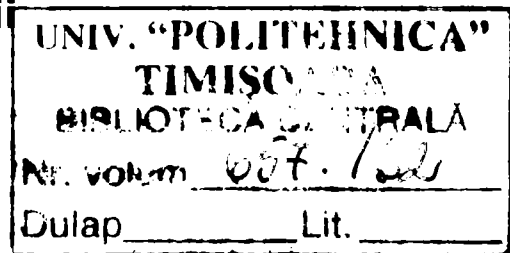


# CERCETĂRI PRIVIND CREȘTEREA PERFORMANȚELOR REȚELELOR NEURONALE ÎNTR-UN MEDIU DE CALCUL DISTRIBUIT

Teză destinată obținerii  
titlului științific de doctor inginer  
la  
Universitatea "Politehnica" din Timișoara  
în domeniul ȘTIINȚA CALCULATOARELOR  
de către

**Ing. Sorin Babii**



Conducător științific:  
Referenți științifici:

prof.univ.dr.ing Vladimir Crețu  
prof.univ.dr.ing Mircea Petrescu  
prof.univ.dr. Viorel Negru  
prof.univ.dr.ing. Ștefan Holban

Ziua susținerii tezei: 7.12.2007

Seriile Teze de doctorat ale UPT sunt:

- |                        |   |
|------------------------|---|
| 1. Automatică          | 7. Inginerie Electronică și Telecomunicații |
| 2. Chimie              | 8. Inginerie Industrială                    |
| 3. Energetică          | 9. Inginerie Mecanică                       |
| 4. Ingineria Chimică   | 10. Știința Calculatoarelor                 |
| 5. Inginerie Civilă    | 11. Știința și Ingineria Materialelor       |
| 6. Inginerie Electrică |   |

Universitatea „Politehnica” din Timișoara a inițiat seriile de mai sus în scopul diseminării expertizei, cunoștințelor și rezultatelor cercetărilor întreprinse în cadrul școlii doctorale a universității. Seriile conțin, potrivit H.B.Ex.S Nr. 14 / 14.07.2006, tezele de doctorat susținute în universitate începând cu 1 octombrie 2006.

Copyright © Editura Politehnica – Timișoara, 2007

Această publicație este supusă prevederilor legii dreptului de autor. Multiplicarea acestei publicații, în mod integral sau în parte, traducerea, tipărirea, reutilizarea ilustrațiilor, expunerea, radiodifuzarea, reproducerea pe microfilme sau în orice altă formă este permisă numai cu respectarea prevederilor Legii române a dreptului de autor în vigoare și permisiunea pentru utilizare obținută în scris din partea Universității „Politehnica” din Timișoara. Toate încălcările acestor drepturi vor fi penalizate potrivit Legii române a drepturilor de autor.

România, 300159 Timișoara, Bd. Republicii 9,  
tel. 0256 403823, fax. 0256 403221  
e-mail: editura@edipol.upt.ro

## Cuvânt înainte

Teza de doctorat a fost elaborată pe parcursul activității mele în cadrul Departamentului de Calculatoare al Universității „Politehnica” din Timișoara.

Mulțumiri deosebite se cuvin conducătorului de doctorat prof.dr.ing. Vladimir Crețu, pentru sprijinul constant, ideile și sugestiile fără de care această teză nu ar fi fost posibilă. De asemenea, doresc să mulțumesc domnului prof.dr.ing. Ștefan Holban, pentru inițierea în fascinantul domeniu al rețelelor neuronale și pentru recomandările făcute pe parcursul elaborării tezei. Nu în ultimul rând, doresc să mulțumesc domnilor profesori Mircea Petrescu și Viorel Negru, pentru răbdarea cu care mi-au lecturat teza și pentru recomandările făcute.

Doresc să adresez mulțumiri conducerii Departamentului MRC din cadrul Alcatel-Lucent S.A și echipei IT din acest departament, pentru sprijinul tehnic acordat pe parcursul măsurărilor experimentale.

Nu în ultimul rând, le mulțumesc tuturor celor care m-au sprijinit pe tot parcursul elaborării tezei, pentru sfaturile și încurajările lor.

Timișoara, Decembrie 2007

Sorin Babii

Babii, Sorin Ioan

**Cercetări privind creșterea performanțelor rețelelor neuronale într-un mediu de calcul distribuit**

Teze de doctorat ale UPT, Seria 10, Nr. 4, Editura Politehnica, 2007, 138 pagini, 49 figuri, 19 tabele.

ISSN: 1842-7707

ISBN: 978-973-625-559-5

Cuvinte cheie:

Rețele neuronale, BackPropagation, antrenare, algoritmi paraleli, algoritmi distribuiți.

Rezumat:

Teza de față este rezultatul original al investigațiilor efectuate în direcția găsirii unei modalități avantajoase de îmbunătățire a performanțelor în etapa de antrenare a rețelelor neuronale, pe baza resurselor computaționale existente în cadrul Departamentului Calculatoare. Cercetările au fost axate pe dezvoltarea unor algoritmi paraleli care să evolueze într-o rețea de calculatoare, care să accelereze procesul de învățare și să depășească limitările în puterea de calcul, pentru a permite investigații ulterioare în domeniul rețelelor neuronale.

Vreme de sute de milioane de ani, creierul viu, creat și perfecționat continuu prin eternul proces evoluționar al selecției naturale, au fost singurele dispozitive capabile să efectueze procesarea informației.

Ființele umane au inventat calculatorul digital, un dispozitiv artificial de procesare a informației, care a creat perspectiva efectuării de calcule arbitrare înafara sistemelor nervoase biologice. Totuși, curând, a devenit evident faptul că, sub anumite aspecte, proprietățile calculatoarelor digitale sunt mult diferite de cele ale creierelor vii.

În primul rând, avem de-a face cu o diferență de structură: calculatoarele digitale (de obicei) posedă o singură unitate de procesare, care este, totuși, destul de puternică. Creierul, dimpotrivă, sunt alcătuite din neuroni intens interconectați, lucrând în paralel, fiecare constituind o unitate simplă (comparativ) de procesare. În privința capacității de procesare a informației, însă, diferența structurală nu este cea mai importantă, întrucât toate calculatoarele posedă abilitatea de a simula alte structuri decât cea proprie.

Mai importantă este diferența în modul de a câștiga abilitatea de a efectua noi funcțiuni: creierul *învață*, în vreme ce calculatoarele trebuie programate. Pentru a efectua o sarcină dată, calculatorul digital necesită *software*, programe care să implementeze în detaliu modul în care acea sarcină poate fi efectuată. Un calculator fără un program nu poate să proceseze informație sau să efectueze calcule. De obicei, o ființă umană trebuie să înțeleagă o funcțiune dată de procesare a informației și să dezvolte un algoritm care s-o implementeze, și abia după aceea calculatorul poate fi programat să efectueze acea funcțiune. Există însă o multitudine de probleme pentru care încă nu există algoritmi formali, sau pentru care este virtual imposibil să se enumere o serie de pași logici care să determine calculatorul să ajungă la răspunsul corect. Asemenea probleme de obicei presupun observarea unui mare număr de reguli complexe, dependente de context, dintre care majoritatea încă nu sunt cunoscute. Exemple de asemenea probleme sunt cele complexe de recunoaștere a tiparelor, întâlnite în recunoașterea vorbirii, identificarea scrisului de mână, recunoașterea fizionomiilor, sau interpretarea spațială a imaginilor bidimensionale.

Adesea, însă, este posibil să specificăm o problemă destul de precis oferind un set mare de exemple, care să precizeze modul în care obiectele din spațiul intrărilor trebuie asociate cu obiecte din spațiul ieșirilor. În mod uzual, oamenii pot învăța să rezolve asemenea probleme, întrucât creierul animalelor superioare au evoluat către o bună generalizare atunci când li se prezintă un număr de exemple. Acest fapt a condus la dezvoltarea de *Rețele Neuronale Artificiale* [74], [77], dispozitive proiectate să modeleze activitatea rețelelor neuronale biologice la un anumit nivel de detaliu, pentru a atinge (câteva dintre) capacitățile creierului uman<sup>1</sup>.

Ca și în realitate, o rețea neuronală artificială este un sistem de elemente simple de procesare, masiv paralel interconectate. În această privință, rețelele neuronale artificiale se bazează pe cunoștințele noastre actuale legate de sistemele nervoase biologice. Trebuie remarcat însă faptul că un creier uman este cu câteva ordine de mărime mai complex decât oricare rețea neuronală artificială existentă. Se estimează [79] că în creierul uman există circa  $10^{10}$ – $10^{11}$  neuroni. Fiecare dintre acești neuroni recepționează semnal de intrare de la mii și chiar zeci de mii de sinapse care îl conectează la alți neuroni, iar rezultatul activității neuronului poate fi transmis prin intermediul altor mii de sinapse către alți neuroni, influențându-le astfel activitatea viitoare.

Inteligența artificială (AI) este un domeniu de cercetare în care se depun eforturi

---

<sup>1</sup> Cercetările în materie mai sunt cunoscute sub denumirea de *Conexționism* (datorită rolului important jucat de conexiunile din rețea) sau Procesare Paralelă Distribuită (PDP) [Rumelhart86], însă denumirea de *rețele neuronale artificiale* pare să aibă o rezonanță deosebită și este un termen larg utilizat.

variate în vederea modelării și/sau re-creării aspectelor inteligenței naturale. Una dintre cele două paradigme luate în considerație de comunitatea AI se bazează pe ideea re-creării comportării inteligente prin imitarea arhitecturii creierelor biologice. Asemenea imitații sunt întâlnite de obicei sub denumirea de rețele neuronale artificiale, iar programele utilizate pentru execuția lor sunt denumite simulatoare de rețele neuronale.

<b>Cuprins</b>	<b>iii</b>
<b>Listă de tabele</b>	<b>vi</b>
<b>Listă de figuri</b>	<b>vii</b>
<b>1 Introducere</b>	<b>1</b>
1.1 Obiective . . . . .	2
1.2 Limitări . . . . .	2
1.3 Conținutul tezei . . . . .	3
<b>2 Modele conexiuniste</b>	<b>5</b>
2.1 Celula nervoasă . . . . .	5
2.2 Rețele neuronale . . . . .	6
2.2.1 Modele neuronale artificiale . . . . .	6
2.2.2 Scurt istoric . . . . .	6
2.2.3 Modelul McCulloch & Pitts . . . . .	7
2.2.4 Modelul Hebb . . . . .	8
2.3 Modelul Rosenblatt — perceptronul . . . . .	10
2.3.1 Regula delta . . . . .	11
2.3.2 Problema XOR . . . . .	14
2.3.3 Perceptronul multistrat . . . . .	14
2.4 Modelul Hopfield . . . . .	16
2.4.1 Modelul fizic Ising . . . . .	16
2.4.2 Definierea modelului Hopfield . . . . .	17
2.5 Rețele cu propagare înapoi ( <i>backpropagation</i> ) . . . . .	20
2.5.1 Generalizarea regulii delta . . . . .	21
2.5.2 Modul de lucru cu <i>backpropagation</i> . . . . .	22
2.5.3 Alte funcții de activare neliniare . . . . .	23
2.5.4 Rata de învățare și momentul . . . . .	25
2.5.5 Învățarea unui tipar . . . . .	25
2.5.6 Deficiențe ale algoritmului <i>backpropagation</i> . . . . .	26
2.6 Îmbunătățiri ale algoritmului <i>backpropagation</i> . . . . .	26
2.6.1 Parametrii de antrenare a rețelei . . . . .	29
2.6.2 Accelerarea convergenței — rată de învățare adaptivă . . . . .	30
2.7 O implementare secvențială . . . . .	31
2.7.1 Analiza simulatorului secvențial . . . . .	33
2.8 Concluzii . . . . .	35
<b>3 Paralelizarea rețelelor neuronale</b>	<b>36</b>
3.1 Strategii de paralelizare . . . . .	36
3.1.1 Paralelizarea la nivelul setului de antrenare . . . . .	36
3.1.2 Paralelizarea la nivelul nodurilor . . . . .	37
3.2 Direcții de acțiune . . . . .	37
3.2.1 Partiționarea datelor . . . . .	37
3.2.2 Partiționarea rețelei . . . . .	38
3.3 Analiza algoritmilor paraleli . . . . .	38
3.4 Principalele obiective ale paralelizării . . . . .	40
3.5 Surse de ineficiență . . . . .	40
3.5.1 Software adițional . . . . .	40
3.5.2 Echilibrarea încărcării . . . . .	41
3.5.3 Comunicația suplimentară . . . . .	41
3.5.4 Porțiuni inerent secvențiale . . . . .	42

