă au fost înglobate în modelul dezvoltat. Utilizarea reprezentării acestor parametrii în virgulă mobilă ar necesita un volum prea mare de resurse al circuitului FPGA. Din acest motiv și acest sistem utilizează reprezentări în virgulă fixă, pe 16 biți.

Arhitectura acestui procesor s-ar putea clasifica, ca unul SIMD (vezi Figura 4.7). Este compus dintr-o matrice de elemente de procesare (Processing Elements – PE), care operează asupra aceleiași instrucțiuni, date de un secvențiator central, utilizând însă date locale. Stimularea procesorului se face prin două module de intrare seriale care pot citi datele de intrare în mod asincron. În mod similar, există și două module de ieșire.

Neuronii și sinapsele sunt implementate în elementele de procesare neuronale (PE's). Procesorul neuronal prezentat în articolul de (Izhikevich, 2004) a fost implementat cu 10 PE, fiecare emulând 120 de neuroni virtuali și 912 sinapse. Perioada de reîmprospătare a procesorului a fost setată la 500  $\mu$ s, frecvență generată de un numărător din modulul de secvențiere (Sequencer).

#### Descrierea modulelor procesorului:

- **modulul de secvențiere:** menține performanța de funcționare în timp real și coordonează activitatea tuturor PE ale procesorului, care operează concurrent. Această coordonare este asigurată de comunicația desfășurată pe un canal de control pe 2 biți și o magistrală de date de 16 biți care conectează toate elementele de procesare.
- **modulul de intrare:** dispune de un port de intrare de 64 de biți (multiplexat pe 384 canale de intrare) care poate fi conectat la conectoare fizice de I/O sau la o interfață internă corespunzătoare, utilizând logica circuitului FPGA gazdă. Linii de handshake facilitează operația asincronă și permit comunicarea între diferite domenii a frecvenței de tact. Cele 384 canale de intrare sunt transmise pe magistrala de date internă de 16 biți pentru a fi stocate ulterior în memoria de stare.
- **modulul de ieșire:** are o structură și funcționare similară cu modulul de intrare, dar conține un bloc de memorie RAM intern circuitului FPGA, pentru stocarea listei ieșirilor rețelei neuronale.

Modelul neuronal utilizat în acest procesor – dezvoltat în Anglia (Pearson, și alții, 2005) – este unul singular (unicompartimental), aparținând clasei de excitabilitate de ordinul I, conform clasificării date în (Izhikevich E. , 2004). Ponderea fiecărei sinapse poate fi supusă influenței unui zgomot multiplicativ cu distribuție Rayleigh. De asemenea, se poate injecta un zgomot cu distribuție Gaussiană, în magnitudinea valorii de prag a potențialului de membrană. Aceste distribuții ale zgomotelor utilizate au fost alese astfel încât ieșirea modelului să aproximeze cât mai fidel datele empirice

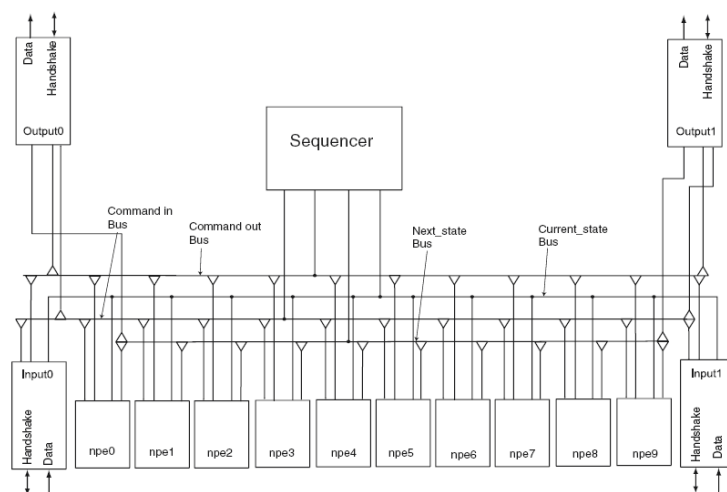


Figura 4.7 Diagrama bloc a topologiei procesorului neuronal SIMD

**Elementele de procesare Neuronale** conțin implementările hardware ale unui neuron și o singură sinapsă (Figura 4.9). Informația contextuală despre cei 120 de neuroni “virtuali” și cele 912 sinapse este stocată local în 4 bancuri de memorii RAM. Această informație este multiplexată secvențial pe hardware-ul implementat cu o viteză foarte mare. O funcționalitate interesantă, implementată de autori, este că o copie a stării întregii rețele este stocată local în fiecare PE, care reprezintă stimulul de intrare pentru neuronii/sinapsele virtuale. Starea reîmprospătată a fiecărui neuron din PE este de asemenea stocată în memoria locală a stării următoare, care va fi transmisă ulterior celorlalte PE-uri.

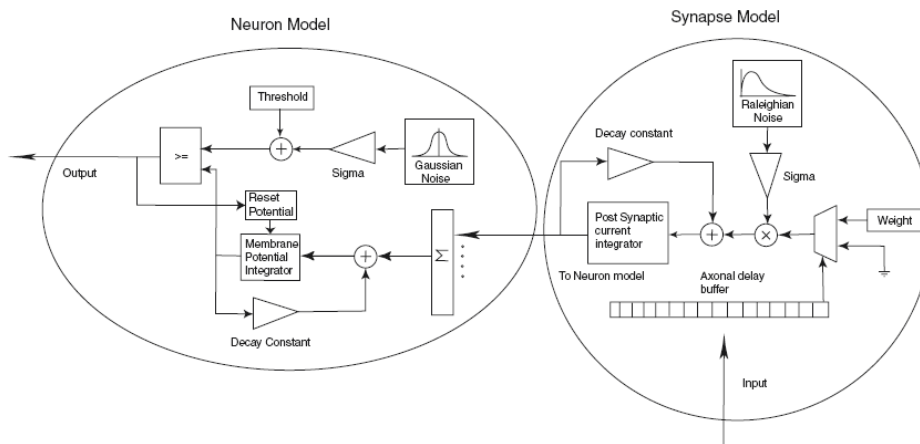


Figura 4.9 Diagramele bloc ale unui neuron și a unei sinapse din SIMD

Circuitul FPGA utilizat a fost Xilinx Virtex-II (XC2V1000-4), cu 1 milion de porți logice echivalente, funcționând la o frecvență de tact de 50 Mhz, situate pe placă de dezvoltare Celoxica RC200, pe care s-au implementat 1200 de neuroni.

În general, pentru a putea estima dacă o rețea neuronală se poate implementa hardware să funcționeze în timp real – ceea ce este o măsură a activității maxime – este necesar a se lua în calcul dimensiunea rețelei, determinându-se apoi rata de împrăștiere necesară. Aceasta mărime poate fi enunțată a fi *cota de activitate a rețelei* adică, numărul maxim de evenimente impuls/neuron/perioadă de reîmprospătare. În cazul de față, dacă ar fi implementați 7000 de neuroni, ei ar genera o activitate pulsatorie de 175 de evenimente pe o perioadă de reîmprospătare, rețeaua putând fi categorizată ca necesitând o cotă de activitate de 0.025, pentru ca modalitatea de procesare implementată să poată funcționa în timp real. Rețeaua implementată, de 1200 de neuroni, are momente de activitate de peste 30 de evenimente de impulsuri pe perioadă de reîmprospătare. Autorii susțin că procesorul prezentat a fost conceput să mențină performanța de funcționare în timp real până la o cotă de activitate de 1, altfel spus, performanța procesorului nu depinde de activitatea rețelei.

#### 4.5. Implementarea cu circuit FPGA a modelului neuronal propriu dezvoltat, validat prin simulare software

Modelul neuronal prezentat în subcapitolul 2.4 a fost conceput în așa fel, încât să poată fi utilizat în aplicații în timp real. Implementarea practică acestui model în hardware, a fost bazată și pe rezultatele simulării software prezentate la subpunctul 2.6.

Acest model exploatează bine posibilitățile oferite de sistemul de dezvoltare FPGA, asigurând o funcționare complet paralelă rețelelor neuronale construite. Totodată, datorită realizării speciale a algoritmului de învățare de tip Hebb, am reușit implementarea de neuroni pulsativi (spiking neurons) cu circuite digitale **fără module multiplicatoare**. Multiplicarea seriilor de impulsuri este efectuată de numărătoare și modulele de control ale acestora.

Inițial, am realizat circuite simple, testând dotările sistemului de dezvoltare FPGA utilizat (Ashenden, 1990) (Cîrstea, Dinu, & Nicula, 2001), cum ar fi numărătoare, unități logice și aritmetice, pe care le-am utilizat ulterior în implementarea modelului neuronal.

Proiectarea neuronului realizat în FPGA a fost divizată în două faze – care vor fi prezentate într-un capitol următor – și anume, proiectarea modului sinaptic și proiectarea corpului celular, adică al somei.



În continuare se va prezenta sistemul de dezvoltare FPGA utilizat și procesul proiectării-implemării hardware.

#### 4.5.1. Descrierea sumară a sistemului de dezvoltare FPGA

##### 4.5.1.1. Structura plăcii de dezvoltare

Părțile componente cele mai importante ale plăcii de dezvoltare FPGA, XESS XSA100, utilizate sunt (Figura 4.10):

- XC2S100 Spartan II FPGA - circuit integrat logic programabil (FPGA XC2S100 Spartan II cu 100 de mii de porți logice echivalente 144 pini QFP)
- XC9572XL CPLD – circuit integrat configurabil care face legătura între portul paralel al calculatorului personal și celelalte componente ale plăcii XSA
- Oscilator Dallas DS1075 – Oscilator programabil, responsabil pentru furnizarea semnalului de tact principal plăcii XSA cu intermediul CPLD
- Flash RAM: 256 Kbyte memorie Flash pentru salvarea de date și program de configurare FPGA
- SDRAM: 16 Mbyte memorie SDRAM, adresabil de către FPGA
- LED - afișor cu șapte segmente, pentru vizualizare variabile în cursul procesului de debug al circuitului realizat în FPGA
- Parallel Port (interfață paralelă): interfața prin care se încarcă biții de configurare FPGA iar apoi utilizabilă pentru eventualele date și/sau semnale de stare/control.
- Port PS/2 – port de conectare mouse sau tastatură PS/2 la placa XSA
- Port VGA – Se poate conecta un monitor VGA la placa XSA

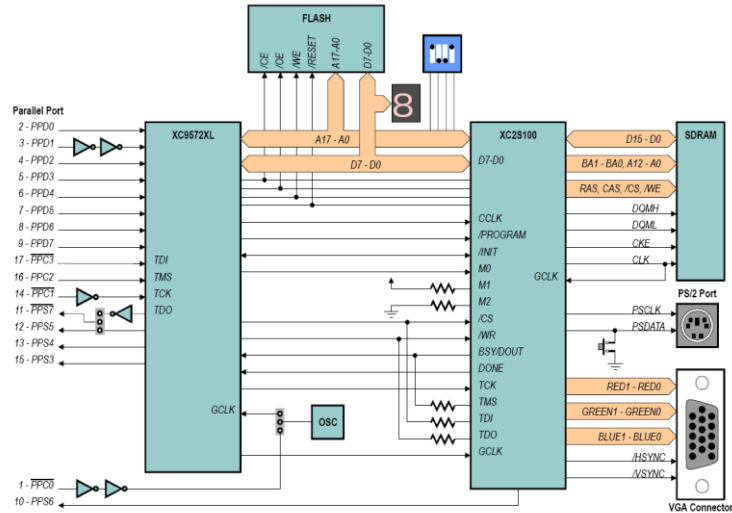


Figura 4.10 Schema arhitecturală a plăcii de dezvoltare FPGA XSA100, produsă de firma XESS (SUA), cu un circuit FPGA Xilinx XC2S100

Oscilatorul programabil Dallas DS1075 furnizează semnalul de tact pentru circuitele integrate FPGA și CPLD. Frecvența maximă a DS1075 este de 100MHz, care se poate diviza, obținând frecvențe de 50 MHz, 33,3 Mhz, 25 MHz.... 48.7KHz. Semnalul de tact este conectat la o intrare dedicată a circuitului CPLD care o transmite spre FGPA.

Comunicația între PC și placa XSA se desfășoară prin portul paralel (PP). Linia de control  $C_0$  a acestui port este direct conectată la oscilatorul DS1075 și este utilizată pentru programarea divizorului acestuia. Bitul de stare  $S_6$  al PP realizează o conexiune directă între circuitul FPGA și calculator. Restul de 15 semnale ale PP sunt controlate de circuitul CPLD. Programarea circuitului CPLD se efectuează prin trei linii de control ale PP ( $C_1$ - $C_3$ ),  $C_1$  furnizează un semnal de ceas de sincronizare, prin  $C_2$  se transmit datele seriale iar prin  $C_3$  starea circuitului CPLD. Starea CPLD-ului se poate citi prin bitul de stare  $S_7$ .

Prin biții registrului de date al PP se pot trimite date spre FPGA, iar din acesta se pot transmite semnale spre calculator prin biții rămași liberi în registrul de stare. Producătorul a furnizat câte un program utilitar pentru testarea și plăcii XSA și încărcarea biților de configurare FPGA respectiv CPLD.



#### 4.5.1.2. Elementul central al plăcii XSA100, circuitul FPGA Xilinx XC2S100

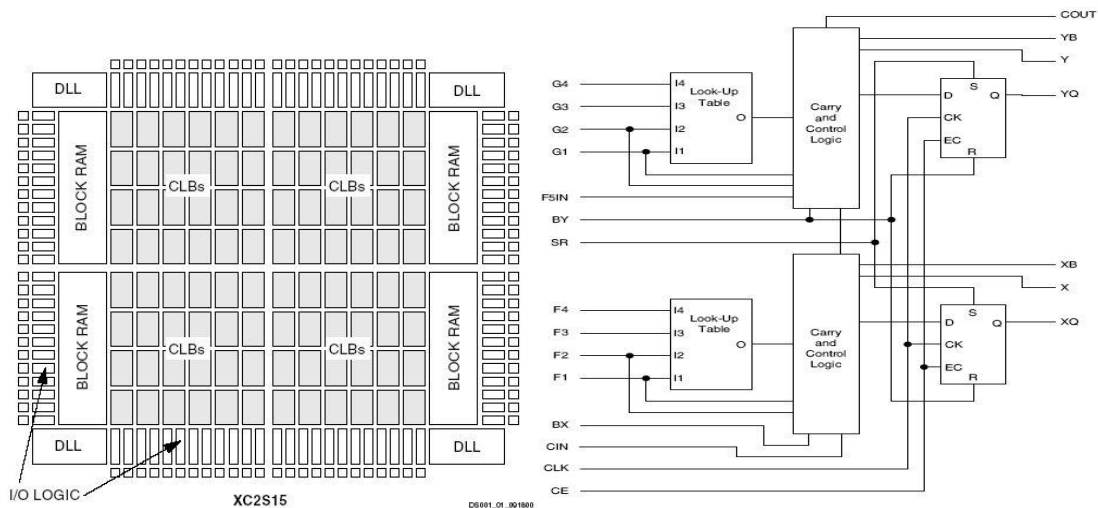


Figura 4.11 Structura internă a circuitelor FPGA din familia Spartan (stânga), părțile componente al unui bloc logic programabil (dreapta) (CLB)

Circuitul Xilinx XC2S100 este membru al familiei Spartan-II FPGA. Circuitele FPGA din această familie au o structură ușor de înțeles și de utilizat, precum se poate urmări și în Figura 4.11.

Această arhitectură de bază se bazează pe *Blocuri Logice Programabile (CLB)*, care sunt înconjurată de Unități Logice Configurabile de Intrare/Ieșire (I/O Logic). Ca elemente arhitecturale se mai pot aminti blocuri RAM interne. Toate aceste elemente constructive sunt conectate printr-o magistrală de control extrem de flexibil programabilă, care este realizată pentru a asigura un număr infinit de procese de reconfigurare. Șirul de biți, care determină configurarea acestor conexiuni între elementele de bază ale FPGA, se poate descărca prin porturi seriale, paralele sau dintr-o memorie PROM externă. Un aspect important este faptul că aceste circuite sunt capabile să funcționeze la frecvențe de până la 200 MHz, ceea ce permite realizarea unor rețele neuronale artificiale foarte rapide cu ajutorul acestora.

#### 4.5.1.3. Mediul de dezvoltare software VHDL Xilinx ISE Webpack

Firma Xilinx, producătorul circuitului FPGA utilizat, pune la dispoziția utilizatorilor acest mediu gratuit de dezvoltare al proiectelor VHDL. Acest proces de dezvoltare, cu ajutorul programului Xilinx Webpack 6.3 se compune din mai mulți pași, după cum arată și Figura 6.3. În primul pas se descrie circuitul logic de realizat cu ajutorul unui limbaj HDL (hardware description language) ales (VHDL sau VERILOG) sau se

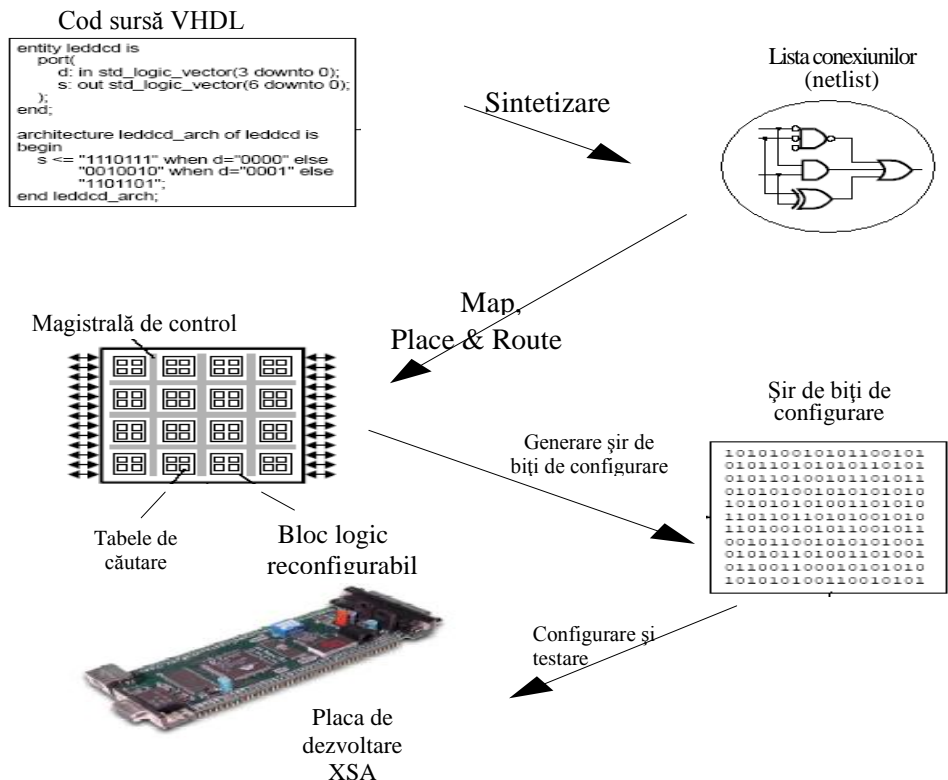


Figura 4.12 Fazele procesului de proiectare FPGA în mediul Xilinx Webpack

realizează o schemă cu ajutorul utilitarului grafic inclus

Programul de sintetizare logică translatează fișierul sursă scris sau schema desenată în formă de “netlist”. Un fișier netlist conține descrierea tipurilor de elemente de bază utilizabile pentru realizarea unui circuit care se comportă conform descrierii date și a conexiunilor dintre aceste elemente.

Utilitarul de implementare (Implementation Tool) realizează o hartă a acestor elemente, pe baza celor aflate efectiv în circuitul FPGA și a conexiunilor posibile între acestea. Blocurile logice programabile din care se compune XC2S100 sunt de fapt formate din subansamble ca tabele look-up (LUT) care pot realiza operații logice. Un utilitar numit Mapping Tool grupează porțile logice descrise de fișierul netlist, în așa fel ca să se poată implementa cu cât mai puține LUT-uri. În fine utilitarul de așezare *Place and Route* asignează aceste grupuri de porți logice unor CLB-uri, deschizând sau închizând “comutatoarele” corespunzătoare în matricea de comandă a magistralei de conexiuni a circuitului FPGA, conectând porțile logice în modul descris în program.

Din această matrice se generează un fișier binar (șirul de biți de configurare).

Odată cu descărcarea acestui fișier în circuit, proiectul realizat începe să funcționeze. Toate aceste utilitare sunt parte integrantă a mediului de dezvoltare Xilinx ISE WebPack 6.3.

#### 4.5.2. Procesul de realizare practică a modelului neuronal, implementarea sinapsei

##### 4.5.2.1. Fazele implementării sinapsei în FPGA

După cum am amintit, construcția unui circuit în integratul FPGA se poate realiza în două feluri:

- Proiectare cu ajutorul circuitelor elementare secvențiale și combinaționale oferite de mediul de dezvoltare.
- Proiectare prin descriere (algoritm) cu ajutorul codului VHDL

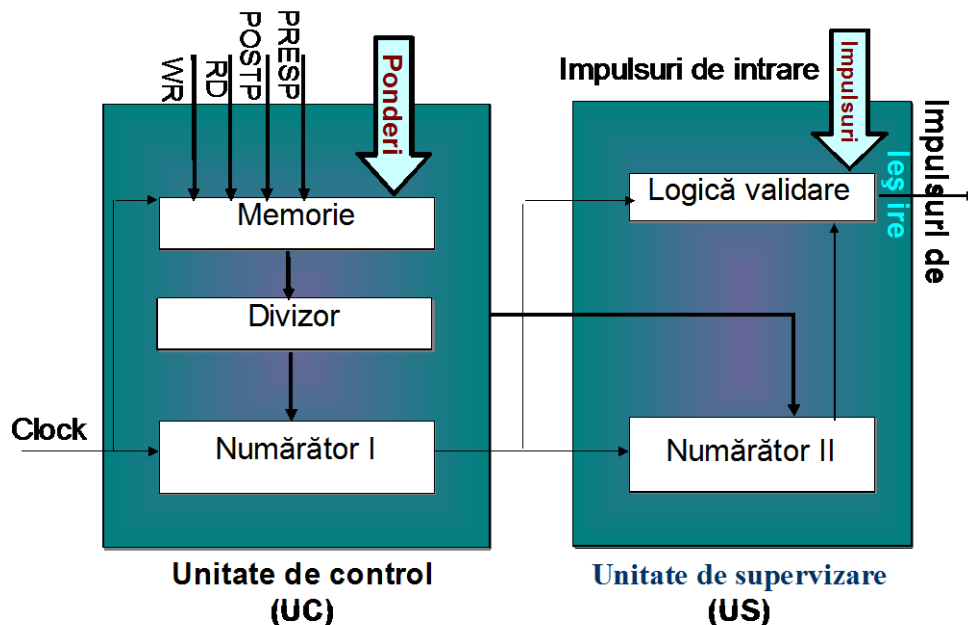


Figura 4.13 Schema structurală a sinapsei

În cursul proiectării a fost necesar împărțirea sinapsei în mai multe circuite funcționale diferite, și definirea exactă a intrărilor și ieșirilor acestora.

În concordanță cu funcționarea sinapsei biologice și a modelului teoretic dezvoltat pentru un neuron pulsativ artificial, sinapsa electronică realizată preia impulsurile pre-sinaptice, le multiplică conform valorii ponderii actuale și le transmite, ca impulsuri post-sinaptice, către corpul neuronal (somă). Impulsurile de intrare sosesc în sinapsă prin intrarea *spike\_in* și cele de ieșire o părăsesc prin ieșirea *spike\_out*. Sinapsa va fi de fapt un circuit logic secvențial, în care este necesar a sincroniza diferite părți componente, cu ajutorul unui semnal de tact.

Valorile ponderale sunt stocate într-un registru de  $N$  biți, ele trebuind inițializate înainte, respectiv citite în timpul și după procesul de învățare. Aceste operații se efectuează cu ajutorul semnalelor RD și WR (Figura 4.13), sincronizate pe semnalul de tact exterior.

Deoarece este vorba despre mai multe sinapse în cadrul aceluiași neuron, pentru a putea executa operațiile menționate anterior am implementat o logică de adresare și o magistrală de ponderi. Selecția unei sinapse se realizează setând semnalul CE la valoare logică "1", ceea ce înseamnă permiterea citirii sau programării ponderii acelei sinapse. În cazul în care CE este 0 logic, nu pot avea loc aceste operații, deoarece neuronul este în stare funcțională (rețeaua din care face parte neuronul respectiv este activă).

Scrierea ponderilor în sinapse se compune din următorii pași:

1. Se selectează sinapsa setând CE=1,
2. Se activează semnalul WR,
3. Se așează o valoare pe magistrala de ponderi (de exemplu prin portul paralel),
4. La frontul crescător al semnalului de ceas această valoare se va înscrie în registrul intern al sinapsei,
5. Se șterge CE, separând sinapsa de la magistrala de ponderi.
6. Operația de citirea a ponderilor din sinapse are următorii pași:
7. Se selectează sinapsa setând CE=1,
8. Se activează semnalul RD,
9. La frontul crescător al semnalului de ceas sinapsa va așeza valoarea ponderii sale pe magistrala de ponderi,
10. Se citește ponderea prin PP al PC,
11. Se șterge CE, separând sinapsa de la magistrala de ponderi.

Conform modelului utilizat, învățarea se efectuează pe baza IA pre- și post-sinaptice. Rezultă, că sinapsa trebuie să fie dotată cu două intrări care să semnaleze valorile acestor impulsuri.

Din punct de vedere funcțional, sinapsa se împarte în:

- ✓ Unitate de control
  - Stochează valorile ponderii într-un registru intern,
  - Asigură citirea și scrierea ponderilor,
  - Recepționează semnalele impulsurilor pre- și post-sinaptice utilizate în învățare,
  - Conține un circuit divizor realizat din porți logice, care calculează partea întregă a câtului divizării frecvenței de tact cu valoarea ponderii
  - Conține un numărător care divizează frecvența de bază (tact).
- ✓ Unitate de supervizare
- ✓ Numărător
- ✓ Logică de validare.

Divizorul calculează valoarea cu care trebuie divizată frecvența de bază, pentru a realiza un număr de impulsuri de ieșire proporțional cu ponderea sinapsei.

Numărătorul unității de supervizare va număra impulsurile de ieșire, iar dacă acest număr a atins valoarea ponderii atunci interzice emiterea altor astfel de impulsuri până la sosirea unui nou impuls pre-sinaptic.

Deoarece resursele hardware ale FPGA sunt finite, ordinul de mărime al unor variabile a fost limitat. În faza inițială a proiectării, valorile ponderilor au fost reprezentate pe patru biți, ceea ce înseamnă, că sosirea unui impuls de intrare poate determina emiterea a maxim cincisprezece impulsuri post-sinaptice.

Multiplicatorul de impulsuri implementat funcționează după ideea următoare: două numărătoare conectate în serie, divizează frecvența de bază cu  $N_1 * N_2$ , în așa fel, încât frecvența semnalului de tact de ieșire să rămână constantă. Valoarea numărătorului  $N_2$  este definită de pondere, deci cunoscând frecvența de bază putem calcula valoarea  $N_1$ .

Semnalul de tact de intrare este de 32 de ori mai rapid decât cel de ieșire. Impulsurile de ieșire nu sunt distribuite în mod absolut liniar. Liniaritatea completă putea fi atinsă, dacă tactul de intrare ar fi cu mult mai mare decât cel de ieșire și ar fi cel mai mic multiplu comun al valorilor 2,3,5,7,9,11,13. Am renunțat la această cale, deoarece astfel, numărătorul  $N_1$  este mai simplu de implementat și utilizează mai puține resurse.

#### 4.5.2.2. Descrierea detaliată a funcționării sinapsei

Valorile ponderilor sunt înscrise în Registrul unității de control. De aici, aceste valori vor ajunge la intrarea numărătorului 1, prin circuitul de divizare. Circuitul divizor realizează o împărțire întreagă a frecvenței de bază (predefinită) cu valorare ponderii. Astfel am aflat la al câtelea front crescător al semnalului de ceas va trebui să genereze impuls de ieșire.

Numărătorul unității de control va număra de la zero până la valoarea dată de divizor. Dacă a atins valoarea aceasta, atunci va emite un impuls de durată egală cu perioada semnalului de tact, se va reseta și va reîncepe numărarea la fronturile pozitive de tact.

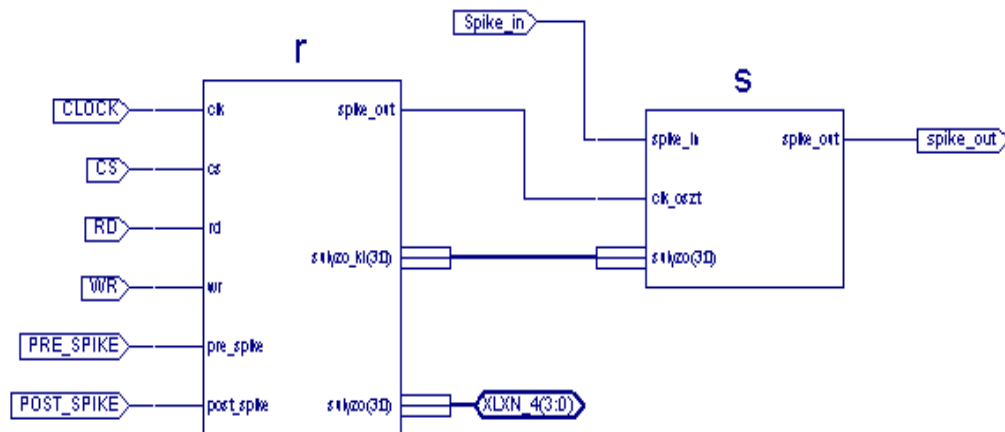


Figura 4.14 Schema de conectare interioară a sinapsei

Numărătorul unității de supervizare numără de la zero, la fiecare front crescător al impulsurilor emise de unitatea de control, până la valoarea ponderii. Acest numărător este resetat și validat de un impuls pre-sinaptic. Dacă acest numărător a atins valoarea dată de pondere, atunci se oprește și prin logica de supervizare împiedică emiterea altor impulsuri post-sinaptice, până la sosirea unui nou impuls pre-sinaptic.

Modificarea ponderii duce la modificarea valorii limită a numărătorului 1. Astfel pentru o pondere mare impulsurile de ieșire din acest numărător vor fi mai dense precum și cele post-sinaptice, iar pentru pondere mică ambele vor fi mai rare.

Unitatea de control (r), precum și cea de supervizarea (s) a sinapsei a fost scrisă în cod VHDL, din care am generat câte un modul, utilizat în schema de conectare din Figura 4.14.

#### 4.5.2.3. Îmbunătățirea sinapsei

După efectuarea de teste cu implementarea sinaptică prezentată, am ajuns la concluzia, că necesită unele modificări, pentru a îmbunătăți capacitatea de învățare a rețelei.

Imaginile din Figura 4.15 și Figura 4.16 prezintă implementarea îmbunătățită a sinapsei. Se poate observa, că a apărut un modul de memorie, care are ca scop menținerea stării semnalului pre-sinaptic de-a lungul procesului de învățare. Procesul de învățare începe cu frontul crescător generat de apariția unui IA pre-sinaptic și durează un număr de cicluri de tact date de valoarea ponderii. Reprezentarea ponderilor a rămas neschimbată, pe patru biți, deoarece experimentele efectuate, după cum se va vedea mai târziu, au arătat că această precizie este suficientă pentru ca RNP construită să poată rezolva problema propusă. Rezultă, că modulul de memorie amintit, va trebui să mențină informația pe o durată de maximum 15 cicluri de tact, timp în care nu poate sosi alt impuls de intrare. După expirarea acestui interval, un semnal de control (Clearmem, Figura 4.15) va șterge conținutul modulului de memorie.

În sinapsă am mai introdus un semnal de control suplimentar, care permite sau nu rularea procesului de învățare, în funcție de nevoia de a învăța RNP sau a-i testa performanțele după învățare.

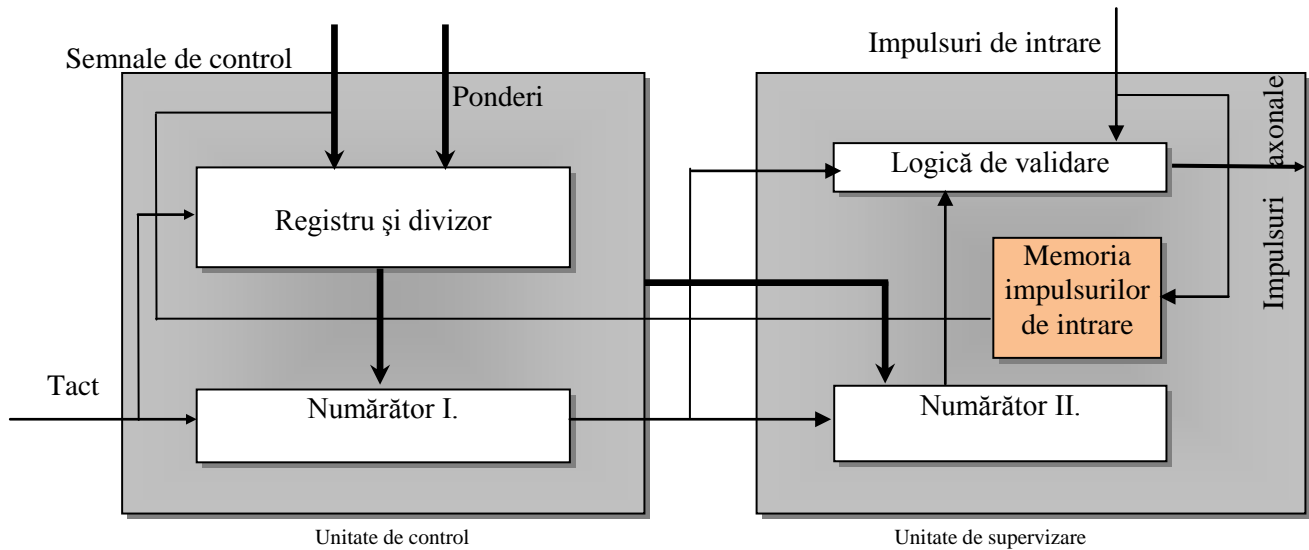


Figura 4.15 Schema de principiu a sinapsei

Figura 4.16 prezintă schema de conectare interioară a sinapsei modificate, unde  $r$  este UC,  $s$  este US, iar FDC modulul de memorie (bistabil D).

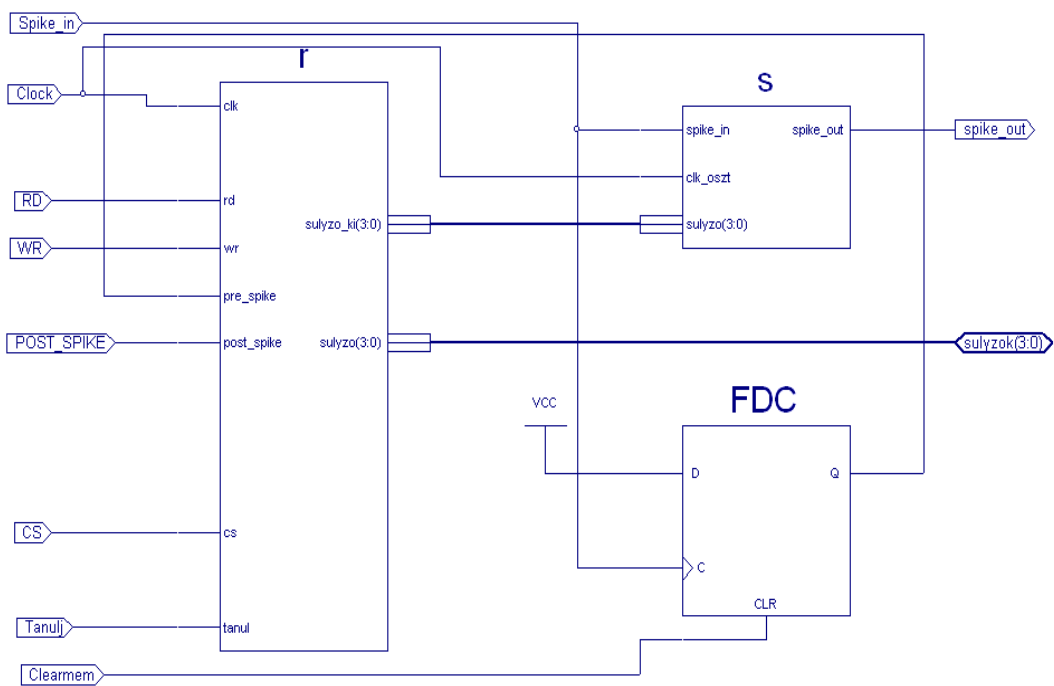


Figura 4.16 Schema de conectare a sinapsei modificate

#### 4.5.2.4. Unitatea serială de intrare date (USID)

Această unitate asigură transferul valorilor ponderilor de la calculator la rețeaua aflată în FPGA, respectiv livrarea impulsurilor de intrare rețelei neuronale.

Unitatea serială de intrare date, citește semnalul de tact al transferului serial (CLOCK), biții de date transmiși serial (DATA), respectiv semnalele de activare a sinapselor pentru citire/scriere (CEN – Chip Enable) și semnalul de validare a impulsurilor de intrare (SPE – Spike Enable) de la portul paralel al calculatorului (Figura 4.17).

Datele sunt înscrise într-un registru de deplasare, prin linia de date la frontul pozitiv al semnalului de tact serial, respectând următorii pași:

- Un bit de date este înscris pe linia de date,
- Este generat un front crescător pe linia CLOCK (Serial\_clock),

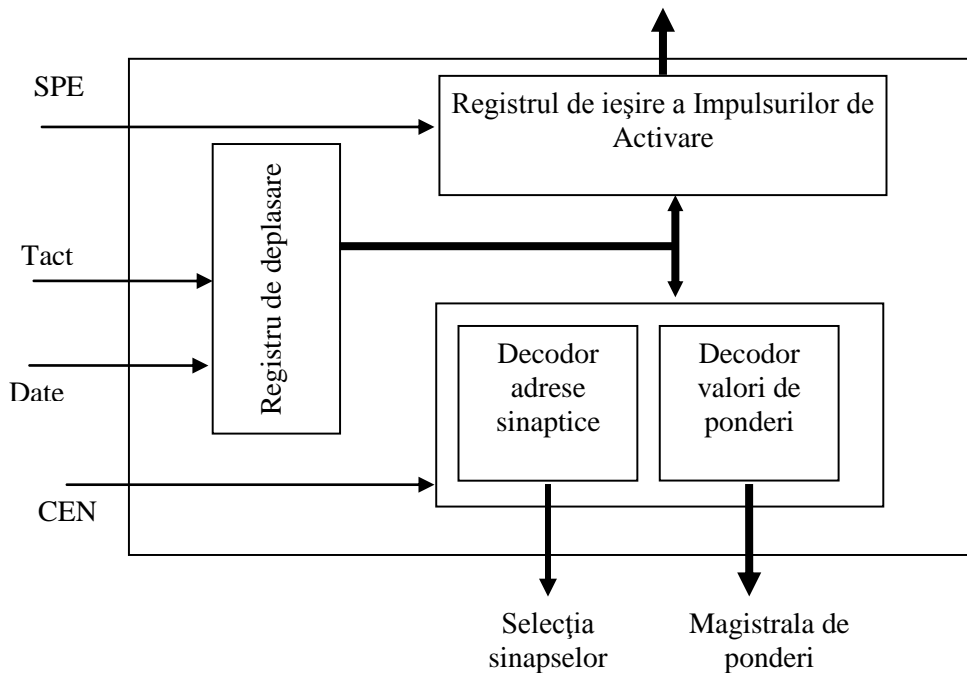


Figura 4.17 Schema constructivă a USID

Transferul unui bit se realizează într-un singur ciclu de tact. În Figura 4.18 se poate vedea o diagramă de timp a unei transmisii de patru pachete de date pe linia serială către registru de deplasare din USID FPGA (ambele semnale sunt reprezentate în formă negată).

Acest proces este identic atât în cazul citirii/scrierii valorilor ponderilor, cât și în cazul transmisiei impulsurilor de intrare și a fost folosit numai pentru experimente în care s-a urmărit evoluția rețelei fără pretenții în ceea ce privește performanțele în viteză de prelucrare.

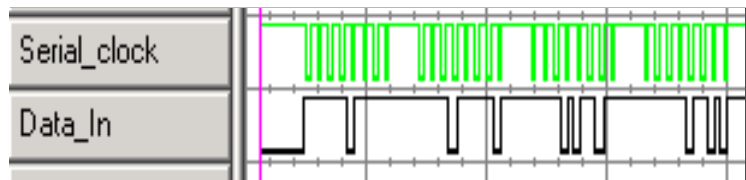


Figura 4.18 Transmisia serială de date către FPGA

În cazul transferului de ponderi,

pe lângă valorile acestora se trimite și adresa fiecărei valori. Astfel biții cei mai puțini semnificativi conțin adresele (reprezentate pe 3 sau 4 biți, în funcție de numărul de sinapse utilizate în neuronul respectiv), urmați de cei patru biți reprezentând ponderile.

O logică din USID decodifică adresele și activează numai o singură sinapsă pentru programarea ponderii, valoarea acesteia din urmă fiind de asemenea decodificată din registru de deplasare în care a intrat și înscrisă pe magistra de ponderi.

Pentru ca ponderile să poată fi citite într-un moment ulterior, USID este deconectată de la magistra de ponderi.

### 4.5.3. Fazele proiectării somei (al corpului neuronal)

Ca și în cazul natural, unde soma este Unitatea Centrală de Prelucrare a neuronului biologic, în cadrul NP artificial, realizat în FPGA, circuitul care implementează soma este cel care efectuează prelucrarea impulsurilor de intrare. Circuitul somei se compune din patru părți componente majore, după cum se poate vedea și în Figura 4.19:

- Unitatea de citire și modulare a impulsurilor de intrare sosite de la sinapse (SYNIN)
- Unitatea de calcul și reîmprospătare a valorii potențialului de membrană (MPOT)
- Registru de stocare a potențialului de prag încărcat (THRH)
- Circuitul comparator de generare a impulsurilor de activare a neuronului (AXON)

Prima versiune testată a neuronului pulsativ artificial a fost implementată cu opt sinapse încorporate, deci un astfel de neuron putea fi conectat la cel mult opt neuroni învecinați prin sinapsele sale. Pe baza unor calcule efectuate, a parametrilor neuronilor naturali și a estimării necesității de stocare a informației a algoritmilor utilizați în viitor, am decis reprezentarea valorii potențialului de membrană pe opt biți. Această decizie a fost influențată de către considerentul economisirii resurselor sistemului de dezvoltare FPGA, în vederea posibilității realizării de rețele mai mari cu ajutorul acestuia.

Impulsurile post-sinaptice, care ajung la somă sunt modulate de unitatea de intrare a acesteia, cu un anumit factor de multiplicare (în primele experimente această valoare a fost setată la cinci sau zece), ele fiind transmise apoi unității MPOT. Pentru ca potențialul de membrană să fie compatibil, adică complet comparabil cu potențialul de prag (threshold potential), și acesta din urmă a fost reprezentat pe opt biți.

Neuronul pulsativ artificial astfel realizat va emite un IA atunci și numai atunci, când circuitul de ieșire al somei (AXON-ul) – un comparator pe 8 biți – va semnaliza, că potențialul de membrană calculat a depășit potențialul de prag.

#### 4.5.3.1. Descrierea detaliată a funcționării unității centrale a somei, adică a modului de calcul a potențialului de membrană

Intrarea acestui modul, scris în VHDL, este de fapt – după cum se vede și în Figura 4.19 – ieșirea modului de citire a impulsurilor post-sinaptice. Modulul MPOT primește ca intrare și semnalul de tact principal, semnalul de reacție a stării axonale (dacă s-a emis IA) și semnalul de autorizare a citirii valorii potențialului de membrană (OE – Output Enable). Ieșirea acestui modul este valoarea nou calculată a potențialului de membrană. Calcularea acestei valori se efectuează în felul următor:

La fiecare front crescător al semnalului de tact, modulul MPOT adună valoarea potențialului de membrană stocat în ciclul anterior cu valorile ce se prezintă la intrare. În primele experimente, valoarea inițială a fost setată la zero, iar potențialul de repaus la valoare de zece.

Dacă potențialul de membrană a depășit cel de prag și s-a emis IA, atunci această adunare nu se mai efectuează, ci se resetează potențialul de membrană la zero, conform fazei de hiperpolarizare, prescrisă de modelul de urmat. Dacă valoare de intrare este egală cu zero, adică nu a fost prezent nici un impuls post-sinaptic, atunci potențialul de membrană scade cu cinci (valoare variată în diferitele experimente) unități în fiecare ciclu de tact.

După cel de-al 32-lea ciclu de tact, dacă semnalul de reacție de la ieșirea axonală este de valoare neschimbată, logic zero, potențialul de membrană va scade până la potențialul de repaus. Acest aspect este necesar pentru a implementa proprietatea neuronilor naturali care spune că ne se poate genera IA numai dacă într-un anumit interval de timp sosesc destule impulsuri post-sinaptice pentru a ridica potențialul de membrană peste cel de prag.

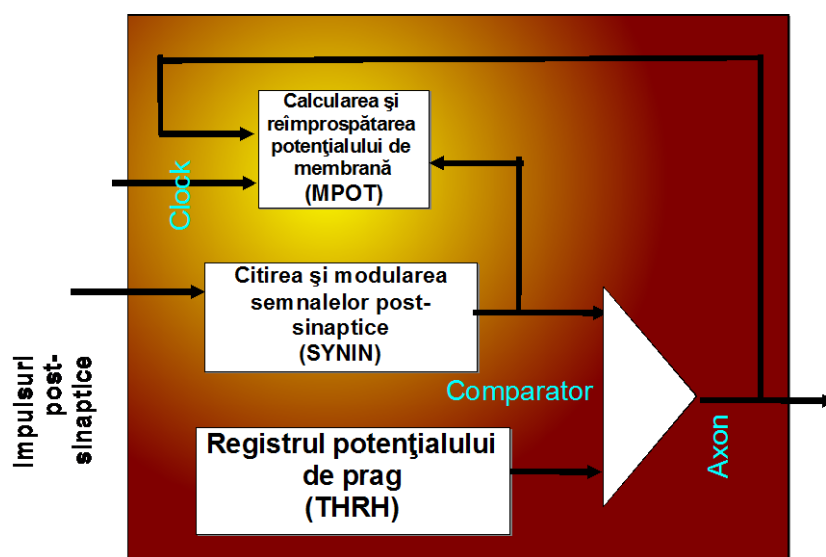


Figura 4.19 Schema constructivă a circuitului ce realizează soma neuronului



Figura 4.20 prezintă schema de conectare internă a somei, într-o fază intermediară a proiectării. În această fază modulul de calculare a potențialului de membrană era compus din două părți (Figura 4.20.) și anume modulele Regiszt și ADD8. Acesta din urmă a fost mai târziu inclus în Regiszt pentru a rezolva unele probleme de sincronizare ivite. Modulul mux\_8\_of\_8x8 este de fapt un multiplexor 8 din 2x8 scris în VHDL, care servește la citirea sau reinițializarea din exterior a valorii actuale a potențialului de membrană. Modulul SR8CE este un registru de deplasare utilizat pentru inițializarea potențialului de prag, iar COMPM8 este circuitul comparator al cărei ieșire coincide cu axonul neuronului.

Figura 4.21 prezintă un circuit de testare a funcționării modulului MPOT în care s-a utilizat o versiune a neuronului cu doar șase sinapse, pentru o transparență mai bună a testului. Tot în acest sens

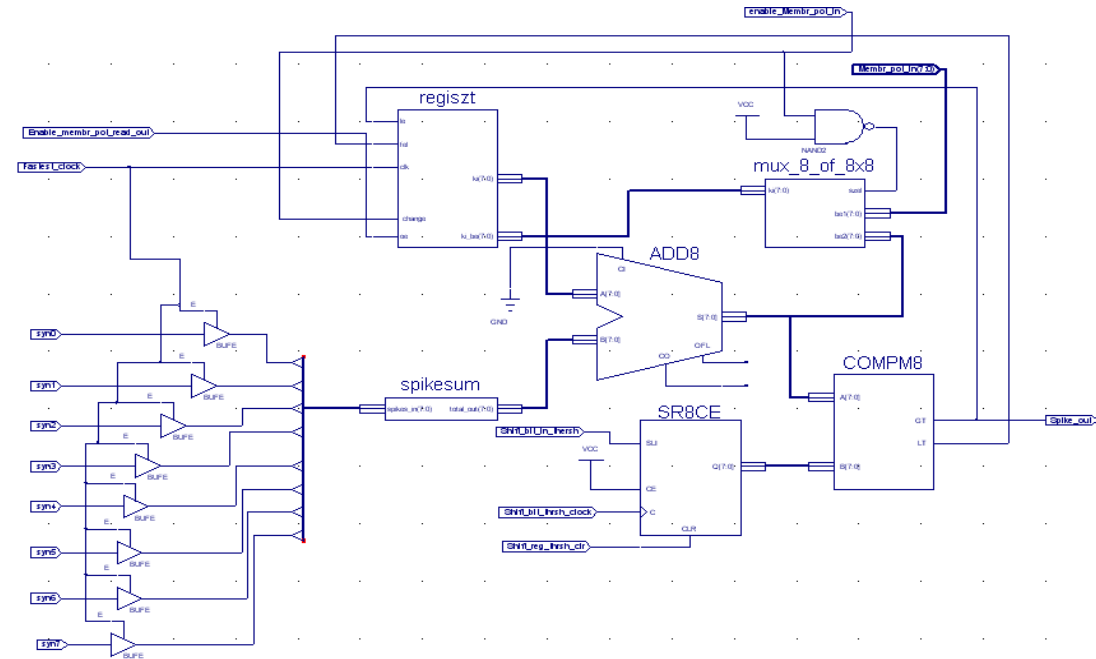


Figura 4.20 Modelul somei implementate, într-o fază intermediară a proiectării

am utilizat un comutator extern – aflat pe placa de dezvoltare – în locul reacției axonale, pentru a verifica corectitudinea efectuării operațiilor propuse de către modulul MPOT. Am procedat la realizarea de astfel de circuite de testare, pentru a avea certitudinea funcționării corecte a subansamblelor modelului neuronal, precauție necesară datorită relativei complexități a implementării finale.

Semnalele din Figura 4.21. au fost

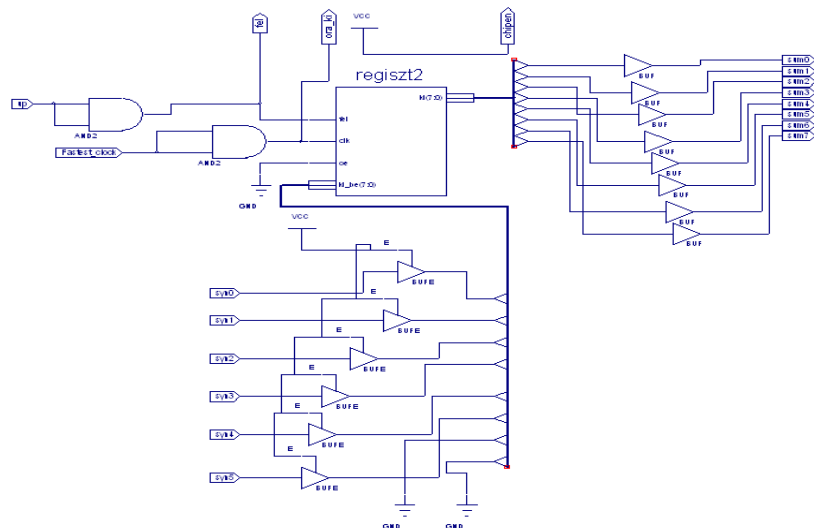


Figura 4.21 Circuitul de testare a modulului central al somei, MPOT

șantionate cu ajutorul unui Analizor Logic Tektronix cu 36 de canale și arată variația potențialului de membrană al unui neuron în cursul unui experiment de testare a funcționării somei. Valoarea



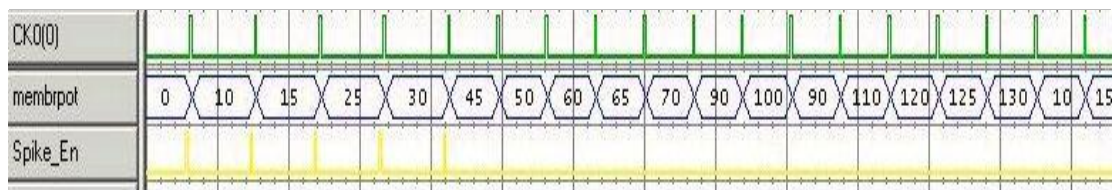


Figura 4.23 Valori ale potențialului de membrană, eșantionate în timpul testării

reprezentată atunci, când sunt prezente impulsuri post-sinaptice iar dacă acestea dispar sau expiră intervalul de 32 de cicluri de tact, potențialul scade treptat respectiv scade brusc la valoarea de repaus.

După efectuarea a numeroase astfel de teste am ales un model mai simplificat al somei (Figura

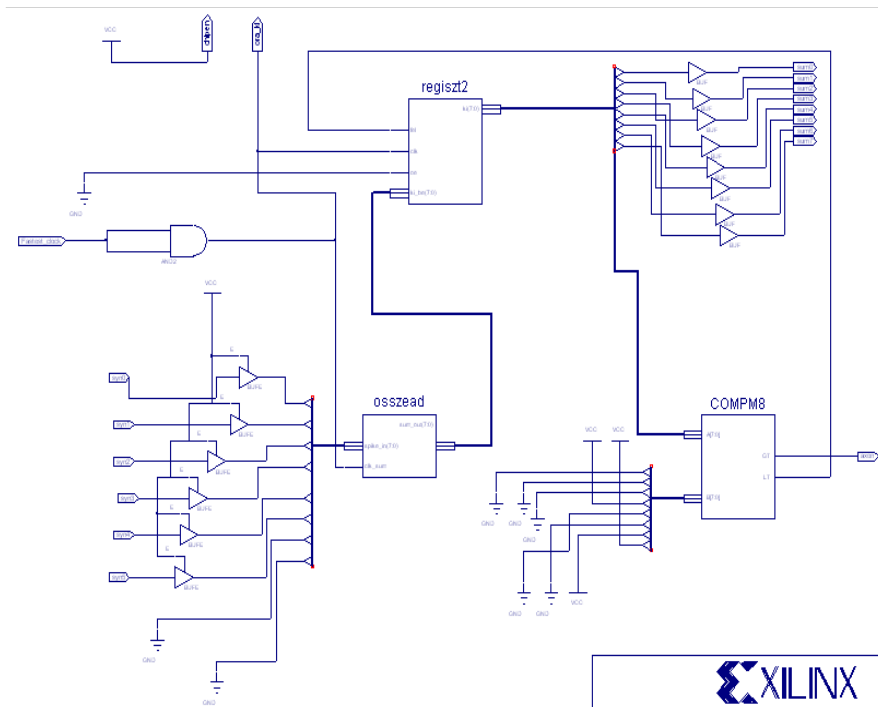


Figura 4.22 Schema somei implementate în neuronul pulsativ complet

4.22.) pe care am utilizat-o la implementarea completă a unui NP.

Diferența față de schema din Figura 4.20. este lipsa multiplexorului de citire a potențialului de membrană, care se va utiliza ulterior, respectiv versiuni optimizate ale modulelor MPOT și SYNIN. Valoarea potențialului de prag a fost fixată la 200.

#### 4.5.4. Stadiul actual de dezvoltare a somei

Forma finală a implementării somei, care a fost utilizată la realizarea RNP prezentate în subcapitolul următor are următoarele componente (Figura 4.24):

- Unitatea de citire și modulare a impulsurilor de intrare sosite de la sinapse (SYNIN)
- Unitatea de calcul și reîmprospătare a valorii potențialului de membrană (MPOT)
- Registrul de stocare a potențialului de prag încărcat (THRH)
- Modulul de stocare a stării ieșirii axonale (AXOUT)
- Circuitul comparator de generare a impulsurilor de activare a neuronului (AXON)

Se poate observa, că acest model s-a îmbogățit cu un modul nou – față de cel din Figura 4.19, și anume modulul de stocare a stării ieșirii axonale (AXOUT). Introducerea acestuia a fost necesară pentru ca intrările POSTSPIKE ale sinapselor, care este de fapt ieșirea axonală a aceluiași neuron, să semnalizeze starea de activare (de aprindere) a neuronului pe întreaga durată a procesului de învățare. Astfel s-a asigurat (urmând regula Hebb) creșterea ponderilor în cazul în care neuronii pre- și post-sinaptice sunt simultan activi.



#### 4.6.2. Câteva rezultate experimentale de testare a NP implementat, cu diferite configurații parametrice

În Figura 4.25. este prezentat un grup de semnale măsurate cu analizorul logic, pe care se pot observa următoarele:

- ✓ Semnalul trasat cu linie verde (Serial-clock), este tactul transmisiei seriale, la al cărui front crescător se înscriu impulsurile de intrare ale neuronului. Se poate vedea, că la acest experiment am transmis impulsuri de intrare simultan pentru mai multe sinapse.
- ✓ Cu linii galbenă s-a trasat semnalul Spike-En, care în stare logică 1 validează înscrierea intrărilor în sinapse.
- ✓ Semnalele de culoare mov (Spikes-out[0÷7]) reprezintă ieșirile celor opt sinapse, emițând un număr de impulsuri post-sinaptice egale cu valoarea ponderii sinapsei respective
- ✓ Cu linii albastră sunt reprezentate valorile potențialului de membrană a neuronului, în diferite momente ale rulării testului

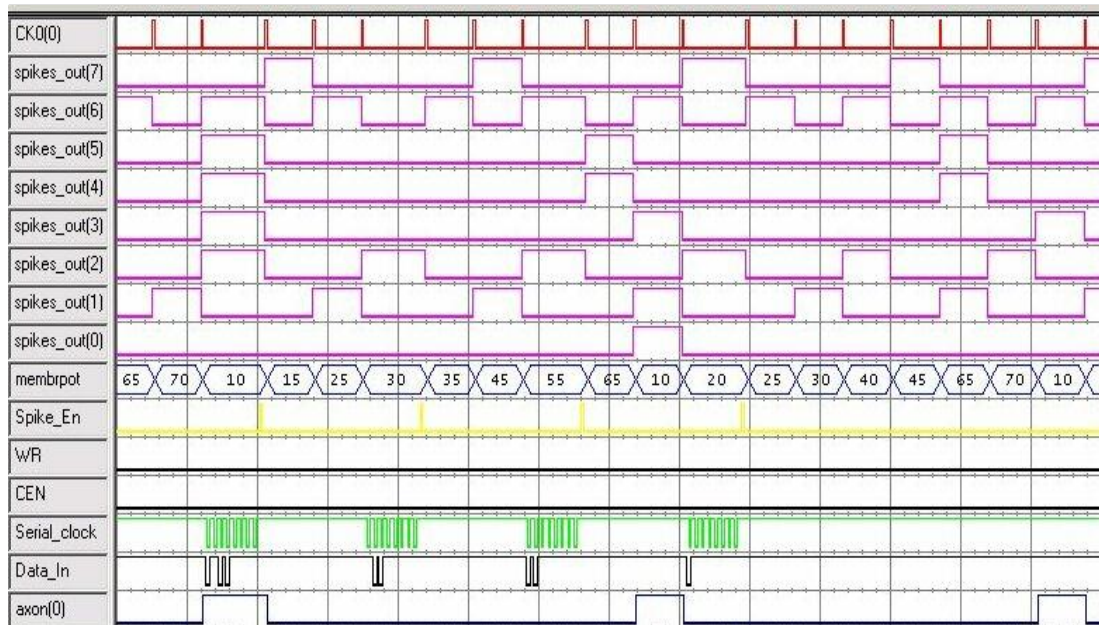


Figura 4.25– Impulsuri axonale într-un experiment cu potențial de prag scăzut

Aici se poate observa, că la frontul descrescător al tactului principal (roșu, CK0) potențialul de membrană scade cu de atâtea ori 5 unități (aceasta fiind valoare setată a factorului de modulare din SYNIN la acest experiment), câte sinapse emit impuls post-sinaptic în același ciclu de tact.

Totodată se pot observa momentele de activare a neuronului (emiterea de impulsuri axonale, semnalul axon(0)). Acestea apar la momentele în care potențialul de membrană depășește cel de prag (setat aici la 80). La frontul crescător al tactului CK0 din ciclul imediat următor, potențialul de membrană scade la 10.

Figura 4.26. respectiv Figura 4.27. prezintă două segmente ale unor semnale eșantionate într-un alt test asemănător, al cărui parametrii au fost setați astfel:

În comparație cu cazul anterior, unde valorile ponderilor au fost inițializate cu valori între 1 și 3, aici am utilizat valori mai mari, aleator alese (ponderi[8]={4,8,14,6,7,5,9,10});

- ✓ Potențialul de prag a fost setat la 143,
- ✓ Potențialul de repaus era 10,
- ✓ Factorul de modulare SYNIN=5.



Pe aceste două figuri se poate urmări un ciclu complet. În Figura 4.26, se poate observa, cum după un IA emis, potențialul de membrană intră în faza de hiperpolarizare, scăzând brusc la zero, sub

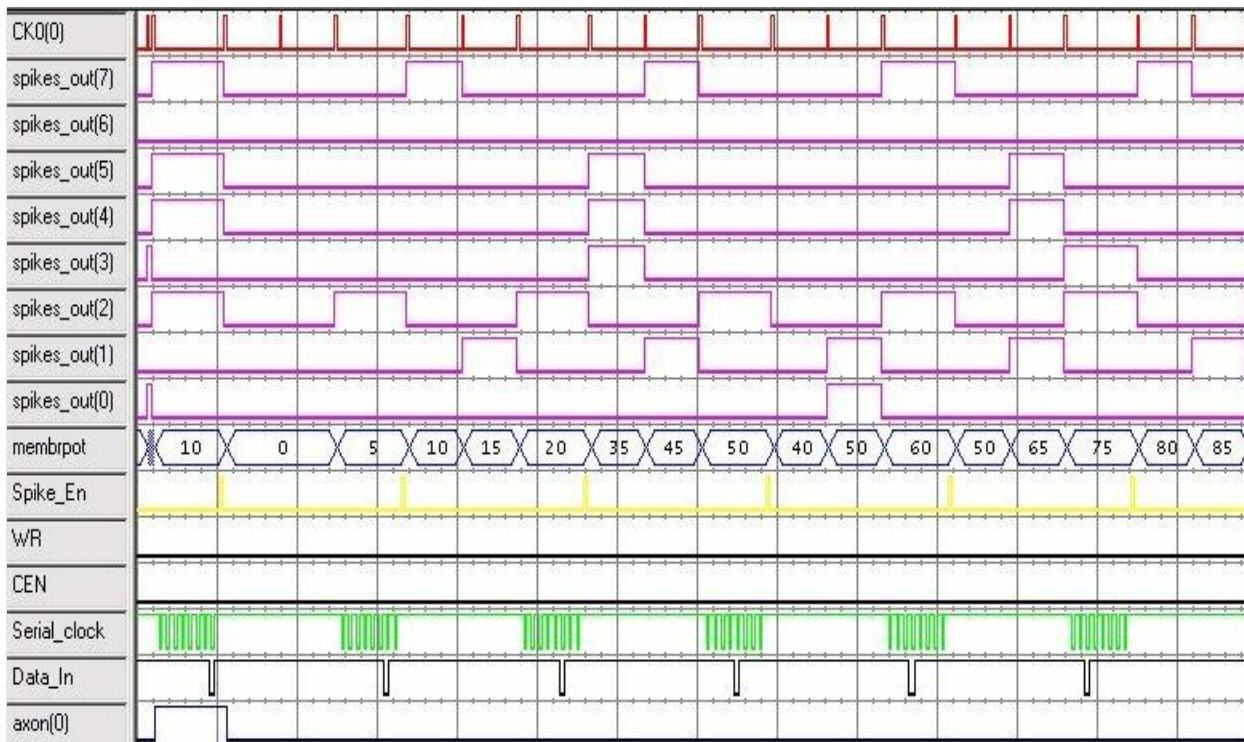


Figura 4.26 IA utilizând valori inițiale mari ale ponderilor și impulsuri de intrare distribuite – prima parte

valoarea de repaus. Apoi, în funcție de numărul sinapselor care emit impulsuri post-sinaptice, potențialul de membrană crește cu cinci sau multiplu de cinci unități. Dacă într-un anumit moment, nici una dintre sinapse nu a fost activă, atunci această valoare scade cu zece unități în fiecare ciclu de ceas.

Partea a doua a măsurătorii este prezentată în Figura 4.27., unde se poate urmări, cum evoluează valoarea potențialului de membrană, până când neuronul se aprinde din nou și emite IA.

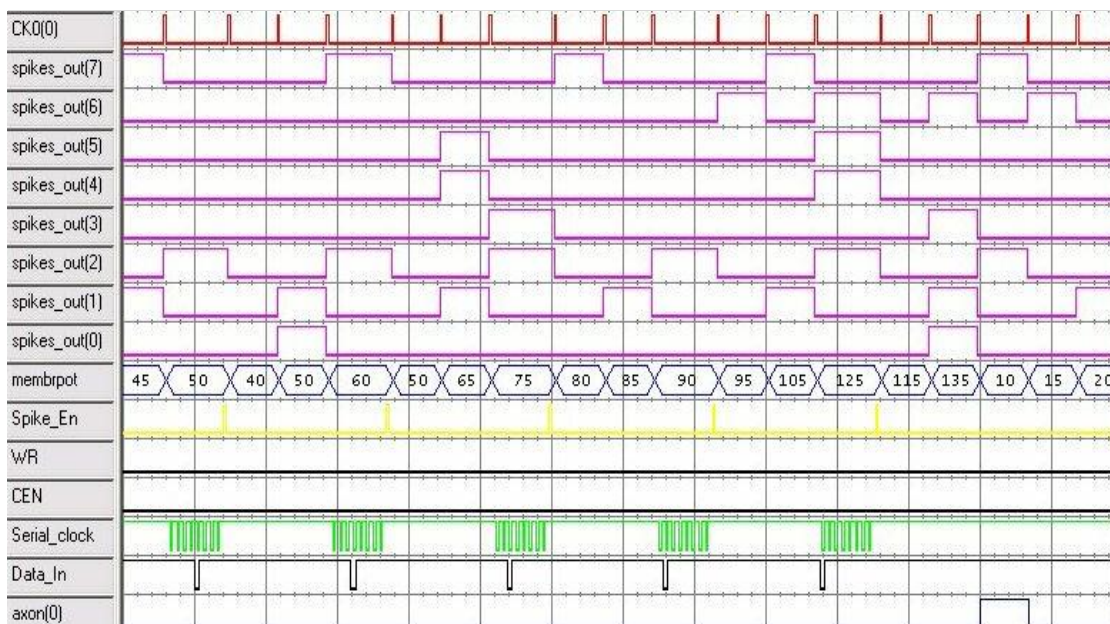


Figura 4.27– IA utilizând valori inițiale mari ale ponderilor și impulsuri de intrare distribuite – partea a doua

## 4.7. Implementarea unei rețele neuronale pulsative artificiale în circuitul FPGA. Rezultate experimentale

### 4.7.1. Implementarea unei RNP în circuit FPGA

După mai multe luni de experimentări, am decis să implementez neuronul pulsativ cu versiune de nouă sinapse. Cu ajutorul a numai doi astfel de neuroni, am reușit să construiesc o rețea neuronală pulsativă artificială, capabilă de a diferenția orice două modele de intrare – chiar și în prezența zgomotului – prezentate într-o matrice de 3x3 puncte. Pentru a utiliza cât mai eficient resursele circuitului FPGA în care s-a efectuat implementarea, fără repercusiuni asupra performanțelor rețelei, modelul neuronal a fost puțin modificat. Astfel, potențialul de membrană al somei și potențialul de prag au fost reprezentate pe șapte biți, ponderile sinapselor rămânând reprezentate neschimbat pe patru biți. Topologia rețelei construite este prezentată de Figura 7.1. Ca și la testarea neuronilor individuali, comunicarea cu calculatorul s-a realizat cu ajutorul portului paralel. Comunicația în sens direct (PC → FPGA) se desfășoară pe șapte biți iar cea inversă pe patru biți. Câte un bit din cele disponibile în registrele de date și stare ale portului paralel nu se pot utiliza ca funcție I/O generală din cauza unor restricții impuse de place de dezvoltare FPGA.

Am utilizat, deci, din nou, USID (Unitatea serială de intrări de date – un circuit inteligent de conversie a comunicării serial-paralele) puțin modificată pentru a comunica cu software-ul de comandă scris în C++ și rulat pe PC. Cel mai puțin semnificativ bit al registrului de date din PP a fost asignat ca linie de transmisie date seriale, bitul următor este ceasul de sincronizare a transmisiei seriale, restul biților fiind semnale de control. Cu ajutorul acestor biți s-a rezolvat atât inițializarea ponderilor sinapselor din rețea, cât și furnizarea modelelor de intrare acesteia. Cei patru biți ai registrului de stare al PP îndeplinesc două sarcini diferite în funcție de poziția unui comutator aflat pe placa de dezvoltare XSA utilizată.

Dacă acest comutator este pe valoare logică 0, atunci acești patru biți vor fi conectați la magistrale de ponderi a rețelei neuronale, astfel putându-se citi valorile tuturor ponderilor în orice moment al testării rețelei. În caz contrar, dacă comutatorul este pe poziție logică 1, atunci doi dintre cei patru biți reprezintă ieșirea axonală a celor doi neuroni pulsativi implementați în rețea. La transmiterea ponderilor un pachet trimis pe linia serială de date  $D_0$  este formată din patru + cinci biți corespunzători valorii

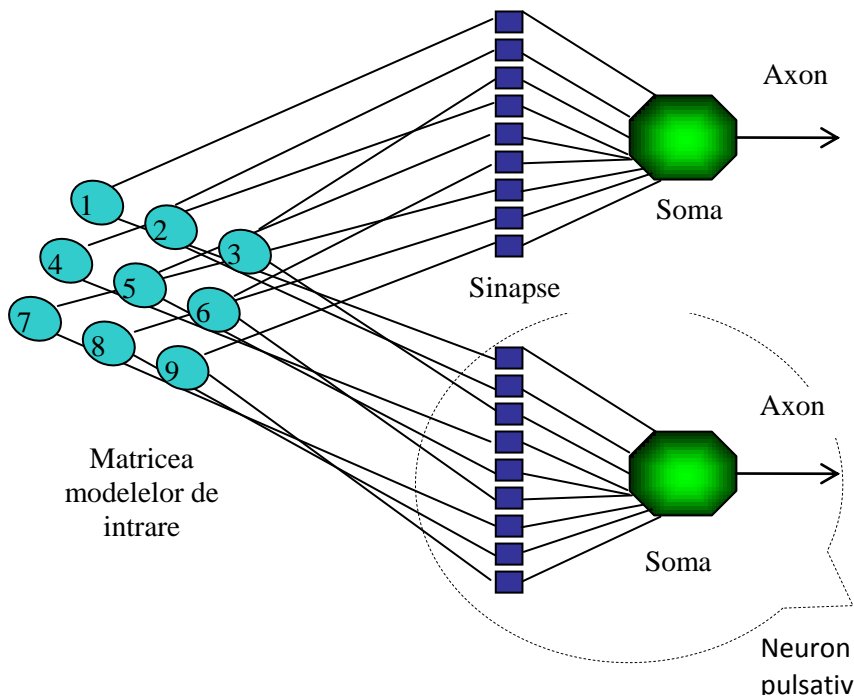


Figura 4.28 Arhitectura rețelei neuronale pulsative implementate

ponderii respectiv adresei sinapsei respective. Aceste pachete sunt decodificate de USID și cu ajutorul semnalelor de comandă CEN și WR se adresează apoi se înscriu valorile în registrele interne ale sinapselor prin magistrala de ponderi a rețelei.

Algoritm de desfășurare a testului rețelei neuronale este următorul:

1. Programul de control (C++) rulat pe PC generează 18 valori aleatoare între 0 și 15 ca valori de inițializare pentru ponderile celor 2x9 sinapse ale neuronilor rețelei. Apoi aceste valori sunt

transmise și înscrise la locul potrivit, timp în care rețeaua încă nu este funcțională, furnizarea semnalului de tact principal fiind sistată.

2. Programul de control calculează valorile binare pentru modelul de intrare care urmează. Fiecare punct întunecat din matricea de intrare corespunde unui bit de valoare logică 1 iar punctele inactive (albe) au valoare logică 0 (Figura 4.29). Acest șir de biți astfel calculat este transmis apoi rețelei prin intermediul PP respectiv USID. Acesta din urmă distribuie valorile decodificate și le plasează la sinapselor corespunzătoare ca impulsuri de intrare. Nici în această fază rețeaua nu funcționează, din cauza lipsei semnalului de tact principal.
3. În acest pas, startul funcționării este dat de primul front crescător al semnalului de tact pornit, care duce la generarea de un număr de impulsuri post-sinaptice egal cu valoarea inițială a ponderii respective de către sinapsele care au primit impuls de intrare. Acestea ajungând la sumă se modulează și se însumează, modificând potențialul de membrană al acesteia.

Pașii de la punctele 2 și 3 se repetă de mai multe ori, în funcție de setările efectuate asupra programului de control în experimentul respectiv. Repetarea modelelor de intrare se face în așa fel încât probabilitatea de apariție a ambelor modele să fie egală (De exemplu fiecare model se repetă la fiecare al doilea pas).

Rularea testului se întrerupe prin oprirea semnalului de tact, pentru a putea citi valorile ponderilor prin registrul de stare al PP, respectiv ieșirile axonale ale neuronilor (pentru a afla dacă au emis IA sau nu). Toate aceste date sunt salvate în fișiere pentru a putea fi prelucrate ulterior.

Testul se termină odată cu închiderea programului de control pe PC.



Figura 4.29 Perechile de modele de intrare, utilizate în cursul învățării (notați ordinea și valorile biților)

Schema de conectare a RNP implementate, împreună cu circuitele auxiliare (convertorul serial-parallel, USID, modulul de sincronizare, unitatea de control a magistralei de ponderi) este prezentată în ANEXA 3- Schema de implementare a unei rețele neuronale pulsative, a acestei teze.

#### 4.7.2. Descrierea algoritmului de învățare supervizată implementat

Perechile de intrare-ieșire utilizate la învățare au fost stocate în microcontroler. Rețeaua neuronală implementată poate funcționa în două regimuri de lucru și anume regim de învățare și regim de testare a reacției la intrări. Regimul de lucru este setat de semnalul *learn\_or\_test* iar selecția algoritmului de învățare este efectuată cu ajutorul semnalului *sel\_learn*. De fapt regulile algoritmilor de învățare sunt aceleași (regulile Hebb) numai interpretarea impulsurilor post-sinaptice și a celor axonale sunt interpretate în mod diferit. Algoritmul Hebb implementat se poate rezuma astfel:

- ✓ Valorile ponderilor sinaptice se măresc în cazul prezenței simultane a impulsurilor pre- și post-sinaptice.
- ✓ Dacă numai un impuls pre- sau numai unul post-sinaptic este prezent, atunci valorile ponderilor sinaptice trebuie să scadă. În cazul învățării nesupervizate, semnalele post-sinaptice proprii ale rețelei sunt conectate ca o reacție în modulul de învățare a sinapselor.

- ✓ În cazul în care nu este prezent nici impuls pre-sinaptic nici unul post-sinaptic, valoarea ponderii sinaptice respective rămâne neschimbată.

În cazul învățării supervizate în loc de ieșirile proprii ale rețelei sunt introduse forțat așa numite IA virtuale din exterior în modulele de învățare ale sinapselor. Aceste IA virtuale vor fi de fapt ieșirile prescrise rețelei pentru modelul de intrare prezentat.

După cum a fost prezentat și mai sus, intrarea este o matrice de 3x3 elemente. Scopul este de a antrena rețeaua în mod supervizat, în așa fel încât să recunoască, de exemplu, literele T și H prin activarea neuronului numărul 1 respectiv 2 la apariția secvențială a celor două litere la intrare.

Conectarea fizică a plăcii de dezvoltare FPGA XILINX XC2S100 cu mini-sistemul conținând microcontrolerul – ca supervisor – PIC16F876 s-a realizat montând ambele pe placa de extensie XST, parte auxiliară a pachetului de dezvoltare FPGA. Realizarea comunicării de date dintre cele două plăci (intrări, ieșiri prescrise, setarea modului de lucru) s-a efectuat conectând porturile B și C ale PIC16F876 la pini corespunzători ai plăcii XSA (Figura 4.30).

Programul în limbaj de asamblare implementat pe microcontroler realizează următorul algoritm. În memoria de date a supervisorului s-au stocat valorile intrărilor pentru literele T și H, respectiv ieșirile impuse pentru acestea. În mod asemănător, tot în memoria de date a microcontrolerului s-a stocat un set de valori de intrare de test, conținând combinațiile repetate ale literelor T și H, alterate sau nu de zgomot. Programul rulat de supervisor selectează modul de lucru al rețelei (învățare) și tipul acesteia (supervizată) cu ajutorul liniilor de control ale circuitului din FPGA, *learn\_or\_test*, respectiv *sel\_learn* (ANEXA 3- Schema de implementare a unei rețele neuronale pulsative). În ciclul de învățare, acest program prezintă rețelei neuronale din FPGA, în mod alternant, combinațiile pentru literele de învățat (T=111010010, H=101111101), respectiv împreună cu acestea, valorile de ieșire impuse fiecăruia (T=10, H=01). La următorul front crescător de tact se setează semnalul *spike enable* (SP) la 1 logic, semnalând rețelei, că impulsurile de intrare sunt pregătite, apoi urmează 15 cicluri de tact în care rețeaua prelucrează aceste intrări, neuronii emițând sau nu IA. Se procedează la fel pentru următorul caracter de prezentat.

### 4.7.3. Măsurători și rezultate experimentale

Experimentele inițiale au fost efectuate utilizând calculatorul prin PP pentru a furniza intrările circuitului din FPGA și pentru a citi date de la aceasta. Scopul primordial în aceste cazuri a fost de a observa evoluția rețelei în timpul procesului de învățare, lucru ușor de realizat dat fiind faptul, că comunicare prin portul PP a rezultat într-o funcționare relativ lentă – cu un tact de ~100 Hz – ușor de urmărit și eșantionat cu ajutorul analizorului logic Tektronix cu 36 de canale. Evident această cale nu a fost cea care să permită o funcționare în timp-real, ea îndeplinindu-și însă țelul de a verifica convergența algoritmului de învățare. Rezultatele fără echivoc pozitive ale acestor experimente au determinat trecerea la pasul următor și anume trecerea la un program de control mai rapid, prin utilizarea microcontrolerului PIC16F876. Acesta poate lucra la o frecvență maximă de 20 MHz. Astfel am reușit să prezint rețelei câte un nou model de intrare la fiecare aproximativ 1,5μs (în doi pași temporali, de 16 cicluri de ceas fiecare), aplicând algoritmul de învățare nesupervizat.

Figura 4.31 respectiv Figura 4.32 prezintă semnale eșantionate cu analizorul logic, care arată reacția celor doi neuroni, impulsurile axonale emise de aceștia, la cele două modele de intrare furnizate rețelei. Semnalul trasat cu culoare roșie, CKO – în Figura 4.31 – este semnalul de tact principal al sistemului. Se poate observa, că înainte de apariția ciclurilor de tact pe acest semnal, semnalele

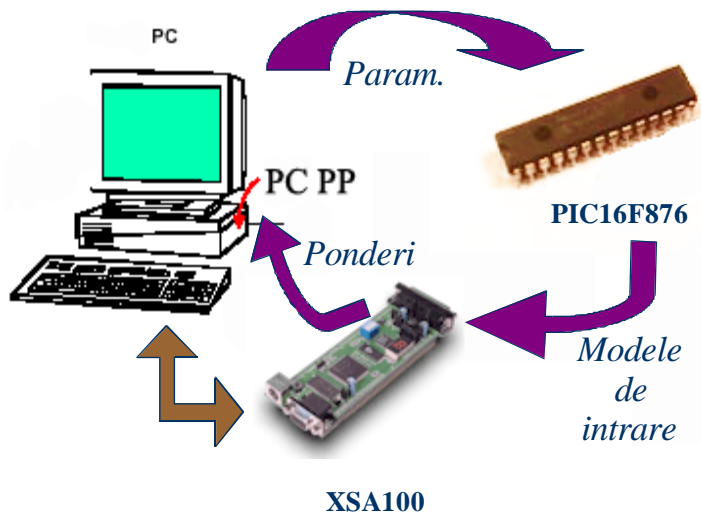


Figura 4.30 Comunicația dintre diferitele echipamente ale sistemului



*Serial\_Clock* și *Data\_In* desfășoară o transmisie serială a valorilor de intrare. Ca efect al acestei intrări, impulsurile post-sinaptice generate, determină schimbări rapide de valori ale potențialului de membrană (aceasta a fiind afișată numai pentru unul dintre neuroni, semnalul *membrpot*). Într-un moment dat potențialul de membrană depășește cel de prag, deoarece neuronii și-au acordat – în cursul antrenării – ponderile sinaptice pe câte una din modelele de intrare, astfel se emite un impuls de activare (*axon2(0)*). Urmează apoi transmiterea serială a următorului model de intrare la care răspunde în mod asemănător celălalt neuron.

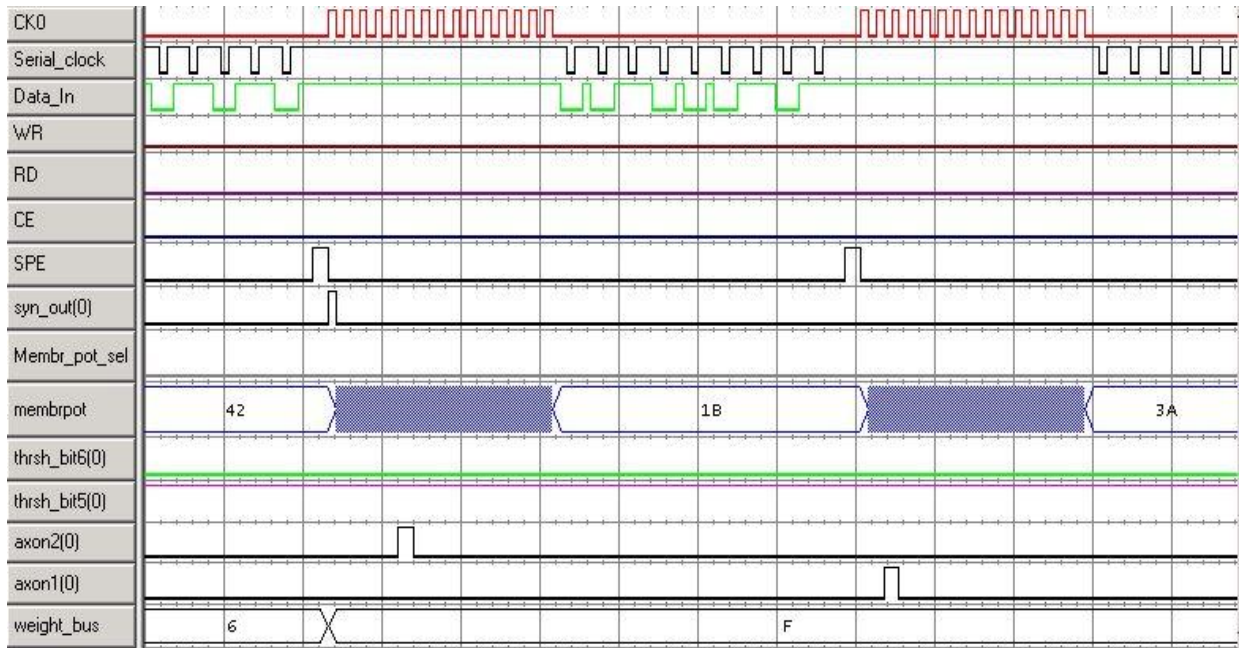


Figura 4.31 Transmisii seriale ale valorilor de intrare cu răspunsurile axonale date de neuroni la acestea

Figura 4.32 marchează un alt moment important al experimentului și anume procedeul de citire a valorilor ponderilor după două intrări transmise rețelei. Aceste valori apar pe magistrala de ponderi (*weight\_bus*) după ce au fost activate semnalele *RD* (semnalarea unei operații de citire) respectiv *CE* (activare sinapselor pentru scriere/citire) și pe linia serială s-a transmis adresa sinapsei citite în fiecare ciclu.

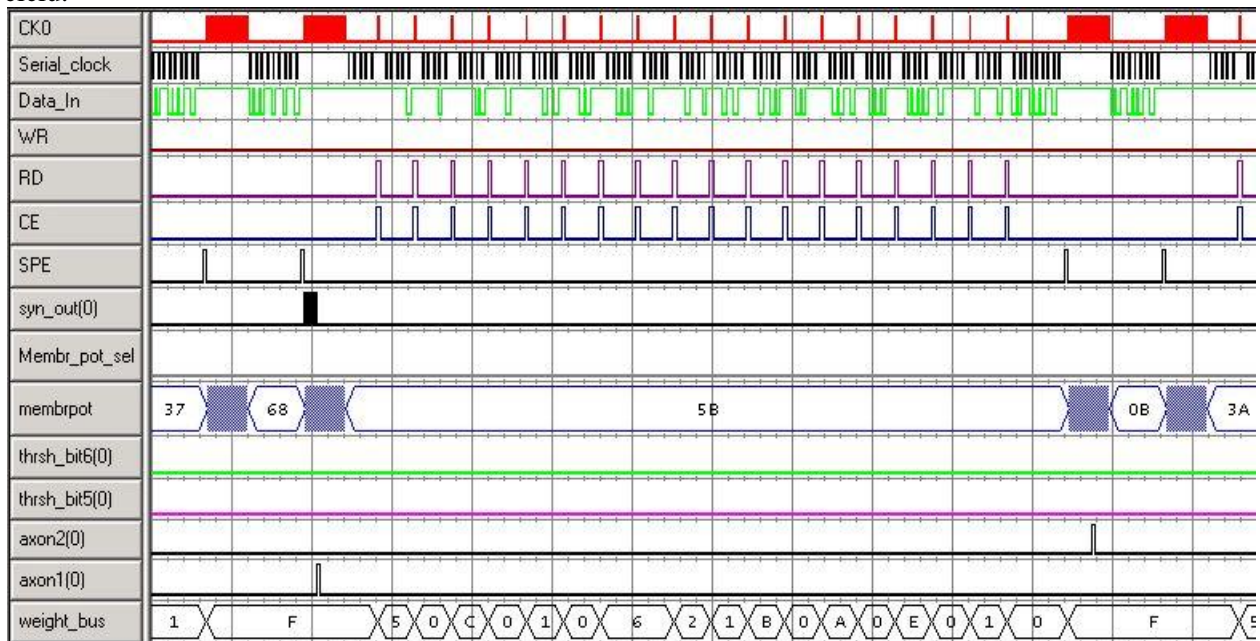


Figura 4.32 Citirea valorilor de pondere sinaptică (hexadecimal) între două faze de transmitere a perechii de modele de intrare



## 4.8. Analiza performanțelor rețelei neuronale neuromorfe implementate

### 4.8.1. Rezultatele obținute de rețeaua neuronală implementată în hardware

Unul din rezultatele majore ale rețelei neuronale neuromorfe pulsative (spiking neural network) implementate în circuitul FPGA este atingerea funcționării în timp real. De asemenea este de notat faptul că rețeaua este capabilă să rezolve foarte rapid și cu o eroare minimă, problema propusă, deși parametrii ei au o precizie relativ mică. Toate acestea sunt adevărate și în cazul unor modele de intrări care suferă distorsiuni, adică apar cu erori introduse de zgomote.

Aceste afirmații sunt susținute de graficele din Figura 4.33 și Figura 4.34 care prezintă variația ponderilor sinaptice ale celor doi neuroni de-a lungul procesului de învățare. Se poate observa, că aceste valori (axa Y) se stabilizează în jurul valorii maxime (15) respectiv minime (0) în funcție de locul de apariție a punctelor întunecate (valori logice 1) din modelul de intrare asimilat de neuronul respectiv. Acest proces se realizează în aproximativ 20 de pași (axa X), unde fiecare pas durează 16 cicluri de tact. Calculând la frecvența de lucru a microcontrolerului PIC utilizat, de 20 MHz, obținem  $16 \times 20 \times 100 \text{ ns} = 32 \mu\text{s}$ , ceea ce reprezintă timpul de învățare a rețelei pentru orice două modele de intrare.

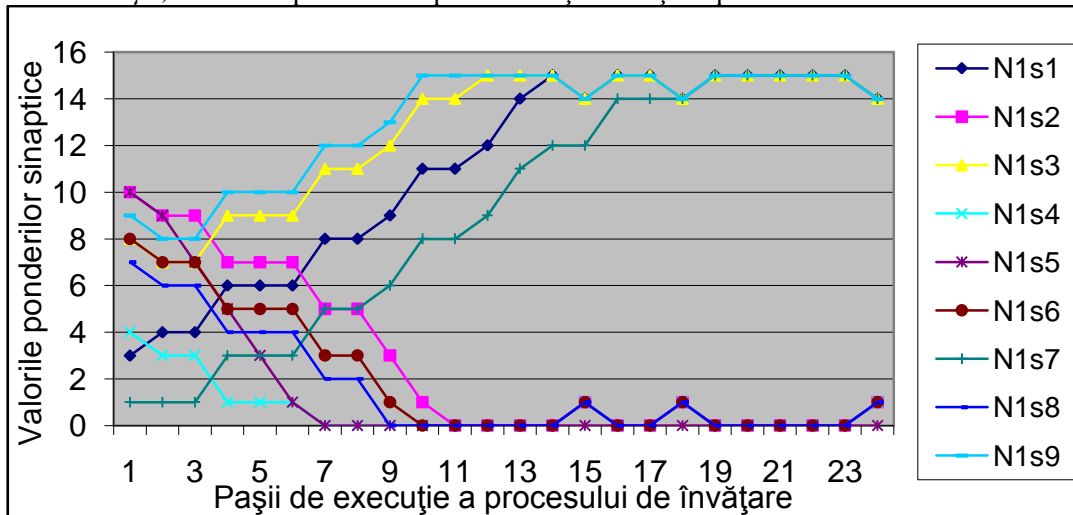


Figura 4.33 Variația ponderilor sinaptice a primului neuron în timpul antrenării

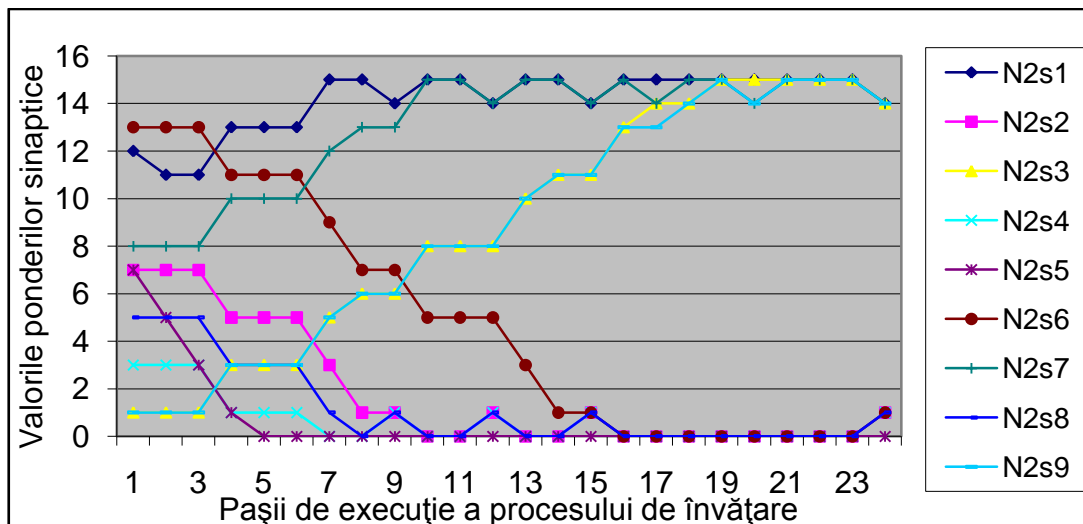


Figura 4.34 Variația ponderilor sinaptice a celui de-al doilea neuron în timpul antrenării

După cum am prezentat într-un capitol anterior, structura circuitului neuronal implementat, permite oprirea procesului de învățare, astfel rolul microcontrolerului dispare. Utilizând însă pentru generarea semnalului de tact, circuitul oscilator programabil aflat pe placa de dezvoltare FPGA XSA100

putem ridica frecvența de funcționare a rețelei la 100MHz. Astfel am ajuns la un timp de identificare a modelelor de intrare de către RNP de numai  $\sim 3,2\mu s$ . Evident această valoare reprezintă un rezultat important, care permite utilizarea de astfel de rețele neuronale în aplicații de control automat în timp real.

#### 4.8.2. Comparație cu o rețea neuronală clasică realizată în software

Cu scopul de a confrunța performanțele obținute prin procedeele noi, descrise în capitolele anterioare, am realizat o simulare de rețea neuronală feed-forward, cu ajutorul mediului Matlab, implementând aceeași structură ca a rețelei din FPGA. Astfel, și neuronii rețelei simulate aveau câte nouă intrări cu ponderile aferente. În acest caz, însă ponderile au fost reprezentate pe 16 biți în virgulă mobilă, cu o plajă de valori cuprinse între  $[-1, 1]$ .

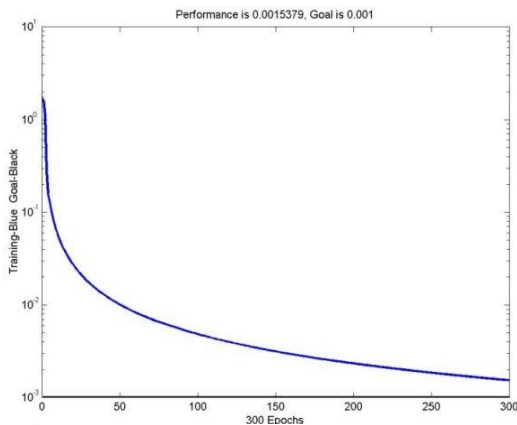


Figura 4.35 Curba de variație a erorii rețelei neuronale feed-forward, realizate cu scop de comparație în Matlab, în timpul procesului de antrenare

Setul de valori de intrare/ieșire ale algoritmului de învățare au fost alese în așa fel, încât să formeze un model pe o matrice de  $3 \times 3$  ca și cel implementat hardware. Ieșirile prescrise rețelei neuronale pentru cele două modele de intrare au fost setate la 1 respectiv  $-1$ . Rețeaua neuronală astfel creată a fost antrenată în 300 de cicluri ( $\sim 300ms$ ) cu algoritmul de învățare backpropagation de tip Levenberg-Marquardt. Rata de învățare a fost setată la 0.1, iar criteriul de oprire a fost atingerea unei erori de 0.001 față de valoarea prescrisă.

Figura 4.35 prezintă graficul variației erorii de-a lungul procesului de antrenare a rețelei neuronale simulate în Matlab.

Comparând această curbă cu cele din Figura 4.37 respectiv Figura 4.36 – care prezintă eroarea medie pătratică a primului și respectiv al celui de-al doilea neuron pulsativ al rețelei implementate în FPGA – putem

nota performanța net superioară a variantei hardware.

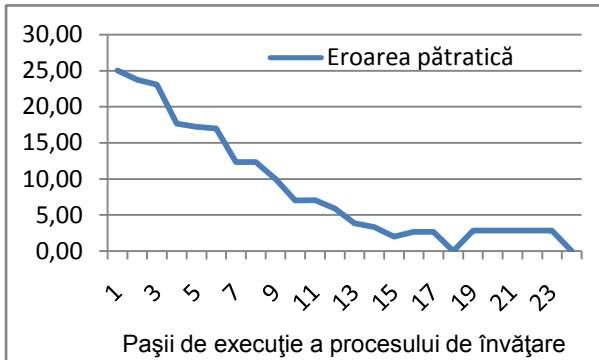


Figura 4.37 Eroarea medie pătratică a primului neuron al RNP implementate hardware

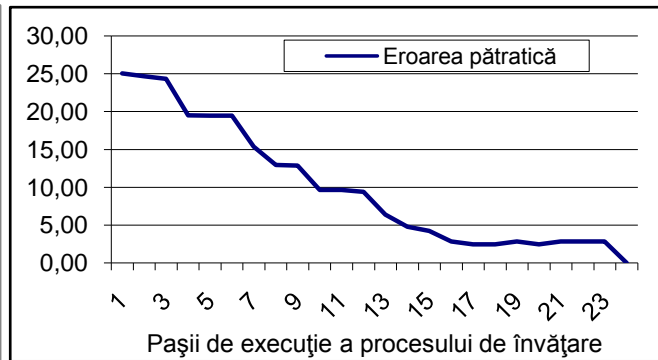


Figura 4.36 Eroarea medie pătratică a celui de-al doilea neuron al RNP implementate hardware

## 4.9. Concluzii

Acest al patrulea capitol al tezei de doctorat începe prin a prezenta caracteristici generale ale sistemelor hardware analogice și digitale în ceea ce privește implementarea rețelelor neuronale artificiale, aducând motivații pro și contra în fiecare caz. De asemenea sunt prezentate cele mai importante implementări hardware de rețele neuronale din literatura de specialitate actuală cu accent asupra realizărilor pe diferite platforme digitale. Se trece apoi la domeniul specific temei acestei lucrări, prezentând implementări software și hardware ale RNP realizate pe calculatoare paralele sau chiar în circuite FPGA.

După această parte de sinteză se prezintă contribuțiile proprii în ceea ce privește implementarea complet paralelă de RNP în circuite FPGA. Inițial se expune implementarea modelului teoretic adaptat realizării hardware, validat prin simulare software în capitolele anterioare. Aceste subcapitole conțin descrieri amănunțite ale proiectării diferitelor subsansamble ale acestor rețele neuronale, cum ar fi sinapsele și modulele de somă sau corp celular. Acestea sunt testate prin simularea funcționării circuitelor proiectate pentru a implementa funcționalitatea specifică acestora.

O aplicație de test a acestor circuite neuronale realizate reprezintă continuarea logică a acestui capitol, aplicația fiind cea clasică de diferențiere a două forme simple cum ar fi literele T și H sau semnele + și X. Se descrie rețeaua neuronală care implementează în circuit FPGA prin resurse hardware dedicate fiecărui subsansamblu al acestuia (implementare complet paralelă), specificând algoritmul de învățare utilizat și enumerând rezultatele măsurătorilor experimentale. În cadrul unei analize a performanțelor ( $32\mu s$  - timpul de învățare a rețelei,  $32\mu s$  timp de recunoaștere a formelor de intrare) acestei RNP se prezintă și o comparație cu o rezolvare pur software (în mediul Matlab) a problemei propuse, rezultând performanțe net superioare ale versiunii hardware, care funcționează în timp real.

Rezultatele și contribuțiile personale prezentate în acest capitol au fost publicate într-un număr de lucrări științifice ale autorului tezei, apărute în publicații ale unor conferințe de specialitate (Bakó & Brassai, Hardware spiking neural networks: parallel implementations using FPGAs, 2006), (Bakó & Brassai, 2005), (Brassai & Bakó, 2007), în reviste de specialitate (Bakó & Brassai, Spiking neural networks built into FPGAs: Fully parallel implementations, 2006) (Bakó, Székely, & Brassai, Development of Advanced Neural Models. Software And Hardware Implementation, 2004), precum și în rapoartele unor granturi de cercetare în care autorul a fost membru.



## 5. IMPLEMENTĂRI PARȚIAL PARALELE ALE RNP UTILIZÂND MICROCONTROLERE ÎNCORPORATE ÎN CIRCUITE FPGA

### 5.1. Prezentarea rezultatelor actuale ale domeniului din literatura de specialitate

O propunere interesantă referitoare la modul de utilizare a rețelelor de neuroni pulsanti a fost făcută în lucrarea *Liquid State Machine*, de către Maass, Natschlagher și Markram (Maass, Natschlagher, & Markram, 2002). Cercetătorii din grupul condus de W. Maass au înțeles că o rețea de neuroni pulsanti conectată aleatoriu poate implementa în mod eficace un filtru temporal complex prin intermediul elementelor complexe ale activității cu reverberație și ale dinamicii sinaptice. Ținând seama de informații primare extinse în timp, precum vorbirea, activitatea colectivă a rețelei poate fi descrisă ca o traiectorie printr-un spațiu de stări multi-dimensional, iar această traiectorie ar trebui să fie recognoscibilă caracteristică pentru intrarea la dispoziție. Un simplu decodor „read-out” ar trebui să fie suficient pentru a clasifica tiparul (modelul) temporal.

În lucrarea „Isolated Word Recognition With The *Liquid State Machine*: A Case Study”, (Verstraeten, Schrauwen, Stroobandt, & Van Campenhout, 2005) autorii studiază recunoașterea limbajului în Liquid State Machine (LSM). Pentru un set lexical standard, se testează o varietate de codificări temporale care sunt introduse într-o LSM și un următor decodor liniar simplu clasifică fiecare tipar în parte. Surprinzător, Verstraeten et al. au găsit doar performanțe (încurajator de bune) pentru LSM atunci când back-end-ul codificator aproximează schema de codificare a urechii interne. Reguli de învățare echivalente celor folosite în rețelele neuronale sigmoide tradiționale au provenit de asemenea din rețele (stratificate) de neuroni pulsanti, de exemplu o regulă error-backpropagation bazată pe un gradient precum Spikeprop (Bohte, Kok, & La Poutré, 2000). Cu o astfel de regulă de învățare, s-a demonstrat că o rețea neuronală pulsantă poate în fapt calcula funcții non-liniare separabile precum rețelele neuronale tradiționale.

Booij și Nguyen prezintă regula de învățare bazată pe propagarea înapoi a erorii (Booij & Nguyen, 2005), care elimină importante limitări ale regulii Spikeprop în lucrarea intitulată „A gradient descent rule for spiking neurons emitting multiple spikes”. Înainte de toate, în derivația lor gradientul erorii este calculat pentru un neuron emițând multiple impulsuri. În plus, simple metode euristice sunt discutate care diminuează efectul disruptiv al discontinuităților în potențialul de membrană pentru algoritmi de minimizare a gradientului. Eficacitatea algoritmului îmbunătățit este demonstrată pe o versiune temporală a clasicei probleme XOR și pe clasificarea transmisiilor (prin angrenaje) de impulsuri Poisson. În mod interesant, pentru cea mai simplă problemă XOR Booij și Nguyen ne arată că există o soluție „hair trigger” într-o rețea fără straturi ascunse.

Un alt domeniu unde neuronii pulsanti sunt considerați a fi deosebit de puternici este cel al *memoriei asociative*. Knoblauch studiază tehnologia de ultimă oră în lucrarea „Neural Associative Memory for Brain Modeling and Information Retrieval” (Knoblauch, 2005). Începând cu o expunere asupra rețelelor neuronale asociative tradiționale sub forma clasicului model Willshaw, soluții folosind neuroni pulsanti sunt sugerate și este discutat faptul dacă memoriile asociative neuronale distribuite au avantaje practice față de memoria locală/limitată.

Multe aplicații folosesc rețele neuronale tradiționale ca aproximatori de funcții când funcția necesară este necunoscută. Având un set de date, rețelele neuronale s-au dovedit a fi foarte bune în a conecta acestea, interpolând valorile dintre ele, apoi aproximând funcția.

Așa cum Iannella și Kindermann subliniază în articolul lor „Finding Iterative Roots with a Spiking Neural Network”, o sarcină mai complexă dar adesea și mai profitabilă va fi cea care va găsi componentele (recurente) care alcătuiesc funcțiile ce urmează a fi approximate. În lucrarea lor (Iannella & Kindermann, 2005), ei demonstrează că progresul stă în învățarea unor astfel de soluții iterative într-o rețea de neuroni pulsanti. Doi algoritmi pentru învățare sunt prezentați, unul semi-supravegheat în care funcția este cunoscută și sarcina fiind aceea de a găsi valorile funcției, iar celălalt ne-având această extra-informație.

O problemă importantă în utilizarea rețelelor neuronale pulsative este faptul că ele sunt în mod tipic computațional mai intensive decât rețelele neuronale tradiționale. A fost deja stabilit că evenimentul bazat pe natura impulsurilor cronometrate reduce în mod drastic randamentul de comunicare între neuroni, permițând în principiu implementări paralele eficiente.

Aceasta încă lasă deschis modul de a calcula eficient în ce moment un neuron se activează având în vedere impulsurile sale de intrare. În „Spiking Neural Nets with Symbolic internal State”, O’Dwyer și Richardson prezintă un algoritm pentru a calcula eficient această activare bazată pe manipularea simbolică și intervalul aritmetic (O’Dwyer & Richardson, 2005).

Sunt de părere că aceste lucrări arată că se face un progres substanțial spre aplicarea rețelelor neuronale pulsative: studiul lui Verstraeten et al. demonstrează că rețelele de neuroni pulsați ar putea avea un real viitor în recunoașterea vorbirii; Booij și Nguyen obțin niște reguli de învățare mai eficiente; Knoblauch tratează atât capacitățile cât și limitările rețelelor neuronale pulsative pentru memoria asociativă; munca lui Ianella și Kindermann are ca scop găsirea soluțiilor funcționale, iar O’Dwyer și Richardson prezintă un algoritm care ajută la calcularea eficientă a activității neîntrerupte într-o rețea de neuroni pulsați.

Pe lângă acest progres - deja foarte palpabil - care duce spre rețele neuronale pulsative practice, trebuie remarcat faptul că diferite lucrări au fost prezentate recent care fac ca rețelele de neuroni pulsați să fie privite cu un și mai mare interes din punct de vedere practic: lucrări de Deneve (Deneve, 2005), Rao (Rao, 2005) și Zemel et al. (Zemel, Huys, Natarajan, & Dayan, 2005) propun diferite moduri de a realiza inferența Bayesiană în rețele neuronale pulsative. Aceste lucrări dezvăluie o perspectivă reală și ar putea fi baza a mai multor aplicații. Desigur, aplicațiile rețelelor de neuroni pulsative prezentate aici reprezintă doar începutul celor ce vor urma.

## **5.2. Clasificarea nesupervizată de mulțimi complexe (complex clusters) cu rețele de neuroni pulsativi**

Întrucât neuronii produc în creier potențiali de acțiune „tot sau nimic”, importanța coordonării acestor impulsuri a fost recent recunoscută ca o mijloc de codare a informației neuronale. Datorită acumulării de probe biologice, s-a arătat teoretic că o astfel de codificare temporală permite procesarea puternică de informații de către neuroni și că ar putea fi esențială în rezolvarea problemelor dinamice combinatorice. Aceste considerente au generat un interes considerabil în domeniul rețelelor neuronale artificiale secvențiale. Hopfield prezintă un model de neuroni pulsați folosiți pentru a descoperi fascicule într-un spațiu de intrare asemănător rețelelor bazate pe funcțiile radiale. Natschläger și Ruf prezintă (Natschläger & Ruf, 1998) un algoritm de învățare pentru rețele neuronale pulsative care efectuează clasificări nesupervizate de fascicule folosind temporizări de impulsuri ca intrare. Acest model codifică șabloanele de intrare în decalaje temporale în propagarea impulsurilor de-a lungul sinapselor și s-a demonstrat că este fiabil în a găsi centre ale zonelor multi-dimensionale, dar este limitat în capacitatea și precizia de clasificare.

Inspirat fiind de domeniile receptive locale ale neuronilor biologici, am codificat variabile de intrare de către neuroni cu profiluri de sensibilitate gradate și suprapuse. Mai mult decât atât, fiecare dimensiune de intrare a unui set de date de dimensiune mare este codificat independent, evitându-se o creștere exponențială a numărului neuronilor de intrare. Aceasta s-a arătat că validează în mod eficient rețeaua pentru a clasifica cu succes un număr de probleme de clasificare de fascicule la costuri relativ joase în ceea ce privește neuronii. Codificarea multiplă are în vedere detectarea fiabilă a fasciculelor asupra unei game remarcabile și flexibile de scale spațiale. Acest aspect este în mod special important în sarcinile de clasificare necontrolată, unde informația de scalare este a-priori necunoscută.

Extinzând rețeaua la straturi multiple, putem demonstra cum anume aspectul secvențial al neuronilor pulsați poate fi exploatat pentru a valida clasificarea corectă a fasciculelor non-sferice sau suprapuse. Într-o rețea RBF multistrat, putem demonstra că neuronii din primul strat sunt specializați pe componente de fascicule. Sincronismul codificării neuronilor pentru componente apropiate este apoi arătat a fi ușor remarcabil de către un strat RBF, rezultând într-o formă de clasificare ierarhică. Adăugarea unor conexiuni excitatoare laterale se spune că validează o astfel de rețea multistrat pentru a clasifica în mod corect structuri integrate complexe, sincronizând codificarea neuronilor pentru elemente ale aceluiași fascicul. Astfel, conexiunile laterale adiționale produc un număr scăzut de neuroni, în același timp măbind complexitatea fasciculelor clasificabile și folosind pentru aceasta doar o regulă de învățare Hebbiană locală.

### 5.3. Implementarea unei rețele de neuroni pulsativi cu întârziere

#### 5.3.1. Structura rețelei

Structura rețelei este formată din neuroni pulsanți conectați în buclă deschisă (feedforward) de tipul integrate-and-fire permeabil (leaky integrate-and-fire type). Cum este ilustrat în Figura 5.1A, conexiunile din rețea constă dintr-un set de sinapse  $k$ , fiecare cu o greutate  $w_{ij}^k$  și o întârziere  $d^k$ . Un impuls de intrare generează un set de potențiale post-sinaptice ponderate (PSP), o funcție modelând impactul unui impuls asupra neuronului țintă ca o funcție temporală). Mărimea ponderii sinaptice determină înălțimea PSP-ului. La fiecare pas temporal al algoritmului, valorile acestor PSP-uri sunt adunate la neuronul țintă, iar când suma depășește pragul  $\theta$ , un impuls de ieșire este generat.