3.5.5	Limitări specifice problemei . . . . .	42
3.6	Considerații specifice rețelelor neuronale . . . . .	43
3.7	Experimente în vederea analizei performanțelor . . . . .	44
3.7.1	Probleme de dimensiune fixă . . . . .	44
3.7.2	Probleme de dimensiune variabilă . . . . .	44
3.8	Concluzii . . . . .	45
<b>4</b>	<b>DataParSim</b>	<b>47</b>
4.1	Strategia de partiționare a datelor . . . . .	47
4.1.1	Configurație inel. Configurație arbore . . . . .	48
4.2	O implementare avansată . . . . .	50
4.2.1	Topologia procesoarelor . . . . .	50
4.2.2	Comunicația . . . . .	50
4.2.3	Pasul înainte . . . . .	51
4.2.4	Pasul înapoi . . . . .	52
4.2.5	O îmbunătățire a eficienței . . . . .	54
4.2.6	Necesarul de memorie . . . . .	55
4.3	Simulatorul DataParSim . . . . .	55
4.3.1	Fișierul de configurare . . . . .	55
4.3.2	Programul administrator . . . . .	56
4.3.3	Programul <i>slave</i> . . . . .	60
4.4	Rezultate experimentale . . . . .	64
4.4.1	Platforme hardware . . . . .	64
4.4.2	Influența dimensiunii variabile a setului de antrenare . . . . .	66
4.4.3	Influența numărului variabil de procesoare . . . . .	67
4.4.4	Set de antrenare variabil cu numărul de procesoare . . . . .	69
4.5	Concluzii . . . . .	71
<b>5</b>	<b>NetParSim</b>	<b>73</b>
5.1	Construcția algoritmului paralel . . . . .	73
5.1.1	Divizarea rețelei . . . . .	73
5.1.2	Topologia procesoarelor . . . . .	74
5.1.3	Notății . . . . .	75
5.1.4	Paralelizarea . . . . .	75
5.1.5	Memorarea tiparelor de intrare și de ieșire . . . . .	79
5.1.6	Îmbunătățiri ale algoritmului . . . . .	80
5.2	Analiza algoritmului . . . . .	80
5.2.1	Necesarul de memorie . . . . .	80
5.2.2	Pasul înainte . . . . .	80
5.2.3	Pasul înapoi . . . . .	81
5.3	Simulatorul NetParSim . . . . .	81
5.3.1	Fișierul de configurare . . . . .	82
5.3.2	Codul aplicației . . . . .	83
5.4	Rezultate experimentale . . . . .	87
5.4.1	Influența dimensiunii variabile a rețelei neuronale . . . . .	88
5.4.2	Efectul numărului variabil de procesoare . . . . .	89
5.4.3	Influența scalării rețelei cu numărul de procesoare . . . . .	90
5.4.4	Modificarea dimensiunii setului de antrenare . . . . .	93
5.5	Concluzii . . . . .	99
<b>6</b>	<b>Experimente cu aplicații uzuale</b>	<b>100</b>
6.1	NETtalk . . . . .	100
6.1.1	Setul de date NETtalk . . . . .	100
6.1.2	Implementare . . . . .	101
6.1.3	Simulări și rezultate . . . . .	102
6.2	Compresie de imagine . . . . .	105
6.2.1	Suportul teoretic . . . . .	105



6.2.2 Simulări și rezultate . . . . .	106
6.3 Recunoașterea țintelor de sonar . . . . .	107
6.3.1 Configurația rețelei . . . . .	108
6.3.2 Simulări și rezultate . . . . .	109
6.4 Concluzii . . . . .	109
<b>7 Concluzii finale și perspective</b>	<b>112</b>
7.1 Concluzii . . . . .	112
7.2 Rezumat al contribuțiilor . . . . .	113
7.3 Concluzii finale . . . . .	114
7.4 Perspective de cercetare . . . . .	114
<b>Bibliografie</b>	<b>116</b>
<b>Lista lucrărilor publicate</b>	<b>123</b>

2.1	Operații în virgulă flotantă necesare. . . . .	34
4.1	Eficiența funcție de dimensiunea setului de antrenare . . . . .	66
4.2	Accelerarea: SunBlade . . . . .	67
4.3	Accelerarea: Linux . . . . .	68
4.4	Set de antrenare fix per procesor: SunBlade . . . . .	69
4.5	Set de antrenare fix per procesor: Linux . . . . .	70
5.1	Dimensiune variabilă a rețelei: SunBlade. . . . .	88
5.2	Dimensiune variabilă a rețelei: Linux. . . . .	90
5.3	Accelerarea – număr variabil de procesoare: SunBlade. . . . .	91
5.4	Accelerarea – număr variabil de procesoare: Linux. . . . .	93
5.5	Scalarea rețelei cu numărul de procesoare: SunBlade. . . . .	95
5.6	Scalarea rețelei cu numărul de procesoare: Linux. . . . .	96
5.7	Eficiența pentru actualizarea la epocă . . . . .	97
5.8	Timpii de execuție – SunBlade. . . . .	98
5.9	Timpii de execuție – Linux. . . . .	98
6.1	Aplicații utilizate pentru măsurarea performanțelor antrenării . . . . .	100
6.2	Performanțele compresiei de imagini: SunBlade150 . . . . .	107
6.3	Performanțele compresiei de imagini: Linux . . . . .	108
7.1	Dimensiuni minime optime. . . . .	113

2.1	Conexiune neuronală . . . . .	5
2.2	Schema modelului matematic McCulloch & Pitts . . . . .	8
2.3	Circuite logice cu McCulloch & Pitts . . . . .	8
2.4	Schema modelului matematic Hebb . . . . .	9
2.5	Structura schematică a perceptronului . . . . .	11
2.6	Schema algoritmului cu LMS pentru corecția ponderilor . . . . .	12
2.7	Perceptron uni-strat cu 2 intrări . . . . .	14
2.8	Spațiul intrărilor pentru perceptronul uni-strat . . . . .	14
2.9	Perceptronul pentru problema XOR . . . . .	15
2.10	Reprezentarea spațială a separării . . . . .	15
2.11	Structura schematică a unui perceptron multistrat. . . . .	16
2.12	Modelul fizic Ising. . . . .	17
2.13	Variația funcției sigmoidă pentru diferite valori ale lui $\beta$ . . . . .	19
2.14	Rețea neuronală cu $l$ straturi ascunse . . . . .	20
2.15	Reprezentarea grafică a funcției tangentă hiperbolică . . . . .	24
2.16	Reprezentarea grafică a funcției Gaussiene. . . . .	24
2.17	Reprezentarea grafică a funcției secantă hiperbolică. . . . .	25
2.18	Variația erorii rețelei funcție de valoarea momentului. . . . .	25
2.19	Coborârea cu pantă mică cu gradient conjugat . . . . .	28
2.20	Variația funcției sigmoidă unipolară și a derivatei sale . . . . .	29
4.1	Însumarea gradientelor componenți. . . . .	48
4.2	Însumarea gradientelor componenți într-un arbore binar. . . . .	49
4.3	Circulația ponderilor în pasul înainte. . . . .	51
4.4	Calculul gradientelor parțiali la pasul înapoi. . . . .	53
4.5	Eficiența funcție de dimensiunea setului de antrenare . . . . .	66
4.6	Accelerarea: SunBlade . . . . .	68
4.7	Accelerarea: Linux . . . . .	68
4.8	Set de antrenare fix per procesor: SunBlade . . . . .	70
4.9	Set de antrenare fix per procesor: Linux . . . . .	70
5.1	Distribuirea unităților și a ponderilor. . . . .	74
5.2	Inel de procesoare. . . . .	74
5.3	Broadcast pentru activarea unităților de intrare. . . . .	76
5.4	Calculul valorilor delta pentru unitățile ascunse. . . . .	79
5.5	Dimensiune variabilă a rețelei – SunBlade. . . . .	89
5.6	Dimensiune variabilă a rețelei – Linux. . . . .	91
5.7	Accelerarea – număr variabil de procesoare: SunBlade. . . . .	92
5.8	Accelerarea – număr variabil de procesoare: Linux. . . . .	94
5.9	Scalarea rețelei cu numărul de procesoare: SunBlade . . . . .	95
5.10	Scalarea rețelei cu numărul de procesoare: Linux . . . . .	96
5.11	Eficiența pentru actualizarea la epocă: SunBlade . . . . .	97
5.12	Eficiența pentru actualizarea la epocă: Linux . . . . .	98
6.1	Fereastra de prezentare utilizată de NETtalk. . . . .	101
6.2	Performanțele NETtalk după prezentarea a 100000 de cuvinte . . . . .	103
6.3	Rețea neuronală pentru compresia de imagine. . . . .	105
6.4	Performanțele compresiei de imagini: SunBlade150 . . . . .	107
6.5	Performanțele compresiei de imagini: Linux . . . . .	108
6.6	Configurația rețelei pentru recunoașterea țintelor de sonar . . . . .	109
6.7	Performanțele aplicației Sonar: SunBlade150. . . . .	110
6.8	Performanțele aplicației Sonar: Linux. . . . .	110



# 1 Introducere

Rețelele neuronale nu sunt programate să îndeplinească o anumită sarcină. Ele trebuie să învețe sarcina dată prin metoda încercărilor și a erorilor, în care un supervisor furnizează rețelei neuronale răspunsurile corecte la toate intrările. Regula de antrenare cea mai răspândită este cea denumită *algoritmul backpropagation*.

Totuși, antrenarea acestor rețele neuronale artificiale este o sarcină foarte costisitoare din punct de vedere computațional, necesitând milioane de înmulțiri în virgulă flotantă, chiar și pentru rețele și probleme de dimensiuni reduse. De asemenea, ele necesită volume mari de memorie. Acești doi factori determină ca rețelele neuronale să fie mari consumatoare de timp și impun limitări efective asupra dimensiunii problemelor care pot fi tratate.

Există mai multe căi de compensare a acestor dezavantaje.

O metodă ar fi reducerea dimensiunii problemei prin *pre-procesarea datelor de intrare*, reducând astfel fie numărul de iterații necesare antrenării rețelei, fie chiar dimensiunea rețelei [103]. Trebuie reținut însă că asemenea soluții sunt aproape întotdeauna specifice problemei, și această abordare nu poate fi generalizată pentru acoperirea tuturor tipurilor de probleme.

O altă posibilitate ar fi îmbunătățirea performanțelor regulii de antrenare cu *backpropagation*, fie prin modificări *ad hoc*, fie aplicând rezultate din teoria optimizării numerice. Lucrări care tratează așa-numitele metode ale *gradientului conjugat* [42, 44, 60] aparțin ultimei categorii. În ultimul timp, o serie de eforturi au fost orientate spre investigarea acestei metode și îmbunătățirea performanțelor acestora: metode hibride [93, 19, 17, 18], metode de evaluare a erorii de învățare [21, 38, 40, 55] sau combinarea metodei gradientului conjugat cu metode de pre-procesare [46, 27, 51, 65]. Tot în această categorie mai figurează și abordări paralele ale metodei gradientului conjugat [11].

O a treia abordare constă în accelerarea algoritmilor existenți, fie prin implementarea lor direct în hardware (utilizând tehnici VLSI [89], optice [1]) sau prin adaptarea lor în vederea executării pe o arhitectură paralelă. Această ultimă abordare a revitalizat interesul spre acest algoritm de antrenare cu *backpropagation* [52, 70, 62, 95, 91, 92, 87, 102], interes care scăzuse anterior datorită limitărilor prezentate mai sus. Majoritatea implementărilor paralele din ultimii ani au vizat arhitecturi dedicate [50, 78, 101, 80, 49]. Creșterea puterii de calcul la stațiile de lucru și a vitezei de interconectare a lor a condus spre implementări distribuite ale algoritmului *backpropagation* [66, 88, 45].