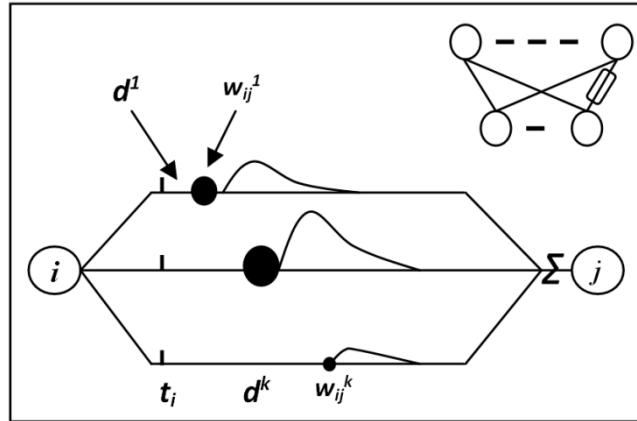


Figura 5.1 Structura de bază a rețelei implementate

Semnalele de intrare pot fi codificate în ponderile sinaptice prin **învățare întârziată**. După învățare, momentul de activare al unui neuron de ieșire reflectă distanța de la mostra de evaluat până la mostra de intrare învățată (realizând astfel un fel de neuroni RBF). Pentru învățarea nesupervizată, ponderile dintre neuronii sursă și cei de ieșire din stratul țintă (Winner-Take-All) sunt modificate folosind o variantă temporală a regulii de învățare Hebbiană. Dacă startul unui PSP la o sinapsă precedă cu puțin un potențial de acțiune în neuronul țintă, ponderea sinapsei crește, deoarece a exercitat o influență semnificativă asupra momentului activării printr-o contribuție relativ mare către potențialul de membrană. Sinapsele anterioare și ulterioare scad în pondere, reflectând impactul lor scăzut asupra momentului de emisie al neuronului țintă. Pentru o pondere cu o întârziere  $d_k$  de la neuronul  $i$  la neuronul  $j$  vom folosi:

$$\Delta w_{ij}^k = \eta L(\Delta t) = \eta(1-b) \exp\left[-\frac{(\Delta t - c)^2}{\beta^2}\right] + b \quad \text{Ec. 65}$$

după (Natschläger & Ruf, 1998) (ilustrat în Figura 5.1B). Ponderile sunt limitate de o valoare minimă și una maximă. O intrare în rețea este codificată de un tipar de timpi de activare în cadrul unui interval de codificare și fiecare neuron de intrare trebuie să se activeze o dată în timpul acestui interval de codificare.

#### 5.3.2. Variabile de intrare continue codificate în decalaje temporale

Pentru a extinde precizia de codificare și capacitatea de clasificare, am investigat folosirea multiplelor domenii receptive locale pentru a distribui o variabilă de intrare asupra mai multor neuroni de intrare. Un astfel de cod de distribuție unde variabilele de intrare sunt codificate cu funcții de activare gradate și integrate, reprezintă o mult studiată metodă pentru a reprezenta parametri cu valori reale (Snippe, 1996). În aceste studii, funcția de activare al unui neuron de intrare este modelată ca un domeniu receptiv local care determină rata de activare. O traducere a acestei paradigme în timpi de activare relativi este simplă: un neuron stimulat la un nivel optim se aprinde la  $t = 0$ , în timp ce o valoare de până la  $t = 9$  este atribuită unor neuroni stimulați la un nivel mai puțin optim (ilustrat în Figura 5.2).

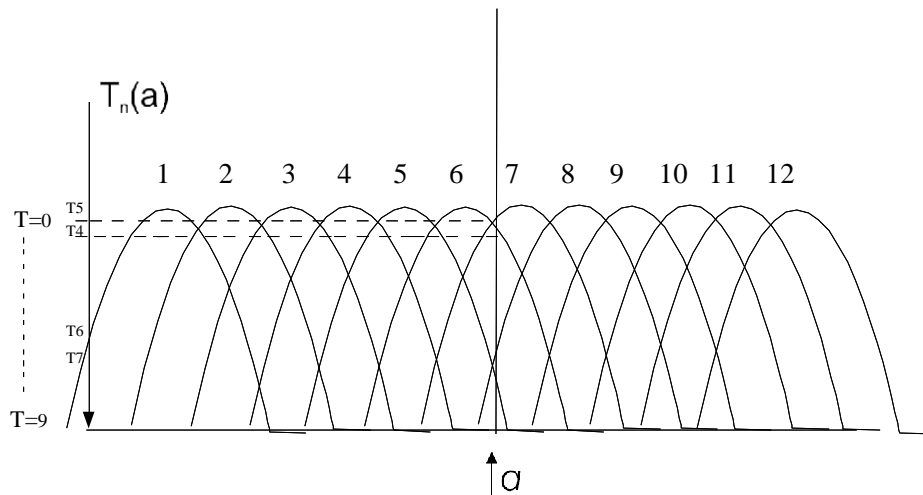


Figura 5.2 Conceptul de codificare a variabilelor de intrare în impulsuri decalate temporale

Pentru a codifica date multidimensionale în maniera descrisă mai sus, o alegere trebuie făcută în ceea ce privește dimensionalitatea domeniilor receptive ale neuronilor. Codificarea cea mai puțin costisitoare din punctul de vedere al neuronilor este aceea în care fiecare intrare este codificată independent cu domenii receptive 1-D. Din moment ce suntem interesați de clasificare multidimensională, această codificare este foarte oportună deoarece se comportă liniar în ceea ce privește numărul de neuroni necesari pe dimensiune și este de asemenea adaptabilă, permițând codificarea unor dimensiuni cu precizie mare fără costuri neuronale excesive.

### 5.3.2.1. Implementarea codării valorilor de intrare în decalaje temporale pe FPGA – Neuronii de intrare, de tip pseudo-RBF

Am experimentat mai multe variante de astfel de codări, luând în calcul resursele disponibile în circuitul FPGA Xilinx Spartan 3 XC3S1000 utilizat. Pentru a utiliza cât mai eficient blocurile logice reconfigurabile din acest circuit, respectiv, pentru a păstra cât mai multe astfel de resurse implementării neuronilor pulsativi, am încercat mai multe variante de stocare a valorilor funcțiilor RBF ale intrărilor în module BlockRAM. Acestea sunt disponibile în mai multe variante de configurație, cu diferite lățimi ale magistralelor de adrese și de date. După cum se vede și din tabelul de mai jos (Figura 5.3), totalul de 54Kbytes (432Kbiți în tabel) de memorie RAM disponibilă este împărțită în 24 de blocuri (Figura 5.4). Dilema proiectantului este în a decide cum să utilizeze aceste resurse, sacrificând – chiar dacă parțial – natura pur paralelă a implementării, sau utilizând mai multe BlockRAM-uri decât ar fi necesare din punct de vedere al cantității de informație ce necesită stocare.

După cum se poate vedea și în Figura 5.5 (generată în Matlab), funcțiile domeniilor receptive locale nu sunt Gaussiene ci triunghiulare, pentru a fi ușor implementate în hardware. Pe axa x a acestei diagrame am reprezentat valorile intrării rețelei, care sunt valori pozitive, întregi și pot varia între 0 și 192.

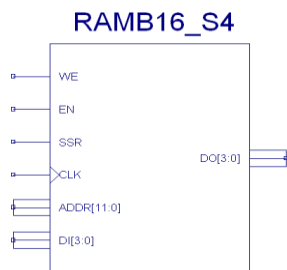


Figura 5.4 Modul BRAM, 1024 cuvinte de 4 biți

Fiecare valoare a acestui interval este distribuită pe mai mulți neuroni de intrare (min. 2 – max. 4) datorită suprapunerii funcțiilor de activare ale acestora. Cu cât este mai mare valoarea unei funcții de activare ale unuia dintre cei 12 neuroni de intrare implementați în acest proiect, cu atât mai repede se va activa acest neuron,

adică decalajul temporal al apariției impulsului de activare față de momentul apariției valorii de intrare va fi cu atât mai mică. Evident, există multe posibilități de a varia aria de suprapunere a triunghiurilor, modificând astfel numărul neuronilor de intrare care se vor activa pentru fiecare valoare de intrare, adică a celor care vor coda această valoare. Totodată, pentru fiecare neuron de intrare există o gamă redusă de valori de intrare pentru care se activează cu decalaj temporal corespunzător. Aceste valori pot fi

Part Number	XC3S50	XC3S200	XC3S400	XC3S1000
System Gates	50K	200K	400K	1000K
Logic Cells	1,728	4,320	8,064	17,280
Dedicated Multipliers	4	12	16	24
Block RAM Blocks	4	12	16	24
Block RAM Bits	72K	216K	288K	432K
Distributed RAM Bits	12K	30K	56K	120K
DCMs	2	4	4	4

Figura 5.3 Parametrii principali al circuitelor Xilinx FPGA din familia Spartan 3



calculate și în timpul funcționării circuitului neuronal, dar acest lucru ar consuma timp și resurse prețioase. Pentru a evita aceste impedimente, am decis calcularea acestor valori în prealabil și stocarea lor în memoriile RAM disponibile în FPGA.

Utilizarea optimă a memoriilor BRAM ar fi fost atinsă, dacă s-ar fi stocat valori ale mai multor neuroni de intrare în același modul de memorie. Astfel însă ar fi fost imposibilă accesare simultană a acestora. Din acest motiv, am ales să utilizez mai multe memorii BRAM, din tipul **RAMB\_16\_S4**, a cărui capacitate este de 1024 cuvinte de 4 biți și se comportă ca o memorie statică, sincronă (Figura 5.4).

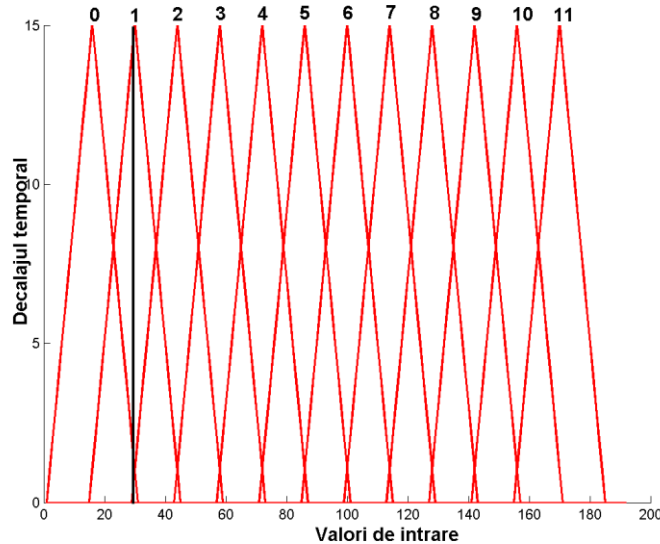


Figura 5.5 Funcțiile triunghiulare ale domeniilor receptive utilizate la implementare

Aceste valori au fost calculate cu ajutorul unui program scris în Matlab (vezi Anexa). Programul permite modificarea ușoară a intervalelor suprapuse și a numărului de neuroni de intrare utilizați. De asemenea, programul Matlab generează codul VHDL necesar inițializării memoriilor BRAM, cu valorile calculate, după cum arată Figura 5.6.

```

INIT_00 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFEDCBA9876543210123456789ABCDEF";
INIT_01 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF";
INIT_02 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF";
INIT_03 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFEDCBA9876543210123456789ABCDEF";
INIT_04 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF";
INIT_05 : bit_vector := x"FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF";
INIT_06 : bit_vector := x"FFFFFFFFEDCBA9876543210123456789ABCDEF";

```

Figura 5.6 Valorile funcțiilor triunghiulare ale domeniilor receptive cu care s-au inițializat memoriile BRAM

## 5.4. Realizarea unei rețele de neuroni pulsativi cu antrenare on-chip și aplicarea sa practică

Inspirat din fenomenele biologice, un nou tip de model de rețea neuronală, rețeaua de neuroni pulsanți, a atras atenția cercetătorilor din domeniul calculului neuronal în ultimii ani. Rețeaua implică neuroni cu funcționare prin impulsuri ca și unități de calcul. Folosește cronometrarea impulsurilor individuale (folosiți în literatura de specialitate sub denumirea de potențial total-sau-nul de acțiune sau spike) sau frecvența instantanee a șirului de impulsuri sau alte caracteristici temporale ale impulsurilor, produse de neuroni pentru codificarea informației.

Această nouă generație a modelelor neuronale, spre deosebire de modelele neuronale convenționale (cu funcții sigmoidale, etc.), poate să ofere explicații în ceea ce privește fenomenul procesării informației în creier printr-o interpretare biologică mult mai realistă, și pe lângă aceasta, pot să atingă performanțe computaționale competitive. Pe lângă acestea, aceste modele neuronale au avantaje ale implementării hardware, datorită faptului că codificarea bazată pe cronometrare oferă o

potențială compactitate și este mai robustă decât implementarea analogică. În cadrul modelelor neurale pulsative, modelul lui Hopfield este cel proiectat ca și un set de detectori de coincidențe pentru impulsuri de intrare pe canale multiple.

Poate fi folosit pentru a rezolva recunoașterea analogică de tipare în mod natural. Hopfield a indicat faptul că aceste modele neuronale funcționează de fapt ca unități de funcții de bază radiale (RBF). Această idee de bază a fost apoi extinsă asupra mai multor scheme sofisticate. Natschläger și Ruf au propus o rețea RBF folosind neuroni pulsanti (Natschläger & Ruf, 1998), care codifică tiparele în întârzierile sinaptice. Fiecare întârziere este tratată ca și o întârziere medie a populației sub-sinaptice, și este ajustată indirect prin ajustarea eficacității populației sub-sinaptice. Momentul de activare al neuronului pulsativ are o dependență regulată de distanța euclidiană dintre vectorul de intrare și vectorul central reprezentat prin întârzierile sinaptice.

Bothe et al. (Bohte, Kok, & La Poutre, 2000) combină arhitectura rețelei neurale pulsative (RNP) cu o metodă de codificare temporală care folosește câmpuri receptivă Gaussiene suprapuse pentru a traduce valorile de intrare analogice în decalaje temporale ale impulsurilor de activare.

Am proiectat un algoritm de antrenare Hebbian modificat pentru această RNP pentru a implementa o aplicație de clustering (clasificare) nesupervizată. În ambele lucrări (Natschläger & Ruf, 1998) și (Bohte, Kok, & La Poutre, 2000), întârzierea sinaptică este tratată ca și o întârziere medie a unei populații sub-sinaptice. În timp ce această arhitectură reflectă realitățile biologice, ea crește capacitatea de calcul a modelului și se delimitează de la algoritmul de antrenare nesupervizată mult mai flexibil.

În această lucrare, am propus realizarea unei rețele neuronale pulsative bazate pe ideea originală a lui Hopfield. Acest model este construit astfel încât timpii de activare ai impulsurilor de ieșire depind de sincronizare și de timpul incidental al sincronizării impulsurilor de intrare. Sincronizarea în rotație se poate optimiza prin coordonarea relației dintre modelul impulsurilor de intrare și modelul întârzierii sinaptice. Pentru a înlesni calculul și realizarea circuitului, vom abstractiza întârzierea sinaptică ca și o variabilă direct ajustabilă în acest model neuronal.

Cele mai multe aplicații ale rețelelor neuronale pulsative și convenționale se implementează prin folosirea unui simulator soft. Pe lângă acestea, calculele în cazul SNN menționat anterior sunt destul de complexe pentru o implementare hardware.

Această constituie o limită a uneia dintre cele mai importante avantaje ale rețelelor neuronale, paralelismul masiv. Fiind și scopul principal al disertației de doctorat, am proiectat circuite hardware pentru rețeaua neuronală propusă folosind circuite Field Programmable Gate Array (FPGA).

Cum circuitul FPGA este un sistem digital, avantajele tehnicii digitale, cum ar fi robustețea la zgomot și flexibilitatea designului, sunt eficiente în mod natural. Dar o tehnică digitală necesită de obicei o suprafață mare atunci când se folosește în cazul circuitelor neuronale paralele. Gătuirea se recunoaște de obicei în apariția unui multiplicator în cazul cadrului calculelor necesare. S-au raportat de asemenea mai multe scheme ale rețelelor neuronale fără multiplicatori. Marchesi et al. (Marchesi, Orlandi, Piazza, & Uncini, 1993) folosește registre mobile pentru implementarea multiplicatorilor de aproximație într-o rețea neuronală de valori reale, precum și în publicațiile proprii (Bakó & Brassai, 2006) (Bakó, Brassai, & Székely, 2006) (Bakó, Székely, & Brassai, 2004) (Brassai & Bakó, 2007) în care am implementat rețele neuronale pulsative bazate pe modele cu prelucrarea ratei impulsurilor realizate fără circuite multiplicatoare.

Totuși, circuitele hardware sunt disponibile doar în cazul calculului forward. Rețeaua antrenată este dirijată folosind un simulator offline. Hikawa (Hikawa, 1999) a propus o schemă hardware bazată pe frecvențe cu antrenare on-chip. Ambele scheme (Marchesi, Orlandi, Piazza, & Uncini, 1993) și (Hikawa, 1999) sunt rețele back-propagation (BP), ai căror algoritmi de antrenare BP ar putea să ajungă într-un minim local. Spre deosebire de acestea, modelul neuronal pulsativ propus în această lucrare evită această problemă, rezultând din esența sa de rețea RBF. Folosind natura de sincronizare a impulsurilor temporale, atât calculele forward cât și cele de antrenare sunt efectuate cu circuite fără multiplicatori.

### 5.4.1. Rețeaua de neuroni pulsativi realizată în FPGA

Arhitectura rețelei constă dintr-o rețea feed-forward de neuroni pulsanți cu terminale sinaptice întârziate multiple Figura 18. Neuronii rețelei generează potențiali de acțiune, sau impulsuri, când variabila de stare internă a neuronului, numit „potențialul membranei”, depășește un prag  $\gamma$ . Relația dintre impulsurile de intrare și variabile stării interne este descrisă prin modelul de răspuns al impulsului (spike response model SRM), introdus de Gerstner (Gerstner W. , Spiking neurons, 2001). În funcție de alegerea funcțiilor de răspuns ale impulsurilor potrivite, se poate adapta acest model pentru a reflecta dinamica unei mari varietăți de neuroni pulsanți.

#### 5.4.1.1. Modelul neuronal RBF pulsativ implementat

Modelul rețelei neuronale pulsative se poate vedea în Figura 5.7. Fiecare neuron, care funcționează ca și o unitate RBF, primește impulsuri de intrare prin conexiuni de la neuronii de intrare. Segmentul intermitent între fiecare neuron de intrare la terminalul sinaptic corespunzător denotă întârzierea, ceea ce este definită ca și diferența dintre momentul de activare al neuronului pre-sinaptic și momentul activării neuronului post-sinaptic.

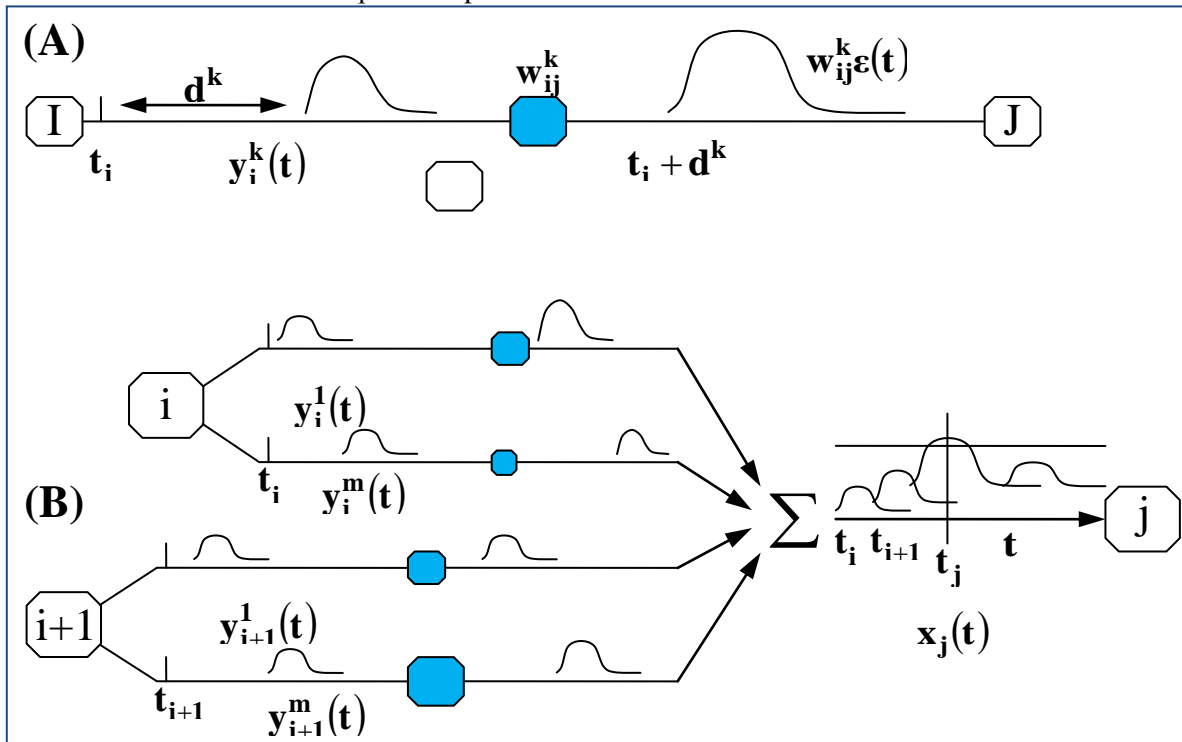


Figura 5.7 Principiul de funcționare al modelul neuronal implementat

Din punct de vedere formal, un neuron  $j$ , având un set  $\Gamma_j$  de predecesori imediați (neuroni pre-sinaptici), primește un set de impulsuri cu timpul de activare  $t_i$ ,  $i \in \Gamma_j$ . Fiecare neuron generează cel mult un impuls în timpul intervalului de simulare, și emite când variabila de stare internă atinge o valoare de prag  $\gamma$ . Dinamica variabilei de stare internă  $x_j(t)$  este determinată de impulsurile incidente, ale căror influență se poate descrie prin funcția-răspuns a impulsului  $\varepsilon(t)$  ponderat cu eficacitatea sinaptică („pondera”) $w_{ij}$ :

$$x_j(t) = \sum_{i \in \Gamma_j} w_{ij} \varepsilon(t - t_i) \quad \text{Ec. 66}$$

Funcția-răspuns a impulsului în Ec. 66 modelează eficient potențialul post-sinaptic ne-ponderat (PSP) al unui singur impuls incident asupra unui neuron. Înălțimea PSP este modulată de către ponderea post-sinaptică  $w_{ij}$  pentru a obține potențialul post-sinaptic efectiv. Funcția-răspuns a impulsului folosită în experimente este definită prin Ec. 68. În rețeaua realizată, o conexiune individuală constă dintr-un număr fix de terminale sinaptice, unde fiecare terminal constituie o sub-conexiune asociată unei întârzieri și unei ponderi diferite (Figura 5.7). Întârzierea  $d^k$  a terminal sinaptic  $k$  se definește prin diferența dintre momentul de activare a neuronului pre-sinaptic și momentul în care potențialul post-sinaptic începe să crească.

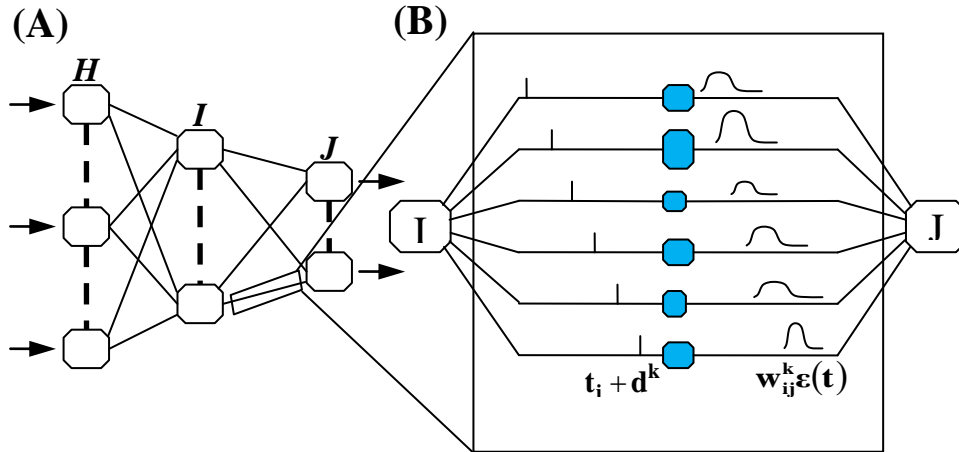


Figura 5.8 Structura logică a rețelei neuronale implementate

Vom descrie un impuls pre-sinaptic la un terminal sinaptic  $k$  ca și un PSP de înălțime standard cu o întârziere  $d^k$ . Contribuția ne-ponderată a unui singur terminal la variabila de stare se dă prin formula:

$$y_i^k(t) = \varepsilon(t - t_i - d^k) \quad \text{Ec. 67}$$

cu o funcție-răspuns a impulsului de forma unui PSP, cu  $\omega(t) = 0$  pentru  $t < 0$ .

Timpul  $t_i$  este timpul de activare a neuronului pre-sinaptic  $i$ , iar  $d_k$  este întârzierea asociată terminalului sinaptic  $k$ . Funcția răspuns a unui impuls care descrie un PSP standard are forma:

$$\varepsilon(t) = \frac{t}{\tau} e^{-\frac{t}{\tau}} \quad \text{Ec. 68}$$

modelând o funcție  $\alpha$  simplă pentru  $t > 0$ , iar  $\tau$  modelează constanta de timp a întârzierii potențialului membranei care determină timpul de creștere și de întârziere a PSP. Extinzând Ec. 66 pentru a include sinapse multiple pe conexiune și inserând Ec. 67, variabila de stare a unui neuron  $x_j$  a unui neuron  $j$  care primește intrări de la toți neuronii  $i$  se poate descrie ca și suma ponderată a contribuțiilor pre-sinaptice:

$$x_j(t) = \sum_{i \in \Gamma_j} \sum_{k=1}^m w_{ij}^k y_i^k(t) \quad \text{Ec. 69}$$

unde  $w_{ij}^k$  înseamnă ponderea asociată terminalului sinaptic  $k$  (Figura 5.8 B). Timpul de activare  $t_j$  a unui neuron  $j$  este determinat ca și primul moment în care variabila de stare depășește pragul  $\mathcal{G}$ :  $x_j(t \geq \mathcal{G})$ . Astfel, timpul de activare  $t_j$  este o funcție neliniară a variabilei de stare  $x_j$ :  $t_j = t_j(x_j)$ . Pragul  $\mathcal{G}$  este constant și egal pentru toți neuronii rețelei.

#### 5.4.2. Arhitectura rețelei implementate pentru clasificarea de fascicule cu domenii receptive

Rețeaua neuronală hibridă RBF-spiking a fost implementată cu două intrări, având un spațiu de intrări bidimensional, format din  $192 \times 192 = 36864$  puncte. Scopul implementării este ca rețeaua să poată identifica punctele din jurul unor centre definite. Numărul centrelor este determinat de numărul neuronilor de ieșire. Fiecare neuron de ieșire va tinde să se activeze atunci, când intrările sunt destul de aproape de centrul fasciculului (cluster) pe care l-a învățat.

Pentru a studia capacitatea și precizia sporită asigurată de codificarea descrisă, am examinat clasificarea unui număr de probleme generate artificial. Când am clasificat intrări constând din două variabile separate, am aflat că o rețea formată din 24 de neuroni de intrare (fiecare variabilă codificată de 12 neuroni) este capabilă să clasifice corect 17 fascicule distribuite în mod egal, prezentând o creștere semnificativă în capacitatea de acumulare a fasciculelor. Astfel de rezultate pozitive au fost obținute și în cazul unor acumulări de fascicule spațiate mai puțin regulat. Am aflat de asemenea că, diminuând

întinderea câmpurilor domeniilor receptive, se pot observa acumulări de fascicule distribuite tot mai compact.

#### 5.4.2.1. Structura proiectului VHDL care implementează rețeaua neuronală

Proiectul VHDL are următoarele elemente componente principale, aranjate conform ierarhiei de mai jos:

- ❖ Rețeaua neuronală
  - Bloc de intrare (Input0block.vhd)
    - Neuronul de intrare (InputNeuron.sch)
      - Blocul funcțiilor domeniilor receptive locale (myBRAMbasefunc.vhd)
      - Modulul funcțional al neuronului de intrare (psRBFInputNeuron.vhd)
  - Modulul sinapselor (Synapse.vhd)
  - Modulul neuronilor de ieșire (RBPoutputSpikingNeuron.vhd)
  - Generatorul pașilor temporali de funcționare (Timestepgenerator.vhd)
- ❖ Modulul de comunicație cu calculatorul (DownloadUnit.vhd)

Figura 5.9 respectiv Figura 5.10 prezintă schema bloc și cea de implementare a întregii rețele neuronale.

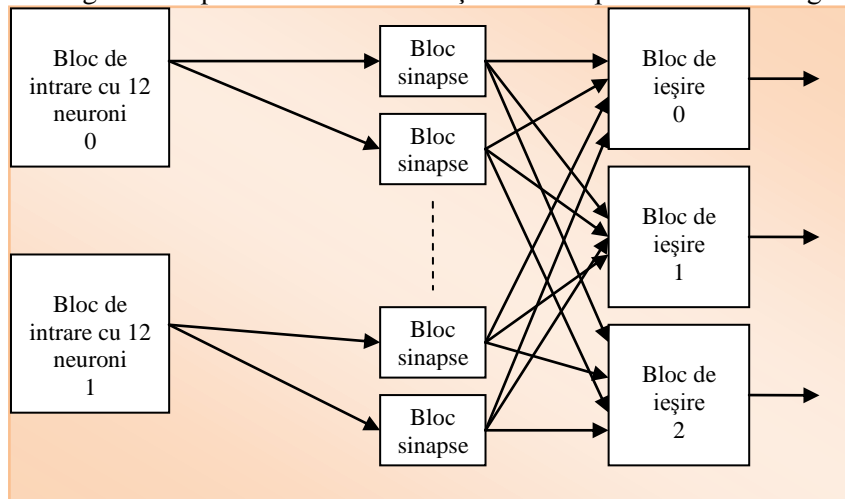


Figura 5.9 Schema bloc a structurii rețelei neuronale implementate

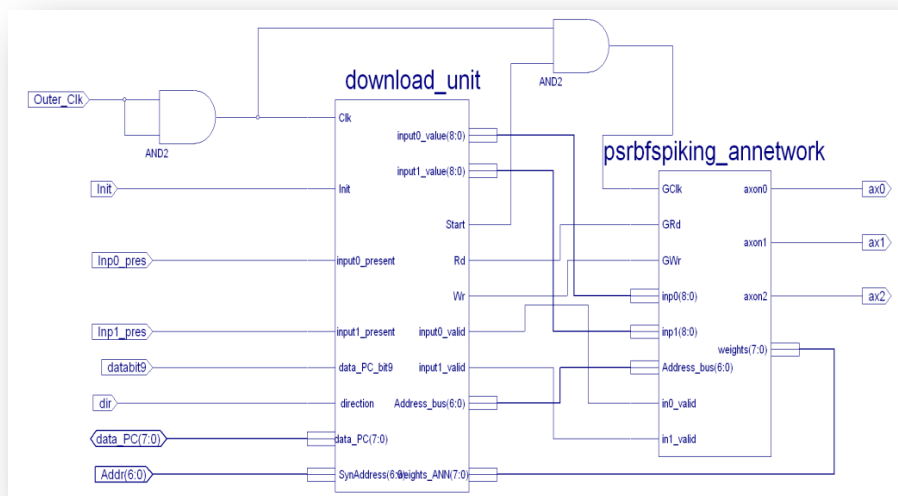


Figura 5.10 Schema de conectare a modulelor principale, în vederea punerii în funcțiune a rețelei neuronale

Simulările au fost executate pentru fiecare modul, pentru care a fost posibil, în cadrul programului Modelsim XE Starter, care este o versiune limitată, dar gratuită al unui pachet de programe de simulare extrem de puternice.

În continuare, se va prezenta fiecare modul în parte, detaliind rolul și funcționare acestora, cu câteva exemple de rezultate de simulare.

#### 5.4.2.2. Neuronul de intrare

Acest modul este generat de două ori, pentru fiecare variabilă de intrare. El conține un modul de memorie BRAM inițializată cu valorile funcțiilor domeniilor receptive locale, respectiv un modul funcțional. Schema din Figura 5.11, este o variantă a acestui modul de neuron de intrare, care a fost folosit pentru testarea funcționării, și conține multe ieșiri, care sunt prezente doar pentru a putea urmări în cadrul simulării variabilele interne.

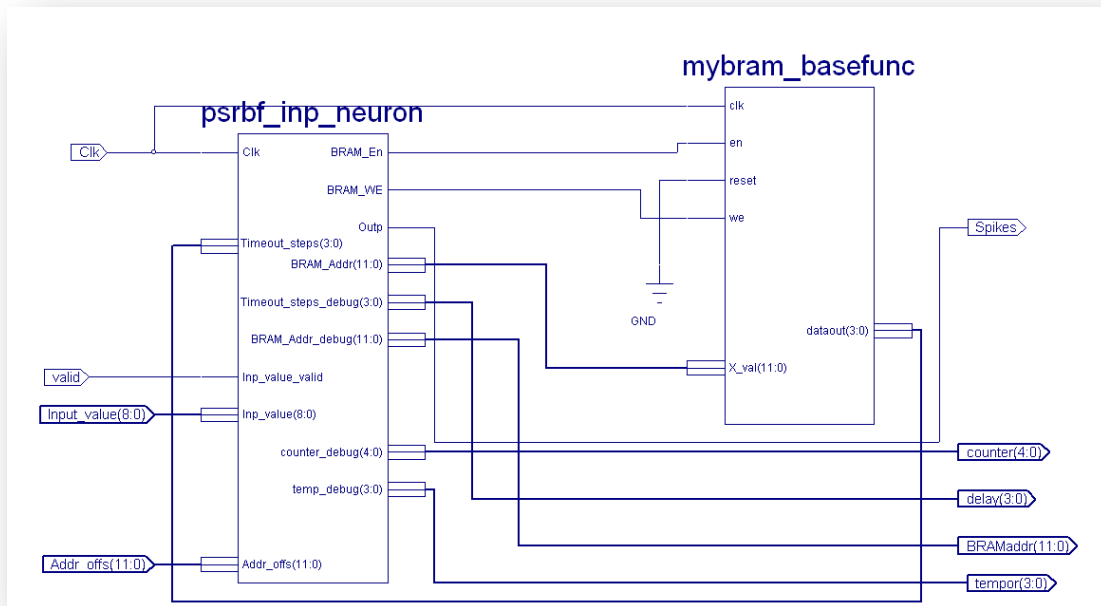


Figura 5.11 Schema de conexiune a unui neuron de intrare

Modulul de control funcțional (Figura 5.12) este un automat cu mai multe stări și execută următoarele operații: pentru valoarea de intrare sosită [Inp\_value (8:0)], adresează memoria BRAM, activând în ordinea necesară și semnalele de control ale acesteia (Rd, En) și extrage valoarea decalajului corespunzător, dat de funcția stocată. Folosind această valoare se generează un impuls pre-sinaptic cu decalaj temporal  $d^k$ . Acest lucru se întâmplă în mod paralel pentru fiecare dintre cei 24 de neuroni de intrare ai celor două blocuri de intrare.

Diagramele din figurile de mai jos prezintă testarea prin simulare – conform stimulilor întocmiți – a subansamblelor neuronilor de intrare. Se poate urmări cum modulul de memorie livrează valorile de decalaj, iar apoi generarea impulsurilor decalate, adică codificarea unei valori de intrare.

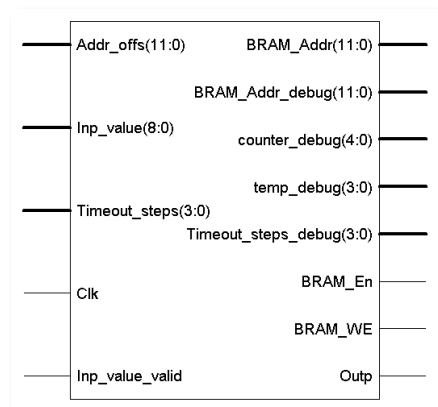


Figura 5.12 Modulul de control funcțional

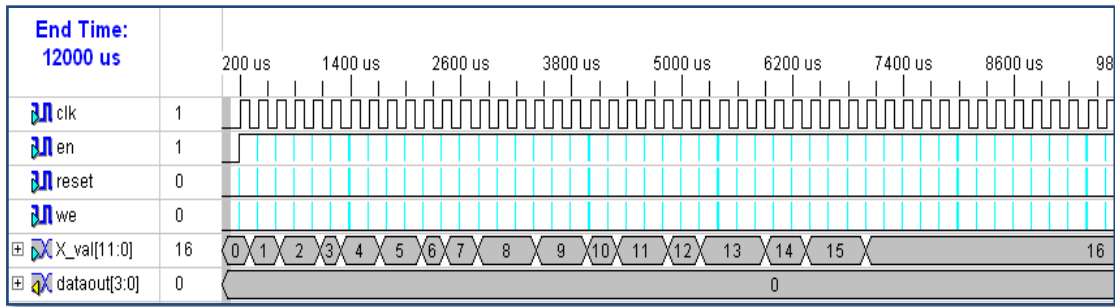


Figura 5.13 Stimul de intrare pentru simularea modului BRAM



Figura 5.14 Rezultat de simulare, valori ale funcțiilor domeniilor receptive locale (semnalul **dataout**) pentru valori crescătoare ale variabilei de intrare (semnalul **xval**)

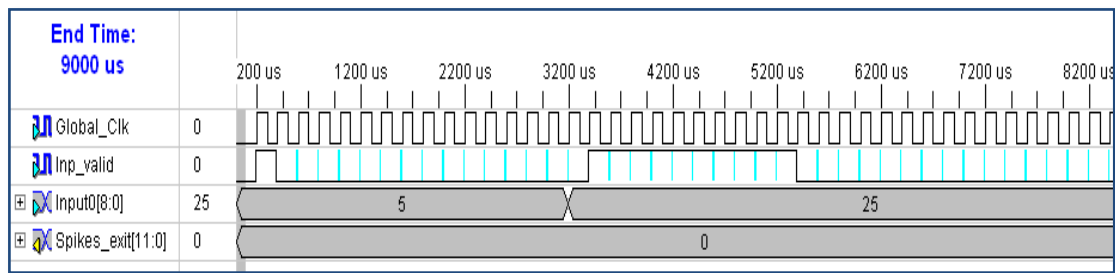


Figura 5.15 Stimul pentru testarea codării intrării în impulsuri decalate temporal

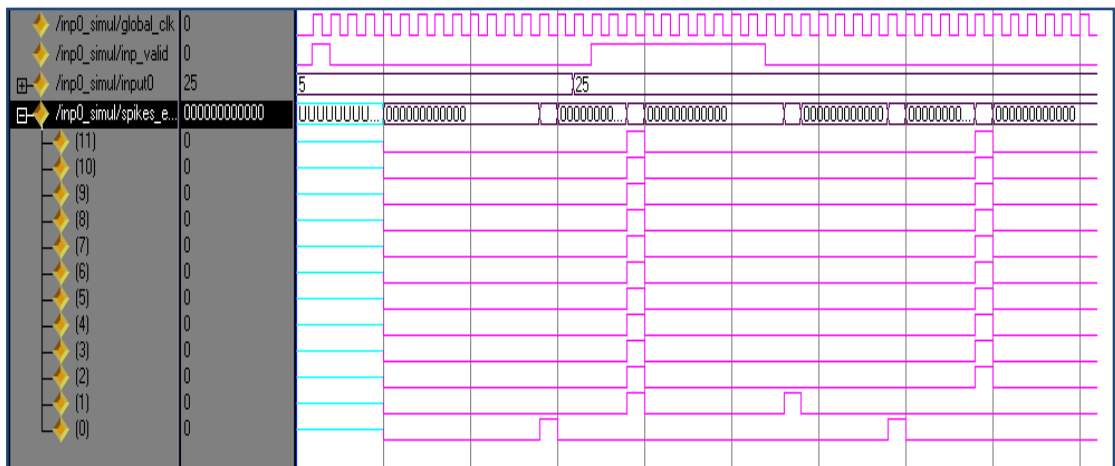


Figura 5.16 Rezultat de simulare al codării intrării în impulsuri. De exemplu, pentru valoarea de intrare 25, se vor activa neuronii de intrare 0 și 1, cu decalaje de 5 respectiv 12 pași (cicluri de tact), iar ceilalți neuroni emit impuls doar la sfârșitul ciclului de funcționare (după 16 pași), generat de un modul dedicat pentru acest scop.



### 5.4.2.3. Modulul sinapsă

Acest modul este responsabil pentru amplificarea sau atenuarea impulsului pre-sinaptic, transmițând o valoare ponderată către soma neuronului. Totodată aici este implementat și algoritmul de învățare. Învățarea se face ne-supervizat, conform unei versiuni adaptate a regulilor Hebb, descrise în capitolele anterioare. Sinapsele verifică eticheta temporală a fiecărui impuls pre-sinaptic, și modifică valoarea ponderii conform regulilor Hebb în varianta modificată. Astfel, de-a lungul unui ciclu de funcțiune de 16 perioade ale semnalului de tact, pentru impulsurile, care au

sosit cu maxim 5 perioade înaintea activării post-sinaptice, fortificare ponderală est drastică. În cazul acelor impulsuri, care au sosit cu un decalaj temporal mai mare de 5 dar mai mic de 10 perioade de tact, ponderea crește moderat, iar pentru restul impulsurilor ponderea scade ușor. Dacă impulsul pre-sinaptic sosește după apariția impulsului axonal, ponderea sinapsei respective descrește abrupt.

Intrările și ieșirile acestui modul sunt semnale ce servesc la controlul procesului de învățare, respectiv la adresare sinapsei în vederea inițializării și citirii valorilor ponderilor.

### 5.4.2.4. Modulul de soma al rețelei neuronale pseudo-RBF-pulsative

Soma este modulul care este responsabil pentru sumarea impulsurilor incidente, decalate temporal și amplificate de către sinapse. Fiecare modul de somă primește 24 de intrări, de câte 8 biți, din care calculează valoare potențialului de membrană. La depășirea valorii de prag, se emite un impuls de activare, iar valoarea potențialului de membrană se resetează la valoarea de hiperpolarizare prestabilită.

### 5.4.2.5. Modulul de comunicație cu calculatorul

Pentru a putea urmări evoluția procesului de învățare, precum și pentru realizarea inițializării rețelei, s-a implementat un modul special, destinat comunicației dintre un calculator și placa de dezvoltare FPGA. Această comunicație se desfășoară paralel pe două căi, și anume prin portul paralel și printr-o placă de achiziții date National Instruments conectabilă pe magistrala USB. Tot prin acest modul se prezintă și variabilele de intrare rețelei neuronale hardware ce funcționează în circuitul FPGA.

## 5.5. Realizări de rețele neuronale neuromorfe utilizând microcontrolere soft-core înglobate în circuitele FPGA

În trecut, sistemele de calcul au fost construite din componente separate cum ar memorii, elemente de I/O, circuite FPGA, etc. Totodată, diferitele programe de proiectare au fost (unele sunt și astăzi) utilizate pentru a specifica funcționalitatea individuală a acestor componente, precum și pentru a le integra să funcționeze ca un sistem complet.

Componenta centrală a acestor sisteme, procesorul, a fost (și mai este) cel mai des o unitate separată. Interfațarea acestuia cu celelalte componente – logică, memorii, I/O – este realizată cu ajutorul pachetelor software mai sus menționate. Partea software a acestor sisteme de calcul a fost (și este) proiectată cu programe specializate pentru dezvoltarea se soft-uri și sisteme de operare. Interfațarea hardware-software este realizată prin utilizarea unei combinații ale acestor unelte menționate și în unele cazuri unelte de co-verificare sau co-simulare.

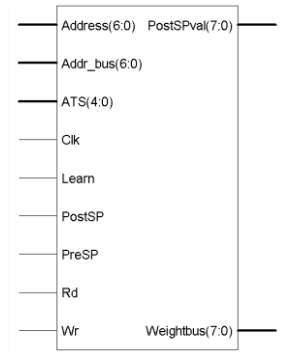


Figura 5.17 Modulul sinapsă

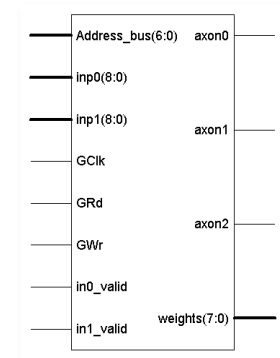


Figura 5.18 Modulul de soma

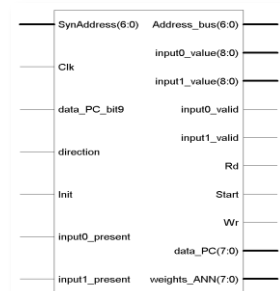


Figura 5.19 Modulul de comunicație cu calculatorul



În a doua parte a anilor 1990 multe funcții logice au fost integrate în dispozitive singulare cum ar fi familiile de circuite 4K sau Virtex™ ale firmei Xilinx. Memoriile de capacitate mare nu sunt încă disponibile în aceste circuite, dar memoriile încorporate se dezvoltă odată cu apariția noilor familii de astfel de circuite și cu dezvoltarea tehnologiei de fabricație a acestora. Aceste progrese (Figura 5.20) au făcut posibilă ca o parte mai importantă a proiectării și implementării logice să se realizeze într-o

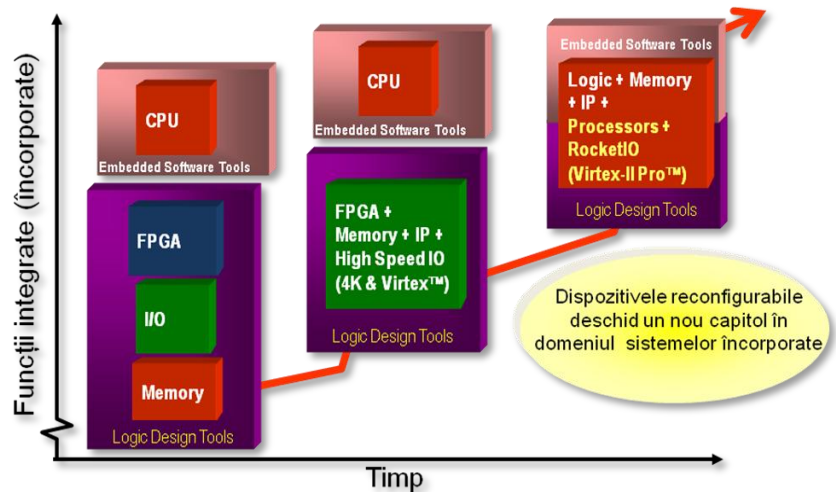


Figura 5.20 Evoluția funcțiilor integrate în circuitele FPGA

singură unitate, eliminând nevoia de interfațare dintre componente. Totuși nu s-a ajuns la schimbări radicale în ceea ce privește dezvoltarea în domeniul unităților centrale de prelucrare (CPU).

Odată cu apariția circuitelor Platform FPGA și a celor SoC ASIC (System-on-Chip Application Specific Integrated Circuit), fuzionarea tuturor acestor funcționalități discrete a devenit realitate. Utilizând potențialul acestor sisteme, care sunt complet reconfigurabile hardware și software, se obține nu numai un produs finit excepțional dar și o platformă de dezvoltare foarte avantajoasă care necesită, însă cooperarea dintre elementele de siliciu, de module predefinite (proiecte de circuite și programe) și îmbinarea coezivă a acestora cu o suită de unelte de proiectare bine puse la punct.

Dezvoltarea și progresele recente din acest domeniu, împreună cu apariția procesoarelor încorporabile în circuitele FPGA, au făcut ca alegerea acestora ca platformă pentru implementarea hardware a RNP parțial serializată să fie evidentă. În continuare se vor prezenta contribuții ale tezei în acest domeniu, și anume implementări de exemple aplicative și teste benchmark clasice utilizând procesoare încorporate pe 8 respectiv pe 32 de biți.

### 5.5.1. Microcontrolerul soft-core Xilinx PicoBlaze

PicoBlaze, cunoscut și ca KCPSM3, este un microcontroler RISC cu o amprentă de siliciu extrem de redusă, de numai 96 de slice-uri dintr-un FPGA din familia Xilinx Spartan3 (optimizat pentru această familie) ceea ce reprezintă aproximativ 1,25% din totalul de 7680 de astfel de structuri logice disponibile de exemplu în circuitul XC3S1000. Cea mai mare frecvență de tact la care poate funcționa este de 87MHz, ceea ce rezultă o capacitate de calcul respectabilă de ~43.5 MIPS.

Toate aceste sunt performanțele unui microcontroler de tip soft-core, care este livrat de furnizorul Xilinx în forma unui cod VHDL, putând fi încorporat în orice proiect destinat implementării pe un circuit FPGA.

După cum se poate urmări și în Figura 5.21, PicoBlaze dispune de o memorie RAM internă de uz general de 64 de octeți și este capabil să comunice cu restul componentelor din proiectul în care este încorporat prin 256 porturi de intrare și 256 porturi de ieșire de câte 8 biți fiecare. Portul accesat este adresat prin ieșirea de 8 biți PORT\_ID. Operațiile de citire porturi sunt validate de impulsuri pe ieșirea READ\_STROBE iar cele de scriere porturi de impulsuri pe ieșirea WRITE\_STROBE. Valorile astfel scrise/citite pot fi valori constante sau conținutul unuia dintre cei 16 regiștrii (de un octet) interni ai procesorului.

PicoBlaze dispune de un sistem de întrerupere ne-vectorizat foarte simplu, cu o singură intrare de semnalizare a acesteia. Toate cele ~50 de instrucțiuni sunt executate întotdeauna în două cicluri de tact și pot utiliza o stivă cu adâncimea de 31.

Programele pentru acest procesor sunt scrise într-un limbaj de asamblare cu instrucțiuni specifice acestuia. Asamblorul livrat de producător generează un cod VHDL care instanțiază o memorie ROM (utilizând modulele BRAM din FPGA) inițializată cu codurile de instrucțiuni și datele care

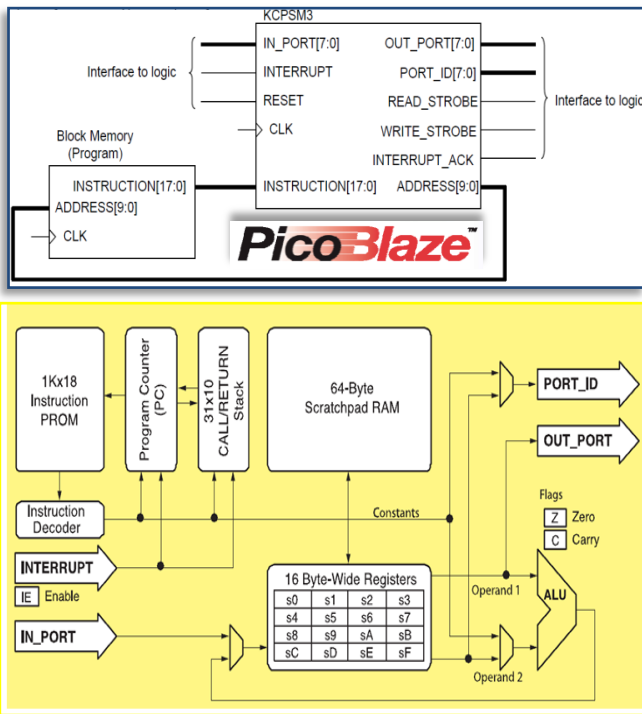


Figura 5.21 Diagrama bloc a microcontrolerului soft-core Xilinx PicoBlaze încorporat în RNP implementate

implementează programul scris. Acest cod VHDL se poate integra cu ușurință în proiectul care conține deja cel puțin un microcontroler PicoBlaze, funcționând ca memorie de program pentru acesta.

PicoBlaze poate adresa 1024 locații a memoriei program, aceasta fiind una din limitările sale principale. Totuși, în multe cazuri este foarte util în a executa algoritmi ce ar necesita multe resurse hardware reconfigurabile. Pentru testarea și simularea funcționării programelor scrise, am utilizat software-ul Mediatronix pBlaze IDE v3.6 care este un debugger simplu dar eficient. Problema majoră a acestui utilitar este faptul, că nu este 100% compatibil din punct de vedere al sintaxei programului cu asamblorul PicoBlaze.

### 5.5.2. Implementare aplicativă: Detector de componente de frecvență

Scopul acestei aplicații este de a demonstra, că rețele neuronale pulsative implementate pot rezolva probleme ingineresti. Rezultatele prezentate în cele ce urmează trebuie văzute ca o demonstrație a capacităților acestor sisteme, și nu ca o expunere a unei rezolvări ieșite din comun a problemei în cauză.

**Formularea problemei:** sistemul dezvoltat trebuie să fie capabil să învețe a separa dintr-un semnal analog (perturbat de zgomot) discretizat, trei componente de frecvență dominante, în jurul a trei valori predeterminate.

#### 5.5.2.1. Rețeaua neuronală pulsativă implementată pe FPGA

S-a implementat o RNP cu două intrări și trei ieșiri, după cum arată și schema bloc din Figura 5.22. Această rețea trebuie să clasifice corect perechile de valori (numerice, întregi, pozitive) care definesc un spațiu de intrări de  $256 \times 256 = 65536$  puncte (Figura 5.23). Fiecare punct al acestui spațiu de intrări reprezintă câte o componentă de o anumită frecvență și putere din componența semnalului analizat. Antrenarea rețelei s-a efectuat în așa fel, încât fiecare ieșire să se activeze atunci, când semnalul respectiv conține o componentă a cărei frecvență și putere se află în vecinătatea valorii prescrise (punct de focus), adică să aproximeze prezența acestei componente cu o anumită histereză.

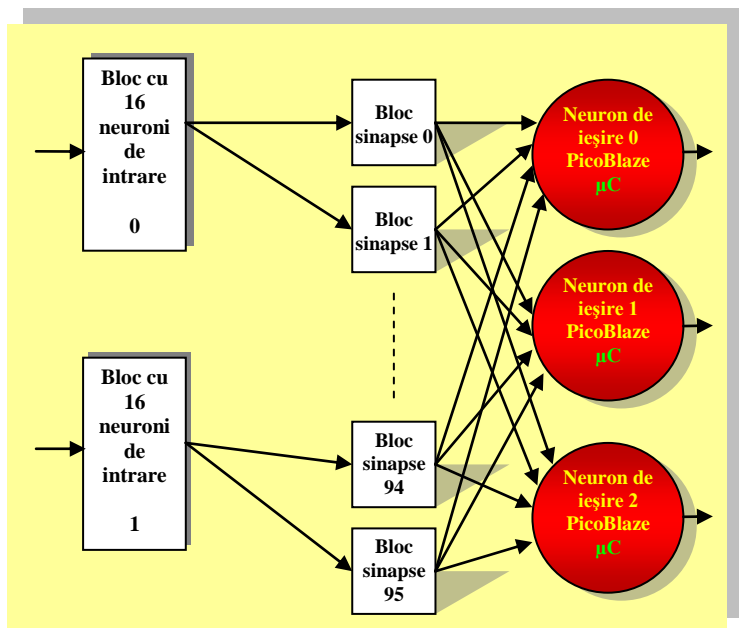


Figura 5.22 Schema bloc a rețelei pseudoRBF-spiking implementată

**Parametrii rețelei neuronale implementate:**

- Numărul intrărilor:  $n = 2$
- Nr. neuronilor de intrare/intrare:  $n_{intr}=12$
- Nr. neuronilor de intrare activi/intrare:  $n_{act} = 3$
- Numărul pașilor dintr-un cadru temporal:  $T=16$

**5.5.2.2. Implementarea neuronilor de intrare – codificarea variabilelor de intrare în impulsuri decalate temporal**

Pentru a putea specifica valorile de intrare rețelei neuronale o decizie a fost luată asupra algoritmului de codificare a acestora în impulsuri. Această codificare putea fi implementată în așa fel încât să fie executată on-chip, dar s-a considerat, că această metodă ar utiliza excesiv de multe resurse reconfigurabile sau ar fi executată prea lent de un procesor încorporat. Așadar codificarea s-a efectuat parțial a-priori, după cum se

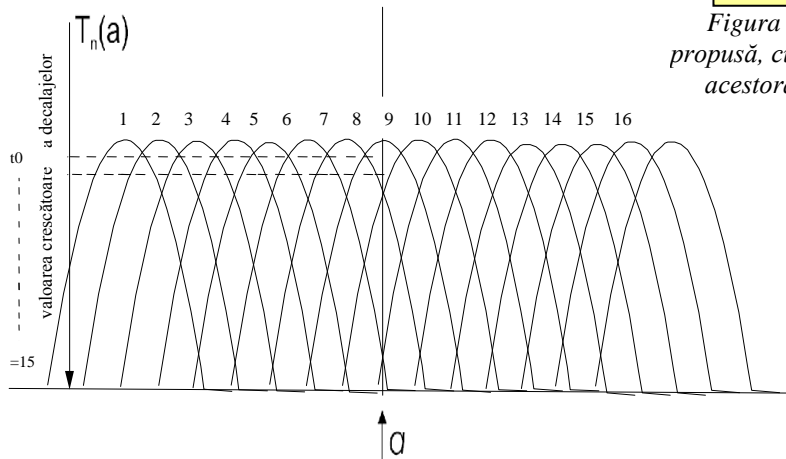


Figura 5.24 Metoda de codificare a valorilor de intrare în impulsuri decalate temporal

pentru a simplifica implementarea am ales o aproximare a acestora cu funcții triunghiulare. Figura 5.25 arată, cum o valoare arbitrară (din domeniul de intrare cu valori întregi între 0-255) activează patru domenii receptive (marcate cu linii punctate), funcția corespunzătoare neuronului 7 fiind activat cu valoarea cea mai mare. Acest lucru va însemna, că acest neuron de intrare va emite primul impuls de activare, iar restul de neuroni activi pentru această valoare se vor activa într-un pas ulterior. Codificarea este realizată într-un cadru temporal de 16 pași, corespunzători cu un

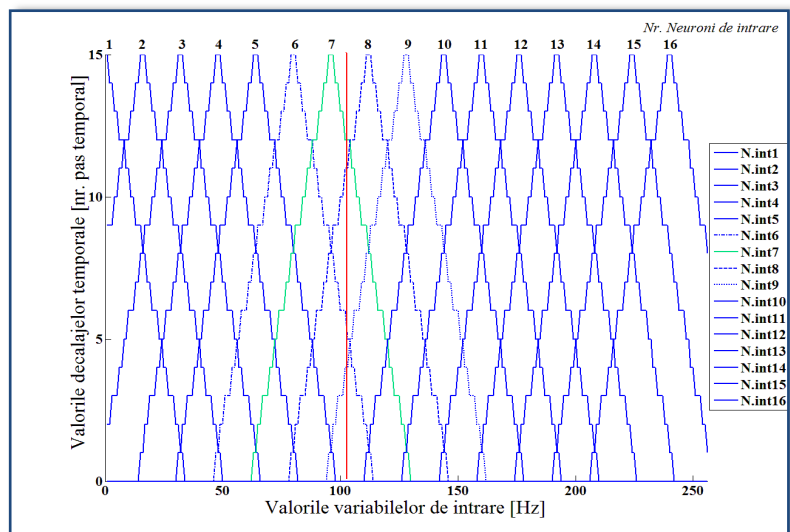


Figura 5.25 Funcțiile triunghiulare ale domeniilor receptive a neuronilor de intrare, cu scalarea decalajelor temporale

Rețeaua se poate considera, ca fiind una hibridă, putând fi numită pseudo-RBF-spiking, dat fiind faptul, că modalitatea de codificare a valorilor de intrare utilizează funcții de bază ca și rețelele neuronale de tip Radial Basis Function (RBF).

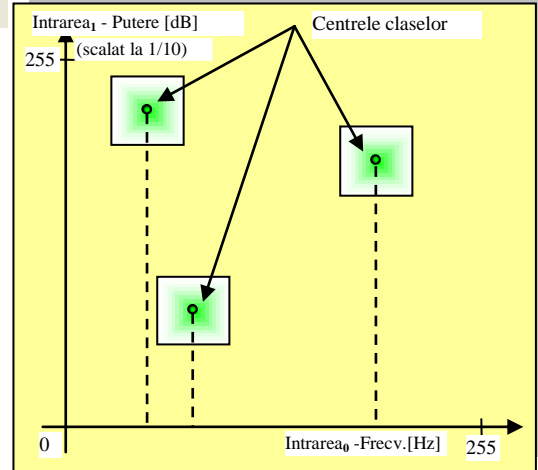


Figura 5.23 Spațiul intrărilor pentru problema propusă, cu exemple de centre de focus și vecinătățile acestora care vor fi învățate de RNP din FPGA

prezintă în continuare.

Scopul este de a găsi o variantă de acoperire a spațiului intrărilor cu domenii receptive (funcții de activare) - variind aria de suprapunere a acestora, în așa fel, încât pentru fiecare valoare de intrare să se activeze un număr redus, dar egal de funcții (de exemplu 4). În cazul ideal, aceste

funcții ar trebui să fie Gaussiene, după cu arată Figura 5.24, dar

număr egal de perioade de tact la care funcționează acest circuit. Astfel, un neuron de intrare activat cu valoarea maximă posibilă a funcției triunghiulare atașate de 15 va emite impuls de ieșire în primul pas temporal al cadrului (pasul 0).

Așadar, această alegere a valorilor de decalaj temporal, rezultă într-un necesar de patru biți/pe valoare pentru stocarea acestora. Valorile au fost calculate cu ajutorul unui program scris în Matlab prin care se poate varia numărul funcțiilor de activare și aria acestora.

Același program Matlab generează și codul VHDL necesar pentru inițializarea modulelor BlockRAM din circuitul FPGA, utilizate pentru stocarea valorilor de decalaj temporal.

După cum a fost prezentat și la subpunctul 3.1.3.1 elementele BlockRAM ale circuitelor FPGA pot fi configurate în mai multe feluri, variind numărul biților de pe magistrala de adrese respectiv de pe cea de date. Pentru a stoca valorile de decalaj temporal am ales componenta cu numele RAMB16\_S36\_S36, care este un modul BRAM cu port dual și o capacitate de 512x32 biți. Arhitectura cu port dual permite ca două locații de memorie să fie accesate simultan. Deoarece valorile stocate sunt inițializate odată cu programarea circuitului FPGA, acest modul va fi utilizat ca și cum ar fi de fapt o memorie ROM, efectuând numai operații de citire, evitând astfel eventualele conflicte de adresare la scriere. Conectând două astfel de module în paralel, devine posibil să se obțină în aceeași perioadă de tact cei 128 biți reprezentând decalajele temporale pentru cele două intrări ale RNP, rezultat obținut conform ecuației următoare:

$$\sum \langle d^{ij} \rangle = \sum_n n_{intr} * \nu \tag{Ec. 70}$$

unde am introdus următoarele notații:

$n$  – este numărul intrărilor RNP, în cazul acestei aplicații 2

$n_{intr}$  – este numărul neuronilor de intrare care codifică valorile unei intrări ale RNP, aici 16

$\nu$  - este numărul de biți pe care este reprezentat un decalaj temporal, în cazul de față 4

$\langle d^{ij} \rangle$  – este numărul de biți necesari pentru reprezentarea tuturor decalajelor temporale ale impulsurilor ce se propagă de la neuronii pre-sinaptici  $i$  alocați unei intrări ale RNP spre neuronul post-sinaptic  $j$ , ( $4 \times 16 = 64$ ).

Valoarea variabilei de intrare este folosită pentru adresarea memoriei, accesând locații diferite pe cele două porturi ale acesteia, utilizând un offset corespunzător (Figura 5.27) la calcularea adresei pe portul B. Am realizat, astfel cea mai eficientă utilizare posibilă a unităților BlockRAM, deoarece pentru a acoperi spațiul valorilor de intrare am folosit 100% din capacitatea acestora.

Aceste valori citite din modulul de memorie - a cărei scheme bloc este prezentată de Figura 5.26 - vor fi folosite de circuitele care implementează neuronii de intrare pentru a codifica valoarea de intrare cu un impuls emis în pasul temporal cu numărul egal cu valoarea decalajului (pași cu numărul 0 la 15).

Pentru a garanta, însă, că pentru fiecare valoare de intrare se va activa un număr egal de neuroni de intrare, a fost nevoie de un artificiu. Odată cu mărirea ariei de suprapunere a funcțiilor, ar fi fost nevoie sau de mărirea numărului acestor funcții pe intrare sau de mărirea numărului pașilor dintr-un

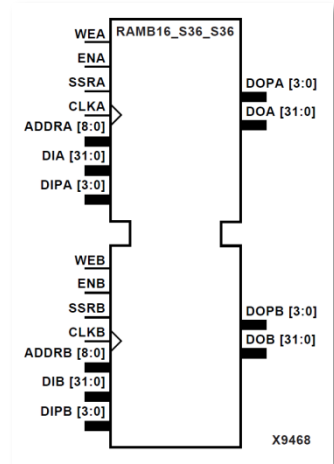


Figura 5.26 Modulul BlockRAM utilizat la stocarea decalajelor temporale

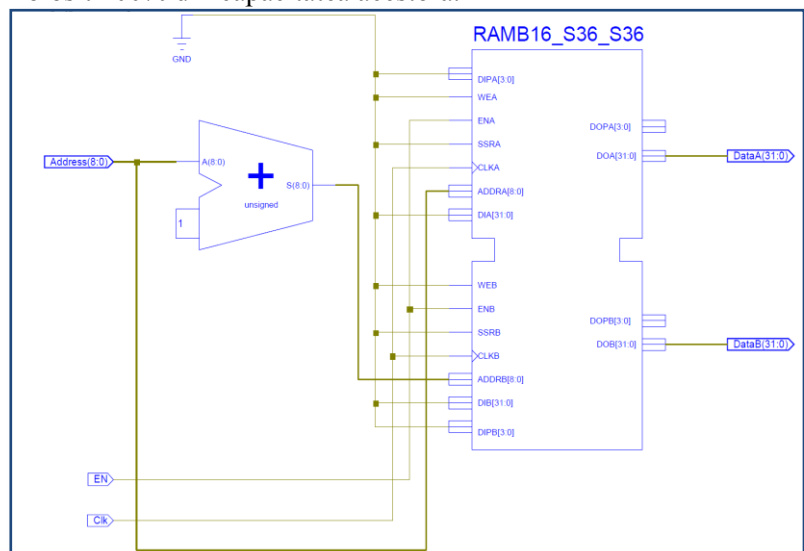


Figura 5.27 Schema de conexiune unui bloc BRAM al unui neuron de intrare

cadru temporal, pentru a menține forma strict triunghiulară a acestora. Aceste modificări, însă ar fi rezultat o situație, în care un modul BRAM nu ar fi fost suficient pentru stocarea tuturor valorilor de decalaj temporal, pentru că valoarea  $\mathcal{U}$  din Ec. 70 ar fi crescut. Așadar, scalarea valorilor de decalaj din intervalul  $[0...35]$  necesar garantării condiției mai sus amintite în intervalul original de  $[0...15]$  s-a dovedit a fi cea mai simplă și eficientă soluție. Aceasta a dus la modificarea ușoară a formei funcțiilor de activare, după cum se poate vedea și în Figura 5.25 ceea ce nu influențează, însă precizia codificării.

În continuare se va prezenta testarea prin simulare a circuitelor ce implementează un neuron de intrare respectiv un întreg bloc de neuroni de intrare pentru o variabilă de intrare a RNP. Simulările au fost realizate cu simulatorul integrat în mediul de dezvoltare Xilinx ISE 10.1.

În Figura 5.28 se pot urmări semnalele de intrare utilizate pentru testarea funcționării unui modul de codificare a intrărilor RNP. Se poate observa, că la intervale de 16 cicluri de tact s-a modificat valoarea semnalului de intrare (valorile 0, 2 și 14) reprezentând decalajul temporal cu care trebuie să apară impulsul de ieșire, valoare citită în modul normal de operare din memoria BlockRAM. A se nota de asemenea prezența semnalului *delay\_valid* care este inactiv pe perioade de câte două cicluri de tact, pentru a semnală că valoarea de decalaj nu este validă pe timpul citirii acesteia din memoria internă.

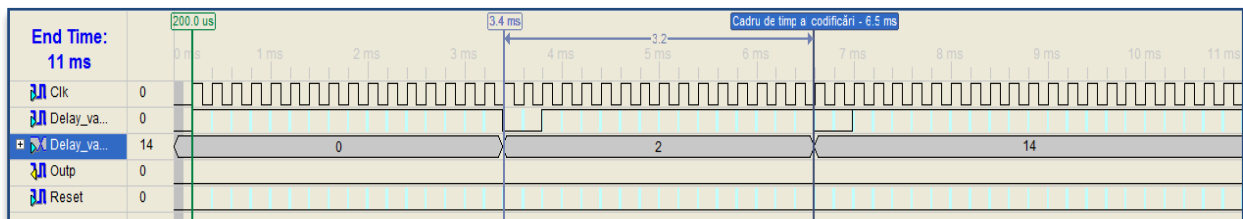


Figura 5.28 Stimul de intrare pentru simularea funcționării unui modul de codificare a intrărilor RNP

Ieșirile circuitului ce implementează modulul de codificare a valorilor de intrare a RNP la stimulul prezentat anterior (Figura 5.28) sunt afișate pe Figura 5.29. Semnalul de notat în această figură este *outp* pe care vor apare impulsurile decalate cu atâtea cicluri de tact cât a arătat valoarea de decalaj specificată pe intrarea *delay\_value*.

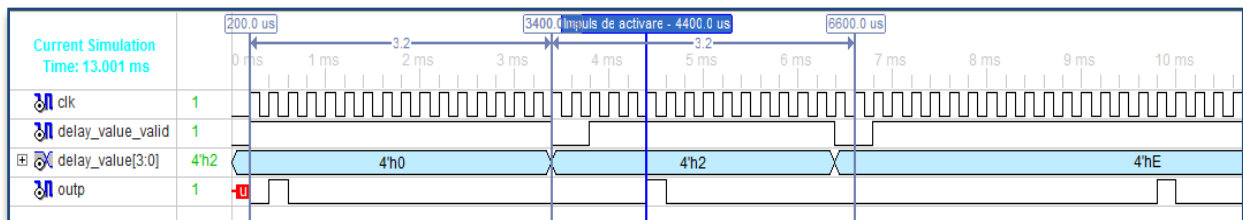


Figura 5.29 Rezultatul simulării unui modul de codificare a valorilor de intrare a RNP

### 5.5.2.3. Blocul neuronilor de intrare

În acest circuit sunt instanțiate câte 16 neuroni de intrare pentru fiecare intrare a RNP. Rețeaua neuronală ce implementează aplicația curentă având două intrări (valori întregi reprezentate pe 8 biți), acest modul a fost generat în două exemplare. Structura acestui circuit (Figura 5.30) înglobează și cele două memorii BRAM care stochează valorile pre+calculate ale funcțiilor de activare a celor  $2 \times 16$  neuroni de intrare.

Modulul funcțional al acestui circuit (Figura 5.31) este de fapt un automat cu mai multe stări, care extrage din memoriile BRAM valorile de decalaj  $d^k$  corespunzătoare funcțiilor de activare iar apoi generează câte un impuls pre-sinaptic decalat în timp cu aceeași valoare (un număr egal de pași temporali din cadrul

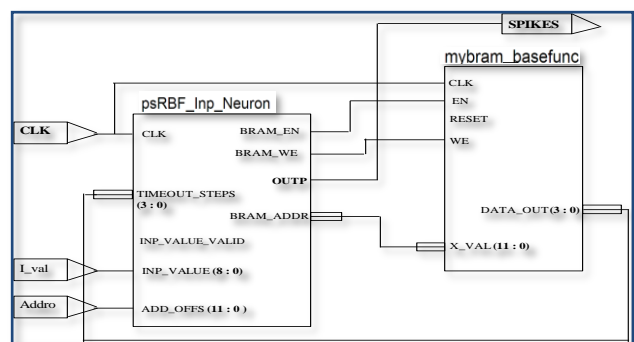


Figura 5.30 Structura blocului neuronilor de intrare



temporal, adică un număr de cicluri de tact) în paralel pentru fiecare dintre cei 16 neuroni de intrare ce codifică variabila de intrare. Pentru a ilustra funcționarea acestui modul, diagramele de timp ce urmează vor arăta rezultatele simulării acesteia. Astfel, în Figura 5.32 se prezintă semnalele de intrare pentru blocul de neuroni de intrare ce codifică una din variabilele de intrare. Se pot observa perechile valorilor de intrare prezentate rețelei neuronale pe intrările *Inp0* și *Inp1* respectiv semnalele de validare ale acestora.

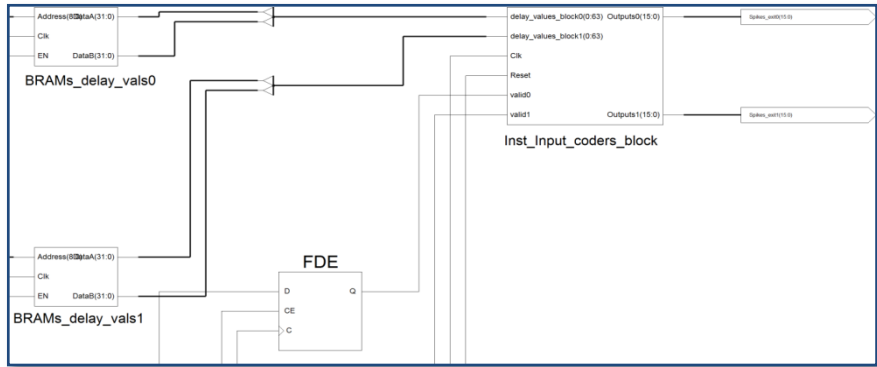


Figura 5.31 Schema de conexiune blocului neuronilor de intrare

Figura 5.33 prezintă rezultatul simulării blocului de neuroni de intrare ce codifică una din variabilele de intrare a RNP, putându-se observa impulsurile de ieșire decalate apărute în urma prezentării valorilor de intrare, precum și offset-ul de adresare a portului B al memorie BRAM, conform funcționării descris la punctul anterior.

Tabelul 11 Utilizarea FPGA a blocului neuronilor de intrare

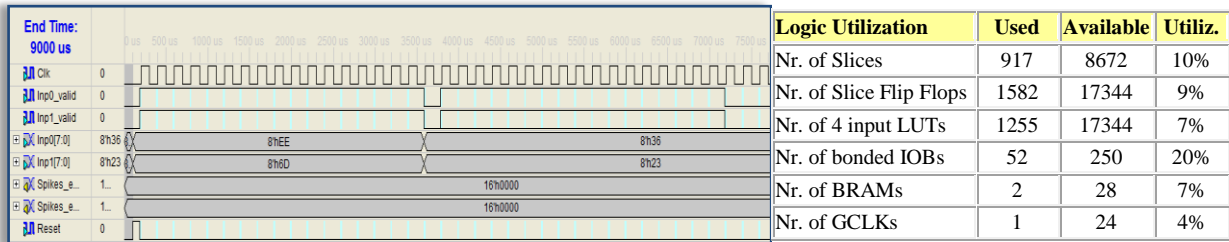


Figura 5.32 Stimul de intrare pentru blocul de neuroni de intrare

Figura 5.33 prezintă rezultatul simulării blocului de neuroni de intrare ce codifică una din variabilele de intrare a RNP, putându-se observa impulsurile de ieșire decalate apărute în urma prezentării valorilor de intrare, precum și offset-ul de adresare a portului B al memorie BRAM, conform funcționării descris la punctul anterior.

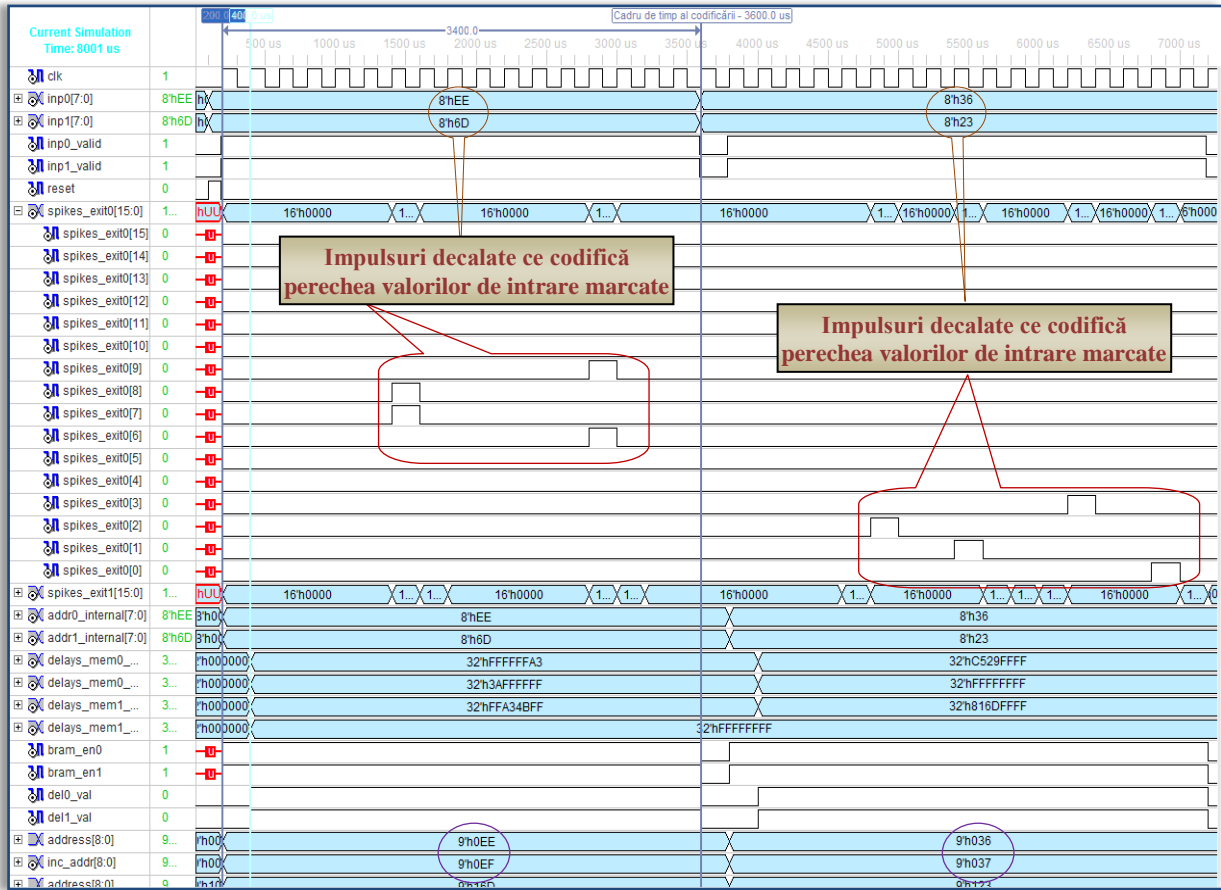


Figura 5.33 Rezultatul simulării blocului de neuroni de intrare ce codifică una din variabilele de intrare a RNP

### 5.5.2.4. Modulul de implementare a sinapselor

Modulul ce implementează sinapsa unui neuron este responsabil pentru amplificarea sau atenuarea impulsurilor pre-sinaptice emise de către neuronul de intrare la care este conectat. Astfel se generează o valoare post-sinaptică ponderată ce va fi transmisă spre soma neuronului de ieșire la care aparține sinapsa respectivă.

Acest modul este de asemenea cel, care implementează regulă nouă de învățare supervizată care este o variantă adaptată a metodei Hebb, bazată pe reguli de decalaje temporale, după cum este prezentat în cele ce urmează. Sinapsele vor testa amprenta temporală (numărul pasului temporal, notat cu *TS* în Tabelul 12) a fiecărui impuls pre-sinaptic și vor ajusta valorile ponderilor stocate în concordanță cu regulile de învățare prezentate în Tabelul 12, luând în calcul și valorile prescrise ale stării (activate sau nu) a neuronilor de ieșire (valoare axonală) pentru valorile variabilelor de intrare actuale.

Așadar, conform tabelului de mai jos, și a condițiilor descrise în Ec. 71, Ec. 72 respectiv Ec.73, în decursul unui cadru de timp de 16 cicluri de tact, pentru acele impulsuri pre-sinaptice care au sosit cu cel mult 5 pași temporali mai devreme decât activarea post-sinaptică (a neuronului de ieșire corespunzător) și dacă valoarea prescrisă este de activare (pentru a corela activarea neuronului de ieșire cu valorile de intrare dintr-o clasă anume), atunci ponderea sinapsei respective este mărită drastic.

#### Parametrii modelului neuronal:

- Valoarea maximă a ponderii:  $\beta_{i \max}^{jk} = 255$
- Potențialul de membrană – valoarea maximă în aplicația curentă:  $\mu_{actmax} = n_{act} * \beta_{i \max}^{jk} * T * n$
- Potențialul de repaus:  $\mu_{rep} = \mu_{actmax} * 0.2$
- Potențialul de prag:  $\theta = \mu_{actmax} * 0.8$
- Potențialul de hiperpolarizare:  
 $\mu_{hip} = \mu_{actmax} * 0.1$
- Potențialul de scurgere:  
 $\mu_{scu} = \mu_{actmax} * 0.05$
- Potențialul de inhibare:  
 $\mu_{inhib} = \mu_{rep} - \mu_{hip} / 2$

Figura 5.34 Parametrii modelului neuronal

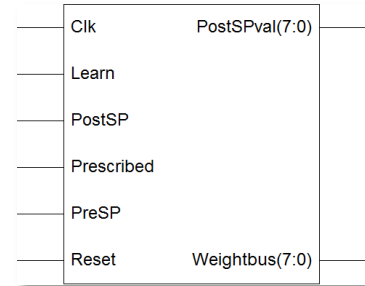


Figura 5.35 Intrările și ieșirile modulului sinapsă

Tabelul 12 Algoritmul de învățare propriu, implementat în hardware

Impuls pre-sinaptic prezent	Valoare axonală - ieșire NP	Reguli temporale	Valoare de ieșire prescrisă	Adaptare pondere	Calculare model somă
0	0	-	0	-	-
0	0	-	1	-	$\Sigma$
0	1	-	0	-	-
0	1	-	1	-	$\Sigma$
1	0	-	0	-	-
1	0	PreTS << PostTS	1	$\uparrow\uparrow$	$\Sigma$
1	0	PreTS < PostTS	1	$\uparrow$	$\Sigma$
1	1	PreTS > PostTS	0	$\downarrow$	-
1	1	PreTS >> PostTS	0	$\downarrow\downarrow$	-
1	1	-	1	-	$\Sigma$

$if \dots \sum \beta_i^{jk} = 0 \dots then \dots \gamma_i = \lambda_i$  Ec. 71

$if \dots \sum \beta_i^{jk} \geq \theta \dots then \dots \gamma_i = \delta_i$  Ec. 72

$if \dots \gamma_{i \pm m}^k = \gamma'_{i \pm m}^k \dots then \dots \gamma_i = \lambda^k$  Ec.73

În cazul impulsurilor care au un decalaj temporal între 5 și 10 cicluri de tact și condiții prescrise identice, ponderea este mărită moderat. Pentru acele impulsuri pre-sinaptice care sosesc după apariția activării axonale și în prezența activării prescrise a acesteia, ponderea sinapsei va fi diminuată proporțional cu diferența dintre aceste două momente.

Intrările și ieșirile acestui modul (Figura 5.35) sunt utilizate pentru a controla procesul de învățare, respectiv pentru a putea citi și inițializa valorile ponderilor. În aplicația de față, proiectul

implementat pe circuit FPGA conține un total de 96 de sinapse, câte una pentru conexiunile dintre cei 32 de neuroni de intrare către cei trei neuroni de ieșire ale RNP, fiecare stocând valoarea ponderii pe un număr de 8 biți. Aceste ponderi sunt inițializate cu valori aleatoare ca apoi să fie adaptate conform regulilor de învățare descrise mai sus în decursul mai multor cicluri de antrenare, în decursul cărora la intrările rețelei neuronale implementate hardware se prezintă perechi de variabile de intrare – de către o Unitate de Monitorizare și Control, aceasta citindu-le dintr-o memorie BRAM separată – reprezentând elementele unui set de antrenare.

#### 5.5.2.5. Modulul neuronilor de ieșire – implementarea somei cu multiple microcontrolere încorporate Xilinx PicoBlaze

Modulul ce implementează soma neuronală este responsabilă pe de o parte pentru sumarea valorilor post-sinaptice ponderate și decalate temporal. Toate cele trei astfel de module ale aplicației de detectare a componentelor de frecvență recepționează câte 32 de impulsuri a reprezentate pe 8 biți fiecare.

Însumând toate aceste valori și urmând același ritm al pașilor temporali ca și modulele de codificare a intrărilor și sinapsele, somele calculează potențialul de membrană (PM) a neuronului artificial, notat cu  $\mu$  în prezentarea parametrilor modelului neuronal din Figura 5.34. Dacă acest PM – inițializat la o valoare de repaus – depășește o valoare de prag predefinită (THS-threshold),

#### Fazele procesării pașilor temporali în RNP implementată

- ✓ **Faza  $\alpha$ :** neuronii de intrare generează impulsurile decalate temporal care codifică valorile de intrare
- ✓ **Faza  $\beta$ :** sinapsele ponderează aceste impulsuri incidente și transmit valorile post-sinaptice modulelor somei
- ✓ **Faza  $\gamma$ :** somele acumulează setul actual de intrări post-sinaptice, apoi rulează algoritmul cornului neuronal

Figura 5.37 Divizarea unui pas temporal în procesarea RNP implementate hardware

neuronul va emite un impuls

axonal (sau impuls post-sinaptic), apoi PM va fi resetat la o valoare mai mică decât valoarea de repaus, plasând neuronul într-o stare specială, numită de hiperpolarizare. În această stare neuronul nu va fi capabil să emită un nou impuls axonal, revenirea la starea de repaus făcându-se într-un număr predefinit de pași temporali.

Dacă nu se recepționează nici o valoare post-sinaptică diferită de zero într-un anumit pas temporal, atunci valoarea PM se va diminua cu valoarea de drenare (starea  $\lambda_k$  conform notației din Figura 5.36), realizând o descreștere liniară a acesteia, aproximând astfel funcționarea modelului neuronal de bază numit *leaky*

*integrate-and-fire*, care prevede o descreștere exponențială.

Implementarea eficientă acestui modul este una din provocările acestei realizări hardware. Însurarea a 32 de variabile reprezentate pe 8 biți precum și implementarea celorlalte funcționalități impuse de modelul neuronal aplicat – care presupun multiple comparații ale valorilor unor registre și ajustarea corespunzătoare a valorilor acestora – realizată într-o variantă complet paralelă, presupune un cost foarte ridicat în ceea ce privește utilizarea resurselor reconfigurabile disponibile în circuitul FPGA. Dacă am fi propus deci, o implementare complet paralelă a acestei RNP, am fi ajuns la nevoia de a restrânge mult dimensiunile acesteia, pentru a nu depăși limitele capacității platformei hardware utilizate, ceea ce ar fi rezultat o scădere sigură a performanțelor și a preciziei de clasificare. În decursul studiului teoretic realizat anterior implementării parțial serializate a acestei aplicații, am prevăzut utilizarea ca platformă un sistem de dezvoltare FPGA cu un circuit Xilinx Virtex4FX12 FPGA, deoarece acesta conține un microprocesor încorporat pe 32 de biți, de tip hard-core (care dispune de arie reconfigurabilă proprie în cadrul chip-ului FPGA), numit PowerPC. Acest PowerPC ar fi putut fi o unitate ce putea procesa cu ușurință toate calculele necesare funcționării tuturor modulelor din proiect. Totuși, am considerat, că nu este justificat îndeajuns a urma această cale de implementare, deoarece PowerPC este mult prea complex și nu ar fi fost exploatat într-o măsură care să contrabalanseze faptul că această metodă ar fi indus o serializare mai profundă a execuției, deci și un timp de execuție mai mare.

#### Notații ale variabilelor și stărilor somei implementate

- $\theta_k$  - Potențialul de prag al somei  $k$
- $\mu_k$  - Potențialul de membrană al somei  $k$
- $\gamma'_k$  - Soma  $k$  este activată
- $\delta_k$  - Soma  $k$  este în stare de hiperpolarizare
- $\lambda_k$  - Soma  $k$  este în stare de drenare curent

Figura 5.36 Notații în descrierea funcționării somei



Cu toate acestea, am găsit o soluție mai fezabilă pentru serializarea parțială a RNP, prin introducerea utilizării microcontrolerului soft-core Xilinx PicoBlaze, pentru implementarea modulelor de somă. Acest procesor pe 8 biți este evident mult mai simplu decât PowerPC, dar are și o amprentă digitală mult mai compactă, utilizând puține resurse reconfigurabile, precum a fost prezentat mai devreme în acest capitol, la punctul 5.5.1. Totuși, pentru a nu pierde în domeniul performanțelor, sistemul dezvoltat trebuie să conțină mai multe astfel de controlere – pentru fiecare somă în parte –, devenind astfel un *sistem de tip multi-core*. Pentru a compensa viteza de execuție mai lentă față de componentele RNP implementate în mod paralel, am conectat procesoarele PicoBlaze încorporate la un semnal de tact mai rapid, obținut utilizând circuitele DCM (Digital Clock Manager) din FPGA.

Implementând prin software rulat (și scris în limbaj de asamblare specific) de procesoarele PicoBlaze – vezi diagrama de execuție în Figura 5.38 – funcționalitatea modulelor somă au fost eliberate importante resurse reconfigurabile, făcând posibilă extinderea RNP pentru a implementa și alte aplicații, cum ar fi testul benchmark de clasificare a setului de date Fischer IRIS, prezentat în subcapitolul 5.5.3. Cu această modificare ar fi posibile și implementații de alte aplicații, de exemplu controlul inteligent și adaptive al unui braț de robot sau a unui robot mobil.

### 5.5.2.6. Dinamismul funcționării modului somă – Rezultate experimentale

Extrem de important și crucial de rezolvat în cadrul implementării acestei aplicații a fost și rezolvarea problemei de temporizare și sincronizare a diferitelor componente ce execută diferiți algoritmi ai RNP (elemente construite din resurse reconfigurabile dedicate respectiv procesoare încorporate). Pentru a putea soluționa această dilemă, a fost introdus un al patrulea procesor PicoBlaze (notat cu PB-UMC, fiind componenta centrală a Unității de Monitorizare și Control) care va genera semnalele de tact celorlalte componente ale sistemului și va sincroniza funcționarea celor trei module de somă implementate tot prin procesoare PicoBlaze. Pe lângă rularea algoritmului somei, procesoarele PicoBlaze-somă transmit prin porturi configurate în acest sens către PB-UMC parametrii principali ai modelului somei (de ex. potențialul de membrană). PB-UMC va citi periodic aceste valori, precum și valorile ponderilor de la cele 96 de sinapse și le va transmite mai departe printr-un circuit USART atașat unui calculator personal la care este conectat. Pe acest calculator datele vor fi stocate și apoi afișate grafic cu ajutorul unui program Matlab.

Pentru a putea trage concluzii asupra funcționării corecte sau pentru a putea efectua corecții implementării, este necesară această vizualizare grafică, precum și o prelucrare numerică prealabilă a fluxului de date ce sosește de la circuitul neuronal implementat în FPGA spre PC.

Diagramele din Figura 5.39 și Figura 5.40 prezintă două seturi de măsurători asupra evoluției potențialelor de membrană în două dintre cele trei some implementate în aplicația curentă, în decursul unui proces de antrenare. În primul experiment (Figura 5.39), am implementat și conexiuni laterale inhibitoare de la o somă la celelalte, pentru a preveni activarea simultană a mai multor some pentru aceleași valori de intrare, diferențiind astfel componentele individuale de frecvență ce trebuie învățate. Eliminând aceste conexiuni inhibitoare laterale, somele se vor comporta diferit (Figura 5.40). Este de notat în figurile de mai jos, că activarea unuia dintre neuronii de ieșire rezultă în resetarea PM al celorlalți neuroni la valoarea de repaus.

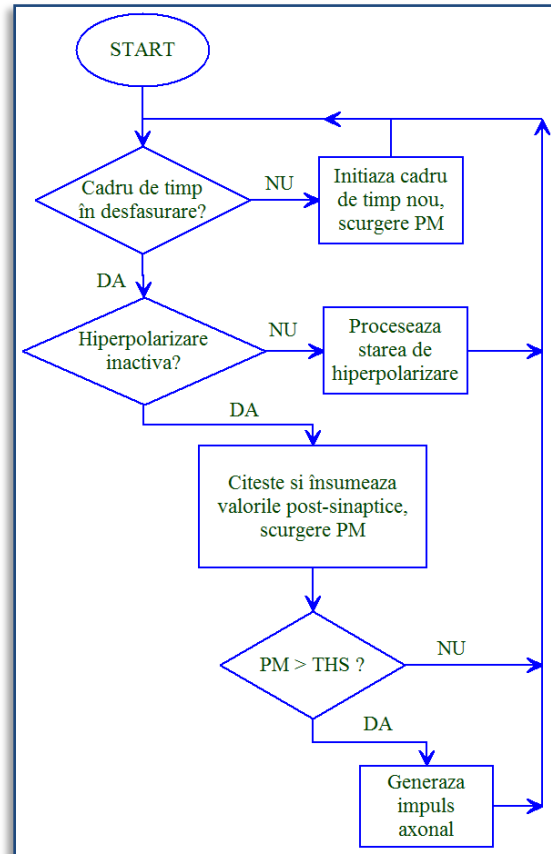


Figura 5.38 Diagrama de execuție a programului în limbaj de asamblare rulat de procesoarele PicoBlaze care implementează modulele de somă

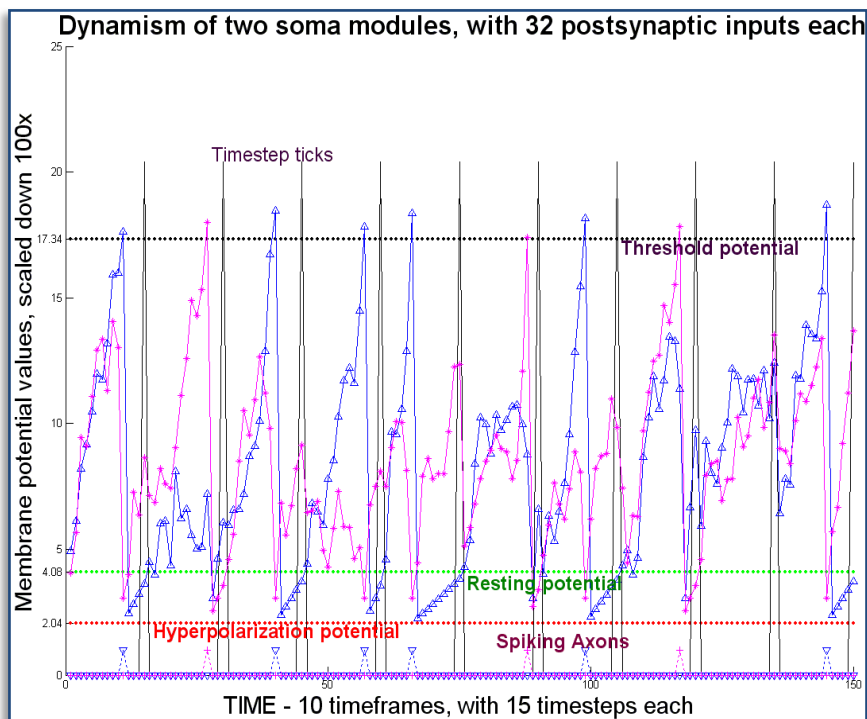


Figura 5.39 Dinamismul somei cu inhibare laterală

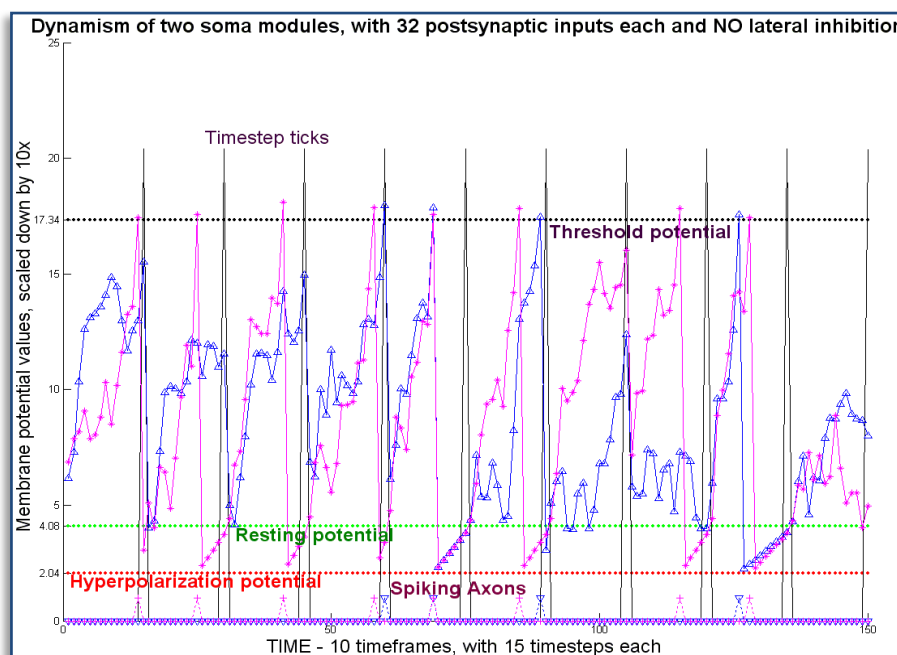


Figura 5.40 Dinamismul somei fără inhibare laterală

După activare, somele vor intra în stare de hiperpolarizare, conform modelului neuronal utilizat, pentru a nu putea emite un nou impuls axonal în același cadru de timp. Pentru a garanta acest lucru, PM va crește liniar, cu valori mici, pentru a ajunge doar la valoarea de repaus până la sfârșitul cadrului temporal.

### 5.5.2.7. Descrierea sistemului test-bench dezvoltat. Prepararea setului de date de antrenare

După cum am afirmat și la formularea problemei, aplicația curentă are ca scop identificarea a

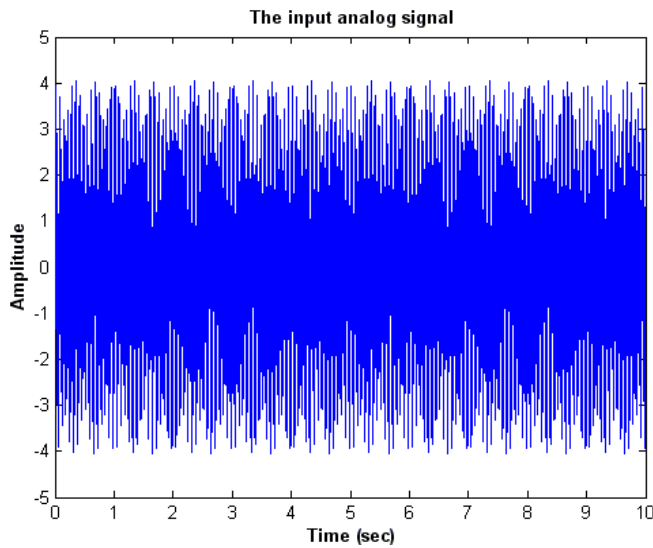


Figura 5.43 Semnalul analog preluat

până la două componente de frecvență dintr-un semnal analogic zgomotos. Un astfel de semnal analog, perturbat de zgomot Gaussian este prezentat în Figura 5.43 și a fost utilizat în experimentele rulate.

Pentru a putea prezenta la intrările RNP perechi de variabile, reprezentând eşantioane din semnalul vizat (valorile fiind: frecvența și puterea componente) a fost nevoie de calcule prealabile asupra acestui semnal. Primul pas a fost efectuarea unei transformări Fourier cu algoritmul FFT, al cărei rezultat se poate vedea în Figura 5.42, urmând ca Figura 5.41 să prezinte secțiunea de interes a acesteia. Valorile din Figura 5.41 reprezintă de fapt spațiul valorilor de intrare

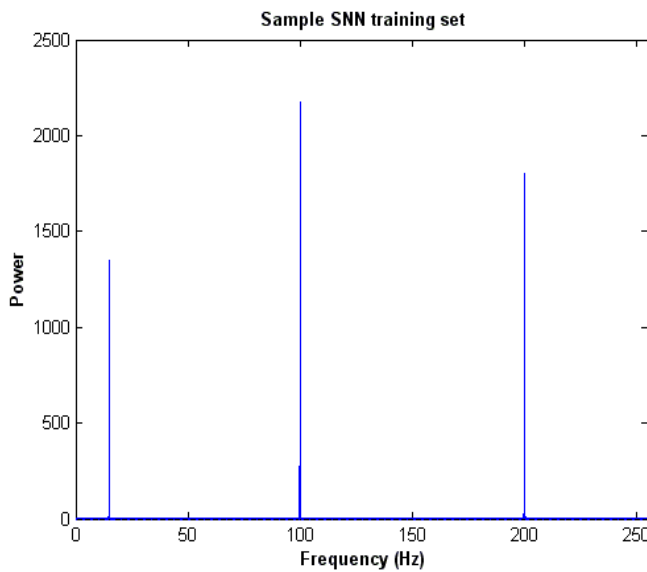


Figura 5.41 Spațiul valorilor de intrare utilizate pentru antrenarea RNP implementate în FPGA

utilizate pentru antrenarea RNP implementate în FPGA.

Graficul din Figura 5.44 arată valorile de antrenare generate în mod aleator, ele fiind concentrare în vecinătatea punctelor de reper (target focus points) care sunt de fapt valorile prescrise în timpul învățării neuronilor de ieșire ai RNP.

Pentru fiecare punct din acest spațiu, valorile coordonatelor reprezintă variabilele de intrare ale RNP,acompaniate de valoarea prescrisă fiecărui neuron de ieșire pentru acel punct, rezultând un algoritm de antrenare

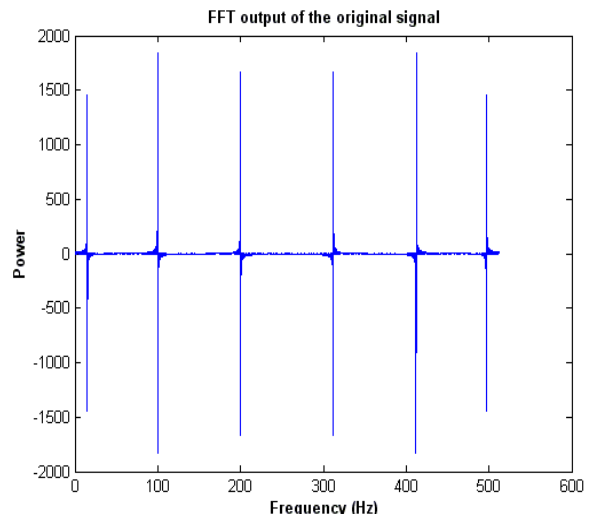


Figura 5.42 Transformata FFT a semnalului analog din Figura 5.43

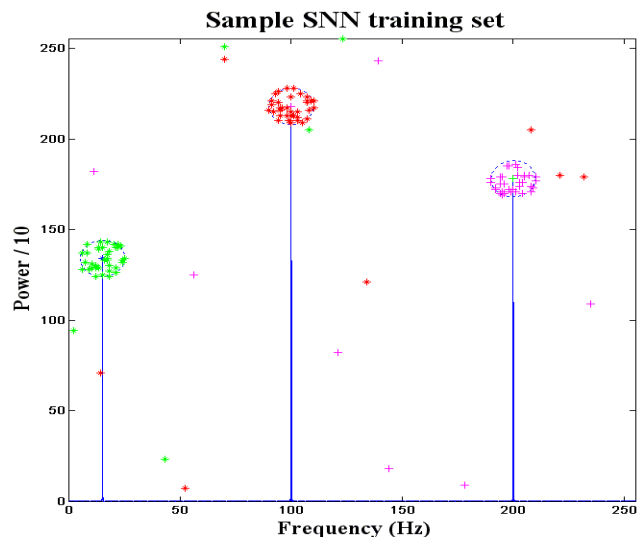


Figura 5.44 Valori de antrenare concentrare în jurul punctelor de focus prescrise, și altele aleator alese ca zgomot

supervizat. În Figura 5.41 se mai pot observa și puncte ce apar în poziții aleatoare în afara vecinătăților punctelor de reper. Acestea au fost introduse pentru a evita blocajul algoritmului în minime locale și a garanta convergența acestuia. Toate valorile setului de antrenare sunt prezentate sistemului neuronal în ordine aleatoare.

Rețeaua neuronală pulsativă hardware a învățat separarea acestor componente de frecvență, activând după antrenare același neuron de ieșire pentru punctele aparținând aceleiași componente. În mod evident, acest lucru se întâmplă cu o anumită histereză, definind o vecinătate detectată a componentelor pentru care s-a efectuat antrenarea.

### 5.5.2.8. Testarea sistemului neuronal implementat

După cum s-a menționat și la subpunctul anterior, întreaga RNP hardware este controlată de cel de-al patrulea procesor PicoBlaze (PM-UMC), componenta principală a Unității de Monitorizare și Control din Figura 5.45 de mai jos. Acest modul controlează funcționarea pas cu pas a sistemului neuronal implementat, așadar și a procesului de antrenare.

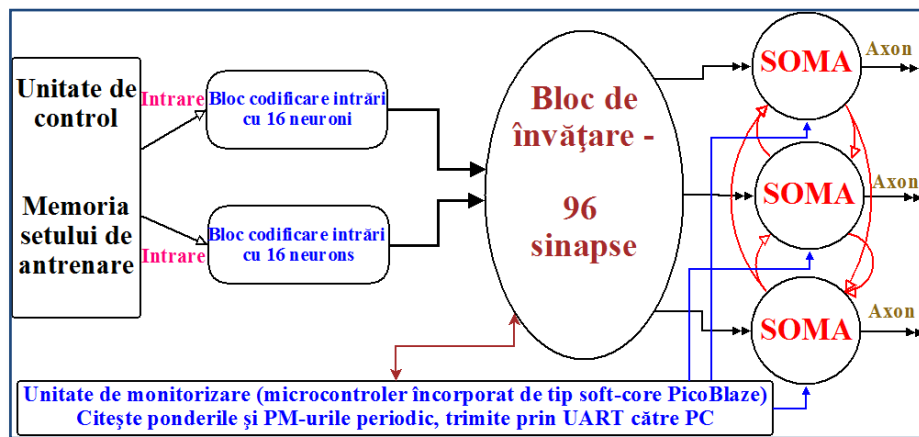


Figura 5.45 Diagrama bloc a sistemului testbench realizat

Diagrama de execuție a algoritmului rulat de PB-UMC, program scris în limbaj de asamblare specific Xilinx PicoBlaze, se prezintă în Figura 5.46 iar fazele și pașii de operare a RNP sunt prezentate pe scurt în Tabelul 13.

Tabelul 13 Pașii de procesare a calculelor RNP

Cicluri de tact/ faze de execuție	a	b	c	d	e	
Nr. ciclu de tact	Control Antrenare	Intrări valide	Tact Bloc Intrări	Tact Sinapse	Tact Some (PBlaze)	Transmisie UART
1	front crescător	NU	fără tact	fără tact	HOLD	
2	front crescător	NU	fără tact	fără tact	HOLD	
3	front crescător	DA	fără tact	fără tact	HOLD	
4	front crescător	DA	fără tact	fără tact	HOLD	valori de intrare
5	fără tact	DA	front crescător	fără tact	HOLD	
6	fără tact	DA	front crescător	front crescător	HOLD	ponderi
7	fără tact	DA	fără tact	fără tact	RULEAZĂ	potențiale de membrană
if TS_nr<16 goto 5					if TS_nr=15 then →	ponderi, valori axonale
if Epoch<129 goto 1						

În primul ciclu de antrenare (epoch) se prezintă pe rând rețelei neuronale fiecare element al setului de antrenare, generând numărul necesar de fronturi ale semnalului de tact către modulul de generare a intrărilor, pentru a permite acestuia să citească aceste valori din memoria BRAM proprie în care sunt stocate. Apoi, Unitatea de Control va livra primul front de tact modulului Blocului Intrărilor, dintre cele 16 ale unui cadru temporal (TS în Tabelul 13) care va începe astfel codificarea valorilor de intrare în impulsuri decalate temporal. La ciclul de tact numărul 7, unele au pe intrările lor primul set

de date post-sinaptice, deci pot fi activate, pentru a rula un ciclu al algoritmului implementat de modulele PicoBlaze.

La momentele corespunzătoare, când valorile respective sunt valide, UMC colecta datele de la componentele sistemului și va trimite pe portul serial cei mai importanți parametri (ponderi, potențiale de membrană, valori axonale) ai RNP spre calculator pentru procesare și vizualizare grafică (Figura 5.39, Figura 5.40). Implementarea acestei rețele neuronale pulsative utilizând procesoare încorporate, s-a dovedit a fi mai eficientă în ceea ce privește necesarul de resurse hardware, decât alte aplicații, implementând RNP cu codarea ratei impulsurilor, realizate complet paralel. Acest lucru a justificat continuarea cercetărilor în această direcție – cu procesare parțial serializată și stocarea variabilelor de stare în memorii BRAM – metodă prezentă și în literatura de specialitate (Maass, Natschlagel, & Markram, 2002) (Schrauwen & Van Campenhout, 2006) (Bakó & Brassai, 2006) (Floreano, Dürr, & Mattiussi, 2008).

Tabelul 14 trece în revistă utilizarea de resurse reconfigurabile FPGA în proiectul implementat, detaliat pentru componentele principale ale acestuia.

Se poate nota, că pentru a optimiza utilizarea resurselor de tip BRAM s-a folosit o component de memorie cu porturi duale, care stochează programul rulat de către procesoarele PicoBlaze ce implementează somele, ea fiind accesată de două astfel de modul simultan, ele rulând același program dar cu date de intrare diferite. Deoarece și această aplicație este de fapt una de clasificare, se poate compara cu aplicații de test clasice (Fischer's Iris, Wisconsin Breast Cancer). Implementarea unui astfel de test se va prezenta în subcapitolul următor, numărul 5.5.3.

Tabelul 14 Utilizarea de resurse reconfigurabile în proiectul implementat

Nume modul	Slices	Slice Reg.	LUTs	LUT RAM	BRAM
[ - ] Spiking NN project	9408	6886	12755	0/281	5/28
Control & Monitoring Unit	54/54	41/41	34/34	0/0	0/1
[ - ] Inst. Spiking ANN	602/8592	43/6845	1106/12323	0/281	0/5
PicoBlaze shared mem. bl.	78	88	49	0	0
Input block	761	693	1277	0	2
Master PicoBlaze	126	69	199	68	1
Synapse Unit (one of 96)	67/6432	59/5664	95/9120	0/0	0/0
SOMA 1 (PicoBlaze uC)	99	67	169	68	2
SOMA 2 (PicoBlaze uC)	97	67	169	68	
SOMA 3 (PicoBlaze uC)	97	67	169	68	

O altă caracteristică valoroasă de notat a RNP implementate în FPGA este faptul că funcționează în timp real, cu toate că procesarea unor subsansamble este serializată parțial prin utilizarea procesoarelor soft-core încorporate. Aceste procesoare PicoBlaze sunt animate de un semnal de tact de 50MHz și execută o instrucțiune în două cicluri de tact. Algoritmul somelor este implementat printr-un program de asamblare de aproximativ 200 de instrucțiuni iar celelalte componente ale RNP funcționează în paralel. Efectuând un calcul simplu, obținem un timp de calcul/ciclu de antrenare de 64μs, și un timp de învățare de 16ms. După terminarea fazei de învățare, circuitul neuronal va putea detecta o component de frecvență într-un timp mai mic de 1μs.

În concluzie, se poate afirma, că această abordare a implementărilor rețelelor neuronale neuromorfe are un potențial foarte promițător în domeniul aplicațiilor în timp real.