În cadrul lucrării de față s-a ales această ultimă variantă: paralelizarea algoritmului de învățare *backpropagation*, utilizând o rețea de calculatoare. După cum se va vedea, este o variantă promițătoare, mai ales prin costul redus al echipamentului.

Teza de față este rezultatul original al investigațiilor efectuate în direcția găsirii unei modalități avantajoase de îmbunătățire a performanțelor în etapa de antrenare a rețelelor neuronale, pe baza resurselor computaționale existente în cadrul Departamentului Calculatoare. Cercetările au fost axate pe dezvoltarea unor algoritmi paraleli care să evolueze într-o rețea de calculatoare, care să accelereze procesul de învățare și să depășească limitările în puterea de calcul, pentru a permite investigații ulterioare în domeniul rețelelor neuronale.

Necesitatea unor asemenea algoritmi a provenit dintr-o problemă concretă: antrenarea unei rețele neuronale, pe o singură mașină, dura circa o săptămână, pentru ca după un asemenea interval de timp, în care trebuia ca nimeni să nu se "atingă" de acea mașină, să se descopere că rețeaua neuronală nu prezenta convergență, întrucât parametrii de învățare nu au fost aleși corect și trebuia reluat procesul, cu alți parametri. Utilizând algoritmi prezentați în lucrare, s-a reușit reducerea timpului de antrenare pentru o rețea neuronală, de la circa o săptămână, la câteva ore, ceea ce a permis efectuarea mai multor experimente, pe diferite seturi de parametri de antrenare, până la obținerea de rezultate satisfăcătoare.

## 1.1 Obiective

Principalul obiectiv al acestei teze a fost construirea unei baze solide pentru procesarea distribuită a fazei de antrenare a rețelelor neuronale într-o rețea de calculatoare. Țelul final a fost implementarea a două simulatoare diferite pentru antrenarea rețelelor neuronale și măsurarea performanțelor lor, comparativ cu varianta secvențială, atât pentru configurații generice de rețele neuronale, cât și testarea lor pentru o serie de aplicații reale de rețele neuronale. Pentru aceasta, pot fi formulate următoarele obiective științifice:

1. Conceperea unor algoritmi care să permită accelerarea procesului de antrenare a unei rețele neuronale, prin distribuirea acestui proces într-o rețea de calculatoare și tratarea în paralel a mai multor tipare de antrenare.

Pornind de la o implementare secvențială a algoritmului de antrenare cu back-propagation, capitolul 4 își propune să rezolve următoarele aspecte:

- proiectarea unei scheme de colectare a gradientilor parțiali generați de fiecare mașină.
  - evaluarea necesarului de memorie.
  - realizarea unei variante avansate, care să elimine redundanța care apare prin memorarea întregii rețele neuronale pe fiecare dintre mașinile care participă la rularea algoritmului.
  - evaluarea performanțelor simulatorului bazat pe această variantă avansată, cu investigații asupra factorilor care pot afecta evoluția algoritmului.
2. Realizarea unor algoritmi care să permită procesarea unor rețele de dimensiuni mari, mai mari decât cele permise de memoria existentă într-un calculator, fără pierderi în performanță. Distribuirea rețelei neuronale se va face prin "secționarea" acesteia și plasarea unui număr de neuroni de pe fiecare strat, pe fiecare mașină din rețeaua de calculatoare. Un simulator bazat pe această strategie este prezentat în capitolul 5, pornind de la o serie de cerințe:
    - simulatorul trebuie să poată fi rulat într-o rețea de calculatoare de dimensiune oarecare, conectate printr-o rețea de tip Ethernet.
    - rescrierea formulelor care intervin în algoritmul backpropagation, în vederea minimizării transmisiei de date între mașini — elementul identificat ca fiind principala cauză pentru care nu se va putea obține accelerarea ideală.
    - se vor efectua o serie de experimente în vederea evaluării performanțelor simulatorului, cu identificarea factorilor care influențează accelerarea și eficiența care se pot obține.

3. Formularea unor recomandări în abordarea unor astfel de probleme.

Simulatoarele realizate vor permite investigații ulterioare în domeniul antrenării rețelelor neuronale — determinarea parametrilor optimi de învățare și a dimensiunii necesare a rețelelor neuronale.

## 1.2 Limitări

Pe parcursul acestui studiu a fost utilizat algoritmul backpropagation clasic de antrenare a rețelelor neuronale feed-forward. Nu au fost luate în considerație tehnici de procesare și alte modalități elaborate în vederea accelerării procesului secvențial de antrenare.

Pentru măsurarea performanțelor simulatoarelor concepute s-au utilizat rețele neuronale generice, fără aplicabilitate practică, de dimensiune convenabilă, iar seturile de tipare utilizate au fost generate aleator. Nu s-au urmărit performanțele de *învățare*

a rețelelor neuronale – care de obicei se măsoară în număr de iterații necesare atingerii unui anumit grad de calitate a învățării – ci s-a dorit să se compare evoluția simulatoarelor distribuite realizate cu cea a unui simulator secvențial, pentru aceeași problemă de rețele neuronale.

### 1.3 Conținutul tezei

Rezultatele obținute sunt sintetizate în teza de față, pe parcursul a trei secțiuni.

Prima secțiune, introductivă, prezintă, pe parcursul capitolului 2, o trecere în revistă a conceptelor fundamentale privitoare la rețelele neuronale, de la modelul biologic la modelele conexioniste, respectiv rețelele cu backpropagation.

Cea de a doua secțiune prezintă două strategii de paralelizare și algoritmi dezvoltati în cadrul cercetărilor întreprinse. Au fost identificate două modalități de distribuire a activității (capitolul 3):

- partiționarea datelor – algoritmul clasic de antrenare se execută în paralel pe mai multe calculatoare, pe seturi diferite de învățare, după care datele sunt colectate și combinate.
- partiționarea rețelei – rețeaua neuronală este distribuită și evoluează pe mai multe calculatoare, comunicarea între neuroni realizându-se prin intermediul rețelei în care sunt interconectate calculatoarele.

Tot în capitolul 3 sunt prezentate elementele de analiză a acestor algoritmi paraleli, sunt enumerate câteva dintre cauzele care ar putea duce la o eficiență scăzută a paralelizării, precum și propuneri de experimente și măsurători de realizat.

Capitolul 4 prezintă în detaliu algoritmul de partiționare a datelor și simulatorul DataParSim. Mai multe mașini execută în paralel și independent același algoritm de antrenare, însă pe seturi de date diferite. O mașină "administrator" se îngrijește de distribuirea inițială a configurației rețelei neuronale și a seturilor de date de antrenare, iar mai târziu colectează rezultatele și le combină într-o nouă configurație a ponderilor din rețeaua neuronală, reluând procesul până la atingerea performanțelor prescrise în privința calității învățării. Sunt prezentate și analizate rezultatele unor experimente în vederea evaluării performanțelor.

Algoritmul pentru varianta cu partiționarea rețelei și simulatorul NetParSim sunt prezentate în capitolul 5. Rețeaua neuronală este divizată astfel încât fiecare dintre calculatoarele utilizate simulează o porțiune din rețea. În acest caz, antrenarea se face pe un singur set de date, iar evoluția mașinilor implicate este sincronă.

Spre deosebire de varianta cu partiționarea datelor, când singura diferență față de varianta secvențială, clasică, este timpul necesar antrenării, partiționarea rețelei permite procesarea unor rețele de dimensiuni mult mai mari, care altfel ar fi imposibil de simulat pe un singur calculator, datorită limitărilor de memorie. Rezultatele experimentelor realizate sunt prezentate și analizate în finalul acestui capitol.

Secțiunea a treia este dedicată prezentării și analizei rezultatelor experimentale pentru o serie de aplicații clasice de rețele neuronale. Testele efectuate au vizat atât studierea performanțelor algoritmilor prezentați, în comparație cu algoritmul secvențial clasic (efectuat pe același tip de mașină), cât și studierea limitărilor impuse ca urmare a paralelizării asupra performanțelor de învățare. Este vorba despre studierea atât a aspectelor cantitative cât și a celor calitative.

În capitolul 7 sunt prezentate concluziile finale și perspectivele de cercetare în continuare.

Testele efectuate au fost realizate pe trei categorii de mașini: o rețea cu până la 12 calculatoare IBM RS6000, într-o rețea ethernet 10 BaseT, respectiv două rețele de până la 16 stații SunBlade 150 și de PC-uri, într-o rețea ethernet 100 BaseT. Desigur, atât algoritmul cu partiționarea datelor cât și cel cu partiționarea rețelei, au fost comparați cu algoritmul secvențial executat pe același tip de mașină.

Capitolul 6 prezintă rezultatele experimentale pentru simulatoare de rețele neuronale "clasice", utilizate de obicei pentru testarea performanțelor: NETtalk, recunoaștere de caractere (fonturi X11), compresie de imagine, recunoaștere de imagini (OCR), recunoașterea țintelor de sonar.

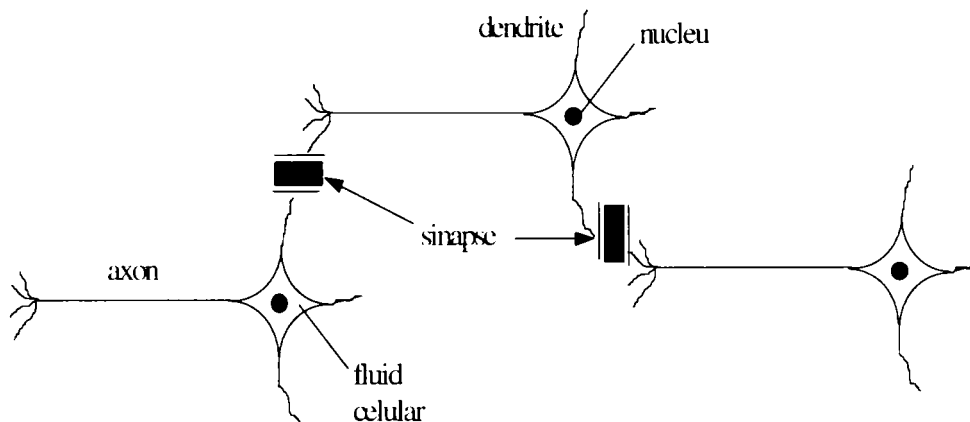
Teza se încheie cu o serie de concluzii desprinse pe parcursul tratării subiectelor propuse în lucrare, rezumatul contribuțiilor, și perspectivele de continuare a cercetărilor în domeniul paralelizării fazei de antrenare a rețelelor neuronale.

Se menționează faptul că testele au fost efectuate pe rețeaua de stații IBM RS6000 din cadrul Laboratorului de Grafică al Departamentului de Ingineria Calculatoarelor și a Programării, Universitatea "Politehnica" din Timișoara, respectiv pe o rețea de stații SunBlade 150 și pe o rețea de PC-uri Hp-Compaq, cu amabilitate puse la dispoziție de Mobile Networks Division din cadrul Alcatel-Lucent SA, România.



### 2.1 Celula nervoasă

Celula nervoasă — sau *neuronul* — este elementul fundamental al sistemului nervos. Schematic, în figura 2.1 este redat simplist un lanț de neuroni, având menționate componentele semnificative din punctul de vedere al utilității în modelele matematice ulterioare [31]:



**Figura 2.1:** Conexiune neuronală

Celula nervoasă are diametrul de ordinul micronilor și este alcătuită din:

- *corpul celular* — conține nucleul neuronului și guvernează transformările biochimice de sinteză a enzimelor necesare funcționalității celulei nervoase;
- *dendritele* — sunt extensii filiforme cu ramificații arborescente ale corpului celular, având diametrul de ordinul zecimilor de microni și lungimi de până la 10 microni. Dendritele sunt principalii receptori de semnal ai neuronului.
- *axonul* — este o prelungire fibroasă, de lungime considerabilă în comparație cu restul celulei nervoase, variind între 0,2-1 centimetru pentru neuronii din cortex și 10 centimetri pentru neuronii periferici. Axonul reprezintă conexiunea de ieșire a semnalului procesat de către neuron.

Modalitatea de interconectare spațială a celulelor nervoase este realizată prin intermediul unor zone cu conductanță electrică îmbunătățită, denumite sinapse, și care constituie practic o legătură între axonul unui neuron și o dendrită a unui alt neuron, cu rol de comunicare între cele două celule.

Măsurătorile au arătat că eficiența energetică a creierului uman este de aproximativ  $10^{-16}$  J/(operație·s) [23], în timp ce cele mai performante calculatoare au o eficiență de  $10^{-6}$  J/(operație·s).

Neuronii au caracteristici comune cu cele ale celulelor vii în general, însă ca organizare și sistem biochimic prezintă individualități ce permit specializarea lor în realizarea unor funcțiuni specifice ca:

- primesc semnale de la neuronii vecini,
- conduc semnalul nervos,
- integrează semnalele primite,

- amplifică sau mediază semnalul nervos rezultat,
- transmit semnalul nervos pe canale preferențiale spre neuronii receptori.

## 2.2 Rețele neuronale

### 2.2.1 Modele neuronale artificiale

În forma cea mai generală - și în accepția inițială - o rețea neuronală este un algoritm proiectat să modeleze modul în care creierul realizează o anumită funcțiune, estimându-se empiric, cu aproximații acceptabile, modalitatea de procesare a informației de către cortex. Acest scop este valabil doar pentru cazul primelor modele. Ulterior modelele neuronale au părăsit ideea inițială de simulare a propagării și procesării semnalului nervos, păstrând însă specificul nomenclaturii biologice pentru desemnarea unor etape. Astfel, aplicațiile tratate în continuare se vor referi la rețele neuronale care își realizează scopul propus printr-un proces inițial de *învățare*. Pentru o bună performanță aplicativă, rețeaua își creează ponderi de interconectare între unitățile de intrare și cele de ieșire, unități care vor fi denumite *neuroni* sau *unități de procesare*. O definiție simplistă a noțiunii actuale de rețea neuronală este dată de Aleksander și Morton, în [3]:

O rețea neuronală este un procesor cu distribuție paralelă care are o mare capacitate de a stoca cunoștințe dobândite din experimente, în scopul utilizării active ale acestora.

Simularea acțiunii cerebrale este elocventă în cele două acțiuni de bază:

- cunoștințele sunt acumulate în timpul unui proces de învățare,
- stocarea cunoștințelor se realizează prin valoarea conexiunilor interneuronale, denumite și ponderi sinaptice.

### 2.2.2 Scurt istoric

- 1943 Warren McCulloch, un neurofiziolog, și Walter Pitts, un tânăr matematician publică un articol [53] despre cum s-ar putea să lucreze neuronii. Au modelat o rețea neuronală simplă prin intermediul unor circuite electrice.
- 1949 Donald Hebb publică o carte [32], *Organization of Behavior*, în care reiterează conceptul de neuroni și modul lor de lucru. El enunță o regulă care afirmă faptul că traseele neuronale se accentuează de fiecare dată când acestea sunt utilizate.
- 1950- Primele modelări (rudimentare), efectuate de Nathaniel Rochester în laboratoarele IBM [71]. Prima încercare a fost nereușită, însă următoarele au fost încununuate de succes. Era perioada înfloririi metodelor tradiționale de calcul și cercetările în domeniul neuronal au trecut în planul secund.
- 1956 *Dartmouth Summer Research Project on Artificial Intelligence* impulsionează cercetările atât în domeniul inteligenței artificiale cât și în cel al rețelilor neuronale. În anii care au urmat, John von Neumann sugerează imitarea funcțiilor simple ale neuronilor prin intermediul releelor telegrafice și al tuburilor electronice. De asemenea, Frank Rosenblatt la Universitatea Cornell, începe studiul Perceptronului. Acesta, implementat hardware, este cea mai veche rețea neuronală încă utilizată în zilele noastre. Din păcate, limitările perceptronului au fost evidențiate destul de curând, ceea ce a dus la o perioadă de stagnare. Abia în 1969, Marvin Minsky și Seymour Papert demonstrează matematic ineficiența lui și oferă și rezolvarea [58].

- 1959 Bernard Widrow și Marcian Hoff de la Stanford dezvoltă modelele denumite de ei ADALINE și MADALINE — primele rețele neuronale aplicate în probleme reale, eliminarea ecoului pe liniile telefonice, care se utilizează și în zilele noastre.
- A urmat o perioadă de stagnare, datorată și lipsei de finanțare, dar și exagerărilor în aprecierea potențialului rețelelor neuronale, care au dus chiar la teama de asemenea "mașini gânditoare". Scrierile lui Isaac Asimov reflectă aceste temeri prin imaginarea unor lumi în care roboții preiau toate sarcinile oamenilor. Alți scriitori imaginează scenarii și mai sinistre, un singur exemplu fiind HAL din *2001: A Space Odyssey*.
- 1982 John Hopfield de la Caltech prezintă o lucrare la Academia Națională de Științe [34]. Abordarea lui este nu doar simpla modelare a creierului, ci crearea de dispozitive utile. El demonstrează clar, analitic, cum ar putea lucra asemenea rețele și ce s-ar putea face cu ele.
- 1982 La Kyoto, în Japonia, are loc *US-Japan Joint Conference on Cooperative/Competitive Neural Networks*. Ca urmare, Japonia anunță efortul către Generația a Cincea, generând o îngrijorare în SUA că americanii ar putea rămâne în urmă. Curând, finanțarea nu a mai fost o problemă.
- 1985 American Institute of Physics organizează ceea ce va deveni o întâlnire anuală — *Neural Networks for Computing*.
- 1987 Institute of Electrical and Electronic Engineer's (IEEE) organizează prima *International Conference on Neural Networks*, la care iau parte mai bine de 1800 de participanți.
- 1989 La conferința *Neural Networks for Defense*, Bernard Widrow afirmă că deja a început al IV-lea Război Mondial, în care câmpul de bătaie sunt producția și comerțul.
- 1990 *US Department of Defense Small Business Innovation Research Program* nominalizează 16 teme care vizau explicit rețelele neuronale, iar alte 13 menționează posibilitatea utilizării rețelelor neuronale.

### 2.2.3 Modelul McCulloch & Pitts

Înainte de descoperirea modalității de funcționare a sistemului nervos celular [33], pentru care Hodgkin, Huxley și Eccles au primit Premiul Nobel în 1963, în 1943, McCulloch și Pitts au elaborat un model simplist al funcționării unui neuron, pornind de la funcționalitatea acestuia [53]. Ei au presupus că neuronul se comportă în modul următor: din sinapse, prin dendrite, acesta acumulează o sarcină electrică ce provine de la axonii neuronilor vecini; în momentul în care valoarea sarcinii depășește o anumită valoare  $T$  — numită *prag* — neuronul devine activ, transmitând prin axonul său un curent constant  $U$ . Această funcționare poate fi descrisă matematic prin:

$$o = \begin{cases} U & \text{pentru } \sum_{i=1}^n x_i \geq T \\ 0 & \text{pentru } \sum_{i=1}^n x_i < T \end{cases} \quad (2.1)$$

unde  $x_i$  = intrarea  $i$ ;  
 $o$  = ieșire.

Schematic, modelul este redat în figura 2.2:

O astfel de unitate este denumită *unitate de procesare*. Considerând valoarea pragului  $T = 1$  și valoarea tensiunii de ieșire  $U = 1$ , se pot construi, cu aceste unități de procesare unitare, circuite logice de tipul AND, NAND, OR, și respectiv celule de memorie, în modul sugestiv prezentat în figura 2.3.

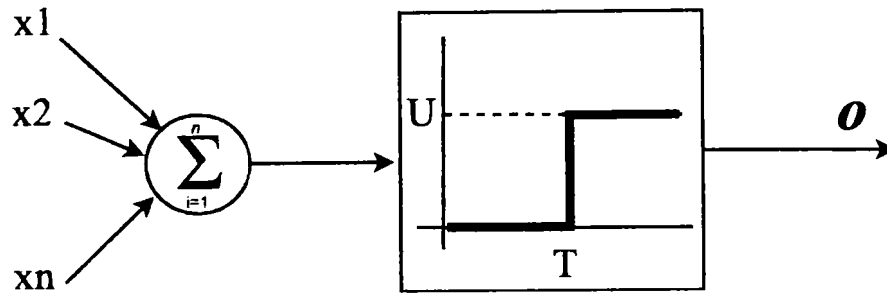


Figura 2.2: Schema modelului matematic McCulloch & Pitts

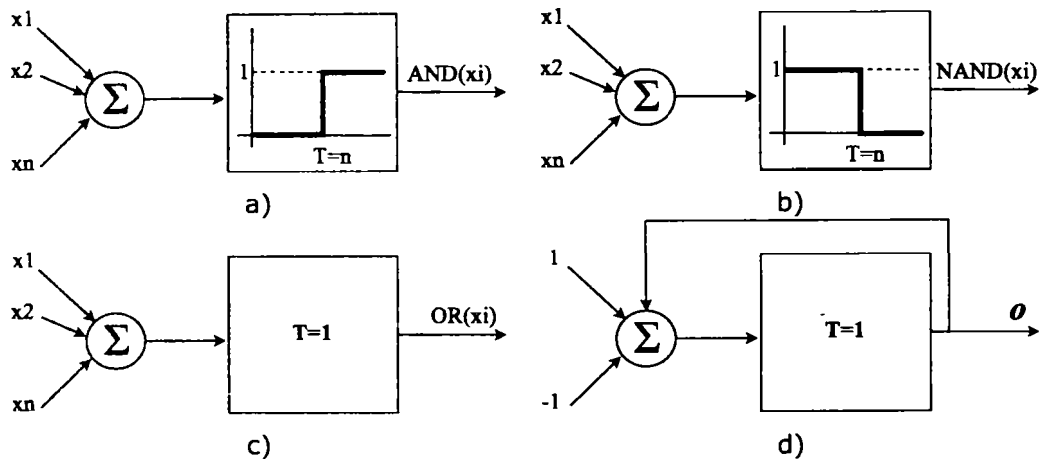


Figura 2.3: Circuite logice realizate cu unități de procesare McCulloch & Pitts

Cei doi matematicieni au sugerat posibilitatea ca și creierul uman să fie o aglomerare de asemenea circuite logice, care ar procesa informația într-un mod similar cu calculatoarele digitale, al căror model era definitivat în aceeași perioadă.

Această prezumție s-a dovedit a fi greșită, însă are valoare intuitivă, deoarece s-a dovedit ulterior că neuronii funcționează oarecum după modelul activării la atingerea unui prag.

Acest model neuronal este caracterizat prin formalism și printr-o definiție matematică elegantă, având un potențial de calcul considerabil. Simplitatea sa presupune:

- stări binare de tip 0 și 1,
- asumarea operării în timp discret,
- sincronicitatea operării tuturor neuronilor.

#### 2.2.4 Modelul Hebb

În 1949, Daniel Hebb, un biolog care se ocupa cu studiul comportamentului primatelor, presupune că prin modificarea sinapselor conectivitatea celulelor nervoase este într-o continuă schimbare, pe măsură ce organismul își dezvoltă diferite aptitudini. Având la

bază prezumțiile lansate de laureatul Nobel Ramon y Cajal încă din 1899, (republicate în 1995, în *Histology of the nervous system* [10], Hebb a enunțat o *regulă a învățării* [32], conform căreia

Valoarea efectivă a variabilei de conexiune între doi neuroni crește prin activarea repetată a conexiunii.

Cartea a avut un răsunet imens printre psihologi, dar a fost neglijată de ingineri, fapt ce a dus la ignorarea acestei reguli de învățare în domeniul rețelelor neuronale. Abia în 1956 a fost utilizată această regulă de învățare, evidențiindu-se fenomenul de *inhibiție* a conexiunii neuronale [71].

Conform regulii de învățare a lui Hebb, sinapsele care conectează neuronii nu transmit axonului neuronului receptor decât o parte a sarcinii electrice primite de la axonul neuronului transmițător. Reprezentarea matematică a acestei reguli se obține utilizând un coeficient subunitar de transmisie, numit *pondere* și notat cu  $w$ .