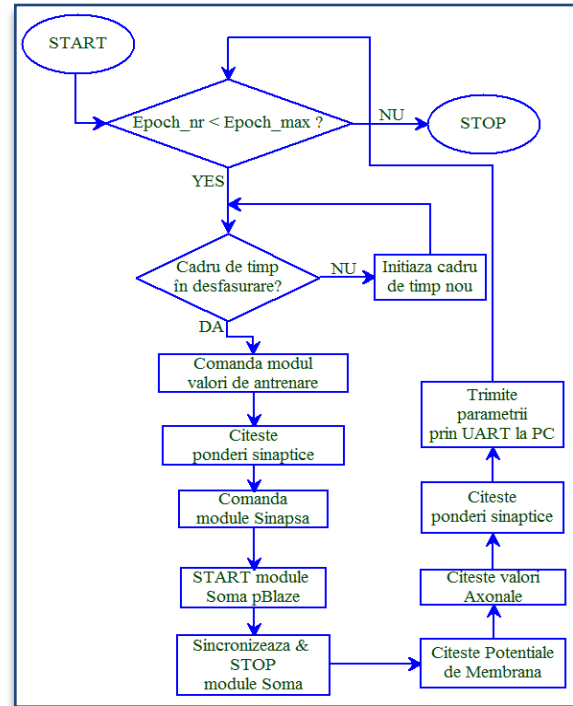


Figura 5.46 Diagrama de execuție a programului în limbaj de asamblare rulat de procesorul PicoBlaze care implementează Unitatea de Monitorizare și Control

### 5.5.3. Implementare test benchmark: Clasificarea setului de date Fisher IRIS

Scopul acestei aplicații este de a testa acest mod de a implementa Rețele Neuronale Pulsative, și a evalua potențialul pentru aplicații similare și puterea de calcul a acestora. Testul reper (benchmark) de clasificare a setului de date Fischer IRIS este probabil cel mai cunoscut în domeniu, fiind astfel o metodă bună de verificare.

#### 5.5.3.1. Descrierea setului de date. Formularea problemei de clasificare

Articolul lui Fisher (Fisher, 1936) este unul de referință în domeniu, și este citat în mod frecvent chiar și în prezent. Acest set de date conține 3 clase a câte 50 de elemente fiecare, unde fiecare clasă se referă la un tip de plantă Iris. Una din clase poate fi separată liniar de celelalte două; celelalte două clase NU pot fi separate liniar între ele. Lungimea și lățimea separelor și petalelor pentru cele trei specii nord-americe de iris. Datele referitoare la *iris setosa canadensis* și *iris versicolor* au fost folosite de către R.A. Fisher pentru a ilustra analiza discriminării liniare. Datele referitoare la *iris virginica* au fost adăugate ca o extensie. Lungimea și lățimea separelor și petalelor sunt doar coordonate convenabile rezultate datorită faptului ca sunt ușor de măsurat.

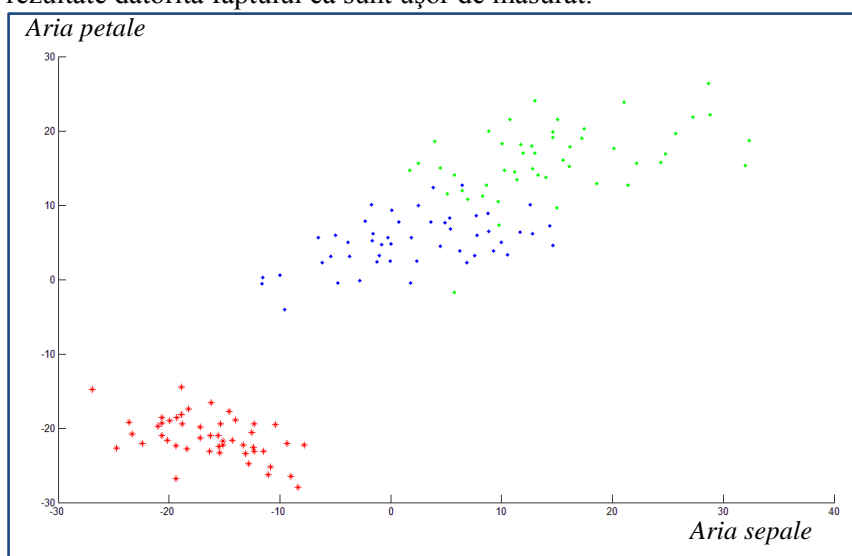


Figura 5.47 Rețea de date Fisher IRIS cu ajutorul proiecției Sammon (vezi subcap.5.6.3.7)  $4D \rightarrow 2D$ .

Valorile afișate pe coordonate sunt de fapt valori calculate de acest algoritm de reprezentare, dar se pot considera ca valorile ariilor celor două tipuri de petale ale florilor IRIS.

Roșu => IRIS-setosa

Albastru => IRIS-versicolor

Verde => IRIS-virginica

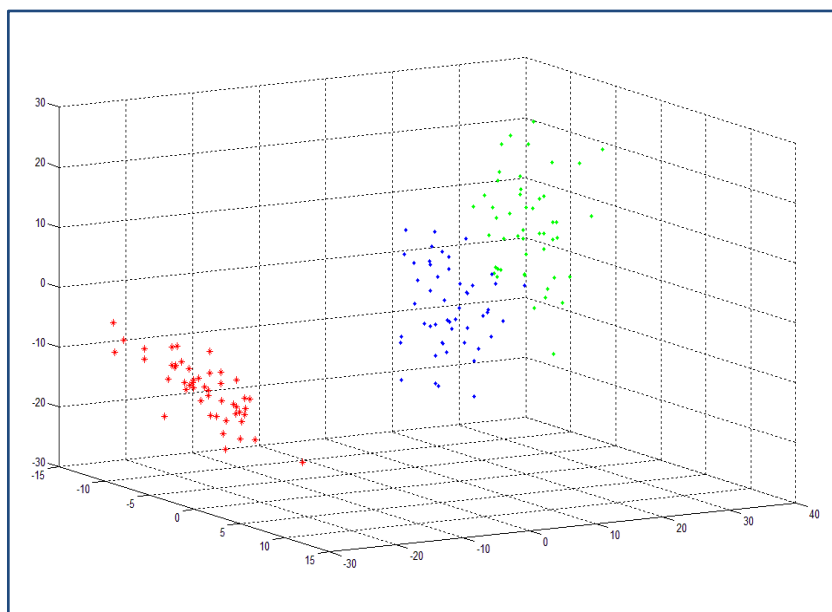


Figura 5.48 Rețea de date Fisher IRIS cu ajutorul proiecției Sammon (vezi subcap.5.6.3.7)  $4D \rightarrow 3D$ .

Roșu => IRIS-setosa

Albastru => IRIS-versicolor

Verde => IRIS-virginica



Folosind lungimea sealei \* lățimea sealei și lungimea petalei \* lățimea petalei ca o aproximare neprelucrată a suprafeței acoperite de seale și petale, obținem o pereche de coordonate care ar putea fi mai relevante pentru botanică. Folosind aceste coordonate, toate cele trei specii din acest set de date pot fi separate, având doar trei clasificări greșite. *Iris setosa* și *iris versicolor* sunt separate perfect. Numărul de atribute: 4 atribute numerice, numerotare și clasa. Distribuția claselor: 33.3% pentru fiecare din cele 3 clase.

### 5.5.3.2. Rețeaua neuronală pulsativă implementată pe FPGA

S-a implementat o RNP cu patru intrări și trei ieșiri, după cum arată și schema bloc din Figura 5.49. Această rețea trebuie să clasifice corect elementele bazei de date Fisher IRIS, care sunt valori cu o zecimală cuprinse între 0 și 8, dar au fost scalate între 0-80 pentru a putea fi

#### Parametrii rețelei neuronale implementate:

- Numărul intrărilor:  $n = 4$
- Nr. neuronilor de intrare/intrare:  $n_{intr}=13$
- Nr. neuronilor de intrare activi/intrare:  $n_{act} = 4$
- Numărul pașilor dintr-un cadru temporal  $T=16$

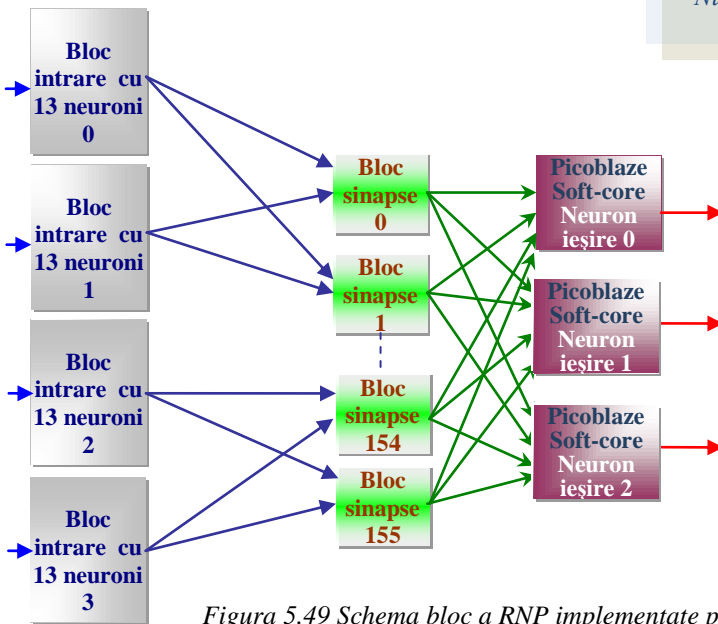


Figura 5.49 Schema bloc a RNP implementate pentru clasificarea setului de date Fisher IRIS

reprezentate ca valori întregi în hardware-ul digital. Antrenarea rețelei s-a efectuat în așa fel, încât fiecare neuron ieșire să se activeze atunci, când valorile de intrare identifică un element care aparține clasei asociate neuronului respectiv. Arhitectura RNP este foarte similară cu ce a aplicației anterioare, dar având un număr diferit de neuroni de intrare pentru codificarea celor patru intrări, neuronii de ieșire fiind implementați și în acest caz cu trei procesoare încorporate de tip soft-core Xilinx PicoBlaze.

### 5.5.3.3. Implementarea neuronilor de intrare – codificarea variabilelor de intrare în impulsuri decalate temporal

Deoarece plaja valorilor din baza de date vizată (IRIS) definește un spațiu al parametrilor de intrare de dimensiuni mai

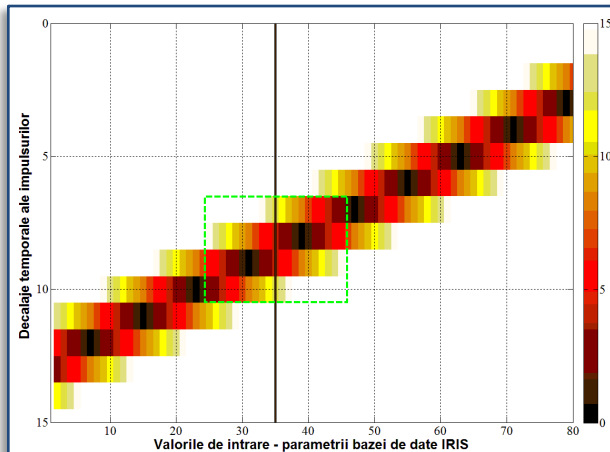


Figura 5.51 Domeniile receptive ale neuronilor de intrare

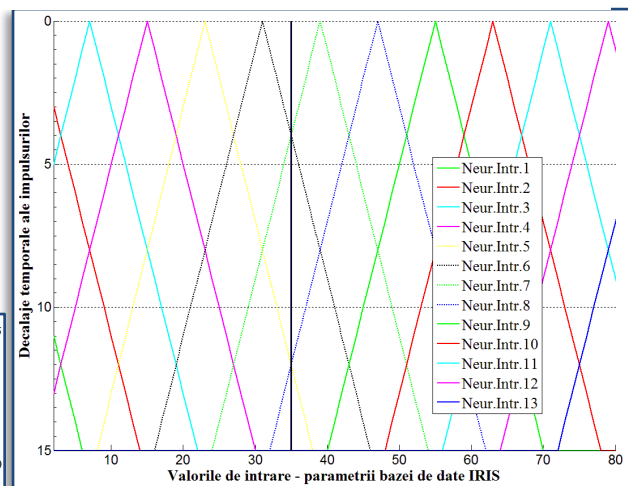


Figura 5.50 Funcțiile triunghiulare asociate neuronilor de intrare și decalajele temporale

reduse decât la aplicația de la punctul 5.5.2, s-a scăzut numărul ( $n=13$ ) neuronilor de intrare, dar s-a menținut constantă ( $n_{act}=4$ ) numărul celor activate într-un cadru de timp. Acoperirea acestui spațiu cu funcții de activare triunghiulare este

reprezentată în Figura 5.50 iar domeniile receptive corespunzătoare acestora în Figura 5.51. Se poate observa în ambele figuri, că valoarea de intrare 35 luată ca exemplu va activa patru neuroni de intrare.

După cum a fost prezentat și la subpunctul 3.1.3.1 elementele BlockRAM ale circuitelor FPGA pot fi configurate în mai multe feluri, variind numărul biților de pe magistrala de adrese respectiv de pe cea de date. Pentru a stoca valorile de decalaj temporal am ales componenta cu numele RAMB16\_S36\_S36, care este un modul BRAM cu port dual și o capacitate de 512x32 biți. Arhitectura cu port dual permite ca două locații de memorie să fie accesate simultan. Deoarece valorile stocate sunt inițializate odată cu programarea circuitului FPGA, acest modul va fi utilizat ca și cum ar fi de fapt o memorie ROM, efectuând numai operații de citire, evitând astfel eventualele conflicte de adresare la scriere. Conectând câte două astfel de module în paralel (4 în total), devine posibil să se obțină în aceeași perioadă de tact cei  $4*4*13=208$  biți reprezentând decalajele temporale pentru cele două intrări ale RNP, utilizând numai o parte a capacității unei locații astfel adresabile de  $4*2*32=256$  biți.

#### 5.5.3.4. Modulul de implementare a sinapselor

Modulul ce implementează sinapsa unui neuron este responsabilă pentru amplificarea sau atenuarea impulsurilor pre-sinaptice emise de către neuronul de intrare la care este conectat. Astfel se generează o valoare post-sinaptică ponderată ce va fi transmisă spre soma neuronului de ieșire la care aparține sinapsa respectivă.

Acest modul este de asemenea cel, care implementează regula nouă de învățare supervizată care este o variantă adaptată a metodei Hebb, bazată pe reguli de decalaje temporale, după cum a fost prezentat în secțiunea 5.5.2.4. Sinapsele vor testa amprenta temporală (numărul pasului temporal, notat cu  $TS$  în Tabelul 12) a fiecărui impuls pre-sinaptic și vor ajusta valorile ponderilor stocate în concordanță cu regulile de învățare prezentate în Tabelul 12, luând în calcul și valorile prescrise ale stării neuronilor de ieșire (valoare axonală) pentru valori variabilelor de intrare actuale.

#### 5.5.3.5. Dinamismul funcționării modului somă – Rezultate experimentale

Utilizarea PicoBlaze pentru implementarea modulelor somă a fost cea mai utilă alegere și în acest caz. Aceste procesoare rulează același algoritm ca și în cazul precedent (Figura 5.38), dar cu parametrii somei adaptate la arhitectura schimbată a RNP. Astfel, numărul de 13 a neuronilor intrare într-un bloc aparținând unei dintre intrări rezultă valori pentru maximul PM, potențialul de prag, valoarea de hiperpolarizare respectiv valoarea de

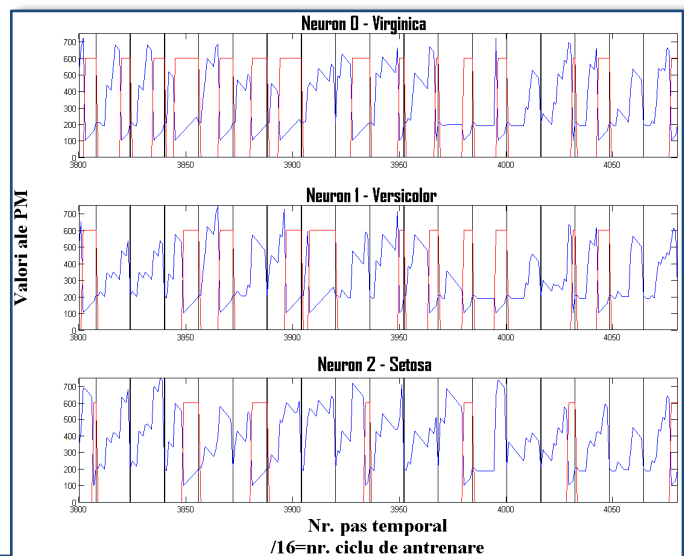


Figura 5.52 Măsurători în timpul antrenării RNP (albastru-variația PM celor trei some, roșu valoarea axonală a acestora)

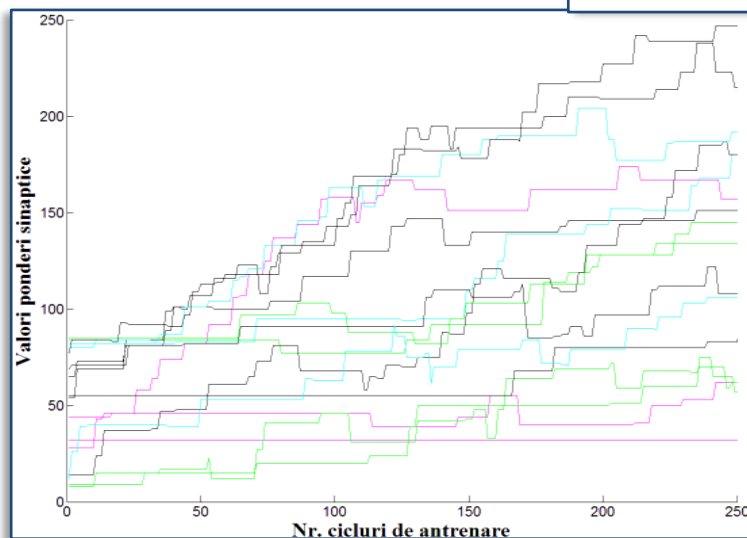


Figura 5.53 Variația ponderilor în timpul antrenării (primele 15)

scurgere.

Se poate nota în Figura 5.53 – în care s-a reprezentat numai un număr redus dintre cele 156 de ponderi sinaptice ale rețelei neuronale pulsative hardware – că unele valori rămân neschimbate. Acest fenomen apare datorită faptului, că procesarea sinapselor depinde de prezența impulsurilor incidente, iar dacă neuronul de intrare la care este



conectată sinapsa respectivă nu s-a activat pentru valoarea de intrare actuală, sinapsa va fi inactivă.

### 5.5.3.6. Descrierea sistemului testbench dezvoltat. Prepararea setului de date de antrenare

După cum am afirmat și la formularea problemei, aplicația curentă are ca scop identificarea clasei de floare IRIS pe baza a patru proprietăți prezentate ca intrări ale rețele neuronale hardware.

S-a menționat și la descrierea aplicației anterioare, că întreaga RNP hardware este controlată de cel de-al patrulea procesor PicoBlaze (PM-UMC) încorporat, componenta principală a Unității de Monitorizare și Control din Figura 5.54 de mai jos. Acest modul controlează funcționarea pas cu pas a sistemului neuronal implementat, așadar și a procesului de antrenare.

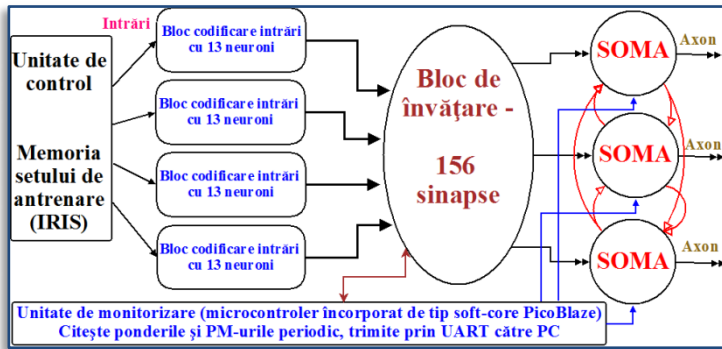


Figura 5.54 Diagrama bloc a sistemului testbench realizat

Tabelul 15 Utilizarea resurselor FPGA a proiectului FPGA IRIS dataset

Device Utilization Summary (XC3S1200E)			
Logic Utilization	Used	Avail.	Utiliz.
Nr. of Slices	6941	8672	80%
Nr. of Slice Flip Flops	5785	17344	33%
Nr. of 4 input LUTs	12806	17344	73%
Nr. of BRAMs	7	28	25%
Nr. of GCLKs	4	24	16%

Din setul de date Fisher IRIS a fost separat o mulțime de antrenare, formată din 60% din toate cele trei clase, alese în mod aleatoriu. Aceste valori de antrenare se prezintă pe rând rețelei neuronale, generând numărul necesar de fronturi ale semnalului de tact către modulul de generare a intrărilor, pentru a permite acestuia să citească aceste valori din memoria BRAM proprie în care sunt stocate. Apoi, Unitatea de Control va livra primul fronturi de tact modulului Blocului Intrărilor, care va începe astfel codificarea valorilor de intrare în impulsuri decalate temporal. La momentele corespunzătoare, când valorile respective sunt valide, UMC colecta datele de la componentele sistemului și va trimite pe portul serial cei mai importanți parametrii (ponderi, potențiale de membrană, valori axonale) ai RNP spre calculator pentru procesare și vizualizare grafică (Figura 5.52 și Figura 5.53).

### 5.5.4. Sistemul de dezvoltare utilizat pentru realizarea proiectelor cu procesoare încorporate

Digilent Nexys2 este o platformă de dezvoltare completă, gata de utilizare, bazată pe un circuit FPGA de tip Xilinx Spartan 3E. Portul USB2 de mare viteză integrat, cei 16Mbytes de memorie RAM și ROM, și mai multe dispozitive și porturi de I/O, fac din Nexys2 o platformă ideală pentru sisteme digitale de toate tipurile, inclusiv sisteme cu procesor integrat bazate pe MicroBlaze de la Xilinx. Portul USB2 asigură plăcii alimentarea și o interfață de programare, astfel încât placa Nexys2 poate fi folosită împreună cu un computer de tip notebook pentru a crea o stație de proiectare cu adevărat portabilă. Nexys2 aduce tehnologii de vârf pe o platformă care poate fi folosită de oricine pentru a câștiga experiență în proiectarea digitală. Poate găzdui nenumărate sisteme digitale bazate pe FPGA, și sistemele proiectate pot ușor crește dincolo de dimensiunile plăcii prin folosirea oricăruia sau a tuturor celor 5 conectoare de expansiune. Patru conectoare de module periferice cu 4 pini pot găzdui până la opt module periferice cu cost redus pentru a adăuga extensii ca control motor, conversie A/D sau D/A, circuite audio, și o gamă de interfețe de senzor și elemente de execuție. Toate semnalele accesibile utilizatorului pe placa Nexys2

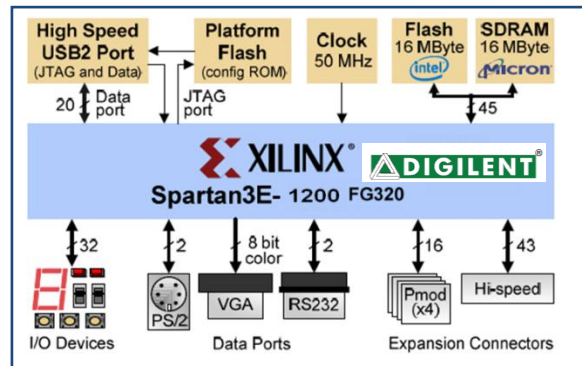


Figura 5.55 Sistemul de dezvoltare Digilent Nexys2 (XC3S1200E)

sunt protejate împotriva perturbărilor electrostatice și a scurtcircuitelor, asigurând un timp îndelungat de funcționare în orice mediu. Elementele principale ale Nexys2 sunt următoarele:

- circuit FPGA de tip Xilinx Spartan 3E cu 500Kgate sau 1200Kgate
- 16MB de PSDRAM Micron și 16MB ROM Intel StrataFlash
- Xilinx Platform Flash pentru configurații FPGA non-volatile
- oscilator de 50MHz plus socket pentru încă un oscilator60 I/O-uri FPGA directate la conectori de expansiune,
- 8 LED-uri, afișaj de 7 segmente și 4 caractere, 4 butoane, 8

## 5.6. Realizări de RNP cu microcontroler pe 32 de biți, încorporat în circuitele FPGA

În acest subcapitol se vor prezenta două aplicații de clasificare care necesită o putere de calcul mai ridicată, ele fiind recunoaștere de caractere respectiv realizarea testului benchmark de clasificare a setului de date Wisconsin breast cancer.

### 5.6.1.1. Microcontrolerul Xilinx MicroBlaze

Nucleul încorporat de tip soft-core MicroBlaze este un procesor cu un set redus de instrucțiuni, optimizat pentru implementările bazate pe circuite FPGA realizate de Xilinx. Figura 5.56 arată o diagramă bloc a nucleului MicroBlaze. Nucleul încorporat de tip soft-core MicroBlaze are următoarele proprietăți:

- 32 de regiștrii pe 32 biți cu destinație generală
- Instrucțiuni pe 32 biți cu trei operanzi și două moduri de adresare
- Magistrale separate pe 32 biți pentru instrucțiuni și pentru date conform specificației OPB (On-Chip Peripheral Bus) de la IBM
- Magistrale separate pe 32 biți pentru instrucțiuni și pentru date cu conexiune directă către blocul de memorie RAM încorporat pe circuitul reprogramabil accesat prin

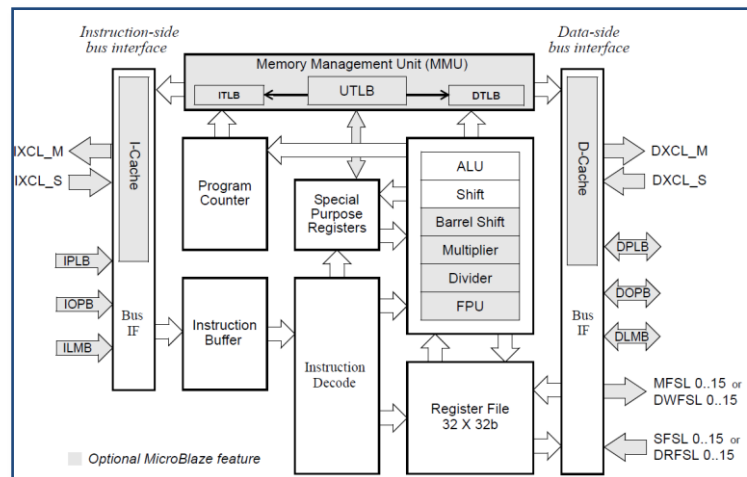


Figura 5.56 Structura logică a procesorului Xilinx MicroBlaze

- LMB (Local Memory Bus)
- Magistrală de adresare pe 32 biți
- Cache pentru instrucțiuni și pentru date
- Elemente logice pentru depanare hardware
- Suportă FSL (Fast Simplex Link)
- Multiplicator hardware (în dispozitivele din familia Virtex-II).

Nucleul MicroBlaze implementează o arhitectură de tip Harvard. Asta înseamnă că are interfețe de magistrală separate pentru acces de date și de instrucțiuni. Fiecare unitate de interfață magistrală se împarte la rândurile ei în Local Memory Bus (LMB) și On-Chip Peripheral Bus (OPB) de la IBM. LMB asigură acces într-un singur ciclu de tact la blocul de memorie RAM pe două porturi încorporat în circuitul reconfigurabil. Interfața OPB asigură conexiunea către perifericele și memoria încorporată sau alte periferice din afara circuitului reconfigurabil. Nucleul MicroBlaze dispune de asemenea de 8 interfețe de intrare și 8 interfețe de ieșire la magistralele de tip Fast Simplex Link (FSL). Magistralele FSB sunt canale de comunicație dedicate unidirecționale fără arbitrar. MicroBlaze utilizează adresare pe 32 de biți, harta memoriei adresabile fiind prezentată în Figura 5.57.

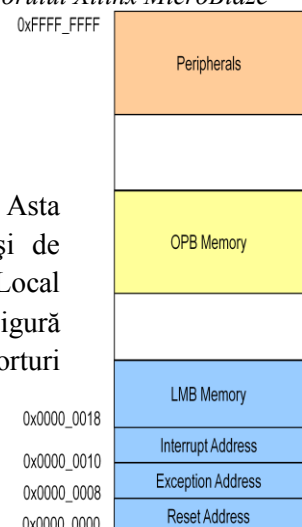


Figura 5.57 Harta memoriei procesorului MicroBlaze

## 5.6.2. Implementare aplicativă: Recunoaștere de caractere cu o rețea neuronală neuromorfă multistrat

### 5.6.2.1. Formularea problemei de recunoaștere de caractere

Scopul acestei aplicații este de a realiza un sistem neuronal bazat pe o RNP implementată într-un mediu mixt hardware-software pentru a evalua performanțele unui astfel de sistem.

În acest sens am construit o RNP cu două straturi care va fi descrisă în cele ce urmează. Caracterele ce vor fi învățate de RNP s-au considerat a fi reprezentate într-o matrice de 6x8, asemănător cu modul de afișare pe un afișaj alfanumeric (Figura 5.58). Fiind vorba de o rețea neuronală bazată pe

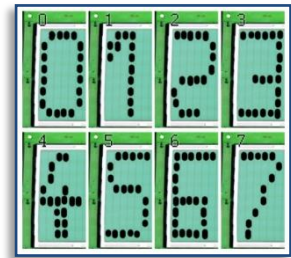


Figura 5.58 Forma caracterelor învățate

impulsuri, valorile de 0 și 1 logic, care reprezintă punctele inactive respectiv active pe afișajul LCD se vor transforma în impulsuri ce sosesc sau nu pe una dintre cele  $6 \times 8 = 48$  intrări ale acesteia.

Valorile prescrise pentru ieșirile rețelei neuronale depind de caracterele prezentate la intrări într-un anumit ciclu de învățare.

### 5.6.2.2. Dezvoltarea componentelor hardware ale sistemului încorporat

În cadrul acestei aplicații am utilizat numai module hardware predefinite ale mediului de dezvoltare Xilinx Platform Studio. Elementul principal al implementării este procesorul Xilinx MicroBlaze, la care am conectat mai multe module de memorie respectiv dispozitive periferice prin magistrala de 32 de biți Processor Local Bus (PLB).

În decursul implementării, s-a constatat, că pentru a acomoda programul scris în limbaj C, rulat pe acest procesor, care implementează de fapt această rețea neuronală pulsativă, capacitatea memoriilor interne de tip BRAM ale circuitului FPGA care sunt folosite ca memorie de instrucțiuni și de date (conectate la magistrala Local Memory Bus - LMB) ale MicroBlaze nu a fost suficientă. Pentru a soluționa această problemă, am conectat la proiectul încorporat un controler de memorie dinamică (MPMC CNTLR în Figura 5.59) cu scopul de a extinde capacitatea memoriei interne cu modulul DDR ce se găsește pe placa de dezvoltare utilizată.

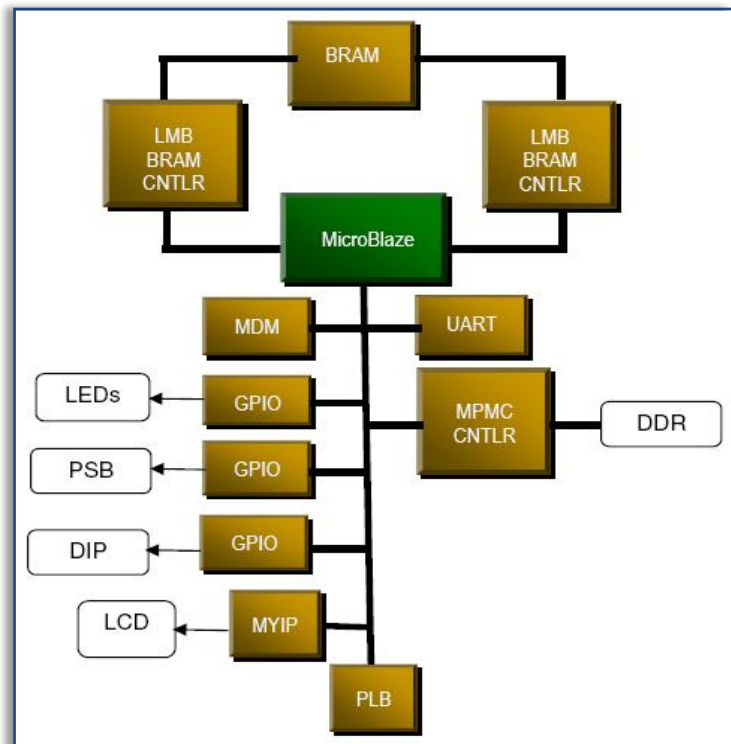
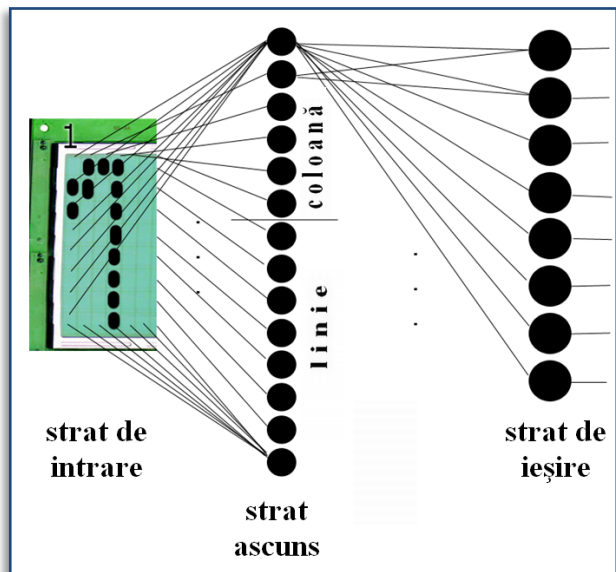


Figura 5.59 Diagrama bloc a sistemului încorporat care implementează RNP pentru recunoașterea de caractere

Dispozitivul standard de intrare/ieșire a sistemului dezvoltat a fost considerat a fi portul serial RS232, atașat proiectului încorporat din biblioteca de dispozitive predefinite a mediului EDK (UART în Figura 5.59). Controlul și setările funcționării sistemului implementat au fost realizate prin butoanele și comutatoarele plăcii de dezvoltare, conectate la MicroBlaze prin modulul standard de accesare a perifericelor generice (General Purpose I/O – GPIO în Figura 5.59). Un alt periferic – implementare proprie – este cel care comandă afișajul LCD al plăcii de dezvoltare (MYIP în Figura 5.59) care este utilizat pentru afișarea stărilor RNP din FPGA, respectiv pentru urmărirea procesului de învățare. Acest din urmă modul a fost implementat în limbaj VHDL.

### 5.6.2.3. Structura rețelei neuronale pulsative dezvoltate

Rețeaua neuronală dezvoltată are o structură cu trei straturi. Primul strat este cel care primește impulsurile de intrare reprezentând caracterele de învățat. Cel de-al doilea strat este cel care va fi antrenat să învețe anumite caracteristici ale acestor caractere, iar stratul de ieșire va conține neuroni pulsativi care se vor antrena să se activeze la apariția unor caractere cu un set de caracteristici bine definite. Astfel, șase dintre cei paisprezece neuroni pulsativi ai stratului ascuns (Figura 5.60) vor fi antrenați cu o regulă care le va acorda să învețe dacă majoritatea punctelor dintr-o coloană a matricii de intrare sunt active. Acești neuroni au câte opt intrări conform numărului de puncte dintr-o coloană a intrării. Restul de opt neuroni pulsativi din stratul ascuns vor fi antrenați să răspundă cu



impuls axonal atunci, când majoritatea celor șase intrări al fiecăruia (conform numărului de puncte dintr-o linie a matricii de intrare) sunt active.

Figura 5.60 Structura RNP pentru recunoașterea de caractere, implementată cu procesor MicroBlaze încorporat în circuit FPGA

Aceste proprietăți sunt apoi propagate spre neuronii din stratul de ieșire, în număr egal cu cel al caracterelor care vor fi învățate de RNP, în aplicația de față aceasta fiind opt. Stratul de ieșire este complet conectat cu stratul ascuns, iar neuronii pulsativ din acesta se vor activa atunci când combinația caracteristicilor raportate de activarea celor paisprezece neuroni din stratul ascuns va corespunde cu caracteristicile prescrise pentru caracterul prezentat la intrarea rețelei neuronale. Așadar, această RNP se poate considera a fi implementată cu codarea ratei impulsurilor deoarece nu sunt considerate decalajele temporale dintre acestea ci doar corelația dintre momentele de apariție pe cele 48 de intrări.

### 5.6.2.4. Modelul neuronal utilizat

Cei mai importanți parametri ai modelului *leaky integrate-and-fire* sunt prezente și în modelul utilizat în această aplicație, adaptat pentru a putea fi implementat în procesorul MicroBlaze încorporat în hardware-ul reconfigurabil FPGA. Astfel, variabila potențialului de membrană (PM) are în continuare un rol primordial în funcționarea neuronului. Fiecărei intrări se asociază o pondere sinaptică care se vor însuma pentru a obține PM. Valoarea de saturație a PM va fi notată în continuare cu THS (ca și în Figura 5.61, care ilustrează modelul neuronal) reprezentând potențialul de prag, care dacă este depășit, neuronul se va activa, emițând un impuls axonal. În urma unui astfel de eveniment, neuronul va trece în faza de hiperpolarizare prin resetarea PM la o valoare mai mică decât cea de repaus. Această din urmă

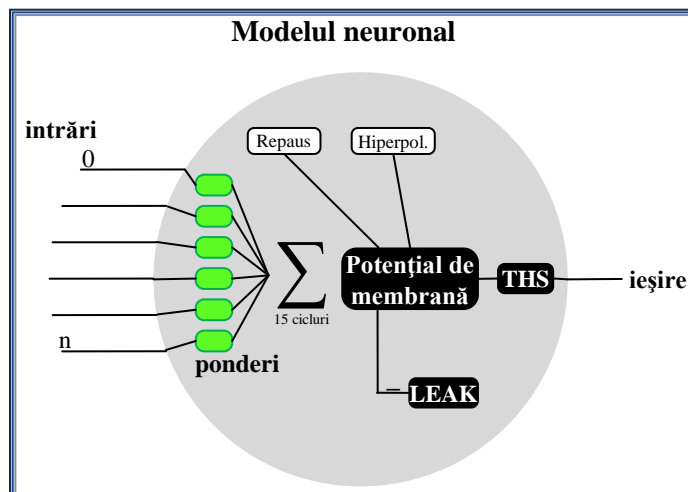


Figura 5.61 Modelul neuronal pulsativ utilizat în aplicația de recunoaștere caractere

valoare este cea până la care va scade PM cu valoarea de scurgere (LEAK) la fiecare pas de procesare în lipsa impulsurilor de intrare.

Funcționarea neuronului (Figura 5.62) este procesată în 15 pași. Potențialul de membrană este inițializat cu valoarea de repaus. Valoarea de hiperpolarizare este setată la aproximativ 1-2% din valoarea THS.

În fiecare pas se evaluează PM conform ecuației de mai jos

$$\mu = \mu + \sum_i x_i * \beta_i^{jk} - \mu_{scu} \quad \text{Ec. 74}$$

cu notații conform celor din Figura 5.34 și Figura 5.36.



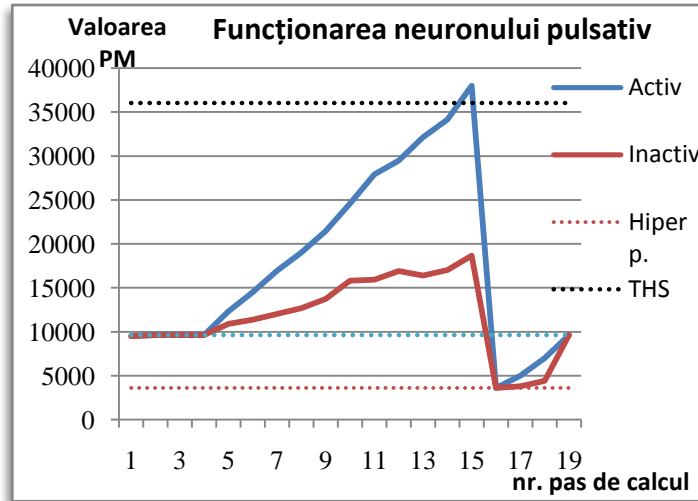


Figura 5.62 Principiul de funcționare a modelului neuronal utilizat în aplicația de recunoaștere caractere

Așadar, PM va crește la fiecare pas cu valoarea sumei ponderilor sinapselor care au recepționat impuls de intrare și va scădea cu valoarea de scurgere. După epuizarea celor 15 pași, se va evalua valoarea PM, comparându-se cu cea de prag (THS) și se va decide dacă neuronul se va activa sau nu. În pasul următor PM se va reseta la valoarea de hiperpolarizare.

#### Structura și funcționarea neuronilor stratului ascuns antrenați pentru coloanele matricii de intrare

Acești neuroni sunt realizați cu opt intrări pre-sinaptice fiecare. Fiecare pondere sinaptică a fost inițializată cu valori aleatoare generate în Matlab ele fiind în domeniul inferior de 10-20% al intervalul valorilor ponderilor care a fost setat între 0 și 1200 și stocat în variabile de tip întreg în procesorul MicroBlaze. Valoarea de prag a acestor neuroni a fost aleasă astfel încât aceștia se vor activa numai atunci, când cel puțin 50% dintre intrările lor primesc impuls pre-sinaptic (4 în cazul de față). Ecuațiile care au fost utilizate pentru calculul valorilor PM de repaus respectiv al THS ( $\theta$ ) sunt date mai jos:

$$\mu_{\text{rep}} = n_{\text{act}} * \frac{\beta_{i \text{ max}}^{jk} - \beta_{i \text{ min}}^{jk}}{i} * 4 = 9600 \quad \text{Ec. 75}$$

$$\theta = n_{\text{act}} * \frac{\beta_{i \text{ max}}^{jk} - \beta_{i \text{ min}}^{jk}}{i} * 15 = 36000 \quad \text{Ec. 76}$$

unde  $n_{\text{act}}=4$ , este numărul intrărilor active.

Valoarea de hiperpolarizare s-a setat la 3200, deoarece această valoare aproximează 1% din valoarea de prag THS iar valoarea de scurgere este 600.

#### Structura și funcționarea neuronilor stratului ascuns antrenați pentru liniile matricii de intrare

Neuronii din această categorie dispun de șase intrări pre-sinaptice. Valorile ponderilor sinaptice sunt inițializate în mod identic cu cazul anterior, având aceeași plajă de valori. Similar, valoarea de prag a fost setată ca neuronii să se activeze dacă cel puțin trei intrări sunt active (50%). În acest caz parametrii modelului sunt calculați astfel:

$$\mu_{\text{rep}} = n_{\text{act}} * \frac{\beta_{i \text{ max}}^{jk} - \beta_{i \text{ min}}^{jk}}{i} * 4 = 7200 \quad \text{Ec. 77}$$

$$\theta = n_{\text{act}} * \frac{\beta_{i \text{ max}}^{jk} - \beta_{i \text{ min}}^{jk}}{i} * 15 = 27000 \quad \text{Ec. 78}$$

unde  $n_{\text{act}}=3$ , este numărul intrărilor active.

Notățiile din ecuațiile de mai sus sunt cele din Figura 5.34 și Figura 5.36. Valoarea de hiperpolarizare s-a setat la 2400, deoarece această valoare aproximează 1% din valoarea de prag THS iar valoarea de scurgere este la fel 600.

### Structura și funcționarea neuronilor stratului de ieșire

Neuronii stratului de ieșire nu diferă mult de cei din stratul ascuns, cu excepția numărului de intrări – paisprezece în cazul celor de ieșire – respectiv a valorilor unor parametri. Valorile ponderilor au fost inițializate și mărginite în mod identic cu cele prezentate anterior.

Totodată s-a considerat, că un caracter dintre cele opt pentru care va fi antrenat RNP – ceea ce a determinat și numărul neuronilor de ieșire, egal cu șapte – va activa în stratul ascuns maxim patru neuroni, adică va deține patru proprietăți detectabile simultan de acest strat. De aici rezultă, că valorile  $\theta$  respectiv  $\mu_{rep}$  vor fi calculate în mod identic cu cele descrise în Ec. 75 respectiv Ec. 76. Valoarea de hiperpolarizare s-a setat la 3200, deoarece această valoare aproximează 1% din valoarea de prag THS iar valoarea de scurgere este la fel 600. Valoarea prescrisă pentru acești neuroni depinde de caracterul prezentat la intrări în ciclul actual de învățare.

#### 5.6.2.5. Prezentarea algoritmului de învățare implementat în RNP

Caracterul 0:	1, 0, 0, 0, 0, 1,	1, 0, 0, 0, 0, 0, 0, 1
Caracterul 1:	0, 0, 0, 1, 0, 0,	1, 1, 0, 0, 0, 0, 0, 0
Caracterul 2:	1, 1, 0, 0, 0, 0,	1, 0, 0, 0, 1, 0, 0, 1
Caracterul 3:	0, 0, 0, 0, 0, 1,	1, 0, 0, 0, 1, 0, 0, 1
Caracterul 4:	0, 0, 1, 1, 0, 0,	0, 0, 0, 0, 1, 0, 0, 0
Caracterul 5:	0, 0, 0, 0, 0, 1,	1, 0, 0, 1, 0, 0, 0, 1
Caracterul 6:	1, 0, 0, 0, 1, 0,	1, 0, 0, 0, 1, 0, 0, 1
Caracterul 7:	0, 1, 0, 1, 0, 1,	1, 0, 0, 0, 0, 0, 0, 0

Figura 5.63 Structura setului de antrenare a stratului ascuns (a se nota, că primele 6 valori sunt pentru neuronii de coloană iar următoarele 8 pentru cele de linie)

ieșirile acestora vor fi propagate spre neuronii de ieșire care vor primi în același timp și valorile prescrise, conform setului de antrenare propriu. Un ciclu de antrenare constă în prezentarea fiecărui caracter la intrările RNP.

În cadrul funcțiilor ce implementează neuronii individuali se pot separa 15 pași de calcul în cadrul cărora se însumează valorile ponderate primite, generate de un caracter prezentat la intrare. Dacă suma rezultantă depășește valoarea de prag, neuronul se va activa. În acest moment are loc adaptarea ponderilor sinaptice conform unui set de reguli bazate pe cele descrise de HEBB, într-o variantă supervizată, după cum se prezintă în Tabelul 16. Aceste condiții prezentate determină dacă și în ce direcție se vor modifica ponderile.

Valoarea cu care crește sau scade o pondere este egală cu 3% din valoarea maximă a ponderilor  $\beta_i^{jk}$ .

Rețeaua neuronală pulsativă implementată pentru această aplicație dispune de 48 de intrări și 8 ieșiri. Neuronii din stratul ascuns respectiv cel de ieșire vor fi antrenați cu algoritmi asemănători. În acest sens am construit seturi de antrenare separate, care conțin valori de intrare și valorile de ieșire prescrise fiecărui neuron. Aceste valori prescrise

sunt ilustrate în Figura 5.63. Aceste valori vor fi prezentate pe rând RNP în decursul ciclurilor de antrenare.

Astfel, se vor prezenta valorile setului de antrenare neuronilor din stratul ascuns, iar apoi

Caracterul 0:	1, 0, 0, 0, 0, 0, 0, 0
Caracterul 1:	0, 1, 0, 0, 0, 0, 0, 0
Caracterul 2:	0, 0, 1, 0, 0, 0, 0, 0
Caracterul 3:	0, 0, 0, 1, 0, 0, 0, 0
Caracterul 4:	0, 0, 0, 0, 1, 0, 0, 0
Caracterul 5:	0, 0, 0, 0, 0, 1, 0, 0
Caracterul 6:	0, 0, 0, 0, 0, 0, 1, 0
Caracterul 7:	0, 0, 0, 0, 0, 0, 0, 1

Figura 5.64 Structura setului de antrenare a stratului de ieșire

Tabelul 16 Algoritmul de învățare supervizat utilizat

Intrare	Ieșire	Valoare prescrisă	Adaptare ponderi
0	0	0	-
0	0	1	-
0	1	0	Scade?
0	1	1	-
1	0	0	-
1	0	1	Crește
1	1	0	Scade
1	1	1	Crește?

În cazul în care valoarea de intrare este 1 iar ieșirea neuronului nu se activează în timp ce valoarea prescrisă ar cere acest lucru, sinapse intrării respective va fi ajustată prin creșterea ponderii

corespunzătoare. Dacă intrarea și ieșirea sunt ambele active dar activarea axonală nu a fost prescrisă, ponderea respectivă va fi diminuată. În celelalte cazuri ponderile vor rămâne neschimbate.

Ar exista, însă, încă două cazuri în care s-ar părea că ar fi nevoie de adaptarea ponderilor. Primul astfel de caz este acela în care atât intrarea, ieșirea cât și valoarea prescrisă este egală cu 1, dar creșterea ponderii în acest caz ar duce la o supra-învățare a rețelei neuronale, deoarece putem considera, că antrenarea și-a atins deja scopul. Al doilea caz asemănător este acela în care ponderea ar trebui să scadă, dacă intrarea este inactivă, ieșirea este activă iar valoarea prescrisă nu cere activare. Aici putem considera, că valoarea respectivă de intrare nu a contribuit la învățarea caracterului respectiv, deci nu va fi luat în considerare. Conform acestor observații în aceste cazuri ponderile nu au fost modificate.

### 5.6.2.6. Proiectarea componentelor software ale sistemului încorporat

Toate funcționalitățile prezentate în secțiunile anterioare au fost implementate într-un program C compilat pentru și rulat pe procesorul Xilinx MicroBlaze. Mediul de dezvoltare utilizat, Xilinx EDK (prezentat în subcapitolul 5.6.4) generează drivere și funcții necesare pentru accesarea perifericelor introduse în proiect.

Faza de analiză și proiectare a unui proiect trebuie să fie gata înainte de realizarea codului, pentru a obține o atenție mărită din partea dezvoltatorilor. Orice dezvoltator recunoaște importanța acestor faze deoarece s-a dovedit că de acestea depinde producerea și re folosirea de software. Pentru analiza și proiectarea programelor s-au creat limbajele de modelare. Unul din aceste limbaje de modelare este limbajul de modelare unificat - UML (The Unified Modeling Language).

UML nu este un simplu limbaj de modelare orientat pe obiecte, ci în prezent, este limbajul universal standard pentru dezvoltatorii software din toată lumea. UML este succesorul propriu-zis al celor mai bune trei limbaje de modelare anterioare orientate pe obiecte (Booch, OMT, și OOSE). UML se constituie din unirea acestor limbaje de modelare și în plus deține o expresivitate care ajută la rezolvarea problemelor de modelare pe care vechile limbaje nu o aveau.

Limbajul de modelare modificat (UML - The Unified Modeling Language) oferă arhitecturi de sisteme ce funcționează pe analiza și proiectarea obiectelor cu un limbaj corespunzător pentru specificarea, vizualizarea, construirea și documentarea artefactelor sistemelor software și de asemenea pentru modelarea în întreprinderi. UML este un limbaj de modelare care oferă o exprimare grafică a structurii și comportamentului software. Pentru această exprimare grafică se utilizează notațiile UML.

Notațiile UML constituie un element esențial al limbajului pentru realizarea propriu-zisă a modelării și anume partea reprezentării grafice pe care se bazează orice limbaj de modelare. Modelarea în acest limbaj se realizează prin combinarea notațiilor UML în cadrul elementelor principale ale acestora denumite diagrame. În cadrul UML-ului descoperim nouă tipuri de diagrame: diagrama cazurilor de utilizare, diagrama de secvență, diagrama de colaborare, diagrama de clase (cea mai utilizată), diagrama de stări, diagrama de componente, diagrama de construcție, diagrama de obiecte, diagrama de activități

Analiza unei aplicații implică realizarea mai multor categorii de modele, dintre care cele mai importante sunt:

- ❖ Modelul de utilizare: realizează modelarea problemelor și a soluțiilor acestora în maniera în care le percepe utilizatorul final al aplicației. Diagramă asociată: **diagramă de cazuri de utilizare**
- ❖ Modelul structural: se realizează pe baza analizei statice a problemei și descrie proprietățile statice ale entităților care compun domeniul problemei. Diagrame asociate: **diagramă de module, diagramă de clase**
- ❖ Modelul comportamental: privește descrierea funcționalităților și a succesiunii în timp a acțiunilor realizate de entitățile domeniului problemei. Diagrame asociate: **diagrama (harta) de stări, diagrama de colaborare, diagrama de interacțiune**

O diagramă oferă utilizatorului un mijloc de vizualizare și de manevrare a elementelor de Modelare. Majoritatea diagramelor se prezintă sub forma unor grafuri, compuse din elemente și arce.

Diagramele pot arăta o parte sau toate caracteristicile elementelor de modelare, conform nivelului de detaliu util în contextul unei diagrame date. Diagramele pot grupa informații interdependente, pentru a arăta, de exemplu caracteristicile moștenite de o clasă.

Diagramele UML sunt:

- diagramele cazurilor de utilizare, care prezintă funcțiile sistemului din punct de vedere al utilizatorului;
- diagrame de clasă, care prezintă structura statică în termeni de clase și asocieri (relații);
- diagrame de colaborare, care sunt reprezentări spațiale ale obiectelor, legăturilor și interacțiunilor;
- diagrame de secvență, care prezintă temporal obiectele și interacțiunile lor;

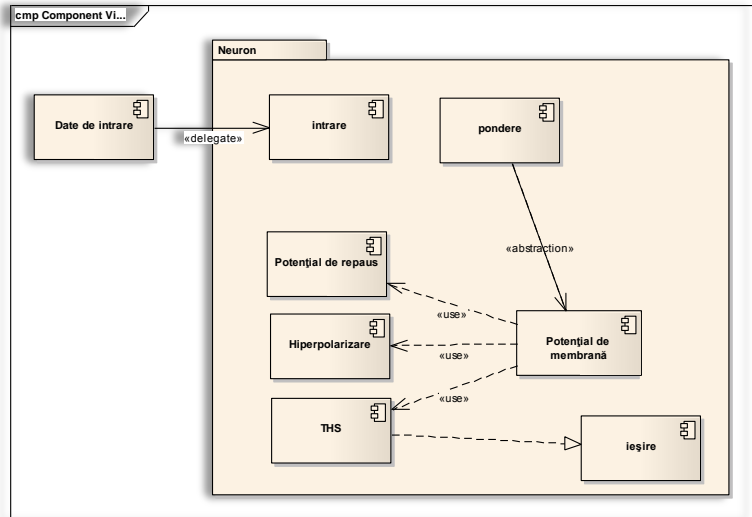


Figura 5.65 Diagrama de componente a modului neuronului pulsativ implementat software

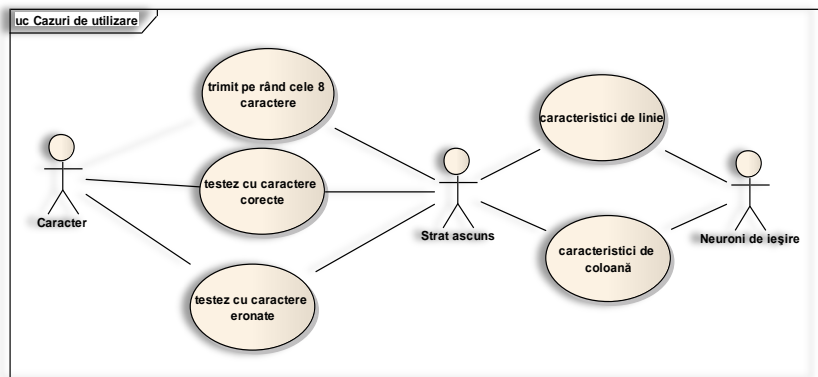


Figura 5.66 Diagrama cazurilor de utilizare (use-case) a RNP implementate pe procesorul MicroBlaze

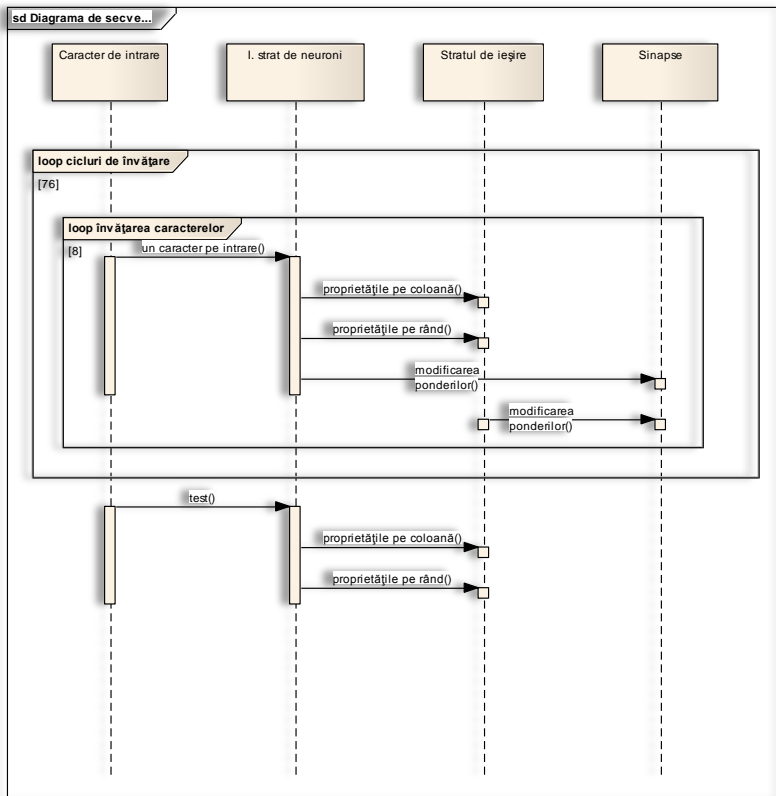


Figura 5.67 Diagrama de secvență a RNP implementate pe procesorul MicroBlaze

▪ diagrame de componente, care prezintă componentele fizice ale unei aplicații;

▪ diagrame de construcție, care prezintă construcția componentelor pe dispozitivele hardware;

▪ diagrame de stări-tranziții, care prezintă comportamentul unei clase în termeni de stări;

▪ diagrame de obiecte, care prezintă obiectele și relațiile lor, fiind niște diagrame de colaborare simplificate, fără reprezentarea mesajelor trimise între obiecte;

▪ diagrame de activități, care reprezintă comportamentul unei operații în termeni de acțiuni.