Ponderile sunt asociate neuronului receptor și își pot modifica valorile în timp conform regulii:

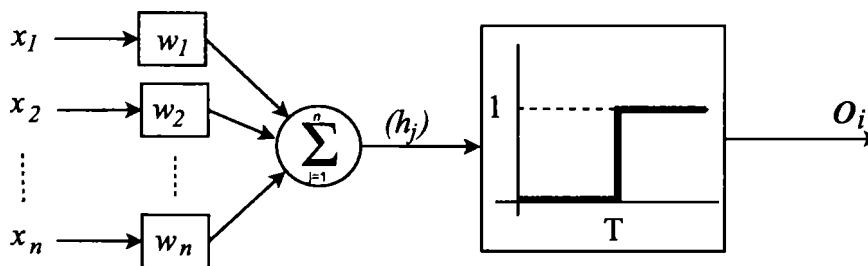
- dacă neuronul transmițător și neuronul receptor sunt ambii activi în același interval de timp, valoarea ponderii asociate conexiunii crește;
- dacă fenomenul de activare concomitentă se repetă de mai multe ori sau dacă intervalul de timp al dublei activări este mare, valoarea ponderii asociate va crește considerabil;
- dacă neuronul transmițător și neuronul receptor nu sunt simultan activi într-o perioadă de timp mare, valoarea ponderii asociate conexiunii scade.

Ulterior, s-a dovedit că într-adevăr, această regulă este valabilă, și stă la baza mecanismului de învățare, respectiv uitare, în creierile naturale.

Plecând de la această regulă, se poate descrie o unitate de procesare, cu ponderi asociate, care este modelul matematic utilizat și în prezent, cu variațiuni nesemnificative. Astfel, considerând o unitate cu indicele  $i$ , cu  $n$  intrări de indice  $j$ , fiecare având asociată o pondere  $w_j$ , funcționarea acesteia se poate descrie prin:

$$o_i = \begin{cases} 1 & \text{pentru } \sum_{j=1}^n x_j \cdot w_j \geq T \\ 0 & \text{pentru } \sum_{j=1}^n x_j \cdot w_j < T \end{cases} \quad (2.2)$$

sau prin schema din figura 2.4.



**Figura 2.4:** Schema modelului matematic Hebb

Pentru simplificare, se va nota suma ponderată care reprezintă intrarea totală în unitatea de procesare cu  $h_i$ :

$$h_i = \sum_{j=1}^n w_j \cdot x_j \quad (2.3)$$

O asemenea unitate poate fi privită ca un procesor de informație foarte primitiv. Dacă se conectează un număr de astfel de procesoare între ele, se obține o *rețea neuronală artificială*, sau, mai corect spus, o *structură de calcul conexionistă*<sup>1</sup>.

Modul de conectare a acestor procesoare reprezintă *topologia rețelei* și determină modul de calcul al ponderilor, și, implicit, funcționalitatea rețelei.

Calculul ponderilor este denumit *învățare* sau *antrenare*, iar metodele iterative de calcul sunt cunoscute sub denumirea de *algoritmi de învățare* sau de *antrenare*.

Ashby, în cartea sa *Design for a brain*, [5], destinată comportamentului adaptiv, presupune și argumentează că un sistem natural nu se naște învățat, ci ajunge la anumite performanțe prin corelarea variabilelor ce descriu evenimentele percepute sistemocentric. Cartea tratează aspectele dinamice ale organismului viu ca sistem abstract, mecanic, și evaluează conceptul de stabilitate.

În 1954, teza de doctorat a lui Minsky [56], se constituie ca prolog al lucrării *Steps toward artificial intelligence* [57], ce constituie și astăzi unul din capitolele fundamentale în domeniul rețelelor neuronale.

Steinbuch, în *Die lernmatrix* [84] continuă ideea lansată de Taylor în *Electrical simulation of some nervous system functional activities* [90] despre o *memorie asociativă* realizabilă prin așa-numitele *matrici de învățare*, care ar fi interpușe între câmpuri de senzori receptori și efectori motori.

Contribuții semnificative la teoria memoriei asociative au adus ulterior și Anderson în *A simple neural network generating an interactive memory* [4], Kohonen în *Correlation matrix memories* [47] și Nakano în *Associatron – a model of associative memory* [61], în același an (1972); independent, au introdus ideea unei matrici corelative de memorie ("correlation matrix memory"), care utilizează o regulă de învățare bazată pe produsul de ieșire al algoritmului.

O figură particulară în contextul evoluției conceptului rețelelor neuronale a fost von Neumann, care a utilizat ideea de redundanță pentru proiectarea unei rețele neuronale privită ca un ansamblu unitar cu un comportament cert, alcătuit din neuroni care pot fi priviți separat ca și componente cu comportament incert. Această idee a fost dezvoltată de Winograd și Cowan, care au utilizat o reprezentare redundant distribuită a rețelelor neuronale [99].

### 2.3 Modelul Rosenblatt – perceptronul

După 15 ani de la enunțarea modelului McCulloch & Pitts, devenit clasic, un prim pas spre rezolvarea problemei recunoașterii tiparelor îl face Rosenblatt, prin lucrarea "*The perceptron: A probabilistic model for information storage in the brain*" [72], în care enunță teorema convergenței perceptronului, și prin lucrările adiacente, în care dezvoltă această teorie [73], [75], [74].

Topologia cea mai simplă este asociată unei structuri conexioniste, denumită *perceptron*, a cărui structură schematică este prezentată în figura 2.5.

În forma cea mai simplă, el este compus dintr-un vector de intrare  $n$ -dimensional,  $x$ , o unitate de procesare și o singură ieșire. Menirea unei astfel de unități este de a învăța o transformare de tipul  $T : \{-1, 1\} \rightarrow \{-1, 1\}$  utilizând seturi de tipare cu intrări de tip  $x$  și ieșirile corespunzătoare  $o = T(x)$ .

Unitatea de ieșire a unui perceptron este un element liniar cu prag, care calculează ieșirea  $o$  conform relației:

$$o = F \left( \sum_{i=1}^n w_i \cdot x_i + \theta \right) \quad (2.4)$$

<sup>1</sup>Expresia *structură de calcul conexionistă* reprezintă de fapt noțiunea descriptivă la care s-a ajuns prin încercarea de a exprima matematic teoriile de moment ce sugerează posibilități de modelare a propagării semnalului prin sistemul neural uman. Este practic o expresie mult mai adecvată atribuită algoritmilor astfel obținuți, care nu își mențin scopul inițial, ci îl adoptă pe cel al aplicațiilor pentru care au fost dedicați.

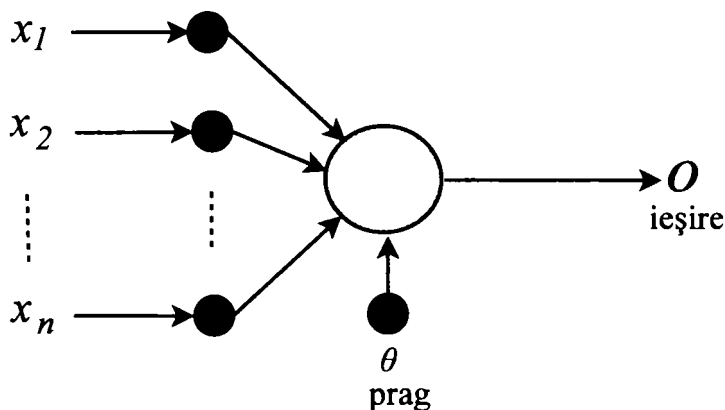


Figura 2.5: Structura schematică a perceptronului

În care funcția  $F$  este funcția Heaviside sau *signum*, iar  $\theta$  poate fi orice funcție sau constantă care ajustează conexiunile astfel încât să se obțină ieșirea dorită. În acest caz, funcția  $F$  este:

$$F(i) = \begin{cases} 1 & \text{pentru } i > 0 \\ -1 & \text{altfel} \end{cases} \quad (2.5)$$

Un tipar de învățare binar pentru un asemenea perceptron constă dintr-o pereche de vectori cu valori binare, unul de lungime  $n$  (numărul de unități de intrare), iar celălalt de lungime  $m$  (numărul de unități de ieșire). Numărul total de tipare alcătuiește o bază de învățare.

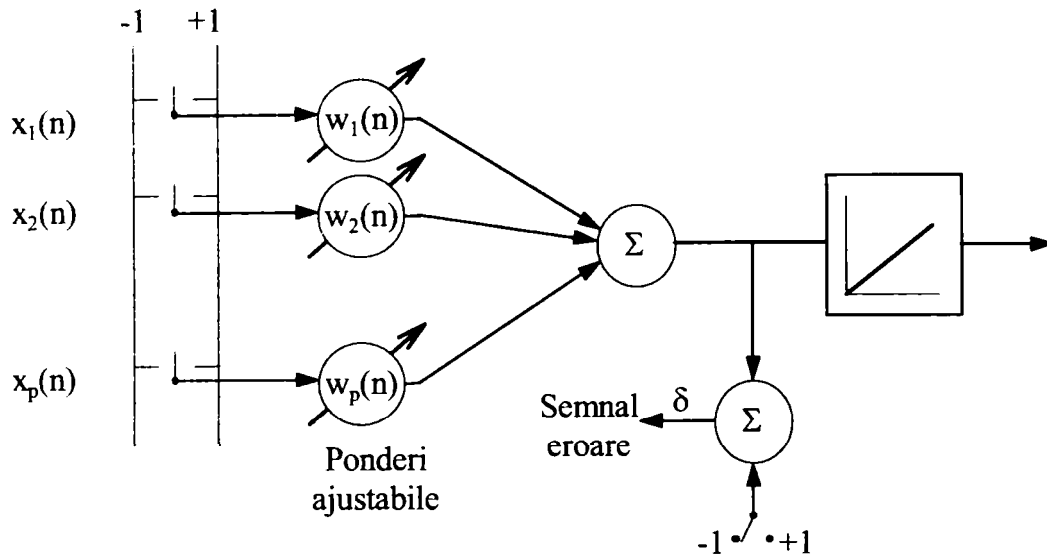
Algoritmul de învățare este următorul:

1. Se prezintă la intrarea rețelei primul vector al unui tipar, activând unitățile corespunzătoare valorii 1.
2. Se prezintă la ieșirea rețelei vectorul 2 al aceluiași tipar, activând unitățile de pe nivelul 2 corespunzător valorilor 1.
3. Se caută perechi de unități care sunt conectate și activate simultan, și se ajustează valoarea ponderii conexiunii dintre acestea cu o valoare mică  $\eta$ . Valoarea lui  $\eta$  se alege în funcție de numărul tiparelor de învățare.
4. Se reiau pașii 1, 2, 3, pentru fiecare tipar – pereche de vectori din baza de învățare.
5. După prezentarea tuturor tiparelor din baza de învățare, ponderile se vor modifica astfel încât prezentând la intrare un vector 1 dintr-un tipar, activând corespunzător unitățile de pe primul nivel și calculând ieșirile din unitățile de pe nivelul 2, se obține vectorul 2 din tiparul de învățare respectiv.

### 2.3.1 Regula delta

În 1960, Widrow și Hoff au introdus *algoritmul celor mai mici pătrate* ("Least Mean Square Algorithm" – LMS), cunoscut și ca *Regula Delta*, ca algoritm de corecție a ponderilor în rețea [98]. Astfel a luat naștere elementul liniar adaptiv, denumit și ADALINE ("ADaptive LINear Element").

Mecanismul de funcționare care utilizează această regulă de corecție a ponderilor este prezentat sugestiv în figura 2.6.



**Figura 2.6:** Schema algoritmului cu LMS pentru corecția ponderilor

Widrow și Hoff realizează practic o generalizare importantă a algoritmului de antrenare a perceptronului, extinzând tehnica intrărilor și ieșirilor continue. *Regula delta* a fost aplicată mai ales unităților cu ieșiri pur liniare.