▪ diagramele de colaborare împreună cu cele de



secvență se numesc diagrame de interacțiune pe când diagramele de stare mai sunt denumite mașini cu stări finite, automate, etc.

Prima fază a implementării software a fost realizarea modelului neuronal în care putem regăsi toate variabilele prezentate, după cum arată și Figura 5.65

Programul ilustrat în diagrama de componente din Figura 5.66 (generată cu software-ul Enterprise Architect) are trei actori.

Primul este actorul caracter care livrează rând pe rând caracterele neuronilor stratului ascuns. Actorul din stratul ascuns analizează aceste informații apoi trimite spre neuronii din stratul de ieșire caracteristicile de linie și coloană. Neuronii de ieșire vor învăța să identifice un anumit caracter.

Scopul principal al diagramei de secvență din Figura 5.67 este ilustrarea interacțiunii dintre obiectele afișate și secvența temporală a desfășurării acestora. În diagramă sunt reprezentate patru obiecte. În primul pas, obiectul *caracter* predă un caracter obiectului *strat ascuns*, asupra căruia se execută operațiile și apoi livrează mai departe caracteristicile de linie și coloană obiectului *neuron de ieșire*. După evaluarea informațiilor vom trece la modificarea ponderilor din obiectele *sinapse*. Acest proces este executat pentru cele opt caractere completând un ciclu de antrenare care este repetat până când neuronii de ieșire identifică corect caracterele prezentate la intrare.

După antrenarea RNP cu setul de antrenare urmează testarea acesteia. În primul pas se trimite un caracter fără erori obiectului *strat ascuns*. Acest caracter propagându-se spre obiectul *neuron de ieșire* obținem rezultatul învățării. După cum se vede în diagramă, în decursul testării ponderile sinaptice nu sunt alterate. Această testare se va repeta apoi și cu caractere de intrare zgomotoase (cu erori de reprezentare).

### 5.6.2.7. Rezultate experimentale

În decursul antrenării am calculat eroarea de aproximare comparând valorile de ieșire și valorile prescrise și am incrementat un numărător dacă acestea nu corespund. Dacă gradientul acestui numărător va fi zero, înseamnă că rețeaua a învățat problema propusă.

În diagrama din Figura 5.68 putem observa că eroarea descrește, având valori mari în primele cicluri, apoi începând cu aproximativ ciclul numărul 5 scade abrupt ca apoi să descrească lent spre zero.

În Figura 5.69 se poate vedea cum variază potențialul de membrană al neuronului nr. 1 în timpul primului ciclu (de 15 pași) de antrenare al caracterului 0, neuron care are prescrist activare pentru acest caracter. Se observă că PM scade din cauza lipsei impulsurilor de intrare. Inflexiunea din pasul 4 se datorează faptului că în experimentul respectiv am inițializat PM cu valoarea de hiperpolarizare, și am prevăzut ca și în lipsa impulsurilor de intrare PM să atingă valoarea de repaus cel târziu în patru pași.

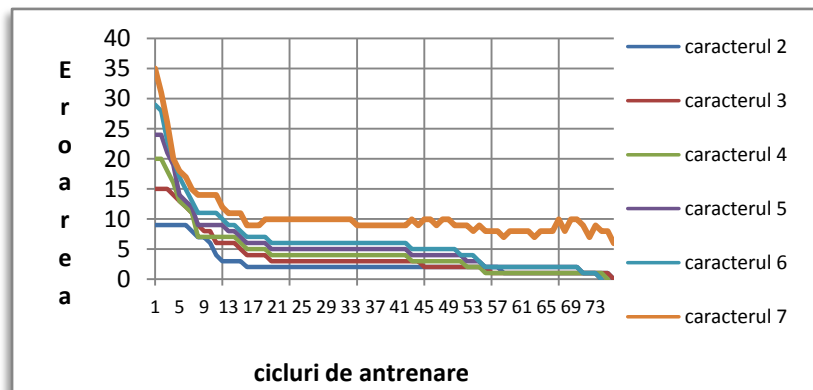


Figura 5.68 Variația erorii în timpul antrenării (pentru caracterele 2,3,4,5,6,7)

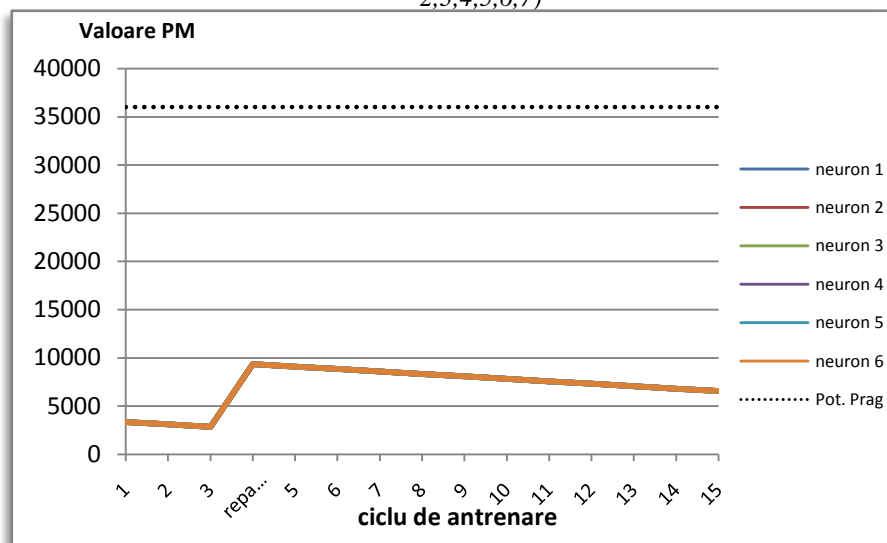


Figura 5.69 Variația potențialelor de membrană a 6 neuroni în timpul primului ciclu de antrenare (al caracterului 0)

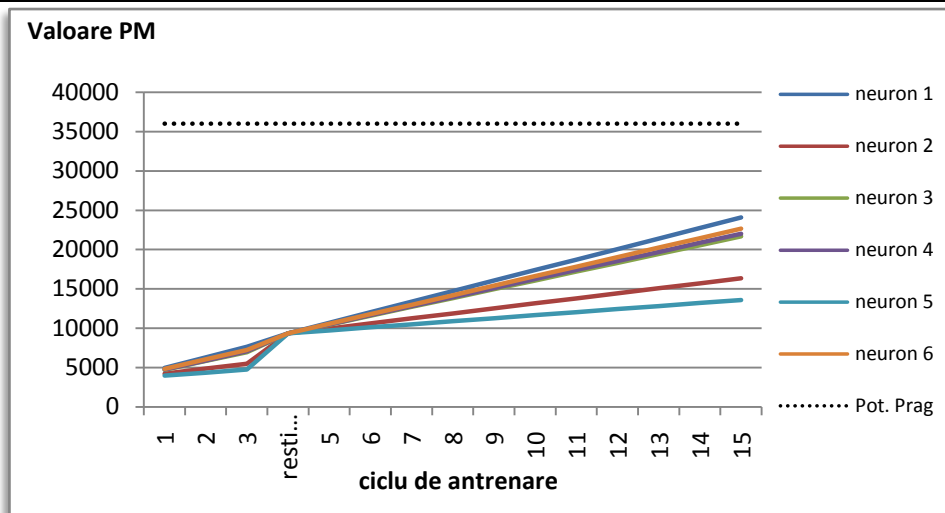


Figura 5.70 Variația potențialelor de membrană în timpul ciclului 32 de antrenare (al caracterului 0)

Se poate vedea în Figura 5.70 că procesul de învățare va ridica treptat valorile potențialelor de membrană, astfel în ciclul nr. 32 PM al neuronului 1, are deja valoarea cea mai mare, dar mai trebuie să crească pentru a atinge valoarea de prag.

Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 0) din Figura 5.71 arată că neuronul 1 se activează, învățând cu succes caracterul 0.

În cele ce urmează se vor prezenta rezultate experimentale similare, obținute prin măsurători efectuate în timpul antrenării celorlalte caractere.

Procesorul MicroBlaze trimite unui calculator prin portul serial aceste date ale RNP implementate, unde sunt salvate și vizualizate de un program Matlab.

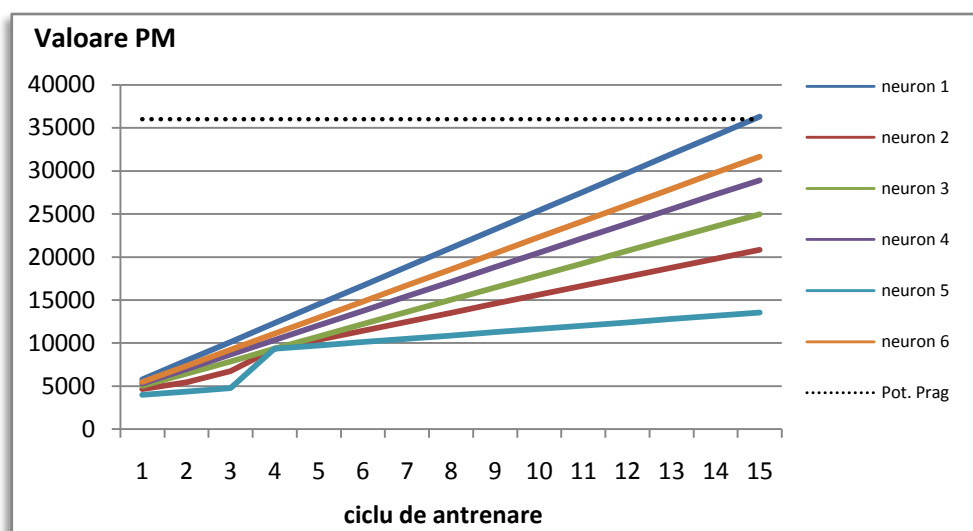


Figura 5.71 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 0), care este învățat de neuronul 1

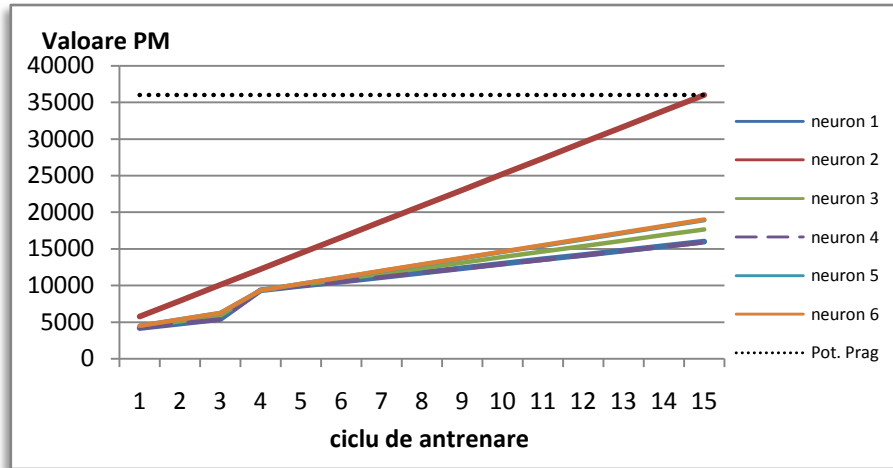


Figura 5.72 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 1), care este învățat de neuronul 2

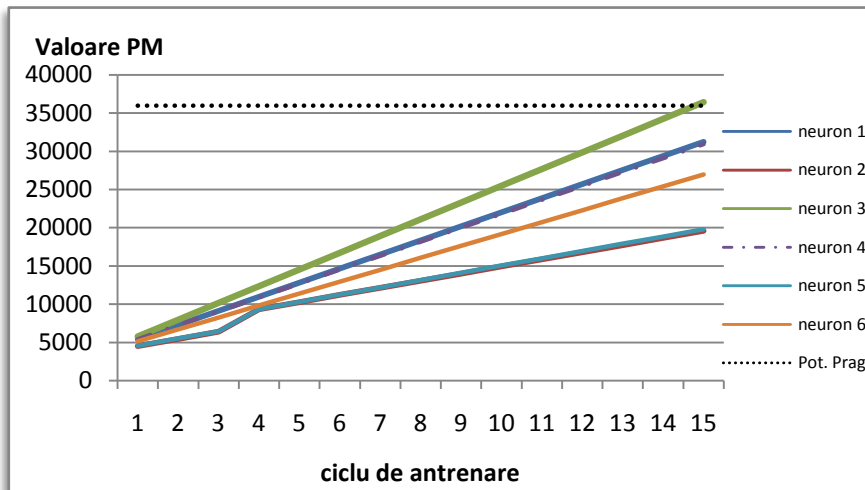


Figura 5.73 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 2), care este învățat de neuronul 3

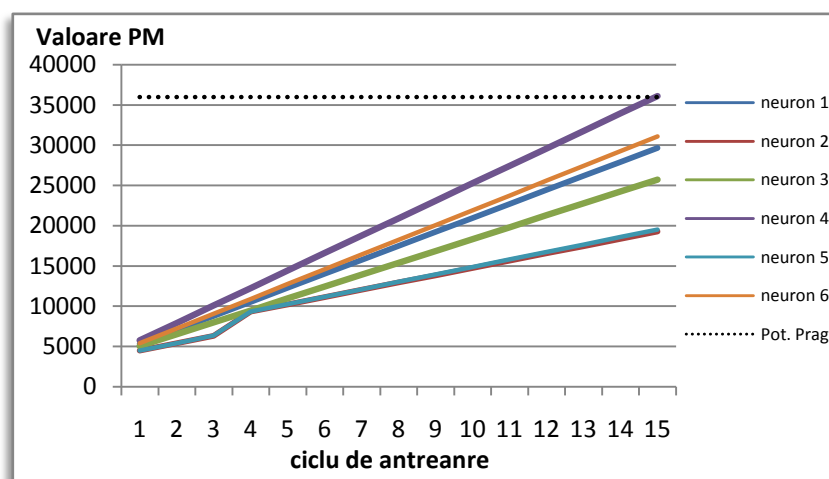


Figura 5.74 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 3), care este învățat de neuronul 4

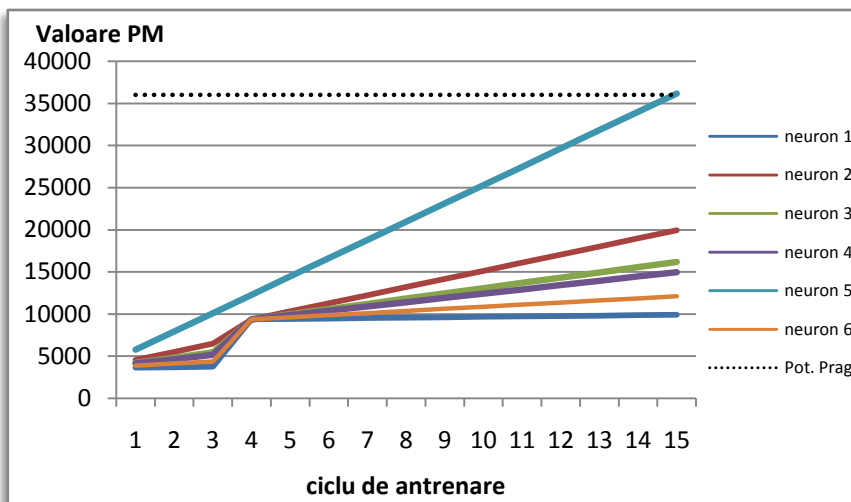


Figura 5.75 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 4), care este învățat de neuronul 5

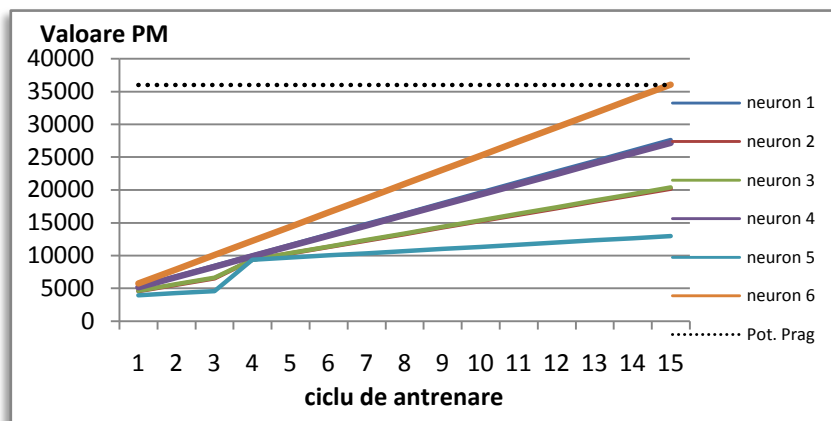
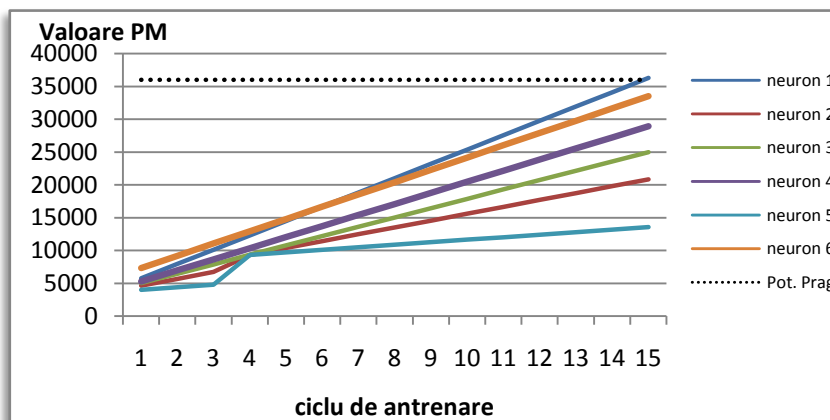


Figura 5.76 Variația potențialelor de membrană în timpul ultimului ciclu de antrenare (al caracterului 5), care este învățat de neuronul 6

În decursul testării RNP după învățare am obținut rezultate de recunoaștere pentru caracterele de la 0-7 corecte în 98% din cazuri. Cele mai multe erori au fost la diferențierea caracterelor 5 și 6.

Pentru a arăta robustețea aplicației am testat RNP implementată în FPGA și cu intrări perturbate de zgomot. În testul ilustrat de Figura 5.77 am prezentat la intrarea RNP un caracter 0 cu erori de reprezentare, dar se vede că sistemul a recunoscut caracterul fără probleme.



0	1	1	1	1	0
0	0	1	0	0	0
1	0	0	0	0	1
1	0	0	0	0	1
1	0	0	1	0	1
1	0	0	0	0	1
1	1	0	0	0	1
0	1	1	1	1	0

Figura 5.77 Testarea RNP cu caracterul 0 perturbat de zgomot

Dacă introducem o diferență mai mare între caracterul original de intrare și cel zgomotos de test, alterând astfel caracteristicile de linie sau coloană ale acesteia, RNP va reuși și în acest caz să identifice corect caracterul ().

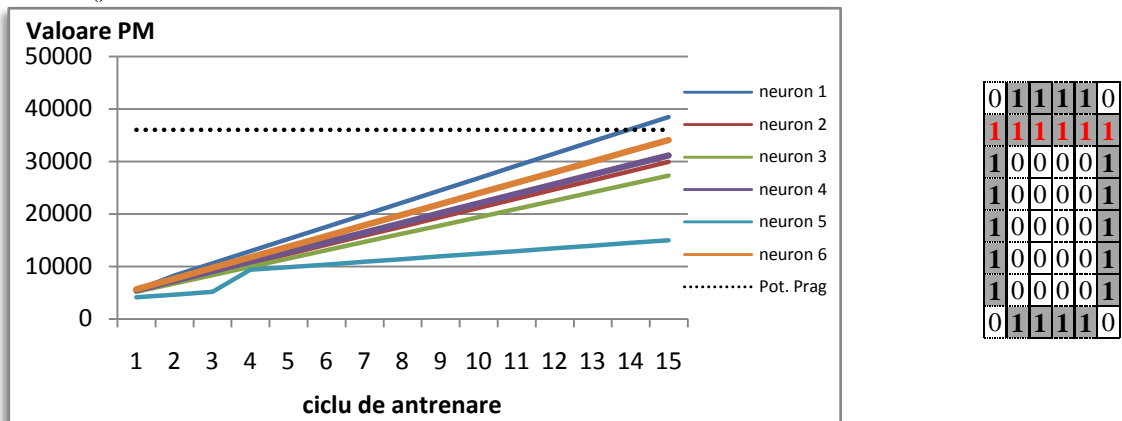


Figura 5.78 Figura 5.79 testarea RNP cu caracterul 0 perturbat de zgomot accentuat

### 5.6.3. Implementare benchmark: Clasificarea setului de date Wisconsin Breast Cancer

#### 5.6.3.1. Formularea problemei de clasificare

Testul utilizează a baza de date Wisconsin Breast Cancer Database (WBCD), care urmărește clasificarea malignă sau benignă a unor țesuturi mamare, pe baza a nouă proprietăți ale nucleelor celulare prelevate. Această bază de date a fost compilată de Dr. William H. Wolberg (Mangasarian & Wolberg, 1990) de la University of Wisconsin Hospitals, Madison, SUA. Conține 699 de elemente cu câte 9 atribute, fiecare reprezentând o caracteristică a țesuturilor cercetate (conform sursei (Open source data visualization and analysis for novice and experts, 2005) respectiv (National Science Foundation, 2007)):

- *Clump thickness*
- *Unif\_Cell\_Size*
- *Unif\_Cell\_Shape*
- *Marginal\_Adhesion*
- *Single\_Cell\_Size*
- *Bare\_Nuclei*
- *Bland\_Chromatine*
- *Normal\_Nucleoli*
- *Mitoses type*

#### 5.6.3.2. Elementele hardware ale sistemului încorporat dezvoltat

În cadrul acestei aplicații am utilizat module hardware predefinite ale mediului de dezvoltare Xilinx Platform Studio precum și module proprii implementate în limbaj VHDL (Figura 5.80). Elementul principal al implementării este și în acest caz procesorul Xilinx MicroBlaze, la care am conectat mai multe module de memorie respectiv dispozitive periferice prin magistrala de 32 de biți Processor Local Bus (PLB).

Perifericul principal care a fost dezvoltat este acela care implementează codificarea valorilor de intrare a rețelei neuronale pulsative care va descrie detaliat în subcapitolul următor.

Acest periferic de generare a impulsurilor de intrare decalate temporal conține 5 registre prin care se realizează comunicare dintre acesta și procesorul

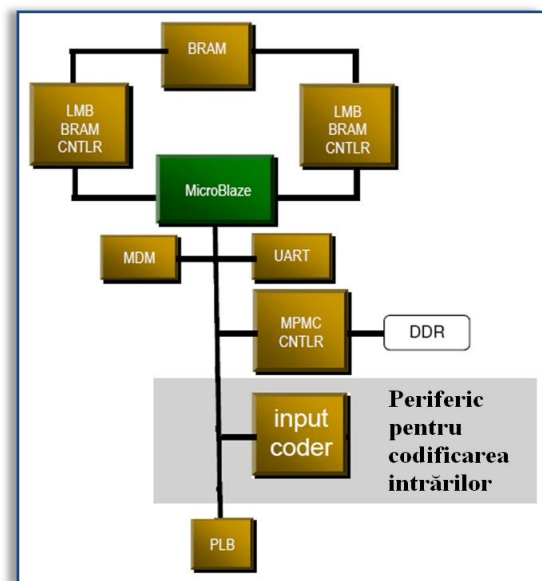


Figura 5.80 Structura hardware a sistemului implementat

MicroBlaze prin magistrala PLB. Pentru a evita coliziuni în scrierea acestor registre dintre acești doi inițiatori de transferuri ai magistralei, am separat rolul registrelor. Astfel, MicroBlaze va putea scrie numai în unul dintre registrele celelalte 4 fiind configurate ca putând fi numai citite de către procesor.

### 5.6.3.3. Principiul de funcționare perifericului de codificare a intrărilor

Scopul acestui circuit este de a codifica valorile de intrare ale rețelei neuronale – care sunt de fapt cele nouă parametri ale fiecărei mostră din setul de date – în impulsuri decalate temporal, pentru a putea fi prelucrate de către RNP cu procesor încorporat implementată.

Elementele bazei de date WBCD, care sunt valori cu două zecimale (scalate, pentru a le putea reprezenta ca numere întregi în FPGA) sunt stocate într-o memorie BRAM încorporată în acest periferic. Procesorul MicroBlaze trimite o valoare de 10 biți într-unul dintre registrele perifericului care va fi utilizată pentru a adresa memoria valorilor setului WBCD. Astfel vom obține o valoare de 64 biți care va reprezenta o mostră, conținând toate cele 9 proprietăți, pe care baza de date le reprezintă în intervalul 0-10 cu o precizie de o zecimală. Pentru a stoca numere întregi în FPGA aceste valori au fost scalate între 0-100 astfel necesitând 7 biți pentru reprezentare, ceea ce rezultă în  $7 \times 9 = 63$  biți. Ultimul bit din cei 64 va fi folosit pentru a stoca clasa în care aparține mostra respectivă (1-malign, 0-belig).

Perifericul este programat să aștepte în stare inactivă sosirea unui front crescător pe semnalul *Slv\_WR\_ack* care va apare pe magistrala PLB atunci, când MicroBlaze va efectua prima operație de scriere la adresa perifericului. În urma acestui eveniment, perifericul va intra în stare de inițializare, ce durează 3 cicluri de tact după care va începe secvența celor 15 pași în care cele 9 valori de intrare se vor coda în impulsuri decalate temporal. Semnalul de intrare *Inp\_valid* al perifericului de codificare (Figura 5.82) trebuie menținut activ timp de 18 cicluri de tact pentru ca în cel de-al patrulea ciclu să obținem la ieșiri

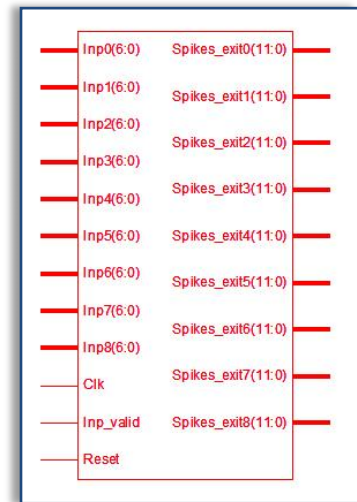


Figura 5.82 Simbolul bloc al perifericului de codificare a intrărilor

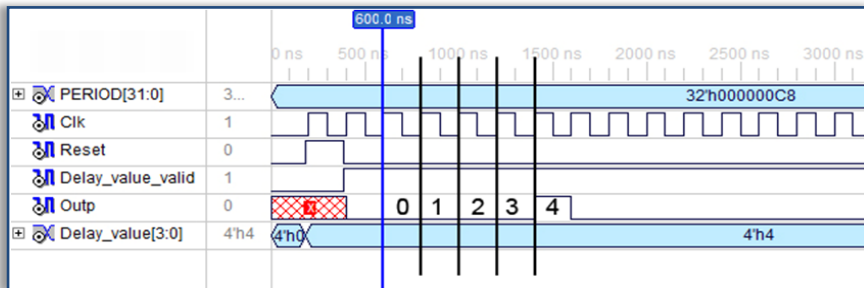


Figura 5.81 Generarea impulsului decalat din valoarea pe 4 biți citită din BRAM

*Spikes\_exitX* primul set de impulsuri.

Modulul de codificare a intrărilor conține 9 memorii BRAM de tipul RAMB16\_S36\_S36 (modul cu port dual pe 32 de biți, astfel putând accesa 64 de biți simultan), care stochează valorile decalajelor conform funcțiilor de activare cu care a fost acoperită spațiul valorilor de intrare.

În cele ce urmează vor fi prezentate două versiuni ale RNP pentru această aplicație, una fără și alta cu strat ascuns. Versiunea mai simplă (Figura 5.86) conține  $12 \times 9 = 108$  neuroni de intrare (Figura 5.84) și având doi neuroni de ieșire

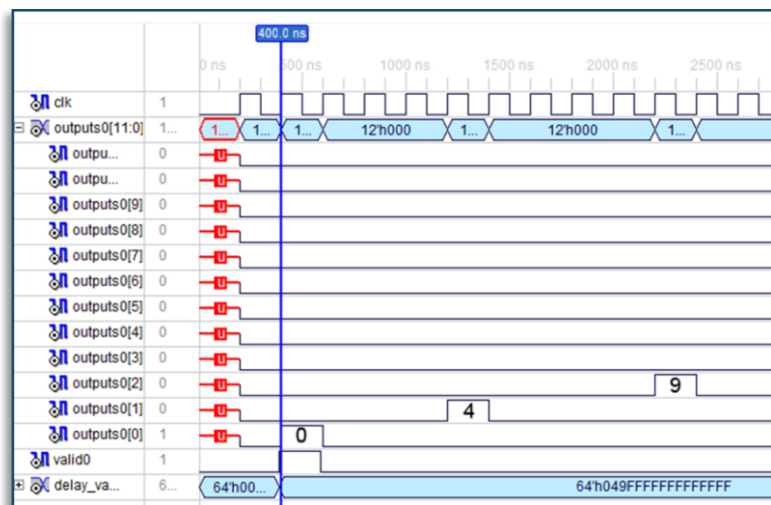


Figura 5.83 Codificarea valorii de intrare în trei impulsuri decalate temporal în decursul unui pas temporal de 15 cicluri de tact



rezultă 216 sinapse. Varianta cu strat ascuns a RNP (Figura 5.87) a fost implementată codificând valorile de intrare cu 16 neuroni de intrare – pentru o precizie ridicată – având 6 neuroni în startul ascuns la care aparțin  $16 \cdot 9 \cdot 6 = 864$  sinapse și doi neuroni de ieșire cu  $6 \cdot 2 = 12$  sinapse.

Rezultatele simulării unui neuron de intrare și a blocului de neuroni de intrare al uneia dintre intrări care codifică valorile sosite pe aceasta în impulsuri decalate temporal sunt prezentate în figurile de mai jos (Figura 5.81 respectiv Figura 5.83).

#### 5.6.3.4. Codificarea valorilor de intrare

Similar cu cele prezentate la aplicațiile cu procesor PicoBlaze în acest caz cadrul temporal de codificare a intrărilor este de 15 cicluri de tact. În figura alăturată (Figura 5.84) am reprezentat acoperirea cu funcții triunghiulare a spațiului intrărilor (valori de la 0 la 100) de către 12 neuroni de intrare. Se poate observa exemplul marcat pe figură, și anume că valoarea 30 va activa trei neuroni de intrare. Dintre acești trei neuroni care codifică valoarea de intrare respectivă, neuronul 3 va emite impuls pre-sinaptic în pasul temporal 1 al cadrului temporal actual, neuronul 2 va emite impuls în pasul 7 iar neuronul de intrare 4 în cel de-al 9-lea ciclu de tact din cadrul. restul neuronilor de intrare nu va emite impuls pentru această valoare.

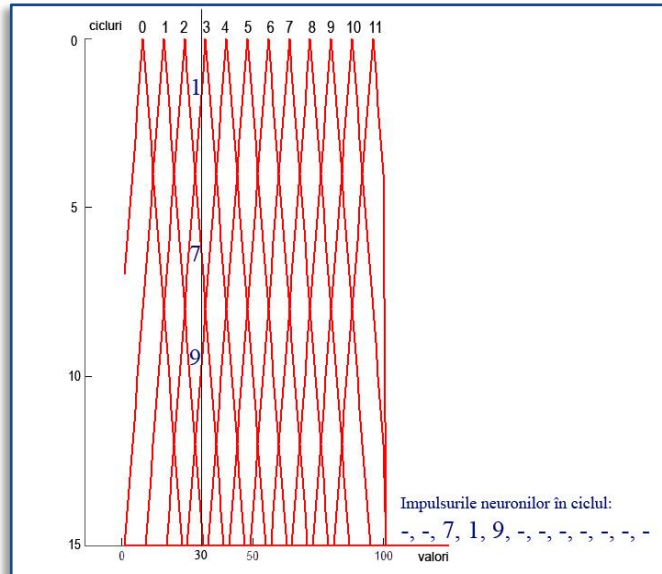


Figura 5.84 Funcțiile triunghiulare ale domeniilor receptive utilizate în cazul RNP fără strat ascuns

Deoarece RNP ce implementează această aplicație de clasificare are nouă intrări, avem nevoie de tot atâtea module BRAM pentru a stoca aceste valori ale momentelor de activare ale neuronilor de intrare. Ca și în cazul aplicațiilor prezentate anterior, aceste module sunt de tip dual-port pentru a putea accesa mai multe valori simultan și a accelera procesarea rețelei neuronale.

Figura 5.85 prezintă o variantă îmbunătățită a acestei codificări, utilizând 16 neuroni pentru o intrare, utilizată în a doua variantă a RNP pentru această aplicație, care implementează și un strat ascuns de neuroni. Această alegere a fost luată pentru a crește precizia codificării deci și a sensibilității clasificării. A rezultat o codificare în care pentru orice valoare de intrare se vor activa patru neuroni în loc de trei (de exemplu în Figura 5.85, valoarea 35 va activa neuronii nr. 6, 7, 8 și 9), activând procesul de învățare în mai multe sinapse.

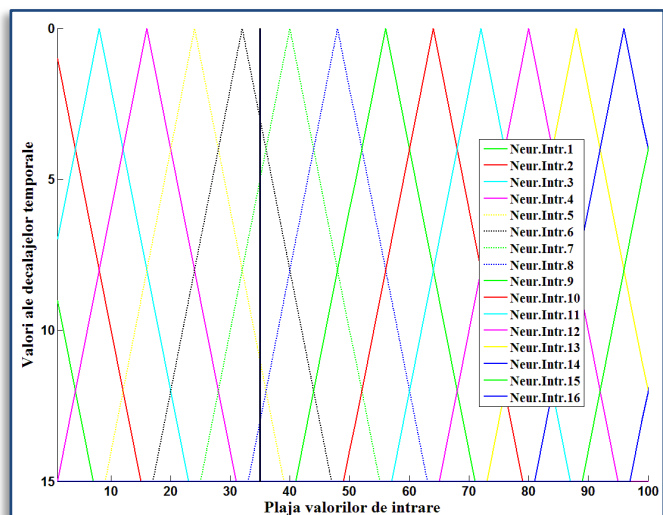


Figura 5.85 Funcțiile triunghiulare ale domeniilor receptive utilizate în cazul RNP cu strat ascuns

#### 5.6.3.5. Structura rețelelor neuronale pulsative implementate pentru rezolvarea clasificării setului de date Wisconsin

Am experimentat realizarea practică a mai multor variante a RNP cu procesor încorporat MicroBlaze pentru aplicația de clasificare a bazei de date WBCD pentru a găsi o soluție care să ofere o rezolvare bună, să funcționeze în timp real și în același timp să nu depășească necesarul de resurse hardware oferit de sistemul de dezvoltare FPGA utilizat (Nexys2-1200). După cum s-a sugerat și în subcapitolele anterioare se va prezenta în continuare structura și funcționarea a două dintre aceste RNP.



**Versiunea 1 – RNP cu două straturi.** Arhitectura acestei rețele neuronale artificiale pulsative implementate hardware este prezentată în Figura 5.86. După cum se vede și în figură, există două straturi de neuroni, cel de intrare, în care găsim nouă blocuri de neuroni de intrare – corespunzătoare celor nouă proprietăți a fiecărei mostre din WBCD, care formează valorile de intrare ale RNP – cu câte 12 neuroni ce codifică valorile de intrare în impulsuri decalate temporal. Fiecărui neuron din acest strat îi corespunde câte o sinapsă ale celor doi neuroni pulsativi din stratul de ieșire. Rezultă, că RNP conține în total  $9 \cdot 12 \cdot 2 = 216$  sinapse, care sunt implementate de funcții software rulate de procesorul MicroBlaze, ponderile asociate fiind stocate în memoria de date a acestuia. Tot aici este implementat și algoritmul somei, funcționând ca două funcții scrise în limbaj C, realizând procesarea celor doi neuroni de ieșire, care se vor activa atunci când mostra de la intrare aparține uneia dintre cele două clase (malign sau benign) asociate acestuia. În contrast cu aceste componente se află neuronii de intrare care sunt implementați utilizând limbajul de descriere hardware VHDL, devenind astfel un periferic separat – conectat la magistrala PLB al MicroBlaze – care funcționează în paralel cu procesorul, fiind construiți din elemente reconfigurabile distincte față de acesta.

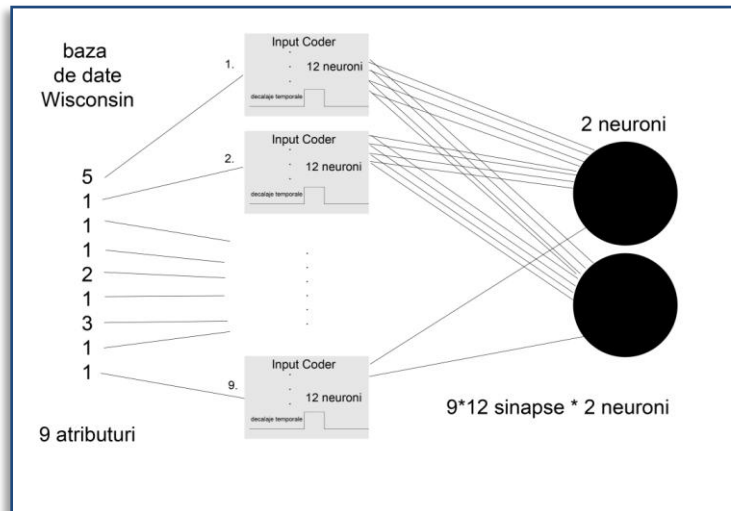


Figura 5.86 Prima versiune a RNP implementate pentru aplicația benchmark de clasificare a setului de date WDBC

Neuronii de ieșire sunt implementați conform aceluiași model de bază, *leaky integrate-and-fire*, adaptat pentru implementare în procesor încorporat soft-core. Calcularea parametrilor acestui model s-a efectuat după cum urmează. La alegerea valorii de prag a potențialului de membrană s-a luat în considerare că în decursul celor 15 cicluri ale cadrului temporal prin cele 108 intrări post-sinaptice ale unui neuron de ieșire poate sosi un număr total de 27 impulsuri, deoarece din cele 9 blocuri de neuroni de intrare fiecare va avea câte trei neuroni activi în acel cadru. Valorile ponderilor au fost mărginite între 0 și 1200 (reprezentabile pe min. 11 biți).

Rezultă, că valoarea  $THS = 27 \cdot 600 = 16200$ , unde 600 este valoarea medie a ponderilor sinaptice. Valoarea de repaus a PM s-a ales a fi 20% din THS, cea de hiperpolarizare 1% iar valoarea de scurgere 3%.

**Versiunea 2 – RNP cu trei straturi.** Diferența majoră față versiunea anterioară constă în introducerea unui strat ascuns de neuroni pulsativi pentru a îmbunătăți performanțele RNP. Acest strat conține 6 neuroni, întreaga rețea fiind complet conectată (Figura 5.87). De asemenea s-a schimbat și acoperirea plajei valorilor de intrare, utilizându-se 16 neuroni de intrare/bloc, dintre care vor fi activi câte patru neuroni într-un cadru de timp de 15 cicluri de tact. Aceste schimbări de parametrii au dus la modificări radicale în structura perifericului de codificare a valorilor de intrare, crescând complexitatea și necesarul de resurse reconfigurabile a acesteia. Un salt important a intervenit și în numărul de sinapse ale rețelei, acest număr devenind (de la 216 în cazul anterior)  $9 \cdot 16 \cdot 6 = 864$  între stratul de intrare și cel ascuns, respectiv  $6 \cdot 2 = 12$  între cel ascuns și neuronii de ieșire, în total 876 sinapse. Acest lucru înseamnă

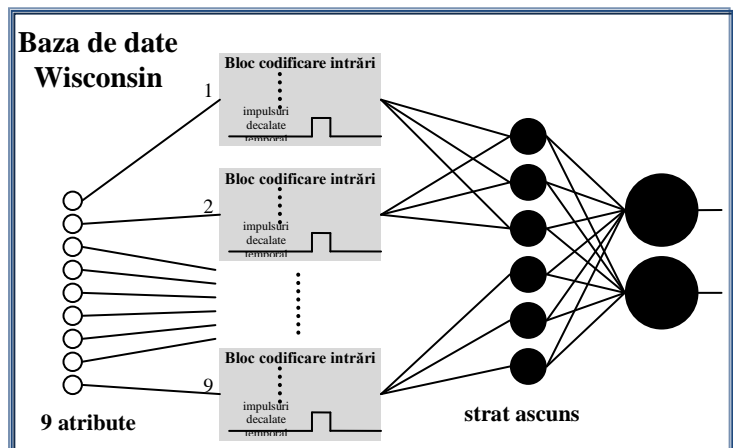


Figura 5.87 Versiunea cu strat ascuns a RNP implementate pentru aplicația benchmark de clasificare a setului de date WDBC

totodată o creștere de 4 ori a numărului de variabile pentru stocarea ponderilor sinaptice, ceea ce a dus la situația în care modulele de memorie BRAM încorporate în circuitul FPGA nu au mai fost suficiente pentru a stoca programul scris pe procesorul MicroBlaze, segmentul de date crescând în dimensiuni. Singura rezolvare a acestei probleme a constat în plasarea acestui segment în modulul de memorie DDRAM extern, prezent pe placa de dezvoltare Nexys2 utilizată. Această modificare, însă implică și introducerea în proiectul dezvoltat în mediul Xilinx EDK (descriș la subpunctul 5.6.4), a unui modul (predefinit în biblioteca de circuite EDK) *memory manager* care implementează protocolul de acces la acest tip de memorie, modul care utilizează la rândul său un quantum important de resurse reconfigurabile. Am ajuns astfel, la o utilizare de peste 95% a capacității circuitului FPGA XC3S1200E.

Neuronii din stratul ascuns au fost implementați pentru a îndeplini rolul de captare a unor proprietăți ale mostrelor de intrare, ale căror caracteristici (nouă la număr) definesc un spațiu de 9 dimensiuni. Reducerea acestui spațiu la 6 dimensiuni a contribuit la capacitatea stratului de ieșire de a învăța segmentarea acestuia în două hiperplanuri. Acest artificiu a dus la o creștere semnificativă a performanțelor sistemului, după cum se prezenta în subcapitolul 5.6.3.1.

### 5.6.3.6. Algoritm de învățare implementat

Neuronii stratului ascuns sunt identici din punct de vedere constructiv cu cei de ieșire în consecință funcționează în mod similar. Acești șase neuroni (Figura 5.87) primesc valori post-sinaptice – de la neuronii de intrare – decalate temporal. Am creat două grupe de câte trei neuroni în acest strat și am alternat valoarea de activare prescrisă acestora conform clasificării date de baza de date WBCD. În mod evident, s-a creat un set de antrenare, utilizând un număr de mostre egal cu 65% din numărul total al elementelor din baza de date alese în mod aleatoriu și restul elementelor au creat setul de test al RNP. Neuronii din stratul de ieșire au fost antrenați în Varianta 2 a RNP, accentuând efectul unei sinapse asupra PM atunci, când cel puțin trei dintre cele șase intrări ale acestor neuroni au fost activate. Valoarea cu care se modifică ponderile sinaptice variază în funcție de pasul temporal la care soma recepționează valoarea post-sinaptică corespunzătoare.

Astfel, dacă această valoare sosește chiar în primul ciclu de tact al cadrului temporal, înseamnă că o funcție de activare a acestuia este centrată exact pe valoarea de intrare actuală, deci ponderea trebuie mărită drastic (pas de învățare inițial, cu valoarea experimentată a fi de 60, la valori ale ponderilor cuprinse între 0-1200). În cazul în care sinapsa trimite valoarea post-sinaptică într-un ciclu ulterior, valoarea cu care crește ponderea lui – în cazul în care și valoarea axonală prescrisă este de activare – cu o valoare dată de formula de mai jos (ponderile se inițializează cu  $\beta_{j \max}^{ik} * 0.3$ ):

$$\Delta\beta_j^{ik} = \beta_{j \max}^{ik} / m \quad \text{Ec. 79,}$$

unde  $m = [1 \div 14]$  este numărul pasului temporal

În cazul în care un neuron a emis impuls axonal atunci când valoarea prescrisă nu cerut acest lucru, sinapsa respectivă este slăbită, conform algoritmului prezentat în Tabelul 12.

### 5.6.3.7. Utilizarea proiecției lui Sammon pentru vizualizarea mulțimilor multidimensionale

Proiecția lui Sammon (Sammon Jr., 1969) este o metodă de proiecție neliniară folosită la cartografierea spațiilor de mari dimensiuni pe spații de dimensiuni mai reduse. Algoritmul cartografiază spațiul original pe spațiul proiectat într-un mod în care distanțele dintre obiectele din spațiul original

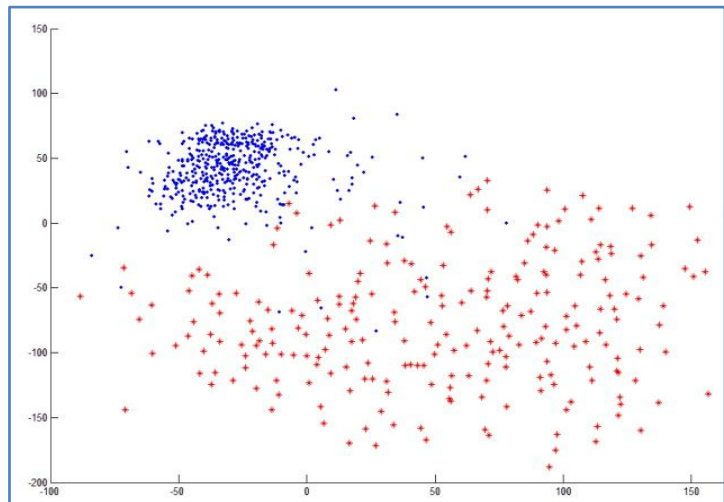


Figura 5.88 Reprezentarea grafică prin proiecție Sammon  $9D \rightarrow 2D$  a clasificării prescrise a bazei de date WBCD cu 699 mostre

sunt conservate în spațiul proiectat, la fel ca în tehnicile Self-Organizing Map (SOM). Totuși, experimentele arată că algoritmul lui Sammon are de obicei performanțe mai bune decât orice alt algoritm de cartografiere (Biswas, Jain, & Dubes, 1981).

Vom nota cu  $d_{ij}^*$  distanța dintre un obiect  $i$  și un obiect  $j$  din spațiul original, și cu  $d_{ij}$  distanța dintre aceleași obiecte în spațiul proiectat. Algoritmul lui Sammon încearcă să minimizeze următoarea funcție de eroare, adică deformarea proiecției:

$$E = \frac{1}{\sum \sum_{i < j} d_{ij}^*} \sum \sum_{i < j} \frac{(d_{ij}^* - d_{ij})^2}{d_{ij}^*} \quad \text{Ec. 80}$$

Această ecuație poate fi minimizată folosind de exemplu o metodă de gradient cum propune și Sammon (Sammon Jr., 1969). Deși proiecția lui Sammon nu este în principiu orientată spre rețele neuronale, ea poate fi modelată folosind o rețea neuronală, cum se arată în (Mao & Jain, 1995).

### 5.6.3.1. Rezultate experimentale

În Figura 5.89 este ilustrată clasificarea originală a elementelor conform valorilor prescrise din baza de date. După 100 de cicluri de antrenare a celor 683 de elemente, rețeaua neuronală pulsativă (Varianta 1 – fără strat ascuns) a clasificat elementele după cum se poate vedea în Figura 5.90.

Punctele colorate cu albastru și roșu din Figura 5.90 – care afișează rezultatul clasificării WBCD cu RNP Varianta 1 (fără strat ascuns) aparțin celor două clase. Punctele notate cu plus (negru) nu au fost clasificate în nici un grup. Elementele bazei de date pot fi încadrate în 9 dimensiuni dar pentru a putea vizualiza grafic am realizat o proiecție din 9D în 2D cu ajutorul algoritmului Sammon.

În concluzie mărimea rețelei neuronale în varianta fără strat ascuns este insuficientă pentru a putea învăța separarea tuturor

elementelor din

baza de date, lăsând aproximativ 20% dintre elemente neclasificate. Totuși este important de remarcat, că chiar și o rețea neuronală pulsativă atât de simplă atinge acest scor apreciabil, lucru greu de imaginat dacă am fi implementat o rețea neuronală clasică de dimensiuni identice, bazată pe modelul perceptron.