Pentru un astfel de perceptron, cu o singură unitate de ieșire, ieșirea este dată de relația:

$$o = \sum_j w_j x_j + \theta \quad (2.6)$$

Funcția erorii, așa cum o indică numele de "cel mai mic pătrat", este suma pătratelor erorilor. Deci eroarea totală  $E$  este definită ca:

$$E = \sum_p E^p = \frac{1}{2} \sum_p (d^p - a^p)^2 \quad (2.7)$$

unde indicele  $p$  evoluează pe setul datelor de intrare, iar  $E^p$  reprezintă eroarea setului indice  $p$ . Variabila  $d^p$  este ieșirea dorită la aplicarea setului  $p$ , iar  $a^p$  este ieșirea actuală pentru acest set. Procedura LMS găsește valorile tuturor ponderilor care minimizează eroarea funcției, prin metoda numită a *gradientului descendent*. Ideea este de a efectua o schimbare în pondere, direct proporțională cu negativul derivatei erorii determinată pe setul curent, în concordanță cu fiecare pondere:

$$\Delta^p w_j^p = -\eta \frac{\partial E^p}{\partial w_j} \quad (2.8)$$

unde  $\eta$  este o constantă de proporționalitate. Derivata este

$$\frac{\partial E^p}{\partial w_j} = \frac{\partial E^p}{\partial a^p} \frac{\partial a^p}{\partial w_j} \quad (2.9)$$

Datorită liniarității unităților,

$$\frac{\partial a^p}{\partial w_j} = x_j \quad \text{și} \quad \frac{\partial E^p}{\partial a^p} = -(d^p - a^p) \quad (2.10)$$



astfel încât

$$\Delta^P u_j^P = \eta \delta^P x_j \quad (2.11)$$

unde

$$\delta^P = d^P - a^P \quad (2.12)$$

este diferența dintre ieșirea dorită și actuala ieșire pentru setul  $p$ . Regula delta modifică ponderea corespunzătoare ieșirii dorite pentru orice polaritate și pentru unități cu intrări-ieșiri continue sau binare. Aceste caracteristici au deschis un bogat câmp de noi aplicații.

Spre deosebire de algoritmul inițial, ce utilizează incrementarea progresivă, acest algoritm ia în considerare și *greșelile* pe care le face rețeaua și încearcă să le corecteze prin schimbarea valorilor ponderilor.

Acest algoritm are rezultate mai bune decât regula lui Hebb simplă din secțiunea 2.2.4, mai ales în cazul prezentării repetate a bazei de învățare, luând de fiecare dată tiparele într-o altă ordine, de preferință aleatoare.

Etapa de traversare o dată a rețelei de către întreaga bază de învățare poartă numele de *epocă*, iar numărul epocilor în care rețeaua învață cu precizia dorită întreaga bază de învățare este un factor de performanță al rețelei.

O deficiență a modelului conexiunist este faptul că precizia de învățare a tiparelor variază de la un tipar la altul, în funcție de modul în care se încadrează tiparul în logica ponderată a bazei.

În cazul în care se dorește ca rețeaua să obțină la ieșire vectorul caracteristic unui vector de intrare care nu face parte din baza de învățare, dar este asemănător ca distribuție cu vectorii de intrare din baza de învățare, este necesară definirea unui parametru ce caracterizează performanța de învățare și anume *eroarea totală* – relativă la toate tiparele:

$$E = \sum_{\mu=1}^p \sum_{i=1}^m \delta_i^\mu \quad (2.13)$$

Eroarea maximă a rețelei este

$$E_{max} = m \cdot p$$

dacă vectorii au valori binare de tipul (0,1), și

$$E_{max} = 2 \cdot m \cdot p$$

dacă valorile sunt de tipul (-1,1).

Raportul

$$\frac{E}{E_{max}} \cdot 100 = E\%$$

este eroarea relativă procentuală a rețelei, și este utilizat ca și criteriu de oprire a etapei de învățare a rețelei.

Datorită avantajului acestui model, de a recunoaște vectorul de ieșire în cazul prezentării la intrare a unui vector care nu face parte din baza de învățare, perceptronul a fost utilizat pe scară largă.

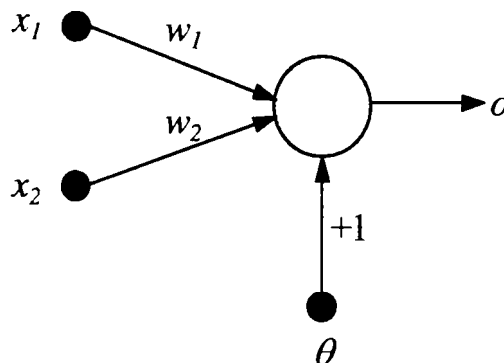
Pe scurt, diferențele principale dintre ADALINE și Perceptron sunt:

- la ADALINE are loc o ponderare liniară, ce se formează adaptiv (printr-o metodă de antrenare supervizată), care produce semnalul de răspuns dorit pentru vectori de intrare diferiți; rezultatul este optim, în sensul celor mai mici pătrate.
- la Perceptron, există o unitate de procesare suplimentară, pentru extracția combinațiilor de semnal din tiparele de intrare. Circuitele de ieșire au ca singur scop sesizarea diferenței dintre vectorii de intrare, fiecare din ei corespunzând unei clase diferite de tipare.

După ADALINE, Widrow propune și realizează pentru prima dată o rețea neuronală stratificată, antrenabilă, cu elemente adaptive multiple, în *Generalization and information storage in networks of adaline "neurons" [97]*.

### 2.3.2 Problema XOR

Euforiei generate de perceptron i-au pus capăt în 1969 Minsky și Papert, care utilizând argumente matematice simple dar elegante, au demonstrat că un perceptron cu un singur strat, numit în continuare *perceptron uni-strat* are limitele sale [58]. În acest context, una dintre problemele care au rămas drept istorice este problema imposibilității unui perceptron uni-strat de a reprezenta o funcție sau-exclusiv (XOR).



**Figura 2.7:** Perceptron uni-strat cu 2 intrări

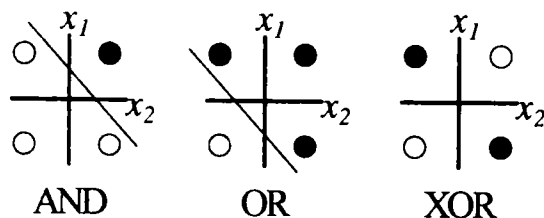
În perceptronul din figura 2.7, valoarea totală a intrării este dată de

$$i = w_1 \cdot x_1 + w_2 \cdot x_2 + \theta \quad (2.14)$$

Pentru o constantă  $\theta$ , ieșirea perceptronului este într-unul din semiplanele delimitate de linia de diviziune definită de ecuația

$$w_1 \cdot x_1 + w_2 \cdot x_2 = -\theta \quad (2.15)$$

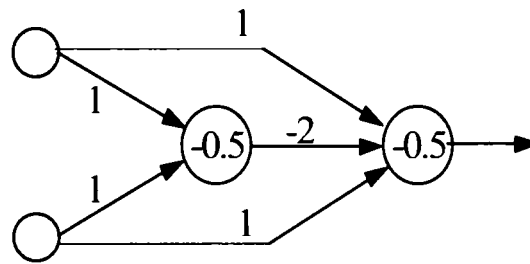
Conform figurii 2.8, spațiul de intrare constă din cele patru puncte (cercuri); cercurile pline, localizate la  $(1, -1)$  și  $(-1, 1)$  nu pot fi separate printr-o linie dreaptă de cercurile goale, localizate la  $(-1, -1)$  și  $(1, 1)$ .



**Figura 2.8:** Reprezentarea geometrică a spațiului intrărilor pentru perceptronul uni-strat

### 2.3.3 Perceptronul multistrat

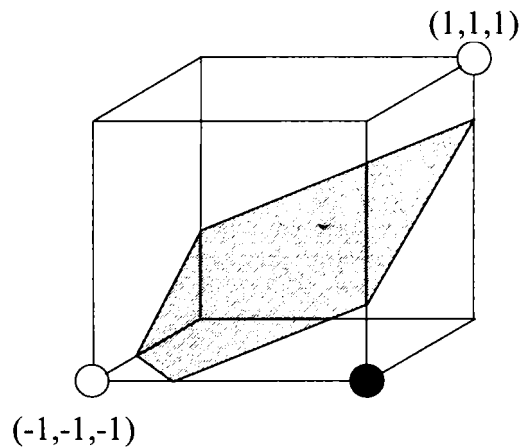
Tot Minsky și Papert au dat și soluția, introducând un nou strat, conectat la toate intrările. Au extins astfel rețeaua la perceptronul multistrat, care poate rezolva problema XOR.



**Figura 2.9:** Perceptronul cu două intrări și unitate ascunsă pentru problema XOR

Rezolvarea problemei XOR este prezentată în figura 2.9.

De data aceasta, cele două grupuri de puncte sunt separate spațial printr-un plan de separare, ca în figura 2.10.



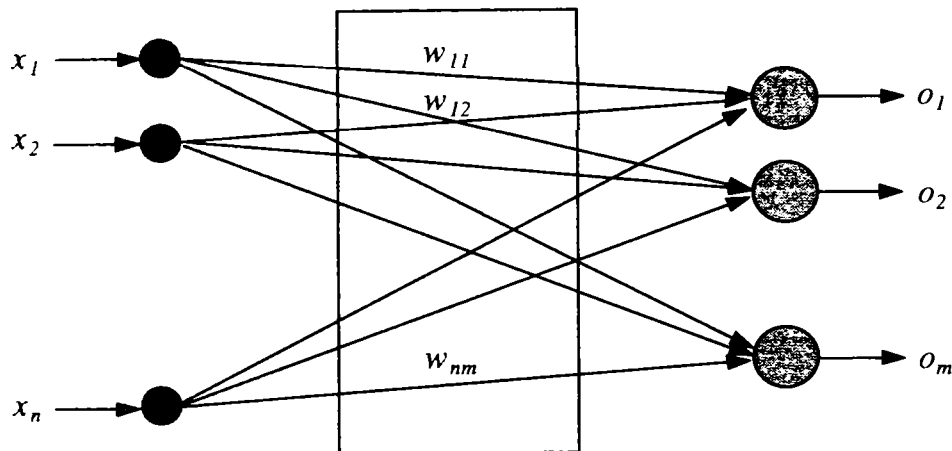
**Figura 2.10:** Reprezentarea spațială a separării celor două grupuri de puncte.

Acest exemplu simplu demonstrează că adăugarea unor unități ascunde crește clasa de aplicații rezolvabile de către perceptronul multistrat.

În figura 2.11 se prezintă schema unui perceptron având la intrare un vector de dimensiune  $n$ , iar la ieșire un vector  $m$ -dimensional. Fiecare intrare  $i$  este multiplicată pentru toate unitățile din rețea, care vor avea  $n$  intrări.

O problemă importantă întâlnită în proiectarea rețelelor tip perceptron multistrat a fost problema *acordării de credit neuronilor din stratul ascuns al rețelei* ("credit assignment theory"). Cel care a introdus pentru prima dată terminologia de *credit assignment* a fost tot Minsky [57], ulterior căruia în anii '60 strădaniile cercetătorilor din domeniu au fost dedicate problemelor legate de această teorie, cu rezultate contradictorii, fără a o rezolva definitiv. A trebuit să vină perioada anilor '80 pentru a găsi soluțiile de ieșire din impas. Printre motivele unei asemenea întârzieri:

- tehnologic — nu existau calculatoare personale sau stații de lucru suficient de performante care să permită experimentarea;



**Figura 2.11:** Structura schematică a unui perceptron multistrat.

- psihologic — monografia [58] nu a fost deloc încurajatoare pentru continuarea studiilor în domeniu;
- financiar — derivat din cel psihologic, care nu încuraja nici o agenție în lansarea unor programe de cercetare în domeniu.

Din nefericire, aceste motive, și altele, au dus la sistarea interesului pentru studiul rețelelor neuronale. Singurele cercetări ale anilor '70, totuși colaterale domeniului, au fost cele în psihologie și neurologie. De aceea, această decadă a fost caracterizată ca fiind "the dormancy decade" (decada de somnolență).

## 2.4 Modelul Hopfield

În 1982, Hopfield a utilizat ideea unei *funcții de energie* [34], pentru a reformula un nou mod de înțelegere a calculelor dintr-o rețea neuronală, folosind conexiuni sinaptice simetrice, și stabilind o analogie între o rețea neuronală recurentă și *modelul Ising* utilizat în fizica statistică. Aceste tipuri de rețele sunt cunoscute sub denumirea de rețele Hopfield.

### 2.4.1 Modelul fizic Ising

Conform acestui model, un material magnetic este compus dintr-o rețea regulată de atomi cu proprietăți magnetice, care au asociată o mărime numită *spin magnetic*, orientat vectorial în sensul dat de compunerea forțelor magnetice la care este supus atomul. Pentru cazul ipotetic în care este permisă orientarea spinului pe o singură direcție, în două sensuri, valoarea acestuia poate fi memorată într-o variabilă binară  $s_i$ . Alegând pentru aceasta variabilă valorile -1 și 1, se obține un model simplificat - modelul Ising (figura 2.12):

Fiecare spin  $s_i$  este influențat de câmpul magnetic  $H_i$ , care este compus din câmpul magnetic exterior  $H_{ext}$  și câmpul produs de atomii învecinați.

Ising a propus pentru calculul lui  $H_i$  relația:

$$H_i = \sum_j p_{ij} \cdot s_j + H_{ext} \quad (2.16)$$

unde  $p_{i,j}$  este un coeficient dat de poziția relativă dintre atomii  $i$  și  $j$ .

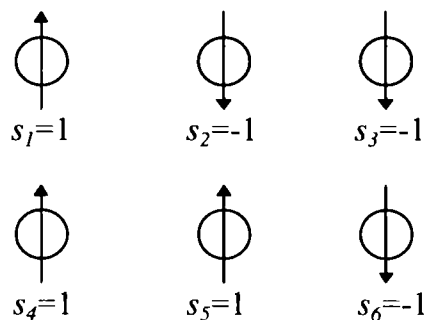


Figura 2.12: Modelul fizic Ising.

Materialul magnetic poate avea înmagazinată o energie potențială a cărei valoare poate fi calculată cu relația:

$$H = -\frac{1}{2} \sum_i \sum_j p_{ij} \cdot s_i \cdot s_j - H_{ext} \sum_i s_i \quad (2.17)$$

Pentru simplificarea discuției în ce privește utilizarea modelului Ising la construcția unei rețele neuronale, în continuare se va considera  $H_{ext} = 0$ .

### 2.4.2 Definirea modelului Hopfield

Trecerea de la modelul Ising la o structură conexiunistă a fost definitivată de Hopfield în 1984 [35]. Rețeaua conexiunistă bazată pe acest model este topologic diferită de perceptron, având fiecare unitate legată de celelalte unități, toate legăturile fiind ponderate. Nu există diferențiere între unitățile de intrare și cele de ieșire, ieșirea fiecărei unități fiind multiplicată ca intrare către toate unitățile, inclusiv ea însăși. Ponderile reprezintă o matrice de dimensiune  $N \times N$ , unde  $N$  este numărul de unități.

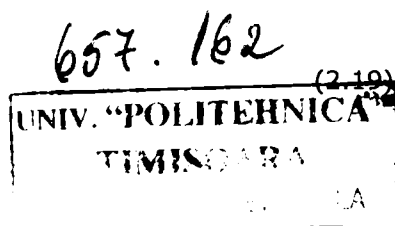
Funcționarea rețelei este, de asemenea, diferită de cea a perceptronului. Astfel, dacă la o rețea structurată pe straturi de intrare, respectiv de ieșire, aplicăm un vector de intrare primului strat și calculăm activările din cel de al doilea strat pentru a determina ieșirea, pentru noua rețea vom considera ca intrare o stare generală, a întregii rețele, reprezentabilă printr-un vector de stare  $S = [S_1, S_2, \dots, S_N]$ .

Datorită câmpului magnetic exercitat asupra lor, atomii dintr-un material magnetic își pot schimba spinul. Acceptând prin convenție că o unitate este activă dacă are spinul egal cu 1 și este inactivă dacă are spinul egal cu -1, în modelul Ising, valoarea spinului se poate calcula cu funcția *signum*. Astfel, dacă semnul valorii câmpului din locația atomului este pozitiv, atunci și spinul este pozitiv, iar dacă este negativ spinul este și el negativ. Considerând că  $h_i = H_i$ , (intrarea într-o unitate este egală cu câmpul) și  $w_{ij} = p_{ij}$  (ponderea unei intrări este egală cu coeficientul de interacțiune dintre atomi), se poate afirma că unitățile acestei rețele conexiuniste au următoarea funcție de transfer:

$$s_i = \text{signum}(h_i) = \begin{cases} 1 & \text{dacă } h_i \geq 0 \\ -1 & \text{dacă } h_i < 0 \end{cases} \quad (2.18)$$

numită *lege dinamică de evoluție*, unde:

$$h_i = \sum_{j=1}^N w_{ij} s_j$$



Un astfel de sistem caută să ajungă întotdeauna într-o stare de energie potențială minimă. Unitățile își vor schimba valoarea de activare până la un moment în care schimbarea activării oricărei unități nu va mai duce la scăderea energiei potențiale a întregului sistem.

Se poate demonstra că un astfel de sistem ajunge întotdeauna într-o stare de minim energetic — stabilă. Considerând energia potențială dată de relația:

$$H = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{ij} s_{ij} \quad (2.20)$$

și  $s'_i = \text{sign}(h'_i)$ , ea fiind starea următoare a lui  $s_i$ , vor putea exista două situații:

1.  $-s'_i = s_i \Rightarrow$  unitatea  $i$  nu își schimbă activarea  $\Rightarrow \Delta H = 0$
2.  $-s'_i = -s_i \Rightarrow$  unitatea își schimbă activarea

În primul caz,  $H' = H$ , deci energia rămâne aceeași. În cel de-al doilea caz,

$$\Delta H = H' - H = -\frac{1}{2} \sum_{i,j} w_{ij} s'_i s_j + \frac{1}{2} \sum_{i,j} w_{ij} s_i s_j \quad (2.21)$$

Înlocuind  $s'_i$  cu  $-s_i$ , se obține:

$$\Delta H = \sum_i s_i \sum_j w_{ij} s_j = \sum_i s_i h'_i \quad (2.22)$$

Având în vedere ca semnul lui  $h'_i$  este și semnul lui  $s'_i$  și că  $s_i = s'_i$ , rezultă că  $\Delta H$  este întotdeauna negativ atunci când o unitate își schimbă activarea, deci sistemul va evolua întotdeauna către o stare de minim energetic.

Pentru un singur tipar de învățare  $\xi$  constând dintr-un vector binar cu componente  $(-1, 1)$  de lungime  $N \times N$ , relația de calcul a ponderilor va fi:

$$w_{ij} = \frac{1}{N} \xi_i \xi_j \quad (i = \overline{1, N}; j = \overline{1, N}) \quad (2.23)$$

deci matricea de ponderi va fi o matrice simetrică.

Pentru mai multe tipare de învățare,  $\xi^\mu$  ( $\mu = \overline{1, p}$ ) relația de calcul a ponderilor va fi:

$$w_{ij} = \frac{1}{N} \sum_{\mu=1}^p \xi_i^\mu \xi_j^\mu \quad (i = \overline{1, N}; j = \overline{1, N}) \quad (2.24)$$

Având în vedere capacitatea limitată de memorare și posibilitatea apariției stărilor mixate ca o combinație liniară a stărilor învățate, acest tip de rețele nu au aplicații practice.

Publicațiile lui Hopfield au generat numeroase controverse. În 1983 Cohen și Grossberg au stabilit un principiu general pentru proiectarea unei memorii cu conținut direct adresabil care a inclus și versiunea rețelei Hopfield ca și caz particular [14]. O componentă distinctivă a convergenței unei astfel de rețele neuronale este modul natural în care timpul ca și dimensiune esențială a învățării se manifestă în dinamica neliniară a rețelei.

În 1983, Kirkpatrick, Gellat și Vecchi au descris o nouă procedură numită algoritmul de *simulare a răcirii treptate* ("simulated annealing algorithm") [43], ideea fiind preluată de la Metropolis, Rosenbluth și Teller [54], și își are rădăcina în statistica termodinamică pentru rezolvarea unor probleme de optimizare combinatorială. Această idee a fost exploatată în 1985 de Ackley, Hinton și Sejnowski [2] în dezvoltarea

algoritmului stohastic de învățare ce utilizează proprietăți ale distribuției Boltzmann, denumit *învățarea Boltzmann*.

Metoda presupune pentru  $S_i$  o dependență de tipul:

$$S_i = \begin{cases} +1 & \text{cu probabilitatea } g(h_i) \\ -1 & \text{cu probabilitatea } 1 - g(h_i) \end{cases} \quad (2.25)$$

unde

$$h_i = \sum_j w_{ij} S_j \quad (2.26)$$

Funcția de probabilitate  $g(h)$  trebuie să îndeplinească condițiile:

1. să ia valori în intervalul  $[0, 1]$
2. să fie dependentă de temperatură.

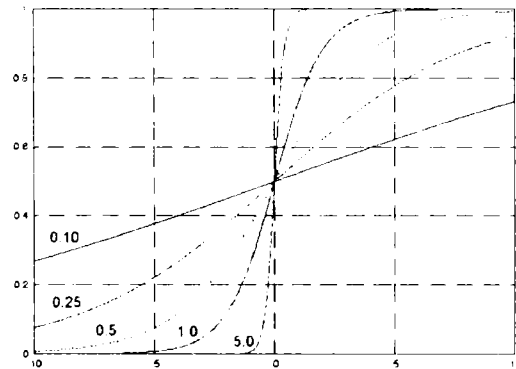
Cea mai utilizată funcție de probabilitate este cea sigmoidă:

$$F(x) = \frac{1}{1 + e^{-x}} \quad F(x) : \mathbb{R} \rightarrow [0, 1] \quad (2.27)$$

Pentru a simula dependența de temperatură a funcției, se utilizează factorul de temperatură  $\beta$  astfel:

$$F(x) = \frac{1}{1 + e^{-\beta x}} \quad \beta = \frac{1}{K_B \cdot T_K} \quad (2.28)$$

unde  $K_B$  este constanta lui Boltzmann, iar  $T_K$  este temperatura absolută în K. Pentru diferite valori ale lui  $\beta$  se obține o variație a funcției de genul celui ilustrat în figura 2.13. Inițial se pornește de la o temperatură ridicată, fapt ce induce o valoare mică a valorii lui  $\beta$ , având ca și consecință o pantă mică a funcției de probabilitate, urmând ca la creșterea valorii temperaturii panta funcției să crească în mod considerabil.



**Figura 2.13:** Variația funcției sigmoidă pentru diferite valori ale lui  $\beta$

Lăsând sistemul să evolueze "la o temperatură ridicată", el va ajunge într-o stare stabilă, care nu este una din stările învățate, dar este mai joasă din punct de vedere energetic. În acest moment, la scăderea temperaturii în mod progresiv, sistemul va ajunge la una din stările învățate sau la negativele acestora. Capacitatea unei astfel de rețele Hopfield stohastice este similară cu cea Hopfield discretă – prezentată anterior, evoluția fiind doar în ce privește lansarea în domeniu a *funcției de activare sigmoide* care a stat la baza trecerii la *rețele cu propagare înapoi*.

Tot în 1984, Braitenberg publică o carte "Experimente în psihologia sintetică" [8], în care descrie diverse mecanisme cu capacitate de autoorganizare, inspirate din comportamentul animal.

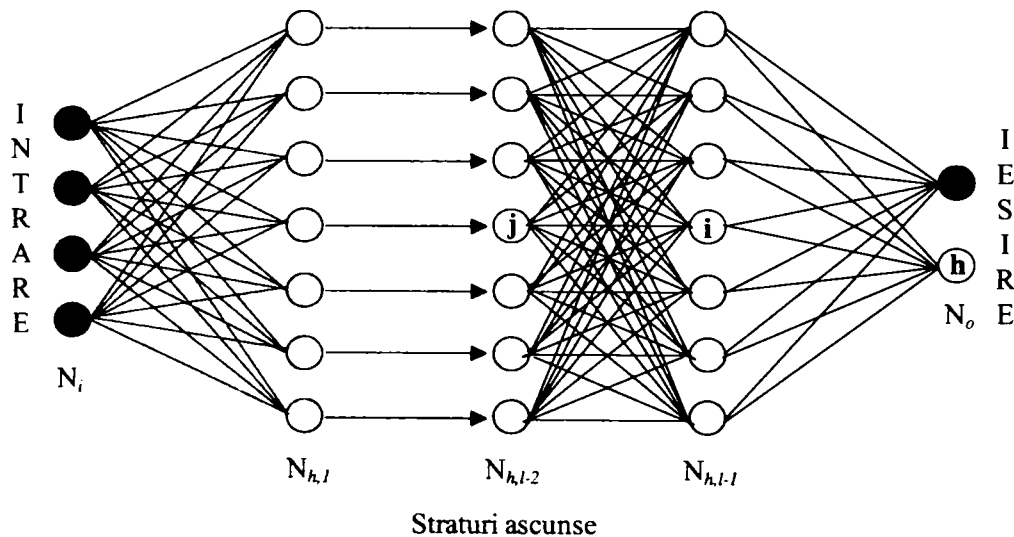
## 2.5 Rețele cu propagare înapoi (backpropagation)

Deși fusese descris de către Bryson și Ho în 1969 [9], și de către Werbos în 1974 [96], abia în 1986 Rumelhart, Hinton și Williams publică studiul algoritmului de propagare înapoi ("the back-propagation algorithm") [77], spre care se îndreaptă atenția cercetătorilor în domeniu.