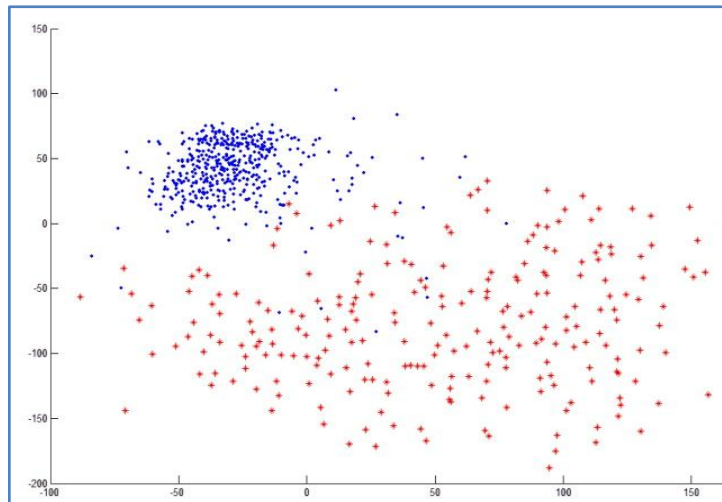


Figura 5.89 Reprezentarea grafică prin proiecție Sammon 9D → 2D a clasificării prescrise a bazei de date WBCD cu 683 mostre

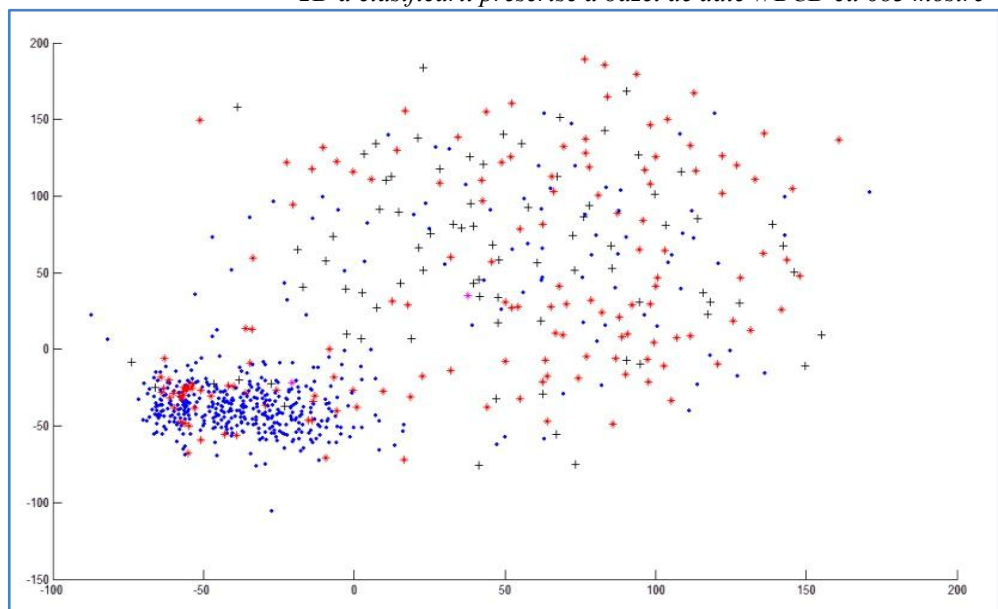


Figura 5.90 Rezultatul clasificării bazei de date WBCD cu RNP fără strat ascuns, implementată hardware, reprezentat grafic prin proiecție Sammon 9D → 2D

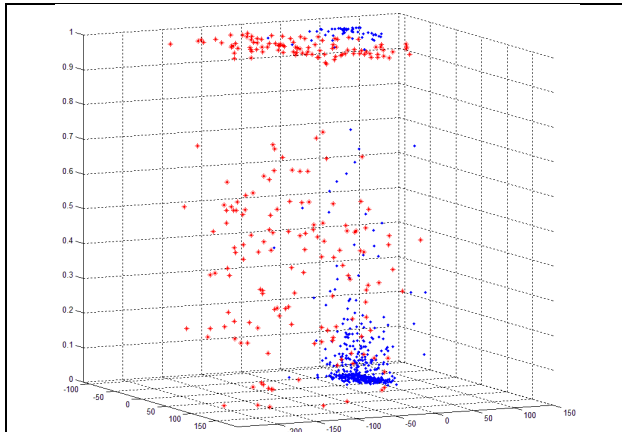


Figura 5.91 Stadiul clasificării într-un ciclu de antrenare intermediar (Varianta 2 a RNP-FPGA)

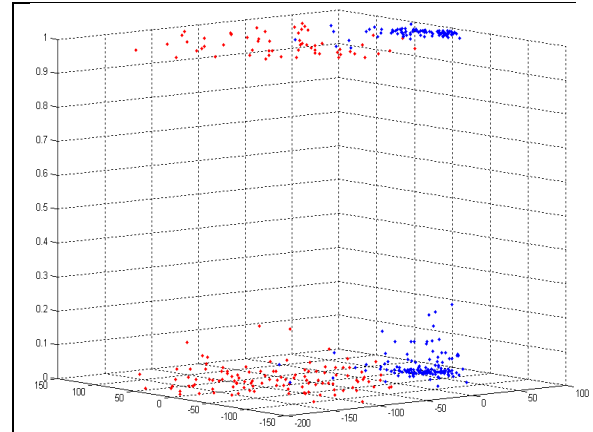
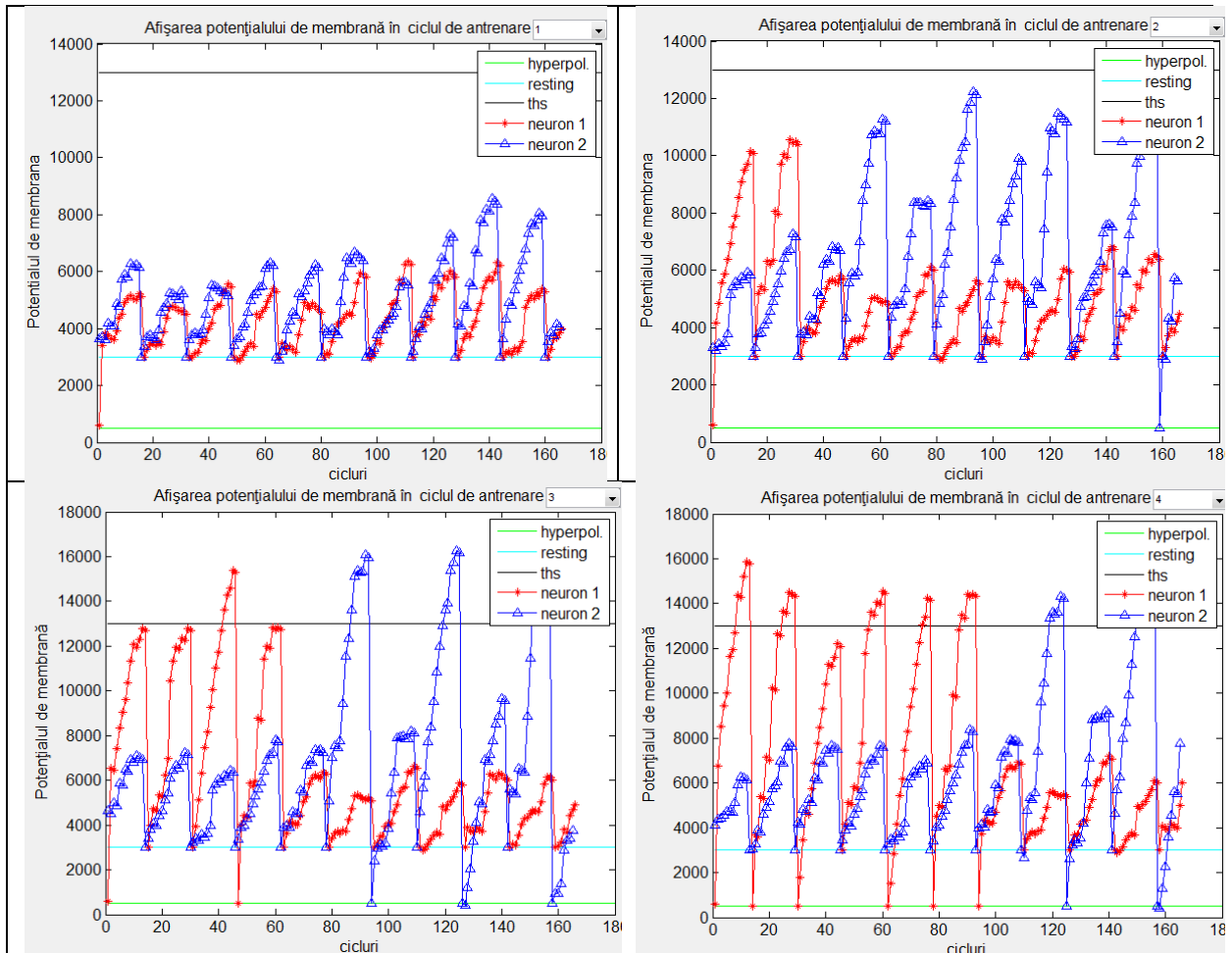


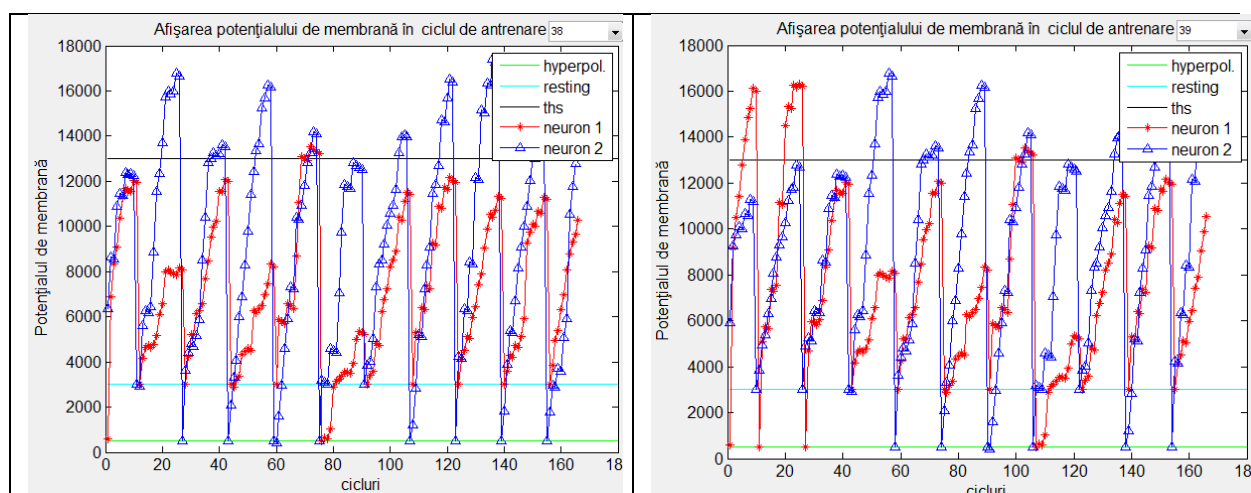
Figura 5.92 Stadiul clasificării în ciclul de antrenare 250 (Varianta 2 a RNP-FPGA)

Varianta 2 a RNP, cu strat ascuns, atinge un scor mult bun, clasificând toate elementele date, greșind clasa în proporție de ~10%. Graficele de mai jos arată – prin proiecție Sammon 9D → 3D – două măsurători în timpul antrenării, în pasul 100 (Figura 5.91) respectiv 250 (Figura 5.92) dintre cele 300 utilizate. Pe axele X, Z găsim valorile scalate generate de algoritmul Sammon (între -150...150 respectiv -200...10), iar pe axa Y a fost reprezentată procente de eroare cu care o anumită mostră a fost clasificată malign (roșu) sau benign (albastru).

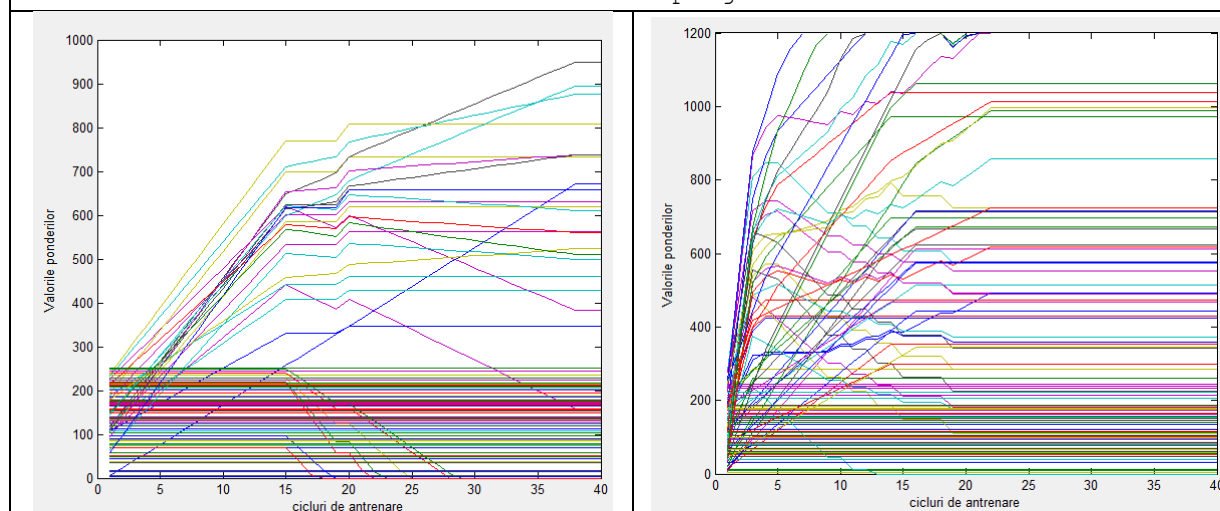
Tabelul 17 Măsurători în timpul procesului de învățare asupra variației potențialelor de membrană respectiv a ponderilor sinaptice







A se nota, că în primele cicluri de învățare PM nu depășește încă THS, iar în ciclurile ulterioare numai unul dintre cei doi neuroni de ieșire se va activa. Se poate observa, de asemenea intrarea neuronilor în starea de hiperpolarizare după fiecare impuls axonal emis la trecerea PM peste nivelul de prag.



În aceste figuri s-a reprezentat variația unui număr redus de ponderi sinaptice (rețeaua neuronală pulsativă implementată hardware conține mai multe sute de sinapse), observându-se că unele sinapse se vor stabiliza după un număr destul de redus de cicluri de antrenare.

Este important de menționat, că rezultatele prezentate mai sus sunt obținute de RNP implementată în hardware în timp real. Procesul de învățare se desfășoară mai lent, încetinit și de funcțiile de comunicare cu calculatorul prin portul serial RS232 al plăcii de dezvoltare FPGA (prin care se trimit datele ce arată funcționarea și comportamentul rețelei). După încheierea procesului de învățare, însă, sistemul neuronal implementat – fapt valabil la toate aplicațiile prezentate în această teză – este capabil să dea un rezultat într-un interval de timp măsurabil în milisecunde (chiar microsecunde la sistemele implementate complet paralel). Este deci o contribuție importantă, realizarea unor astfel de sisteme neuronale hardware cu procesare în timp real.

#### 5.6.4. Mediul de dezvoltare utilizat pentru sisteme încorporate

Embedded Development Kit este un pachet software dezvoltat de Xilinx pentru proiectarea sistemelor programabile încorporate. Pachetul include toate instrumentele, documentația, și module de circuite predefinite de care utilizatorul are nevoie pentru a proiecta sisteme care încorporează nuclee de procesor tip hard-core IBM PowerPC™, și/sau nuclee de procesor tip soft-core Xilinx MicroBlaze™.

Deși pare complex la o vedere în ansamblu, fluxul de proiectare a sistemului combină pur și simplu fluxul hardware standard folosit pentru a crea date de configurare a circuitelor FPGA și fluxul software standard folosit pentru a crea fișiere executabile de procesor ELF. De fapt, exceptând cazul în care resursele de memorie încorporate în circuitul de reconfigurare sunt folosite pentru a stoca imaginea softului, Embedded Development Kit poate fi considerat ca fiind nimic mai mult decât o extensie a instrumentului de generare nuclee CoreGen de la Xilinx.

Primul pas (Figura 5.93) constă în crearea unui „System Netlist” folosind Embedded Development Kit și instanțierea acestui „netlist” în codul HDL al proiectării. Proiectarea hardware este deci sintetizată, unită și implementată folosind exact același flux de proiectare ca cel folosit pentru orice alt nucleu „black box”. Deși se obișnuiește să se includă o porțiune deja creată a imaginii software în circuitul FPGA folosind blocul RAM, fișierul „Compiled BIT” creat în timpul acestei faze de dezvoltare conține doar descrierea hardware a sistemului.

Al doilea pas constă în crearea pachetului „Board Support Package” (BSP) folosind Embedded Development Kit (EDK) și includerea driverelor necesare în codul sursă C al sistemului. Codul sursă este apoi compilat și legat cu diversele funcții disponibile în BSP, el fiind identic cu al oricărui alt sistem procesor.

Deoarece sistemul încorporat este construit folosind resursele FPGA, BSP-ul este personalizat pentru acest set particular de periferice incluse în „System Netlist”. Spre deosebire de un procesor cu destinație generală larg răspândit, sistemele încorporate cu Virtex-II Pro pot include orice combinație de periferice produse de Xilinx sau create de utilizatori. Asta înseamnă că fiecare BSP are potențialul de a fi unic și astfel EDK are sarcina de a personaliza un set generic de drivere după nevoie pentru a asigura funcționarea optimă a „sistemului procesor arbitrar”.

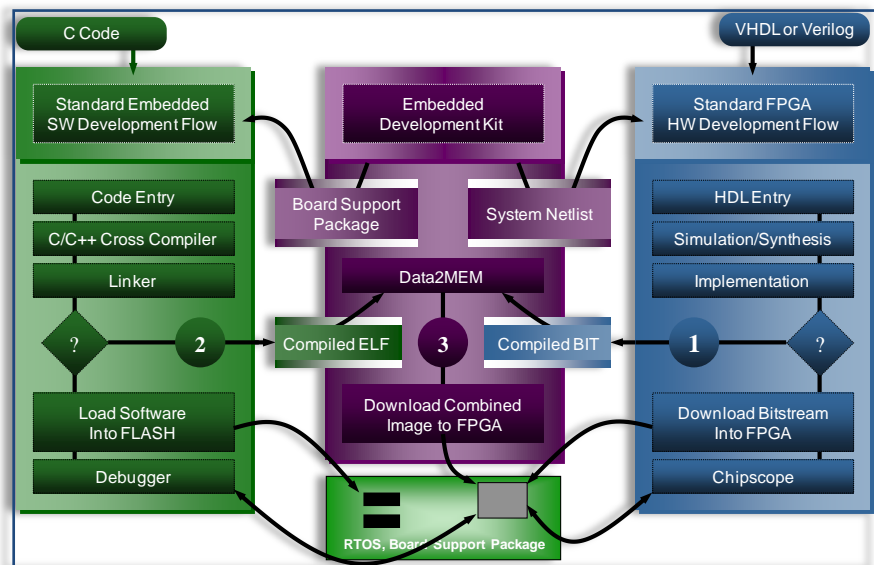


Figura 5.93 Pașii proiectării sistemelor încorporate în mediul Xilinx EDK

Odată ce setul final de periferice și structuri de magistrală au fost bine stabilite, fluxurile de hardware și software pot rula în mod independent. Chiar dacă o parte sau întreaga imagine software este stocată folosind blocul de memorie RAM încorporat în circuitul de reconfigurare, fluxul software nu necesită repornirea fluxului hardware de la zero atunci când se fac schimbări în software. Doar dacă se fac schimbări în „System Netlist”-ul instanțiat atunci hardware-ul trebuie implementat din nou înainte ca o imagine software nouă care se bazează pe schimbările arhitecturale operate să poată fi încărcată și rulată.

Dacă imaginea software este stocată complet extern, configurarea circuitului FPGA și încărcarea dispozitivului de stocare extern se efectuează la fel ca în cazul unei soluții tipice bazată pe două chip-uri. Dacă o parte sau întreaga imagine software este stocată folosind blocul de memorie RAM încorporat în circuitul reconfigurabil și este astfel inclusă în datele de configurare a circuitului FPGA, este nevoie de un pas în plus înainte ca circuitul FPGA să poată fi configurat. EDK pune la dispoziție un instrument denumit Data2MEM care unește secțiunile fișierului „Compiled ELF” (Figura

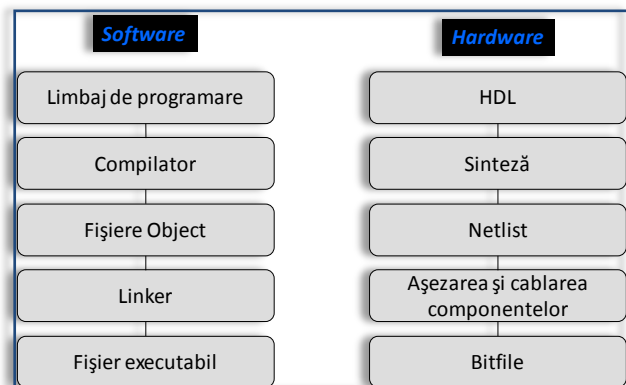


Figura 5.94 Fazele dezvoltării comune hardware-software

5.94) cu secțiunile potrivite din fișierul „Compiled BIT”. Fișierul BIT rezultat este creat de obicei în câteva secunde și poate fi apoi folosit pentru a configura circuitul FPGA. Când întreaga imagine software este stocată în interiorul circuitului FPGA, doar fișierul BIT este necesar pentru a configura sistemul și pentru a încărca imaginea software. Dacă doar o porțiune a imaginii software, cum ar fi segmentul de instrucțiuni, este stocată în interiorul circuitului FPGA, atunci se rulează Data2MEM pentru a crea fișierul BIT combinat și sistemul este gata pentru configurare / încărcare la fel ca orice altă soluție bazată pe două chip-uri, folosind secțiunile care nu au fost unite din fișierul ELF și secțiunile din fișierul BIT combinat.

Deplanarea softului care rulează pe sistem se realizează la fel ca în cazul oricărui procesor cu destinație generală. Se folosește GDB sau alt instrument de depanare pentru conectarea la platformă, asigurarea controlului de execuție al platformei și chiar pentru încărcarea unor noi imagini când se dorește acest lucru. Spre deosebire de procesoarele cu destinație generală, sistemele fizice pot fi testate folosind module ChipScope. Acest lucru asigură un nivel de transparență ridicat al modului de operare al sistemului, care nu poate fi egalat de procesoarele externe.

## 5.7. Concluzii

Acest capitol a prezentat cele mai recente rezultate ale cercetărilor autorului în domeniul implementărilor hardware de rețele neuronale neuromorfe. Pornind de la rezultatele expuse în capitolul anterior, s-a ajuns la concluzia, că implementarea cu succes a unor aplicații complexe bazate pe rețele neuronale pulsative nu este fezabilă în mod complet paralel din cauza necesarului excesiv de resurse reconfigurabile, care nu sunt disponibile în majoritatea circuitelor FPGA de azi. De aceea a fost nevoie de schimbare de direcție ușoară în continuarea cercetărilor, schimbare ce s-a materializat prin abordarea temei de serializare parțială a acestor implementări hardware. Conceptul de serializare paralelă în acest caz înseamnă introducerea unor elemente de execuție a unor algoritmi pentru a înlocui modulele care ar necesita un număr mare de resurse reconfigurabile, adică utilizarea procesoarelor încorporate de tip soft-core. S-au utilizat două astfel de procesoare (care sunt de fapt un set de programe VHDL sau Verilog ce implementează funcționalitatea acestora și pot fi încorporate în orice proiect) și anume Xilinx PicoBlaze și Xilinx MicroBlaze. PicoBlaze, cunoscut și ca KCPSM3, este un microcontroler RISC pe 8 biți, cu o amprență de siliciu extrem de redusă, de numai 96 de slice-uri dintr-un FPGA din familia Xilinx Spartan3 (optimizat pentru această familie) ceea ce reprezintă aproximativ 1,25% din totalul de 7680 de astfel de structuri logice disponibile de exemplu în circuitul XC3S1000. Acest necesar scăzut de resurse face ca PicoBlaze să fie candidatul ideal al implementărilor multi-core, fapt ce a fost exploatat în două dintre implementările din acest capitol. Pe de altă parte, nucleul MicroBlaze (pe 32 de biți) implementează o arhitectură de tip Harvard, cu interfețe de magistrală separate pentru acces de date și de instrucțiuni, mult mai complexă, fiind astfel corespunzător în cazul unor aplicații cu necesar de putere de calcul mai mare.

Realizările practice prezentate în acest capitol sunt de două feluri, aplicative și de test reper (benchmark) și au ca scop determinarea și ilustrarea capacității acestor sisteme inteligente inovative implementate hardware. Toate aplicațiile din acest capitol pot fi incluse în clasa generală a clasificărilor supervizate, cele aplicative fiind detectarea componentelor de frecvență a unui semnal analogic zgomotos (implementat utilizând multiple procesoare încorporate PicoBlaze) respectiv o aplicație de recunoaștere de caractere (realizată ca un sistem încorporat având ca procesor central un nucleu MicroBlaze). S-au implementat două teste benchmark clasice și anume testul de clasificare a setului de date Fisher Iris (multi-core PicoBlaze) și testul de clasificare a setului de date Wisconsin Breast Cancer Database (MicroBlaze).

În aceste aplicații s-a utilizat o metodă originală de codare a intrărilor numerice a rețelelor neuronale implementate prin codificarea acestora în impulsuri decalate temporal, codificare realizată printr-un model neuronal pseudo-RBF. Aceste module de codificare au fost implementate în fiecare aplicație utilizând resurse hardware reconfigurabile dedicate, calculele inerente fiind executate astfel în paralel, ajungându-se la un timp de codificare extrem de redus (~15 cicluri de tact la ~100MHz, adică ~150ns).

Pentru a putea evalua performanțele atinse de sistemele implementate, am studiat literatura de specialitate selectând rezultate de reper cu care se pot efectua comparații relevante. Există o serie de metode de clasificare des utilizate, dintre care se vor prezenta în continuare câteva – cu scop ilustrativ – pentru a înlesni analiza rezultatelor expuse ulterior.



**Algoritmul K-means clustering:** K-medii este unul din cei mai simplii algoritmi de învățare nesupravegheați care rezolvă bine cunoscuta problemă de clustering. Procedeele urmează o cale simplă și ușoară pentru a clasifica setul de date de intrare într-un număr de K grupe (clusters) (K fixat a priori). Ideea de bază este de a defini K centre de greutate, câte unul pentru fiecare grupă. Aceste centre de greutate trebuie fixate rațional, deoarece locații diferite conduc la rezultate diferite. Alegerea cea mai bună este să le fixăm cât este posibil mai depărtate unele de altele. Următorul pas este să luăm fiecare element al setului de date și să-l asociem celui mai apropiat centru de greutate. Prima etapă a grupării se termină când nu mai există elemente negrupate. În acest moment trebuie să calculăm K noi centre ale grupelor rezultate din pasul anterior. Procedeele continuă până în momentul în care pozițiile noilor centre nu se mai modifică semnificativ.

**Algoritmul K-medoids:** K-medoids este un algoritm de clustering înrudit cu algoritmul *k*-means și cu algoritmul medoidshift. Ambii algoritmi *k*-means și *k*-medoids sunt algoritmi de partiționare (împărțind setul de date în grupuri) și ambii încearcă să minimizeze eroarea pătratică, distanța dintre puncte marcate ca făcând parte dintr-un grup (cluster) și un punct desemnat ca fiind centru al aceluși grup. Spre deosebire de algoritmul *k*-means, algoritmul *k*-medoids alege elemente ale setului de date ca centre (medoids sau exemplare). *k*-medoid este o tehnică de partiționare clasică care grupează setul de date de *n* obiecte în *k* grupuri fixate *a priori*. Un instrument folosit pentru determinarea lui *k* este metoda *silhouette*. Este mai robust în condiții de zgomot și valori extreme în comparație cu *k*-means. Un medoid poate fi definit ca un obiect al unui grup, a cărui diferență medie față de toate obiectele aceluși grup este minimă, adică este punctul situat cel mai central din setul de date respectiv.

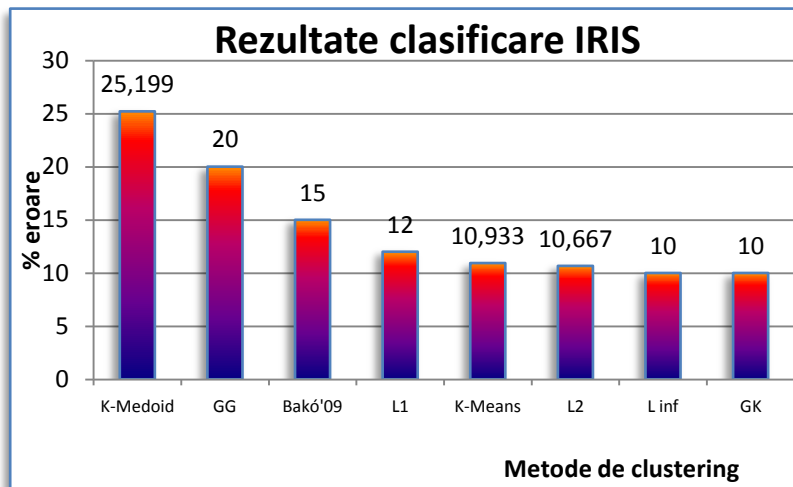


Figura 5.95 Comparatie a rezultatelor clasificării setului de date Fisher IRIS cu diferite metode similare

**Algoritmul L1 și adâncimea relativă a datelor:** În (Jörnsten,

Vardi, & Zhang, 2002) s-a introdus adâncimea relativă a datelor, ReD, ca un instrument de validare a grupurilor folosit împreună cu un algoritm K-median exact. Adâncimea datelor L1 a unei observații relative la un grup de date indică cât de reprezentativă este această observație pentru grupul respectiv. Astfel, pentru date cu etichetă de clasă ne așteptăm ca adâncimea datelor L1 să fie maximizată față de grupul de date corespunzător etichetei de clasă corecte. O regulă simplă de clasificare este să clasificăm observațiile cu etichete necunoscute după maximul adâncimii

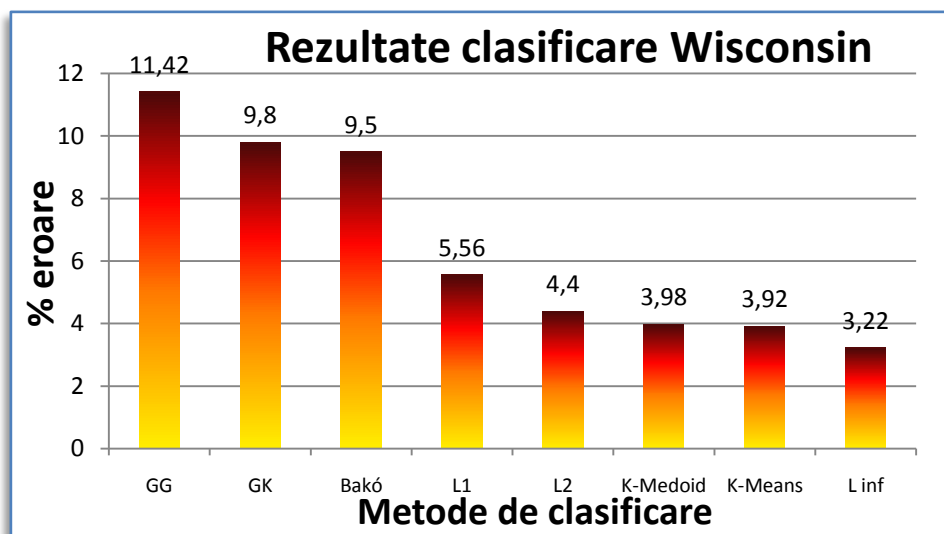


Figura 5.96 Comparatie a rezultatelor clasificării setului de date WBCD cu diferite metode similare

datelor față de datele etichetate. Ne referim la datele etichetate ca fiind setul de antrenare, iar datele neetichetate ca fiind setul de testare.

Graficul comparativ din Figura 5.95 prezintă performanțele implementării benchmark a clasificării setului de date Fisher IRIS, iar Figura 5.96 rezultatele obținute cu circuitul neuronal pulsativ ce implementează clasificarea setului de date Wisconsin Breast Cancer Database. Deși rezultatele afișate nu excelează față de celelalte metode în mod special, trebuie luat în considerare faptul extrem de important, că toate sistemele hardware realizate și prezentate în această teză funcționează cu învățare on-chip și procesează calculele necesare în timp real.

*Tabelul 18 Comparatie a mai multor algoritmi de clasificare rulate în software cu algoritmul implementat hardware pentru clasificare seturilor de date benchmark*

<b>Setul de date Fisher IRIS</b>						
					<b>precizie</b>	
<b>algoritm</b>	<b>intrări</b>	<b>ascuns</b>	<b>ieșiri</b>	<b>iterații</b>	<b>set antrenare</b>	<b>set de test</b>
<i>SpikeProp</i>	50	10	3	1000	96,4%	95,3%
<i>Matlab BP</i>	50	10	3	2.6*10 <sup>6</sup>	97,9%	95,8%
<i>Matlab LM</i>	50	10	3	3750	98,8%	95,9%
<i>RNP<sub>Bakó</sub></i>	4	0	3	300	85%	83,4%
<b>Setul de date Wisconsin Breast Cancer</b>						
					<b>precizie</b>	
<b>algoritm</b>	<b>intrări</b>	<b>ascuns</b>	<b>ieșiri</b>	<b>iterații</b>	<b>set antrenare</b>	<b>set de test</b>
<i>SpikeProp</i>	64	15	2	1500	97,6%	97,3%
<i>Matlab BP</i>	64	15	2	9.2*10 <sup>6</sup>	97,8%	96,2%
<i>Matlab LM</i>	64	15	2	3500	98,1%	96,5%
<i>RNP<sub>Bakó</sub></i>	9	6	2	250	90,2%	89,5%

Rezultatele prezentate în Tabelul 18 au fost obținute prin măsurători proprii (Matlab, RNP) respectiv din literatura de specialitate (Bohte, Kok, & La Poutré, Spike-prop: errorbackpropagation in multi-layer networks of spiking neurons, 2000). Ambele seturi de date au fost împărțite în set de antrenare și set de test. Algoritmul Matlab LM a fost rulat pe 50 de epoci cu 1500 de cicluri fiecare. Simularea software prin algoritmul de propagare înapoi a erorii a fost implementată în mediul Matlab, utilizând rutinele „TRAINLM” (LM) respectiv „TRAINGD” (BP).

Implementările descrise în acest capitol au fost publicate într-un număr de lucrări științifice (Bakó L. , MACRo 2009, 2009), (Bakó L. , 2009), (Bakó & Brassai, 2009), (Bakó & Székely, Challenges for implementations of delay-coded neuromorphic neural networks on embedded digital hardware, 2009), (Bakó & Székely, 2009), (Bakó L. , Brassai, Székely, & Baczó, 2008), ale autorului, una din ele aducând și un premiu la o conferință internațională (“*The Best Presenter in Information Technology*” la Conferința The 4<sup>th</sup> International PhD, DLA Symposium, organizat de University of Pécs, Pollack Mihály Faculty of Engineering, Pécs, Ungaria, 20-21 Octombrie 2008.).

## 6. CONCLUZII FINALE ȘI CONTRIBUȚII ORIGINALE

### 6.1. Concluzii finale

Rezultatele cercetărilor neurobiologice au condus la apariția celei de-a treia generații de rețele neuronale artificiale, care se străduiesc să construiască modele cât mai fidele ale rețelelor neuronale naturale. Dintre posibilitățile oferite de aceste modele, am ales cele ce se bazează pe proprietatea pulsativă a neuronilor naturali. Studiarea acestor neuroni și realizarea de modele teoretice și modele practice cât mai fidele ale acestora, care sunt în același timp implementabile în mod eficient pe platforme digitale reconfigurabile a reprezentat scopul cercetării efectuate de-a lungul desfășurării programului de doctorat. Teza de față încearcă să surprindă cele mai importante și relevante teorii și aplicații existente în acest domeniu de cercetare – care stă de fapt la confluența științelor sistemelor de calcul, electronicii digitale și a biologiei neuronale – descrise de literatura de specialitate. Deși bazele teoretice au fost puse și în cazul acestui domeniu cu mulți ani în urmă (mai ales în ceea ce privește modelarea matematică a proceselor neurobiologice), după cum reiese și din sintezele prezentate în teză, lucrările aplicative cu implementări software și/sau hardware bazate pe aceste modele sunt foarte recente.

Se poate concluziona, deci, că studiul rețelelor neuronale neuromorfe, și implementarea de sisteme hardware cu ajutorul acestora este un domeniu de cercetare foarte proaspăt, cu puține rezultate răsunătoare și cu multe direcții de evoluare abordate de diferitele grupuri de oameni de știință, rezultând posibilități foarte importante de a realiza contribuții originale semnificative. De fapt, prin studierea constantă în ultimii șase ani a lucrărilor apărute în acest sens pe plan mondial, am găsit, că există numai câteva centre de cercetare sunt specializate pe asemenea sisteme, cum ar fi École Polytechnique Fédérale de Lausanne din Elveția, sau alte grupuri de cercetători din Mexic, Argentina, SUA sau Irlanda.

În acest context autorul acestei teze a încercat o abordare originală în acest domeniu, prin dezvoltarea de modele neuronale proprii, cât mai compatibile cu implementare pe circuite FPGA și demonstrarea prin aplicații proprii – aplicative și de tip benchmark – capacitatea de calcul și de asimilare de informație extraordinară a acestor sisteme, evidențiind totodată viteza excepțional de rapidă de care sunt capabile acestea.

Pentru a avea o privire de ansamblu asupra conținutului tezei de doctorat se prezintă în continuare cele mai importante aspecte abordate de fiecare capitol a lucrării în parte.

**Primul capitol** *introdactiv* prezintă aspecte generale ale domeniului mai larg de cercetare în care se încadrează teza, cea a inteligenței artificiale și a mașinilor cu autoinstruire. Fundamentele biologice ale cercetării sunt detaliate de asemenea în acest capitol, urmat de o descriere sumară a uneia din cele mai avansate domenii de aplicare a inteligenței artificiale, și anume sistemele de control inteligent, domeniu în care autorul tezei a publicat o serie de lucrări științifice (Brassai & Bakó, Visual trajectory control of a mobile robot using FPGA implemented neural network, 2008), (Brassai & Bakó, 2008), (Brassai S. , Bakó, Székely, & Dan, 2008), (Brassai, Márton, Dávid, & Bakó, 2008), (Brassai, Gidró, Bakó, & Csernáth, 2008), (Brassai & Bakó, 2007), (Brassai, Bakó, & Dan, 2007), (Brassai, Dávid, & Bakó, 2004).

În prima parte a **Capitolului 2**, intitulat „*Modelarea și simularea rețelelor neuronale neuromorfe*” s-a discutat modelarea rețelelor neuronale neuromorfe, expunând comportamentul intrinsec temporal al acestora, utilizarea impulsurilor în comunicarea datelor în cadrul acestora și evidențiind capacitatea de calcul mai ridicată decât în cazul rețelelor neuronale bazate pe modele clasice.

Studiul neuronilor biologici a rezultat în apariția mai multor ipoteze pentru codarea valorilor reale în trenuri de impulsuri. Acestea pot apărea ca o codare a ratei impulsurilor, care ia în considerare numai valoarea medie a timpului de apariție a unui impuls sau ca codarea exactă a momentelor temporale de apariție a acestor impulsuri.

Când vorbim despre RNP există de fapt o supraabundență de diferite modele, de la cele simple de tipul integrează și activează prin neuroni cu diferite modele sinaptice până la modele comportamentale complete descrise prin ecuații diferențiale. Există câteva studii teoretice (Maass, 1996)

care arată că unele dintre aceste modele sunt mai eficiente decât celelalte, dar aceste studii se limitează doar la un set redus de modele și scheme de codare. O explicație generală a capacității de calcul a diferitelor modele neuronale pulsative și a schemelor de codare conexe reprezintă și astăzi o temă deschisă cercetărilor actuale.

În orice studiu legat de dinamica unei rețele neuronale pulsative, sunt două probleme specifice, și anume ce model descrie dinamica pulsației fiecărui neuron respectiv cum sunt conectați neuronii. O alegere greșită a modelului sau o conexiune greșită, poate duce la rezultate complet eronate.

În cadrul acestui capitol s-a elaborat prima problemă, comparând neuronii pulsatori. Se vor prezenta diferite modele de neuroni pulsatori precum și o clasificare a acestora.

În continuare am prezentat cele mai des utilizate metode de antrenare (ne-supervizate și supervizate) utilizate în învățarea rețelelor neuronale neuromorfe. Astfel, în această parte a capitolului se regăsesc discuții asupra realizării învățării Hebbiene competitive precum și asupra modalităților specifice de aplicare a regulii back-propagation în rețele neuronale pulsative codificate temporal (algoritmul SpikeProp).

După acest studiu amănunțit al acestui domeniu de inteligență artificială se prezintă modelul neuronal teoretic ales pentru implementările hardware ce urmează a fi prezentate în capitolele următoare.

A doua contribuție importantă a acestui capitol este prezentarea generală a strategiilor de simulare software a RNP împreună cu descrierea celor mai importante astfel de medii existente la ora actuală, dar și a simulatorului neuronal implementare proprie, cu care am testat funcționarea modelului dezvoltat. Se detaliază încorporarea algoritmilor de învățare în procesul de simulare software și se expun rezultatele experimentale obținute pe această cale, care au stat la baza primelor publicații ale autorului (Bakó, Székely, & Brassai, Development of Advanced Neural Models. Software And Hardware Implementation, 2004) (Bakó, Analiza și simularea sistemelor cu inteligență artificială neuromorfă, 2005) (Bakó, Székely, Dávid, & Brassai, 2004) în acest domeniu de cercetare.

**Capitolul 3**, cu titlul „*Dispozitive FPGA și implementarea rețelelor neuronale artificiale*” este o sinteză cuprinzătoare a evoluției dispozitivelor reconfigurabile, realizând clasificări și evaluări de performanțe ale acestora, dar descriind și funcționarea celor mai importante elemente constructive ale arhitecturilor de circuite FPGA.

Creșterea continuă în densitate (număr porți/suprafață) a sistemelor FPGA a făcut posibilă realizarea unor proiecte de sisteme on-chip cu o complexitate mai mare de un milion de porți și RAM intern. De aceea, sistemele FPGA s-au dovedit în prezent a fi o platformă hardware atractivă pentru algoritmi RNA care necesită mult spațiu. Un alt avantaj al acestui dispozitiv este capacitatea de a combina programabilitatea cu viteza crescută a operațiilor asociate cu soluțiile hardware paralele.

Arhitecturile bazate pe circuite FPGA pot fi exploatate bine pentru realizarea de rețele neuronale artificiale, deoarece utilizând capacitatea lor de reconfigurare rapidă și concurentă se poate ajunge la performanțe bune în adaptabilitatea ponderală și chiar topologică. Cu toate acestea, luând în calcul că prima astfel de aplicație (Cox & Blanz, 1992) a apărut numai cu aproximativ un deceniu în urmă, realizarea de rețele neuronale de dimensiuni mari, cu mulți neuroni în circuite FPGA este o mare provocare chiar și în zilele noastre. Acest lucru se datorează faptului că majoritatea modelelor neuronale clasice necesită un număr mare de operații de multiplicare, ceea ce duce la realizări hardware foarte costisitoare (multe circuite multiplicatoare complexe). În pofida acestei probleme, există tehnologii care exploatează proprietatea de reconfigurabilitate a acestor circuite și reușesc să implementeze ieftin și eficient RNA în acestea.

Orice realizare de rețea neuronală cu FPGA trebuie să se străduiască să utilizeze într-un anumit fel configurabilitatea acestor circuite, implicit optimizarea lor pentru problema propusă. Această proprietate care se poate aplica într-un timp scurt și de nenumărate ori, este de folos în special în dezvoltarea de *prototipuri și simulări*. Această flexibilitate permite realizarea rapidă de rețele neuronale cu diferite structuri și strategii de învățare, dar și demonstrarea unor principii prin simulări primordiale.

Procedurile de *creștere a densității* realizează o mai bună funcționalitate raportată la unitatea de arie FPGA, utilizând reconfigurarea. Acest lucru se poate valorifica dacă dispunem de un circuit FPGA capabil de reconfigurare – cel puțin parțială – în timpul funcționării (run-time/partial reconfigurability).

Primul tip de astfel de procedură este *multiplexarea în timp*, ceea ce înseamnă că simularea RNA este descompusă în mai multe faze secvențiale, care corespund fiecare unei configurări diferite ale aceluiași circuit FPGA. Astfel se poate realiza un sistem în care circuitul este întotdeauna optimizat pentru faza care este în execuție.

Al doilea tip de procedură pentru creșterea densității de implementare se bazează pe execuția de operații cu multiplicatori constanți și se numește *schimbare dinamică de constante*. Deoarece aceste două proceduri au avut ca scop primordial creșterea densității de implementare, în ceea ce privește o performanță superioară, nu putem avea pretenții față de aceste sisteme, numai în cazul în care timpul de reconfigurare a circuitului este neglijabil în comparație cu timpul utilizat pentru efectuarea calculului.

În circuitele FPGA reconfigurabile dinamic, există posibilitatea de implementare a rețelelor neuronale cu *adaptare topologică*. Altfel spus, putem construi RNA care suferă modificări arhitecturale succesive. În faza de învățare, pe lângă topologia rețelei avem posibilitatea de a acorda și precizia calculului, conform funcției de criteriu utilizate de algoritmul de învățare.

Neurohardware-ul digital diferă de asemenea în reprezentarea numerică a operațiilor aritmetice. Implementările în virgulă flotantă oferă o precizie mai bună, dar mai complexă și cu necesități de suprafață mai mare. Aritmetica în virgulă fixă este mai puțin complexă, cu necesități de suprafețe mai reduse, dar oferă o precizie scăzută față de implementarea în virgulă flotantă. Aritmetica în virgulă fixă poate fi utilă în cazul în care un algoritm RNA nu necesită precizie ridicată pentru implementarea unor aplicații simple.

**Al patrulea Capitol** al tezei de doctorat, „*Posibilități de implementare total paralelă a modelelor neuronale pulsative cu circuite FPGA*” începe prin a prezenta caracteristici generale ale sistemelor hardware analogice și digitale în ceea ce privește implementarea rețelelor neuronale artificiale, aducând motivații pro și contra în fiecare caz. De asemenea sunt prezentate cele mai importante implementări hardware de rețele neuronale din literatura de specialitate actuală cu accent asupra realizărilor pe diferite platforme digitale. Se trece apoi la domeniul specific temei acestei lucrări, prezentând implementări software și hardware ale RNP realizate pe calculatoare paralele sau chiar în circuite FPGA.

După această parte de sinteză se prezintă contribuțiile proprii în ceea ce privește implementarea complet paralelă de RNP în circuite FPGA. Inițial se expune implementarea modelului teoretic adaptat realizării hardware, validat prin simulare software în capitolele anterioare. Aceste subcapitole conțin descrieri amănunțite ale proiectării diferitelor subansamble ale acestor rețele neuronale, cum ar fi sinapsele și modulele de somă sau corp celular. Acestea sunt testate prin simularea funcționării circuitelor proiectate pentru a implementa funcționalitatea specifică acestora.

O aplicație de test a acestor circuite neuronale realizate reprezintă continuarea logică a acestui capitol, aplicația fiind cea clasică de diferențiere a două forme simple cum ar fi literele **T** și **H** sau semnele + și X. Se descrie rețeaua neuronală care implementează în circuit FPGA prin resurse hardware dedicate fiecărui subansamblu al acestuia (implementare complet paralelă), specificând algoritmul de învățare utilizat și enumerând rezultatele măsurătorilor experimentale. În cadrul unei analize a performanțelor ( $32\mu s$  - timpul de învățare a rețelei,  $3,2\mu s$  timp de recunoaștere a formelor de intrare) acestei RNP se prezintă și o comparație cu o rezolvare pur software (în mediul Matlab) a problemei propuse, rezultând performanțe net superioare ale versiunii hardware, care funcționează în timp real.

Rezultatele și contribuțiile personale prezentate în acest capitol au fost publicate într-un număr de lucrări științifice ale autorului tezei, apărute în publicații ale unor conferințe de specialitate (Bakó & Brassai, Hardware spiking neural networks: parallel implementations using FPGAs, 2006), (Bakó & Brassai, 2005), (Brassai & Bakó, 2007), în reviste de specialitate (Bakó & Brassai, Spiking neural networks built into FPGAs: Fully parallel implementations, 2006) (Bakó, Székely, & Brassai, Development of Advanced Neural Models. Software And Hardware Implementation, 2004), precum și în rapoartele unor granturi de cercetare în care autorul a fost membru.

**Capitolul 5** intitulat „*Implementări parțial paralele ale RNP utilizând microcontrolere încorporate în circuite FPGA*” a prezentat cele mai recente rezultate ale cercetărilor autorului în domeniul implementărilor hardware de rețele neuronale neuromorfe. Pornind de la rezultatele expuse în capitolul anterior, s-a ajuns la concluzia, că implementarea cu succes a unor aplicații complexe bazate pe rețele neuronale pulsative nu este fezabilă în mod complet paralel din cauza necesarului excesiv de resurse reconfigurabile, care nu sunt disponibile în majoritatea circuitelor FPGA de azi. De aceea a fost nevoie de schimbare de direcție ușoară în continuarea cercetărilor, schimbare ce s-a materializat prin abordarea temei de serializare parțială a acestor implementări hardware. Conceptul de serializare paralelă în acest caz înseamnă introducerea unor elemente de execuție a unor algoritmi care ar necesita un număr mare de resurse reconfigurabile, adică utilizarea procesoarelor încorporate de tip soft-core. S-au utilizat două astfel de procesoare (care sunt de fapt un set de programe VHDL sau Verilog ce implementează funcționalitatea acestora și pot fi încorporate în orice proiect) și anume Xilinx PicoBlaze și Xilinx MicroBlaze. PicoBlaze, cunoscut și ca KCPSM3, este un microcontroler RISC pe 8 biți, cu o amprență de siliciu extrem de redusă, de numai 96 de slice-uri dintr-un FPGA din familia Xilinx Spartan3 (optimizat pentru această familie) ceea ce reprezintă aproximativ 1,25% din totalul de 7680 de astfel de structuri logice disponibile de exemplu în circuitul XC3S1000. Acest necesar scăzut de resurse face ca PicoBlaze să fie candidatul ideal al implementărilor multi-core, fapt ce a fost exploatat în două dintre implementările din acest capitol. Pe de altă parte, nucleul MicroBlaze (pe 32 de biți) implementează o arhitectură de tip Harvard, cu interfețe de magistrală separate pentru acces de date și de instrucțiuni, mult mai complexă, fiind astfel corespunzător în cazul unor aplicații cu necesar de putere de calcul mai mare.

Realizările practice prezentate în acest capitol sunt de două feluri, aplicative și de test reper (benchmark) și au ca scop determinarea și ilustrarea capacității acestor sisteme inteligente inovative implementate hardware. Toate aplicațiile din acest capitol pot fi incluse în clasa generală a clasificărilor supervizate, cele aplicative fiind detectarea componentelor de frecvență a unui semnal analogic zgomotos (implementat utilizând multiple procesoare încorporate PicoBlaze) respectiv o aplicație de recunoaștere de caractere (realizată ca un sistem încorporat având ca procesor central un nucleu MicroBlaze). S-au implementat două teste benchmark clasice și anume testul de clasificare a setului de date Fisher Iris (multi-core PicoBlaze) și testul de clasificare a setului de date Wisconsin Breast Cancer Database (MicroBlaze).

În fiecare dintre aceste aplicații s-a utilizat o metodă originală de codificare a intrărilor numerice a rețelelor neuronale implementate prin transformarea acestora în impulsuri decalate temporal, codificare realizată printr-un model neuronal pseudo-RBF. Aceste module de codificare au fost implementate în fiecare aplicație utilizând resurse hardware reconfigurabile dedicate, calculele inerente fiind executate astfel în paralel, ajungându-se la un timp de codificare extrem de redus (~15 cicluri de tact la ~100MHz, adică ~150ns).

Graficele comparative din finalul acestui capitol prezintă performanțele implementării benchmark a clasificării setului de date Fisher IRIS respectiv a circuitului neuronal pulsativ ce implementează clasificarea setului de date Wisconsin Breast Cancer Database. Deși rezultatele afișate nu excelează față de celelalte metode în mod special, trebuie luat în considerare faptul extrem de important, că toate sistemele hardware realizate și prezentate în această teză funcționează cu învățare on-chip și procesează calculele necesare în timp real.

Rezultatele prezentate în Tabelul 18 au fost obținute prin măsurători proprii (Matlab, RNP) respectiv din literatura de specialitate (Bohte, Kok, & La Poutré, Spike-prop: errorbackpropagation in multi-layer networks of spiking neurons, 2000). Ambele seturi de date au fost împărțite în set de antrenare și set de test. Algoritmul Matlab LM a fost rulat pe 50 de epoci cu 1500 de cicluri fiecare. Simularea software prin algoritmul de propagare înapoi a erorii a fost implementată în mediul Matlab, utilizând rutinele „TRAINLM” (LM) respectiv „TRAINGD” (BP).

Implementările descrise în acest capitol au fost publicate într-un număr de lucrări științifice (Bakó L. , MACRo 2009, 2009), (Bakó L. , 2009), (Bakó & Brassai, 2009), (Bakó & Székely, Challenges for implementations of delay-coded neuromorphic neural networks on embedded digital hardware, 2009), (Bakó & Székely, 2009), (Bakó L. , Brassai, Székely, & Baczó, 2008), ale autorului.

## 6.2. Contribuții originale

În acest ultim subcapitol al acestei lucrări se vor trece în revistă contribuțiile originale ale tezei de doctorat.

- Am realizat un studiu de sinteză privind aspecte generale ale domeniului mai larg de cercetare în care se încadrează teza, cea a inteligenței artificiale și a mașinilor cu autoinstruire. Fundamentele biologice ale cercetării sunt detaliate și se descrie sumar una din cele mai avansate domenii de aplicare a inteligenței artificiale, și anume sistemele de control inteligent.
- Am realizat un studiu cu privire la alegerea optimă a modelelor pentru rețele neuronale artificiale neuromorfe.
- Un model neuronal pulsativ propriu a fost dezvoltat care se pretează în mod intrinsec implementărilor hardware digitale.
- S-a efectuat o evaluare a celor mai importante reguli de învățare nesupervizate aplicabile în implementări hardware de rețele neuronale precum și a singurei metode supervizate descrise în literatura de specialitate, SpikeProp, care s-a considerat a nu fi fezabilă din același punct de vedere.
- Sinteza realizată asupra simulatoarelor neuronale ale rețelelor bazate pe modele pulsative prezintă strategiile ce trebuie abordate în astfel de sisteme, metode pentru simulări sinaptice rapide precum și motivele pro și contra în alegerea tipului de algoritm ales (sincron – bazat pe tact sau asincron – bazat pe eveniment). Această sinteză înglobează și câte o prezentare generală a celor mai semnificative medii de simulare software de rețele neuronale pulsative.
- Am realizat un simulator software propriu bazat pe un pachet de funcții open-source (rulate pe sistem de operare Linux) pentru validarea modelului neuronal dezvoltat. S-au prezentat ideile de bază ale acestui program, algoritmul simulării, metoda de învățare bazată pe valori de prag utilizată dar și rezultatele experimentale obținute cu acest sistem.
- S-a realizat un studiu cu privire la implementarea rețelelor neuronale artificiale pe sisteme FPGA:
  - s-au prezentat diferite moduri de abordare de cuplare și configurări pentru FPGA
  - s-a prezentat maparea diferiților algoritmi de rețele neuronale pe sisteme FPGA
  - s-au sintetizat modurile de implementare a funcțiilor de activare,
  - s-au prezentat diferite aritmetici de reprezentare a datelor.
- Într-un subcapitol s-au discutat avantajele și dezavantajele diferitelor platforme de implementare hardware ale rețelelor neuronale, prezentând o serie de astfel de implementări existente.
- S-au expus separat și exemple de implementări hardware și software a rețelelor neuronale pulsative.
- Implementarea cu circuit FPGA a modelului neuronal propriu dezvoltat, validat prin simulare software este de asemenea o contribuție importantă a tezei. Se descriu în detaliu procesele de proiectare premergătoare și cele de implementare pe circuit FPGA. Fiecare componentă a fost testată în simulator hardware, urmat de asamblarea primului circuit digital propriu ce realizează un neuron pulsativ.
- Pe baza acestui concept dezvoltat de realizare a unui neuron neuromorf digital s-a prezentat o primă aplicație de test (diferențierea unor modele de intrare simple, de ex. **T** și **H**) realizată cu o RNP simplă, implementată complet paralel în FPGA.
- Utilizând rezultatele obținute cu RNP implementată complet paralel s-a realizat o comparație cu o soluție pur software (rulat pe PC) a aceleiași probleme.
- Am dezvoltat un concept de realizare a unor circuite digitale pentru codificarea valorilor de intrare a RNP în impulsuri decalate temporal, pentru a putea fi interpretate și prelucrate de aceste rețele neuronale speciale. Procedul este foarte flexibil, putând fi aplicat la o gamă largă de aplicații cu spații de intrare de diferite dimensiuni (spații definite de plaja valorilor de intrare).
- Am realizat un studiu asupra fezabilității implementărilor parțial paralele în hardware digital a rețelelor neuronale pulsative prin introducerea utilizării procesoarelor încorporate de tip soft-core.



- În cadrul realizărilor practice parțial paralele am implementat o RNP pentru detectarea componentelor de frecvență a unui semnal analogic perturbat de zgomot pe baza spectrului acestui semnal prezentat la intrările acestei rețele neuronale, realizare cu some neuronale implementate prin procesoare încorporate PicoBlaze.
- Prima aplicație benchmark implementată a fost cea de clasificare a setului de date Fisher IRIS. Ca și în cazul aplicației de detectare a frecvențelor, sistemul neuronal rezultat este unul de tip multi-core, deoarece conține patru instanțe ale procesorului PicoBlaze.
- Trecând la procesoare încorporate soft-core cu putere de calcul mai ridicată (MicroBlaze) am implementat o aplicație de recunoaștere de caractere, prezentând rezultatele experimentale obținute.
- A doua aplicație benchmark este clasificarea setului de date Wisconsin Breast Cancer, care a fost realizată fuzionând circuitul de codificare a valorilor de intrare în impulsuri decalate temporal cu procesorul MicroBlaze, rezultând un sistem foarte rapid și cu performanțe în clasificare comparabile cu cele ce se regăsesc în literatura de specialitate.
- Am realizat o comparație între performanțele măsurate prin experimente ale sistemelor neuronale pulsative realizate cu utilizarea procesoarelor încorporate și a performanțelor unor metode cunoscute.

## BIBLIOGRAFIE

1. Abbott, L. F., & Nelson, S. B. (2000). Synaptic plasticity: taming the beast. *Nature Neuroscience*, 3 (Suppl), 1178-1283.
2. Abramson, D., Smith, K., Logothetis, P., & Duke, D. (1998). FPGA based implementation of a Hopfield neural network for solving constraint satisfaction problem. *Proceedings of 24th euromicro workshop on computational intelligence*, (pg. 688-693). Sweden.
3. Agis, R., Díaz, J., Ros, E., Carrillo, R., & Ortigosa, E. (2006). Event-Driven Simulation Engine for Spiking Neural Networks on a Chip. In *Book Chapter of Book Series Lecture Notes in Computer Science* (Vol. 3985). SpringerLink.
4. Alippi, C., & Storti-Gajani, G. (1991). Simple Approximation of sigmoidal functions: Realistic designs of digital neural networks capable of learning. *IEEE International Symposium on Circuits and Systems*, (pg. 1505-1508). Singapore.
5. Amdahl, G. (1967). Validity of single-processor approach to achieving large-scale computing capability. *Proceedings of AFIPS Conference*, (pg. 483-485). Reston, VA.
6. Asanovic, e. a. (2006, Dec 18). The Landscape of Parallel Computing Research: A View from Berkeley. *Tech. Report No. UCB/EECS-2006-183*.
7. Ashenden, P. (1990). *The VHDL Cookbook (First Edition)*. Adelaide, South Australia: Dept. Computer Science, University of Adelaide.
8. Badoual, M., Rudolph, M., Piwkowska, Z., Destexhe, A., & Bal, T. (2005). High discharge variability in neurons driven by current noise. *Neurocomputing*, 65, 493-498.
9. **Bakó, L. (2005)**. *Analiza și simularea sistemelor cu inteligență artificială neuromorfă*. Brasov, Ro: Referat Nr. 1, din cadrul programului de doctorat, Conducător științific: Prof.dr.ing. Iuliu Székely.
10. **Bakó, L. (2009)**. Partially Serialized Computation in Networks of Pulse-based Artificial Neurons. *1st International Conference on Recent Achievements in Mechatronics, Automation, Computer Science and Robotics, MACRo 2009* (p. 19). Tîrgu Mures, Romania: Sapientia University, Department Of Electrical Engineering, Department Of Mechanical Engineering.
11. **Bakó, L. (2009)**. Real-time clustering of datasets with hardware embedded neuromorphic neural networks. *HiBi 2009 (High performance computational systems Biology) Workshop*. Trento, Italy: COSBi (Microsoft Research - University of Trento Centre for Computational and Systems Biology), In press.
12. **Bakó, L., & Brassai, S. (2006)**. Spiking neural networks built into FPGAs: Fully parallel implementations. *WSEAS Transactions on Circuits and Systems*, 5 (3), 346-353.
13. **Bakó, L., & Brassai, S. T. (2008)**. Embedded neural controllers based on spiking neuron models. In M. Iványi (Ed.), *Fourth International PhD, DLA Symposium, University of Pécs, Hungary, Pollack Mihály Faculty of Engineering*. Komló, Hungary: Rotari Press.
14. **Bakó, L., & Brassai, S. T. (2009)**. Embedded neural controllers based on spiking neuron models. *Pollack Periodica, An International Journal for Engineering and Information Sciences*, 4 (3), xx-xx.
15. **Bakó, L., & Brassai, S. T. (2004)**. Fejlett neuronmodellek szimulációja és megvalósítása. *Számokt 2004 – Cluj-Napoca* (pg. 98-107). Cluj-Napoca: EMT.
16. **Bakó, L., & Brassai, S. T. (2006)**. Hardware spiking neural networks: parallel implementations using FPGAs. *Proceedings of the 8th WSEAS Int. Conference on Automatic Control, Modeling and Simulation, Prague, Czech Republic*, (pg. 261-266). Prague.
17. **Bakó, L., & Brassai, S. T. (2005)**. Természetazonos felépítésű mesterséges neurális hálózatok hardware megvalósítása. *Számokt 2005 Kolozsvár* (pg. 219-230). Cluj-Napoca: EMT.
18. **Bakó, L., & Székely, I. (2009)**. Challenges for implementations of delay-coded neuromorphic neural networks on embedded digital hardware. *2nd INCF Congress of Neuroinformatics*. Pilsen, Czech Republic: In press.
19. **Bakó, L., & Székely, I. (2009)**. Challenges for implementations of delay-coded neuromorphic neural networks on embedded digital hardware. *Frontiers in Neuroinformatics. Conference Abstract: 2nd INCF Congress of Neuroinformatics*.
20. **Bakó, L., Brassai, S., & Székely, I. (2006)**. Fully Parallel Implementation of Spiking Neural Networks on FPGA. *Proceedings of the 10th International Conference on Optimisation of Electrical and Electronic Equipment (OPTIM '06)*. III, pg. 135-142. Braşov (Moeciu): Transilvania University Press.
21. **Bakó, L., Brassai, S., Székely, I., & Baczó, M. (2008)**. Hardware Implementation of Delay-coded Spiking-RBF Neural Network for Unsupervised Clustering. *Proceedings of the 11th International Conference on Optimisation of Electrical and Electronic Equipment (OPTIM '08)* (pg. 51-56). Braşov: Transilvania University of Braşov.
22. **Bakó, L., Székely, I., & Brassai, S. (2004)**. Development of Advanced Neural Models. Software And Hardware Implementation. *Transaction on Electronics and communication, Scientific buletin of the „Politehnica” University of Timișoara*, 214-219.
23. **Bakó, L., Székely, I., Dávid, L., & Brassai, S. T. (2004)**. Simulation of Spiking Neural Networks. *Proceedings of the 9th International Conference on Optimisation of Electrical and Electronic Equipment (OPTIM '04)* (pg. 179-184). Braşov: Transilvania University Press.
24. Banitt, Y., Martin, K. A., & Segev, I. (2005). Depressed responses of facilitatory synapses. *Journal of Neurophysiology*, 94, 865-870.
25. Basterretxea, K., & Tarella, J. (2004). Approximation of sigmoid function and the derivative for hardware implementation of artificial neurons. *IEEE proceedings on Circuits, devices and Systems*. Athens.
26. Beeman, D. (2005). GENESIS Modeling Tutorial. *Brains, Minds, and Media*.
27. Beiu, V., Peperstraete, J., & Vandewalle, J. (1994). Close approximation of sigmoid functions by sum of steps for VLSI implementation of Neural networks. *The Scientific Annals, section: Informatics*, 40 (1).
28. Bernard, C., Ge, Y. C., Stockley, E., Willis, J. B., & Wheal, H. V. (1994). Synaptic integration of NMDA and non-NMDA receptors in large neuronal network models solved by means of differential equations. *Biological Cybernetics*, 70 (3), 267-73.

29. Beuchat, J., Haenni, J., & Sanchez, E. (1998). Hardware Reconfigurable Neural Networks. *Proc. of the 5th Reconfigurable Architectures Workshop*. Orlando, Florida, USA.
30. Bhalla, U. S. (2004). Signaling in small subcellular volumes: II. Stochastic and diffusion effects on synaptic network properties. *Biophysical Journal*, 87, 745-753.
31. Bhalla, U. S., & Iyengar, R. (1999). Emergent properties of networks of biological signaling pathways. *Science*, 283, 381-387.
32. Bishop, C. (1995). *Neural networks for pattern recognition*. Oxford: Clarendon Press.
33. Biswas, G., Jain, A. K., & Dubes, R. C. (1981). Evaluation of projection algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 3 (6), 701-708.
34. Blum, K., & Abbott, L. (1996). A model of spatial map formation in the hippocampus of the rat. *Neural Computation* (8), 85-93.
35. Bohte, S., Kok, J., & La Poutré, H. (2002). Error-backpropagation in temporally encoded networks of spiking neurons. *Neurocomputing* (48), 17-37.
36. Bohte, S., Kok, J., & La Poutré, H. (2000). Spike-prop: errorbackpropagation in multi-layer networks of spiking neurons. În V. M. (Ed.), *Proceedings of the European Symposium on Artificial Neural Networks (ESANN)*, (pg. 419-425).
37. Bohte, S., Kok, J., & La Poutré, H. (2000). Unsupervised classification in a layered network of spiking neurons. *Proceedings of IJCNN'2000, IV*, pg. 279-285.
38. Bohte, S., Kok, J., & La Poutré, H. (2002). Unsupervised classification in a layered RBF network of spiking neurons. *IEEE Trans. Neural Networks*, 426-435.
39. Bondalpati, K. (2001, August). Modelling and mapping for dynamically reconfigurable hybrid architecture. *PhD thesis*. Computer Engineering Department, University of South California.
40. Booiij, O. (2004). Temporal pattern classification using spiking neural networks. *Master's thesis*. University of Amsterdam.
41. Booiij, O., & Nguyen, H. (2005). A gradient descent rule for spiking neurons emitting multiple spikes. *Information Processing Letters*, 95 (6), 552-558.
42. Bower, J. M. (1995). Reverse engineering the nervous system: An in vivo, in vitro, and in computo approach to understanding the mammalian olfactory system. În S. F. Zornetzer, J. L. Davis, & C. Lau (Ed.), *An introduction to neural and electronic networks, second edn* (pg. 3-28). New York: Academic Press.
43. Bower, J. M., & Beeman, D. (1998). *The book of GENESIS: Exploring realistic neural models with the General Neural Simulation System, second edn*. New York: Springer.
44. Bower, J., & Beeman, D. (1998). *The Book of GENESIS: Exploring Realistic Neural Models with the General Simulation System, Second ed*. New York: Springer.
45. Brassai, S. T., & Bakó, L. (2008). Mobilis robot mesterséges idegsejt hálóval való szabályzása pályakövetési feladatokra. *Enelko-SzámOkt 2008* (pg. 116-121). Sumuleu-Ciuc: EMT Cluj-Napoca.
46. Brassai, S. T., & Bakó, L. (2008). Visual trajectory control of a mobile robot using FPGA implemented neural network. În M. Iványi (Ed.), *Fourth International PhD, DLA Symposium, University of Pécs, Hungary, Pollack Mihály Faculty of Engineering*. Komló, Hungary: Rotari Press.
47. Brassai, S. T., Bakó, L., & Dan, S. (2007). FPGA Parallel Implementation of CMAC Type Neural Network with on Chip Learning. *SACI 2007* (pg. 111-115). Budapest: Budapest Tech, Hungary.
48. Brassai, S. T., Dávid, L., & Bakó, L. (2004). Hardware Implementation of CMAC based artificial network with process control application. *Transaction on Electronics and communication, Scientific buletin of the „Politehnica” University of Timișoara*, 209-213.
49. Brassai, S. T., Gidró, L., Bakó, L., & Csernáth, G. (2008). Practical Implementation of an Embedded Intelligent Control System. *Proceedings of the International Symposium for Design and Technology of Electronic Packages, Faculty Of Electrical Engineering And Computer Science*. Predeal: Department Of Electronics And Computers, "Transilvania" University Of Brasov and Center For Technological Electronics And Interconnection Techniques "Politehnica" University Bucharest.
50. Brassai, S. T., Márton, L., Dávid, L., & Bakó, L. (2008). Hardware implemented neural network based mobile robot control. *Proceedings of the International Symposium for Design and Technology of Electronic Packages, Faculty Of Electrical Engineering And Computer Science*. Predeal: Department Of Electronics And Computers, "Transilvania" University Of Brasov and Center For Technological Electronics And Interconnection Techniques "Politehnica" University Bucharest.
51. Brassai, S., & Bakó, L. (2007). Hardware Implementation of CMAC Type Neural Network on FPGA for Command Surface Approximation. *Acta Polytechnica Hungarica*, 4 (3), 5-16.
52. Brassai, S., Bakó, L., Székely, I., & Dan, S. (2008). Neural Control Based on RBF Network implemented on FPGA. *Proceedings of the 11th International Conference on Optimisation of Electrical and Electronic Equipment (OPTIM '08)* (pg. 41-46). Brașov: Transilvania University of Brașov.
53. Brown, D., Francis, R., Rose, J., & Vranesic, Z. (1993). *Field-Programmable Gate Arrays*. Kluwer Academic Publishers, USA.
54. Brown, R. (1988). Calendar queues: A fast 0(1) priority queue implementation for the simulation event set problem. *Journal of Communication ACM*, 31 (10), 1220-1227.
55. Buonomano, D., & Merzenich, M. (1999). A neural network model of temporal code generation and position-invariant pattern recognition. *Neural Computation*, 1 (11), 103-116.
56. Carnevale, N. T., & Hines, M. L. (2006). *The NEURON book*. Cambridge: Cambridge University Press.
57. Chakrabarti, S., Roy, S., & Soundalgekar, M. (2001). Fast and Accurate Text Classification via Multiple Linear Discriminant Projections. *Proceedings of International Conference on Very Large Data Bases* (pg. 658-669). Hong Kong: Morgan Kaufmann.
58. Cîrstea, M., Dinu, A., & Nicula, D. (2001). *A practical guide to VHDL design*. Bucharest: Editura Technică.
59. Claverol, E., Brown, A., & Chad, J. (2002). Discrete simulation of large aggregates of neurons. *Neurocomputing*, 47, 277-297.

60. Cloutier, J., Pigeon, S., & Boyer, F. (1996). VIP: An FPGA-based Processor for image processing and neural networks. *5th international conference on Microelectronics for neural networks and fuzzy systems*, (pg. 330-336). Switzerland.
61. Compton, K., & Hauck, S. (2000). An Introduction to Reconfigurable Computing. *Technical Report*. Evanston, IL, USA: Northwestern University, Department of Electrical and Computer Engineering.
62. Connolly, C., Marian, I., & Reilly, R. (2003). Approaches to efficient simulation with spiking neural networks. In *WSPC*.
63. Cormen, T., Leiserson, C., Rivest, R., & Stein, C. (2001). *Introduction to algorithms, second edn*. Cambridge: MIT Press.
64. Cox, C., & Blanz, E. (1992). GangLion – a fast field-programmable gate array implementation of a connectionist classifier. *IEEE Journal of Solid-State Circuits*, 23 (1), 288-299.
65. Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems* (2), 303-314.
66. Day, M., Carr, D. B., Ulrich, S., Ilijic, E., Tkatch, T., & Surmeier, D. J. (2005). Dendritic excitability of mouse frontal cortex pyramidal neurons is shaped by the interaction among HCN, Kir2, and k(leak) channels. *Journal of Neuroscience*, 25, 8776-8787.
67. Dayan, P., & Abbott, L. (2001). *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*. Cambridge, MA: MIT Press.
68. de Garis, H., & Korkin, M. (2002). The cam-brain machine (cbm) an FPGA based hardware tool which evolves a 1000 neuron net circuit module in seconds and updates a 75 million neuron artificial brain for real time robot control. *Neurocomputing*, 42 (1-4).
69. de Garis, H., Gers, F., & Korkin, M. (1997). A simplified cellular automatabased neuron model. *Artificial Evolution Conference (AE97)*. Nimes, France.
70. De Schutter, E., & Bower, J. M. (1994). An active membrane model of the cerebellar Purkinje cell. I. Simulation of current clamps in slice. *Journal of Neurophysiology*, 71, 375-400.
71. Dehon, A. (2000). The Density Advantage of Configurable Computing. *IEEE Computer*, 5 (33), 41-49.
72. Dehon, A. (2000). The Density Advantage of Configurable Computing. *IEEE Computer*, 33 (5), 41-49.
73. Delorme, A., & Thorpe, S. (2001). Face identification using one spike per neuron: resistance to image degradations. *Neural Networks*, 795-804.
74. Delorme, A., & Thorpe, S. (2003). SpikeNET: an event-driven simulation package for modelling large networks of spiking neurons. *Network: Comput. Neural Systems*, 4 (14), 613-628.
75. Deneve, S. (2005). Bayesian inference in spiking neurons. (S. K. Lawrence, Y. Weiss, & L. Bottou, Ed.) *Advances in Neural Information Processing Systems* (17), 353-360.
76. Destexhe, A., & Sejnowski, T. J. (2001). *Thalamocortical assemblies*. New York: Oxford University Press.
77. Destexhe, A., Mainen, Z., & Sejnowski, T. (1994). An efficient method for computing synaptic conductances based on a kinetic model of receptor binding. *Neural Computation*, 6, 14-18.
78. Destexhe, A., Mainen, Z., & Sejnowski, T. (1994). Synthesis of models for excitable membranes, synaptic transmission and neuromodulation using a common kinetic formalism. *Journal of Computational Neuroscience*, 1, 195-230.
79. Diesmann, M., & Gewaltig, M. -O. (2002). NEST: An environment for neural systems simulations. In T. Plesser, & V. Macho (Ed.), *Forschung und wissenschaftliches Rechnen, Beitrage zum Heinz-Billing-Preis 2001, Volume 58 of GWDG-Bericht* (pg. 43-70). Gottingen: Ges. fur Wiss. Datenverarbeitung.
80. Duda, R. O., & Hart, P. E. (1973). *Pattern Classification and Scene Analysis*. John Wiley & Sons.
81. Eldredge, J. (1994). FPGA Density Enhancement of a Neural Network Through Run-Time Reconfiguration. *Master's thesis*. Department of Electrical and Computer Engineering, Brigham Young University.
82. Eldredge, J., & Hutchings, B. (1994). Density enhancement of a neural network using FPGAs and run-time reconfiguration. *Proceedings of IEEE Workshop on Field-Programmable Custom Computing Machines*, (pg. 180-188).
83. Eldredge, J., & Hutchings, B. (1994). RRANN: A Hardware Implementation of the Backpropagation Algorithm Using Reconfigurable FPGAs. *IEEE World Conference on Computational Intelligence*. Orlando, FL.
84. Ermentrout, B. (2004). *Simulating, analyzing, and animating dynamical systems: A guide to XPPAUT for researchers and students*. Philadelphia: SIAM.
85. Ferrucci, A. (1994). ACME: A Field-Programmable Gate Array Implementation of a Self-Adapting and Scalable Connectionist Network. University of California, Santa Cruz.
86. Ferscha, A. (1996). Parallel and distributed simulation of discrete event systems. In A. Y. Zomaya (Ed.), *Parallel and Distributed Computing Handbook* (pg. 1003-1041). New York: McGraw-Hill.
87. Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annual Eugenics*, 7, 179-188.
88. Floreano, D., Dür, P., & Mattiussi, C. (2008). Neuroevolution: from architectures to learning. *Evolutionary Intelligence*, 1.
89. Földiák, P. (1990). Forming sparse representations by local anti-hebbian learning. *Biological Cybernetics*, 64, 165-170.
90. Frank, G., & Hartmann, G. (1995). An Artificial Neural Network Accelerator for Puls-Coded Model-Neurons. *Proceedings of ICNN'95*. Perth, Australia.
91. Fujimoto, R. M. (2000). *Parallel and distributed simulation systems*. New York: Wiley.
92. Gerstner, W. (2001). Spiking neurons. In W. Maass, & C. Bishop (Ed.), *Pulsed Neural Networks*. Cambridge, MA: MIT Press.
93. Gerstner, W., & Kistler, W. M. (2002). Mathematical formulations of hebbian learning. *Biological Cybernetics*, 87, 404-415.
94. Gerstner, W., & Kistler, W. M. (2002). *Spiking Neuron Models. Single Neurons, Populations, Plasticity*. Cambridge: Cambridge University Press.
95. Gerstner, W., Kempter, R., van Hemmen, J., & Wagner, H. (1996). A neuronal learning rule for sub-millisecond temporal coding. *Nature* (383), 76-78.
96. Ghosh-Dastidar, S., & Adeli, H. (2007). Improved spiking neural networks for EEG classification and epilepsy and seizure detection. *Integrated Computer-Aided Engineering* (14), 1-26.
97. Gibson, J., Belerlein, M., & Connors, B. (1999). Two networks of electrically coupled inhibitory neurons in neocortex. *Nature*, 402, 75-79.

98. Girau, B., & Tisserand, A. (1996). OnLine Arithmetic-Based Reprogrammable Hardware Implementation of Multilayer Perceptron Back-Propagation. *Proceedings of the Fifth International Conference on Microelectronics for Neural Networks and Fuzzy Systems* (pg. 168-175). IEEE Computer Society Press.
99. Giugliano, M. (2000). Synthesis of generalized algorithms for the fast computation of synaptic conductances with markov kinetic models in large network simulations. *Neural Computation*, 12, 903–931.
100. Giugliano, M., Bove, M., & Grattarola, M. (1999). Fast calculation of short-term depressing synaptic conductances. *Neural Computation*, 11, 1413–1426.
101. Glackin, B., McGinnity, T., Maguire, L., Wu, Q., & Belatreche, A. (2005). A Novel Approach for the Implementation of Large Scale Spiking Neural Networks on FPGA Hardware. In *Computational Intelligence and Bioinspired Systems* (Vol. 3512). SpringerLink.
102. Grill, W. M., Simmons, A. M., Cooper, S. E., Miocinovic, S., Montgomery, E. B., Baker, K. B., și alții. (2005). Temporal excitation properties of parenthesis evoked by thalamic microstimulation. *Clinical Neurophysiology*, 116, 1227–1234.
103. Gustafson, J. (1988). Reevaluating Amdahl's Law. *CACM*, 5 (31), 532-533.
104. Gütiğ, R., & Sompolinsky, H. (2006). The tempotron: A neuron that learns spike timing-based decisions. *Nature Neuroscience*, 9, 420–428.
105. Hammarlund, P., & Ekeberg, Ö. (1998). Large neural network simulations on multiple hardware platforms. *Journal of Computational Neuroscience*, 5, 443-459.
106. Hauck, S. (1995). Multi-FPGA Systems. *PhD Thesis*. University of Washington.
107. Hereld, M., Stevens, R. L., Teller, J., & van Drongelen, W. (2005). Large neural simulations on large parallel computers. *International Journal of Bioelectromagnetism*, 7, 44-46.
108. Hikawa, H. (1999). Frequency-based multiplayer neural network with on-chip learning and enhanced neuron characteristics. *IEEE Trans. Neural Netw.*, 10, 545–553.
109. Hines, M. (1984). Efficient computation of branched nerve equations. *International Journal of Bio-Medical Computing*, 15, 69-76.
110. Hines, M., & Carnevale, N. T. (1997). The neuron simulation environment. *Neural Computation*, 9, 1179-1209.
111. Hines, M., & Carnevale, N. (1997). The NEURON simulation environment. *Neural Computation* (9), 1179-1209.
112. Hirsch, M., & Smale, S. (1974). *Differential equations, dynamical systems, and linear algebra. Pure and applied mathematics*. New York: Academic Press.
113. Hodgkin, A. L., & Huxley, A. F. (1952). A quantitative description of membrane current and its application to conduction and excitation in nerve. *Journal of Physiology*, 117 (4), 500-544.
114. Hodgkin, A., & Huxley, A. (1952). A Quantitative Description of Membrane Current and its Application to Conduction and Excitation in Nerve. *Journal of Physiology* (117), 500-544.
115. Hopfield, J. (1995). Pattern recognition computation using action potential timing for stimulus representation. *Nature* (376), 33-36.
116. Houweling, A. R., Bazhenov, M., Timofeev, I., Steriade, M., & Sejnowski, T. J. (2005). Homeostatic synaptic plasticity can explain post-traumatic epileptogenesis in chronically isolated neocortex. *Cerebral Cortex*, 15, 834-845.
117. Hutchings, B., & Writhlin, M. (1995). Implementation approaches for reconfigurable logic applications. *5th international workshop on FPGA*. Oxford, England.
118. Iannella, N., & Kindermann, L. (2005). Finding iterative roots with a spiking neural network. *Information Processing Letters*, 6 (95), 545–551.
119. Izhikevich, E. M. (2003). Simple model of spiking neurons. *IEEE Transactions on Neural Networks*, 14, 1569-1572.
120. Izhikevich, E. (2003). Simple model of spiking neurons. *IEEE Transactions on Neural Networks*, 14, 1569–1572.
121. Izhikevich, E. (2004). Which model to use for cortical spiking neurons? *IEEE Transactions on neural networks*, 15 (5), 1063-1070.
122. Izhikevich, E., Desai, N., Walcott, E., & Hoppensteadt, F. (2003). Bursts as a unit of neural information: Selective communication via resonance. *Trends in Neuroscience*, 26, 161–167.
123. Jahnke, A., Roth, U., & Klar, H. (1995). Towards Efficient Hardware for Spike-Processing Neural Networks. *Proc. World Congress on Neural Networks*, (pg. 460-463).
124. Jahnke, A., Roth, U., & Schoenauer, T. (1998). Digital simulation of spiking neural networks. În W. Maass, & C. M. Bishop (Ed.), *Pulsed neural networks*. Cambridge: MIT Press.
125. James-Roxby, P., & Blodget, B. (2000). Adapting constant multipliers in a neural network implementation. *Proceedings of IEEE Symposium on Field-Programmable Custom Computing Machines*, (pg. 335-336).
126. Kanold, P. O., & Manis, P. B. (2005). Encoding the timing of inhibitory inputs. *Journal of Neurophysiology*, 93, 2887-2897.
127. Knoblauch, A. (2005). Neural associative memory for brain modeling and information retrieval. *Information Processing Letters*, 95 (6), 537-544.
128. Koch, C. (1999). *Biophysics of Computation: Information Processing in Single Neurons*. Oxford, UK: Oxford University Press.
129. Kohn, J., & Wörgötter, F. (1998). Employing the Z-transform to optimize the calculation of the synaptic conductance of NMDA and other synaptic channels in network simulations. *Neural Computation*, 10, 1639–1651.
130. Lazarro, J., & Wawrzynek, J. (1993). Silicon Auditory Processors as Computer Peripherals. *NIPS* (5), 820-827.
131. Lee, G., & Farhat, N. H. (2001). The double queue method: A numerical method for integrate-and-fire neuron networks. *Neural Networks*, 14, 921-932.
132. Lysaght, J., Stockwood, J., Law, & Girma, D. (1994). Artificial Neural Network Implementation on a Fine-Grained FPGA. În R. Hartenstein, & M. Servit (Ed.), *Field-Programmable Logic: Architectures, Synthesis and Applications. 4th International Workshop on Field-Programmable Logic and Applications* (pg. 421-431). Prague, Czech Republic: Springer-Verlag.
133. Lysaght, P., Stockwood, J., Law, J., & Girma, D. (1994). Artificial Neural Network Implementation on a Fine-Grained FPGA. În R. Hartenstein, & M. Servit (Ed.), *Field-Programmable Logic: Architectures, Synthesis and*

- Applications. 4th International Workshop on Field-Programmable Logic and Applications* (pg. 421-431). Prague, Czech Republic: Springer-Verlag.
134. Lytton, W. W. (1996). Optimizing synaptic conductance calculation for network simulations. *Neural Computation*, 8, 501-509.
135. Maas, W. (1998). A simple model for neural computation with firing rates and firing correlations. *Network: Computation in Neural Systems* (9), 1-17.
136. Maas, W. (1997). Fast sigmoidal networks via spiking neurons. *Neural Computation*, 2 (9), 279-304.
137. Maas, W. (1997). Networks of spiking neurons: the third generation of neural network models. *Neural Networks*, 9 (10), 1659-1671.
138. Maas, W. (1997). Noisy spiking neurons with temporal coding have more computational power than sigmoidal neurons. In M. Mozer, M. Jordan, & T. Petsche, *Advances in Neural Information Processing Systems* (Vol. 9, pg. 211-217). Cambridge: MIT Press.
139. Maas, W., & Bishop, C. M. (1999). *Pulsed Neural networks*. Cambridge, MA: MIT Press.
140. Maas, W., & Natschläger, T. (1997). Networks of spiking neurons can emulate arbitrary Hopfield nets in temporal coding. *Network: Computation in Neural Systems*, 4 (8), 355-372.
141. Maass, W. (1996). Lower bounds for the computational power of networks of spiking neurons. *Neural Computation*, 1 (8), 1-40.
142. Maass, W., Natschläger, T., & Markram, H. (2002). Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural Computation*, 14 (11), 2531-2560.
143. Makino, T. (2003). A Discrete-Event Neural Network Simulator for General Neuron Models. *Neural Comput. & Applic.* (11), 210-223.
144. Mangasarian, O., & Wolberg, W. (1990). Cancer diagnosis via linear programming. *SIAM News*, 23 (5), 1-18.
145. Mantas, C., Pucho, J., & Mantas, J. (2006). Extraction of similarity based fuzzy rules from artificial neural networks. *International Journal of Approximate Reasoning* (43), 202-221.
146. Mao, J., & Jain, A. K. (1995). Artificial neural networks for feature extraction and multivariate data projection. *IEEE Transactions on Neural Networks*, 6 (2), 296-317.
147. Marcelo, H., & Martin, A. (1994). Reconfigurable Hardware Accelerator for Back-Propagation Connectionist Classifiers. University of California, Santa Cruz.
148. Marchesi, M., Orlandi, G., Piazza, F., & Uncini, A. (1993). Fast neural networks without multipliers. *IEEE Trans. Neural Netw.*, 4 (1), 53-62.
149. Markaki, M., Orphanoudakis, S., & Poirazi, P. (2005). Modelling reduced excitability in aged CA1 neurons as a calcium-dependent process. *Neurocomputing*, 65, 305-314.
150. Mattia, M., & Del Giudice, P. (2000). Efficient event-driven simulation of large networks of spiking neurons and dynamical synapses. *Neural Computation*, 10 (12), 2305-2329.
151. Mattia, M., & Del Giudice, P. (2000). Efficient event-driven simulation of large networks of spiking neurons and dynamical synapses. *Neural Computation*, 12, 2305-2329.
152. Mayrhofer, R., Affenzeller, M., Prähofer, H., Hfer, G., & Fried, A. (2002). Devs simulation of spiking neural networks. In *Proceedings of Cybernetics and Systems (EMCSR)* (Vol. 2, pg. 573-578). Austrian Society for Cybernetic Studies.
153. McCator, H. (1991). Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip. In R. L. al. (Ed.), *Proc. of NIPS-3, "Advances in Neural Information Processing Systems 3"* (pg. 1028-1031). Morgan Kaufmann Pub.
154. McCulloch, W., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics* (5), 115-133.
155. Mehrtash, N., Jung, D., Hellmich, H., Schoenauer, T., Lu, V., & Klar, H. (2003). Synaptic Plasticity in Spiking Neural Networks (SP2INN): A System Approach. *IEEE Transactions on Neural Networks*, 14 (5).
156. Migliore, M., Hines, M. L., & Shepherd, G. M. (2005). The role of distal dendritic gap junctions in synchronization of mitral cell axonal output. *Journal of Computational Neuroscience*, 18, 151-161.
157. Moffitt, M. A., & McIntyre, C. C. (2005). Model-based analysis of cortical recording with silicon microelectrodes. *Clinical Neurophysiology*, 116, 2240-2250.
158. Molz, R., Molz, P., Moraes, F., Torres, L., & Robert, M. (2000). Codesign of fully parallel neural network for a classification problem. *5th world multi-conference on systematics, cybernetics and informatics*. Orlando, Florida.
159. Morrison, A., Mehring, C., Geisel, T., Aertsen, A., & Diesmann, M. (2005). Advancing the boundaries of high connectivity network simulation with distributed computing. *Neural Computation*, 17, 1776-1801.
160. Motomura, M., Aimoto, Y., Shibayama, A., Yabe, Y., & Yamashina, M. (1998). An Embedded DRAM-FPGA Chip with Instantaneous Logic Reconfiguration. *Proc. of IEEE Symp. on FPGAs for Custom Computing Machines* (pg. 264-266). IEEE Computer Society.
161. Muller, J. (1997). *Elementary Functions: Algorithms and Implementation*. Boston: Birkhauser.
162. Murtagh, P., & Tsoi, A. (1992). Implementation issues of sigmoid function and its derivative for VLSI neural networks.
163. Myers, D., & Hutchison, R. (1989). Efficient implementation of piecewise linear activation function for digital VLSI neural networks. *Electronic Letters*, 1662-1663.
164. Nagami, K., Oguri, K., Shiozawa, T., Ito, H., & Konishi, R. (1998). Plastic Cell Architecture: Towards Reconfigurable Computing for General Purpose. *Proc. of IEEE Symp. on FPGAs for Custom Computing Machines* (pg. 68-77). Los Alamitos, CA: IEEE Computer Society.
165. National Science Foundation. (2007). *UC Irvine Machine Learning Repository*. Preluat de pe <http://archive.ics.uci.edu/ml/>
166. Natschläger, T., & Ruf, B. (1998). Spatial and temporal pattern analysis via spiking neurons. *Network: Computational Neural Systems*, 3 (9), 319-332.

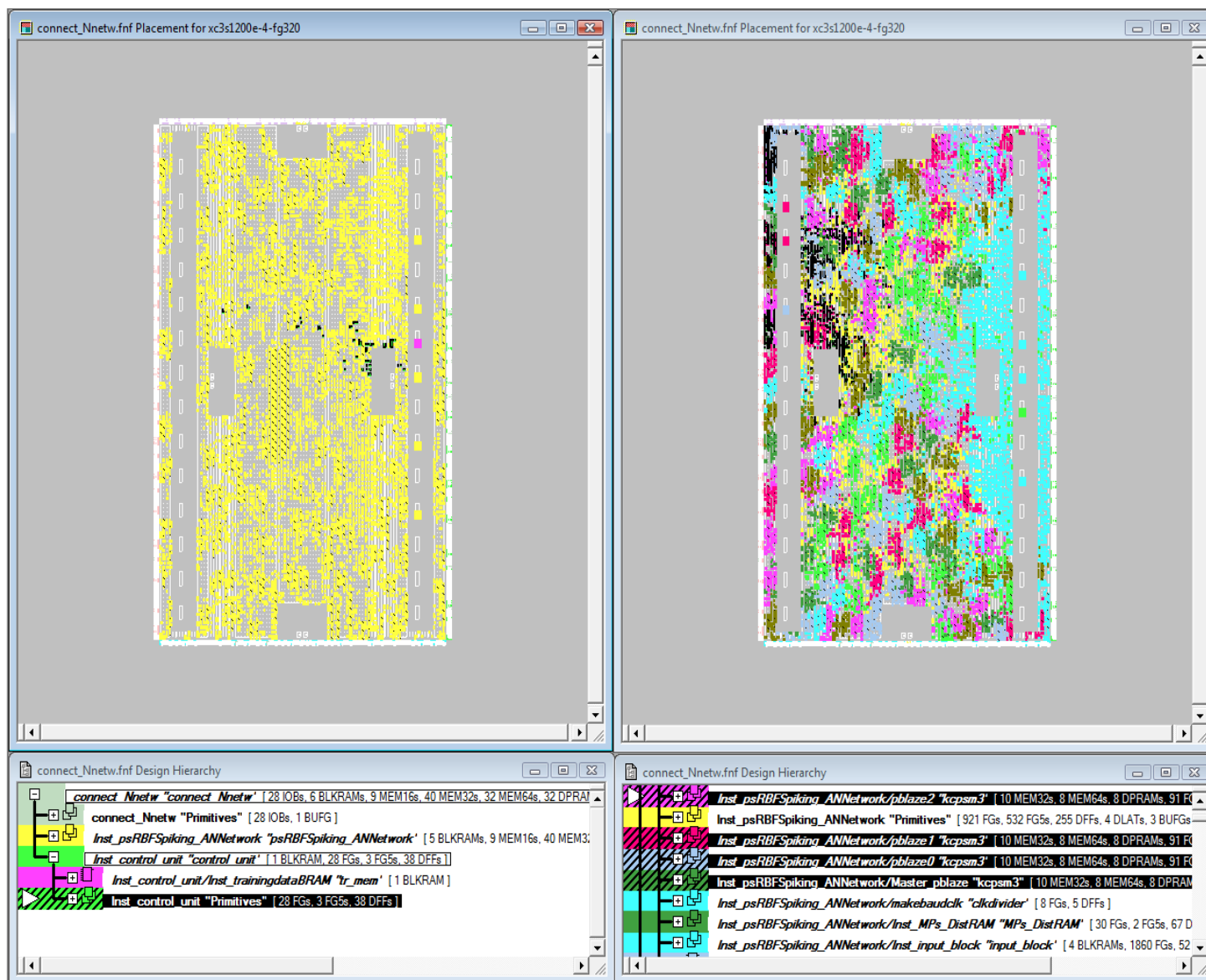
167. Nenadic, Z., Ghosh, B. K., & Ulinski, P. (2003). Propagating waves in visual cortex: A large scale model of turtle visual cortex. *Journal of Computational Neuroscience*, 14, 161-184.
168. Nichols, K. (2003, December). A Reconfigurable Computing Architecture For Implementing Artificial Neural Networks On FPGA. *A Thesis Presented to The Faculty of Graduate Studies of The University of Guelph*.
169. Nordstrom, T. (1995, March). Highly Parallel Computers for Artificial Neural Networks. *Ph.d.thesis (1995:162 f)*. Sweden: Division of Computer Science and Engineering, Lulea University of Technology.
170. O'Dwyer, C., & Richardson, D. (2005). Spiking neural nets with symbolic internal state. *Information Processing Letters - Special issue on applications of spiking neural networks*, 95 (6), 529 - 536.
171. Oldfield, J., & Dorf, R. (1995). *Field Programmable Arrays: Reconfigurable Logic for Rapid Prototyping and Implementation of Digital Systems*. John Wiley and Sons.
172. Olshausen, B. A., & Field, D. J. (2005). How close are we to understanding V1? *Neural Computation*, 17, 1665-1699.
173. Omondi, A. (1994). *Computer Arithmetic Systems: Algorithms, Architecture, and Implementations*. Prentice-Hall, UK.
174. *Open source data visualization and analysis for novice and experts*. (2005). Preluat de pe <http://www.ailab.si/orange/doc/datasets/breast-cancer-wisconsin-cont.tab>
175. Pandya, V. (2005, August). A Handel-C Implementation Of The Back-Propagation Algorithm On Field Programmable Gate Arrays. *A Thesis Presented to The Faculty of Graduate Studies Of The University of Guelph*.
176. Pearson, M., Gilhespy, I., Gurney, K., Melhuish, C., Mitchinson, B., Nibouche, M., și alții. (2005). A Real-Time, FPGA Based, Biologically Plausible Neural Network Processor. În *Book Chapter of Book Series Lecture Notes in Computer Science, Category: Issues in Hardware Implementation* (Vol. 3697). SpringerLink.
177. Pérez-Uribe, A. (1999, Oct.). Structure-Adaptable Digital Neural Networks. *Ph.d. thesis*. Logic Systems Laboratory, Computer Science Department, Swiss Federal Institute of Technology-Lausanne.
178. Pérez-Uribe, A., & Sanchez, E. (1996). FPGA Implementation of an Adaptable Size Neural Network. *VI international conference on ANN ICANN'96*, (pg. 383-1388). Bochum, Germany.
179. Pérez-Uribe, A., & Sanchez, E. (1996). FPGA implementation of an Adaptable-Size Neural Network. În C. Von der Malsburg, W. Von Seelen, J. Vorbrüggen, & B. Sendhoff (Ed.), *Proceedings of the International Conference on Neural Networks (ICANN96). 1112 of Lecture Notes in Computer Science*, pg. 383-388. Heidelberg: Springer-Verlag.
180. Pérez-Uribe, A., & Sanchez, E. (1997). Speeding-up adaptive heuristic critic learning with fpga-based unsupervised clustering. *Proceedings of the IEEE International Conference on Evolutionary Computation ICEC'97*, (pg. 685-689). Indianapolis.
181. Poormann, M., Witkowski, U., Kalte, H., & Ruckert, U. (2002). Implementation of ANN on a reconfigurable hardware accelerator. *Euromicro workshop on parallel, distributed and network based processing*, (pg. 243-250). Spain.
182. Prescott, S. A., & De Koninck, Y. (2005). Integration time in a subset of spinal lamina I neurons is lengthened by sodium and calcium currents acting synergistically to prolong subthreshold depolarization. *Journal of Neuroscience*, 25, 4743-4754.
183. Press, W. H., Flannery, B. P., Teukolsky, S. A., & Vetterling, W. T. (1993). *Numerical recipes in C: The art of scientific computing*. Cambridge: Cambridge University Press.
184. Rao, R. P. (2005). Hierarchical bayesian inference in networks of spiking neurons. (S. K. Lawrence, Y. Weiss, & L. Bottou, Ed.) *Advances in Neural Information Processing Systems* (17), 1113-1120.
185. Reutimann, J., Guigliano, M., & Fusi, S. (2003). Event-driven simulation of spiking neurons with stochastic dynamics. *Neural Computation* (15), 811-830.
186. Rochel, O., & Martinez, D. (2003). An event-driven framework for the simulation of networks of spiking neurons. În *Proceedings of the 11th European Symposium on Artificial Neural Networks - ESANN'2003* (pg. 295-300). Bruges.
187. Ros, E., Carrillo, R., Ortigosa, E., Barbour, B., & Agis, R. (2005). Event-driven Simulation Scheme for Spiking Neural Models based on Characterization Look-up Tables. *Neural Computation*.
188. Roth, U., Jahnke, A., & Klar, H. (1995). Hardware Requirements for Spike-Processing Neural Networks. *Proc. IWANN95*, (pg. 720-727).
189. Rotter, S., & Diesmann, M. (1999). Exact digital simulation of time-invariant linear systems with applications to neuronal modeling. *Biological Cybernetics*, 81, 381-402.
190. Rumelhart, D., Hinton, G., & Williams, R. (1986). Learning representations by back-propagating errors. *Nature* (325), 533-536.
191. Saghatelian, A., Roux, P., Migliore, M., Rochefort, C., Desmaisons, D., Charneau, P., și alții. (2005). Activity-dependent adjustments of the inhibitory network in the olfactory bulb following early postnatal deprivation. *Neuron*, 46, 103-116.
192. Sahin, I., Gloster, C., & Doss, C. (2000). Feasibility of floating point arithmetic in reconfigurable computing systems. *MAPLD International Conference*. DC, USA.
193. Sammon Jr., J. W. (1969). A nonlinear mapping for data structure analysis. *IEEE Transactions on Computers*, C-18, 401-409.
194. Sammut, K., & Jones, S. (1991). Implementing non-linear activation functions in neural network emulators. *Electronic Letters*, 27 (12).
195. Sanchez-Montanez, M. (2001). Strategies for the optimization of large scale networks of integrate and fire neurons. În J. Mira, & A. Prieto (Ed.), *IWANN, Volume 2084/2001 of Lecture Notes in Computer Science*. NewYork: Springer-Verlag.
196. Schoenauer, T., Atasoy, S., Mehrtash, N., & Klar, H. (2002). NeuroPipe-Chip: A Digital Neuro-Processor for Spiking Neural Networks. *IEEE Trans. Neural Networks*, 13 (1), 205-213.
197. Schrauwen, B., & Van Campenhout, J. (2006). Parallel hardware implementation of a broad class of spiking neurons using serial arithmetic. *ESANN'2006 proceedings - European Symposium on Artificial Neural Networks*. Bruges, Belgium: ISBN 2-930307-06-4.
198. Simard, P., & Graf, H. (1993). Backpropagation without Multiplication. *Advances in Neural Information Processing Systems*, 232-239.



199. Skrbek, M. (1999). Fast Neural Network Implementation. *Neural Network World* , 9 (5), 375-391.
200. Sloom, A., Kaandorp, J. A., Hoekstra, G., & Overeinder, B. J. (1999). Distributed simulation with cellular automata: Architecture and applications. In J. Pavelka, G. Tel, & M. Bartosek (Ed.), *SOFSEM'99, LNCS* (pg. 203-248). Berlin: Springer-Verlag.
201. Snippe, H. (1996). Parameter extraction from population codes: A critical assessment. *Neural Computation* , 511-529.
202. Sobieski, J., & Storaasli, O. (2004, Oct). Computing at the Speed of Thought. *Aerospace America* , 35-38.
203. Song, S., & Abbott, L. (2001). Column and Map Development and Cortical Re-Mapping Through Spike-Timing Dependent Plasticity. *Neuron* (32), 339-350.
204. Song, S., Miller, K. D., & Abbott, L. F. (2000). Competitive hebbian learning through spike-timing-dependent synaptic plasticity. *Nature Neuroscience* , 3, 919-926.
205. Song, S., Miller, K., & Abbott, L. (2000). Competitive Hebbian learning though spike-timing dependent synaptic plasticity. *Nature Neuroscience* (3), 919-926.
206. Strenski, D. (2007, Jan 12). FPGA Floating Point Performance. *HPCWire* .
207. Stricanne, B., & Bower, J. M. (1998). A network model of the somatosensory system cerebellum, exploring recovery from peripheral lesions at various developmental stages in rats (abstract). *Society of Neuroscience Abstracts* , 24, 669.
208. Tang, K., & Kak, S. (2002). Fast Classification Networks for Signal Processing. *Circuits Systems Signal Processing* , 2 (21), 207-224.
209. Trappenberg, T. P. (2002). *Fundamentals of Computational Neuroscience*. Oxford: Oxford University Press.
210. Triesch, J. (2007). Synergies between intrinsic and synaptic plasticity mechanisms. *Neural Computation* (19), 885-909.
211. Upegui, A., Peña-Reyes, C., & Sánchez, E. (2004). A Hardware Implementation of a Network of Functional Spiking Neurons with Hebbian Learning. In *Book Chapter of Book Series Lecture Notes in Computer Science* (Vol. 3141). SpringerLink.
212. Upegui, A., Peña-Reyes, C., & Sánchez, E. (2003). A methodology for evolving spiking neuralnetwork topologies on line using partial dynamic reconfiguration. *Proceedings of II- International Congress on Computational Intelligence (CIIC'03)*. Medellín, Colombia.
213. van Emde Boas, P., Kaas, R., & Zijlstra, E. (1976). Design and implementation of an efficient priority queue. *Theory of Computing Systems* , 10, 99-127.
214. Verstraeten, D., Schrauwen, B., Stroobandt, D., & Van Campenhout, J. (2005). Isolated word recognition with the Liquid State Machine: a case study. *Information Processing Letters* (95), 521-528.
215. Vitko, I., Chen, Y. C., Arias, J. M., Shen, Y., Wu, X. R., & Perez-Reyes, E. (2005). Functional characterization and neuronal modeling of the effects of childhood absence epilepsy variants of CACNA1H, a T-type calcium channel. *Journal of Neuroscience* , 25, 4844-4855.
216. Watts, L. (1994). Event-driven simulation of networks of spiking neurons. In *Advances in neural information processing systems* (pg. 927-934).
217. Wilson, M. A., & Bower, J. M. (1989). The simulation of largescale neural networks. In C. Koch, & I. Segev (Ed.), *Methods in neuronal modeling: From synapses to networks* (pg. 291-333). Cambridge: MIT Press.
218. Wolf, J. A., Moyer, J. T., Lazarewicz, M. T., Contreras, D., Benoit-Marand, M., O'Donnell, P., și alții. (2005). NMDA/AMPA ratio impacts state transitions and entrainment to oscillations. *Journal of Neuroscience* , 25, 9080-9095.
219. Xilinx, DS099. (2006, April 26). *Spartan-3 FPGA Family: Complete Data Sheet* .
220. Xilinx, XAPP463. (2005, March 1). Using Block RAM in Spartan-3 Generation FPGAs. *XAPP463 (v2.0)* .
221. Xilinx, XAPP464. (2005, March 1). Using Look-Up Tables as Distributed RAM in Spartan-3 Generation FPGAs. *XAPP464 (v2.0)* .
222. Xin, J., & Embrechts, M. (2001). Supervised learning with spiking neural networks. *Proceedings of International Joint Conference on Neural Networks (IJCNN)*, (pg. 1772-1777).
223. Zeigler, B. P., & Vahie, S. (1993). DEVS formalism and methodology: Unity of conception/diversity of application. *Proceedings of the 1993 Winter Simulation Conference*, (pg. 573-579). Los Angeles, December 12-15.
224. Zeigler, B., Praehofer, H., & Kim, T. (2000). *Theory of modeling and simulation, second edn. Integrating discrete event and continuous complex dynamic systems*. New York: Academic Press.
225. Zemel, R. S., Huys, Q. J., Natarajan, R., & Dayan, P. (2005). Probabilistic computation in spiking populations. (S. K. Lawrence, Y. Weiss, & L. Bottou, Ed.) *Advances in Neural Information Processing Systems* (17), 1609-1616.
226. Zhang, M., Vassiliadis, S., & Fris, J. (1996). Sigmoid generators for neural computing using piece-wise approximation. *IEEE transactions on Computers* , 45 (9).
227. Zhu, J., & Gunther, B. (1999). Towards an FPGA based reconfigurable computing environment for neural network implementations. *roceedings of Ninth International Conference on Artificial Neural Networks*, 2, pg. 661-666.



## ANEXA 1 – DISTRIBUȚIA COMPONENTELOR UTILIZATE PE CIRCUITELE FPGA (DETECTOR FRECVENȚE VS. CLASIFICARE IRIS DATASET)

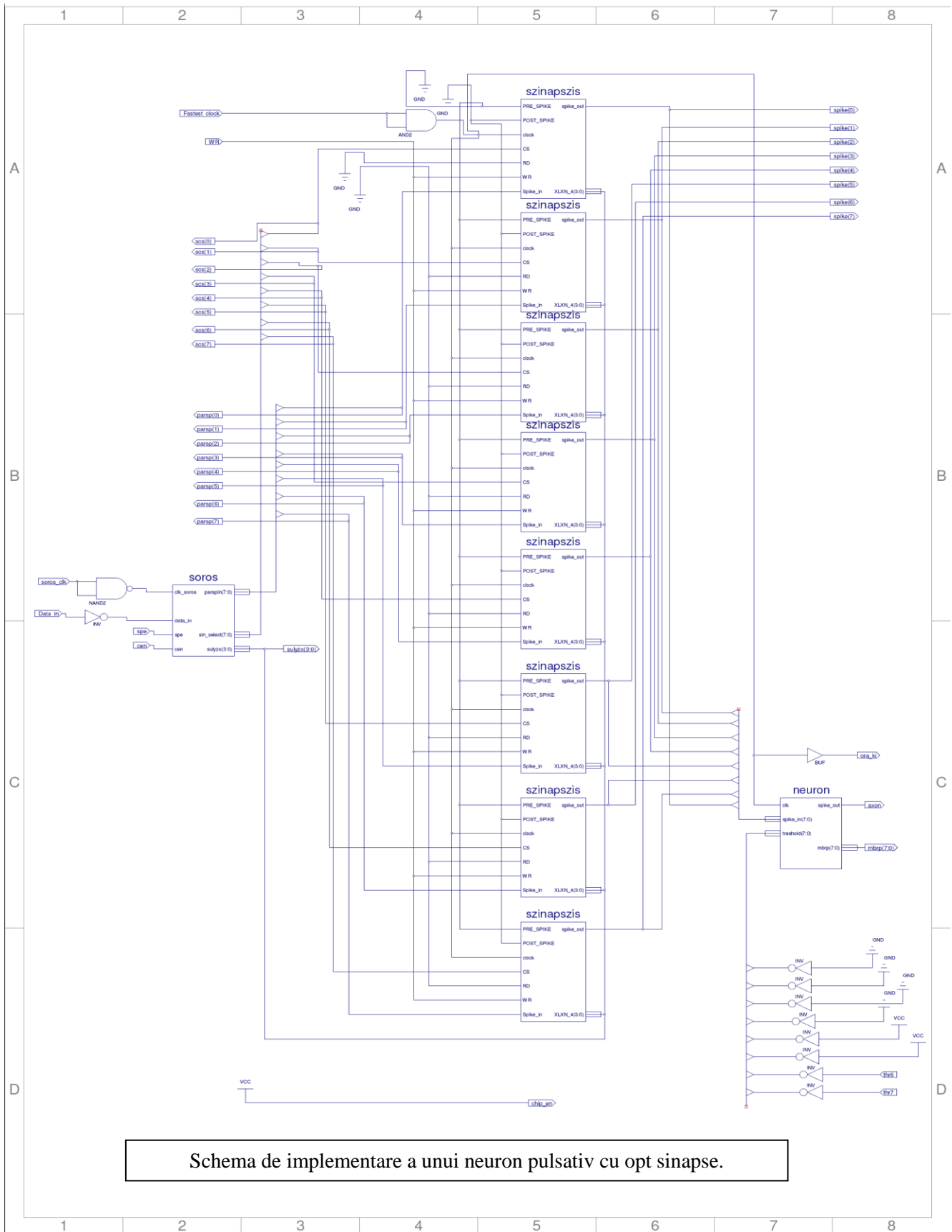


În partea stângă a figurii se poate vedea utilizarea resurselor FPGA de către circuitul ce implementează aplicația de detectarea a componentelor de frecvențe, prezentată în subcapitolul 5.5.2 și pe care se poate distinge, aria de circuit utilizată de modulul de control și antrenare (control unit) cu culoare neagră respectiv RNP reprezentată cu galben.

Imaginea din dreapta figurii prezintă utilizarea de resurse reconfigurabile de către proiectul ce implementează aplicația benchmark de clasificare a setului de date Fischer IRIS, expusă la subpunctul 5.5.3. A se nota faptul, că cele patru microcontrolere PicoBlaze din acest proiect (marcate cu negru pe figură) utilizează o suprafață relativ mică. Suprafața în albastru deschis (dreapta-centru) este mult mai mare, e reprezentând blocul neuronilor de intrare ce codifică valorile de intrare.



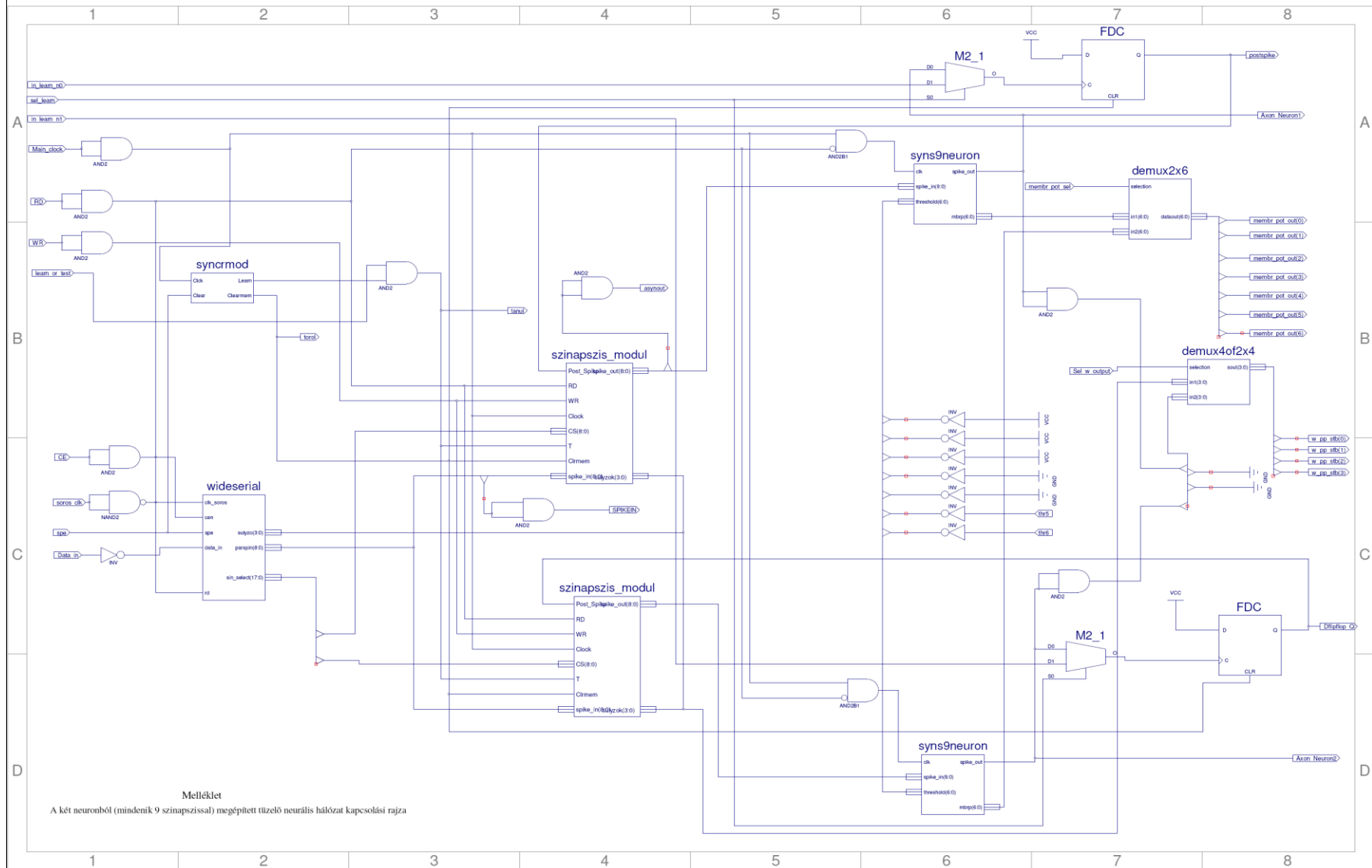
# ANEXA 2- SCHEMA DE IMPLEMENTARE A UNUI NEURON PULSATIV



Schema de implementare a unui neuron pulsativ cu opt sinapse.



# ANEXA 3- SCHEMA DE IMPLEMENTARE A UNEI REȚELE NEURONALE PULSATIVE







# ANEXA 4- SCHEMA DE IMPLEMENTARE A REȚELE NEURONALE PULSATIVE PENTRU CLASIFICARE SETULUI DE DATE IRIS