În același an, Rumelhart și McClelland publică în două volume cartea "Procesare paralelă și distribuită: Explorări în microstructurile cunoașterii" ce se va constitui în perioada următoare ca referință de bază în studiul antrenării perceptronului multistrat [76]. Aceași idee, dar cu o rezonanță mai mică, o utilizează și Parker în 1985–1987 [63], [64] și Le Cun în 1985 [15].

Ideea centrală a acestui model este că erorile pentru unitățile de pe stratul intermediar sunt determinate prin *propagare înapoi*<sup>2</sup> a erorilor de pe stratul de ieșire. Propagarea înapoi este o generalizare a *regulii delta* pentru funcții de activare neliniare aplicate rețelelor multistrat.

O rețea neuronală cu propagare înapoi (figura 2.14) are o structură multistrat, fiecare strat primind intrările de pe stratul anterior și trimițând ieșirile stratului imediat următor. În cadrul unui strat nu există conexiuni între unitățile proprii. Cele  $N_i$  intrări alimentează primul strat, cu  $N_{h,l}$  unități ascunse. Unitățile de intrare sunt doar unități generatoare, în ele neavând loc nici un proces de prelucrare a datelor. Activarea unităților ascunse se face prin aplicarea unei funcții  $F$  asupra ponderilor care le preced și adăugarea unei constante denumită *bias* care nu este obligatorie. Ieșirea unităților ascunse este distribuită stratului imediat următor de unități ascunse  $N_{h,2}$  care, dacă este ultimul strat de unități ascunse, alimentează cele  $N_o$  unități de pe stratul de ieșire.



**Figura 2.14:** Rețea neuronală cu  $l$  straturi ascunse

Algoritmul de propagare înapoi poate fi aplicat oricărui tip de rețele, cu oricât de multe straturi. Pentru cazul rețelelor cu unități binare, s-a demonstrat, de către Hornik, Stinchcombe și White în 1989 [36], [37] și de către Hartman în 1990 [30], că

<sup>2</sup>by back-propagation"



un singur strat de unități ascunse este suficient pentru a aproxima orice funcție cu un număr mare, dar finit de discontinuități, pentru funcții de activare neliniare aplicate pentru straturile ascunse.

Pentru a utiliza funcții neliniare de activare, este necesară *generalizarea regulii delta* care este specifică funcțiilor liniare.

### 2.5.1 Generalizarea regulii delta

Funcția de activare este o funcție diferentiabilă a intrării dată de relația:

$$a_i^p = F(i_i^p) \quad (2.29)$$

în care  $a_i^p$  este ieșirea unității  $i$  obținută pentru intrarea  $i_i^p$  în cazul tiparului  $p$ . Valoarea de intrare este obținută prin relația:

$$i_i^p = \sum_j w_{i,j} a_j^p + \theta_i \quad (2.30)$$

în care  $\theta_i$  este o valoare tip *bias*.

Pentru obținerea unei generalizări corecte a *regulii delta*, se admite:

$$\Delta^p w_{i,j} = -\eta \frac{\partial E^p}{\partial w_{i,j}} \quad (2.31)$$

Mărima  $E^p$  este definită ca eroarea pătratică totală pentru tiparul  $p$  la ieșire:

$$E^p = \frac{1}{2} \sum_p^{N_{i,p}} (d_i^p - a_i^p)^2, \quad (2.32)$$

expresie ce denotă termenul care trebuie minimizat pe parcursul evoluției rețelei. Coeficientul  $\frac{1}{2}$  din fața expresiei este un coeficient de conveniență pentru simplificarea expresiei gradientului obținut prin derivare, acesta neafectând localizarea minimului de eroare și nici procesul de minimizare în sine.

Termenul erorii totale a rețelei la un anumit moment, este dat de suma erorilor tuturor tiparelor:

$$E^T = \sum_p^{N_{i,p}} E^p \quad (2.33)$$

Se poate scrie:

$$\frac{\partial E^p}{\partial w_{i,j}} = \frac{\partial E^p}{\partial i_i^p} \frac{\partial i_i^p}{\partial w_{i,j}} \quad (2.34)$$

Din relația 2.34 se observă că cel de al doilea factor este:

$$\frac{\partial i_i^p}{\partial w_{i,j}} = a_j^p \quad (2.35)$$

Definind:

$$\delta_i^p = \frac{\partial E^p}{\partial i_i^p} \quad (2.36)$$

se obține o regulă nouă, care este echivalentă cu regula delta, constând dintr-un gradient descendent al suprafeței de eroare, dacă schimbările ponderilor se vor face în concordanță cu relația:

$$\Delta^p w_{i,j} = \eta \delta_i^p a_j^p \quad (2.37)$$

Ceea ce rămâne de demonstrat este ce reprezintă  $\delta_i^p$  pentru fiecare unitate  $i$  din rețea. Pentru calculul lui  $\delta_i^p$  se va aplica regula derivării înlănțuite pentru a scrie derivata parțială ca un produs de doi factori, unul reflectând variația erorii ca funcție de ieșire, iar celălalt reflectând schimbarea ieșirii ca funcție de schimbarea intrării. Astfel, rezultă:

$$\delta_i^p = -\frac{\partial E^p}{\partial i_i^p} = -\frac{\partial E^p}{\partial a_i^p} \frac{\partial a_i^p}{\partial i_i^p} \quad (2.38)$$

Din relația 2.38 se observă că:

$$\frac{\partial a_i^p}{\partial i_i^p} = F'(i_i^p), \quad (2.39)$$

expresie care este derivata funcției pentru unitatea  $i$ , evaluată la intrarea  $i_i^p$  a acestei unități. Pentru a calcula primul factor din relația 2.38, vor fi considerate două cazuri:

1. **Unitatea  $i$  este o unitate de ieșire a rețelei.** În acest caz, din definiția erorii tiparului  $E^p$ , se obține:

$$\delta_i^p = (d_i^p - a_i^p)F'(i_i^p) \quad (2.40)$$

pentru o unitate de ieșire  $i$ .

2. **Unitatea  $i$  nu este o unitate de ieșire.** În acest caz, nu se cunoaște contribuția unității la eroarea de ieșire a rețelei. Deci, valoarea erorii poate fi scrisă ca funcție de intrările din stratul ascuns spre stratul de ieșire:

$$E^p = E^p(i_1^p, i_2^p, \dots, i_i^p, \dots) \quad (2.41)$$

Utilizând aceeași regulă de derivare înlănțuită, rezultă:

$$\frac{\partial E^p}{\partial a_i^p} = \sum_{h=1}^{N_o} \frac{\partial E^p}{\partial i_h^p} \frac{\partial i_h^p}{\partial a_i^p} = \sum_{h=1}^{N_o} \frac{\partial E^p}{\partial i_h^p} \frac{\partial}{\partial a_i^p} \sum_{k=1}^{N_h} w_{hk} a_k^p = \sum_{h=1}^{N_o} \frac{\partial E^p}{\partial i_h^p} w_{h,i} = -\sum_{h=1}^{N_o} \delta_h^p w_{h,i} \quad (2.42)$$

Substituind rezultatul în ecuația 2.38, rezultă:

$$\delta_i^p = F'(i_i^p) \sum_{h=1}^{N_o} \delta_h^p w_{h,i} \quad (2.43)$$

Ecuațiile 2.40 și 2.43 redau o procedură de calcul a valorii lui  $\delta$  pentru toate unitățile din rețea, care la rândul lor sunt utilizate pentru a calcula schimbările în valoare ale ponderilor în concordanță cu ecuația 2.37. Aceasta procedură constituie *generalizarea regulii delta* pentru o rețea cu unități neliniare.

## 2.5.2 Modul de lucru cu *backpropagation*

Aplicarea *regulii delta generalizată* implică două faze:

1. În decursul primei faze, se aplică intrarea și se propagă de-a lungul rețelei pentru calculul valorilor  $a_i^p$  pentru fiecare termen al unității de ieșire; această valoare este comparată cu valoarea dorită, rezultând un semnal de eroare  $\delta_i^p$  pentru fiecare termen al unității de ieșire;
2. Cea de-a doua fază implică o trecere înapoi prin rețea, în decursul căreia semnalul de eroare este transmis fiecărui element din rețea, calculându-se schimbările de pondere adecvate. Ponderea conexiunii este ajustată cu o cantitate proporțională cu produsul semnalului de eroare  $\delta$  pe unitatea  $i$ , primind intrarea, iar ieșirea unității  $j$  trimițând acest semnal de-a lungul conexiunii:

$$\Delta^p w_{i,j} = \eta_{i,j} \delta_i^p a_j^p \quad (2.44)$$

unde

$w_{ij}$  ponderea conexiunii dintre unitatea  $i$  și unitatea  $j$

$\Delta^p w_{ij}$  termenul de ajustare a ponderii conexiunii

$\eta_{ij}$  rata de învățare asociată cu ponderea  $w_{ij}$

$\delta_i^p$  semnalul de eroare asociat elementului  $i$  pentru vectorul  $p$  aplicat la intrarea rețelei

$a_i^p$  valoarea de activare a elementului  $i$  al rețelei pentru vectorul  $p$  aplicat la intrarea rețelei

Dacă unitatea este situată la ieșirea din rețea, semnalul de eroare este dat de:

$$\delta_k^p = (d_k^p - a_k^p)F'(z_k^p) \quad (2.45)$$

unde

$z_k^p$  intrarea netă a neuronului  $k$  pentru un vector de intrare  $p$

$d_k^p$  al  $k$ -lea element al vectorului de ieșire dorit al rețelei pentru vectorul  $p$  aplicat la intrarea rețelei

Definită ca funcția de activare sigmoidă, funcția  $F$  este:

$$a_j^p = F(y_j^p) = \frac{1}{1 + e^{-y_j^p}} \quad (2.46)$$

În cazul utilizării funcției de activare temperată  $a_j^p$  devine:

$$a_j^p = F(y_j^p) = \frac{1}{1 + e^{-\beta y_j^p}} \quad (2.47)$$

unde:

$$\beta = \frac{1}{K_B \cdot T_K},$$

$K_B$  este constanta lui Boltzmann

$T_K$  este temperatura absolută în grade Kelvin

Derivata acestei funcții este:

$$F'(y_j^p) = a_j^p(1 - a_j^p) \quad (2.48)$$

$$F'(i_i^p) = \frac{\partial}{\partial i_i^p} \frac{1}{1 + e^{-i_i^p}} = \frac{1}{(1 + e^{-i_i^p})^2} (-e^{-i_i^p}) = \frac{1}{1 + e^{-i_i^p}} \frac{-e^{-i_i^p}}{1 + e^{-i_i^p}} = a_i^p(1 - a_i^p) \quad (2.49)$$

astfel încât semnalul de eroare pentru o unitate de ieșire poate fi scris ca:

$$\delta_j^p = (d_j^p - a_j^p)a_j^p(1 - a_j^p) \quad (2.50)$$

Semnalul de eroare pentru o unitate ascunsă este determinat recursiv în termenii semnalelor de eroare ale unității cu care este direct conectată și ponderile acestor conexiuni. Pentru funcția de activare sigmoidă:

$$\delta_j^p = F'(y_j^p) \sum_{k=1}^N \delta_k^p w_{j,k} = a_j^p(1 - a_j^p) \sum_{k=1}^N \delta_k^p w_{j,k} \quad (2.51)$$

### 2.5.3 Alte funcții de activare neliniare

Practic, numărul funcțiilor care pot fi utilizate ca funcții de activare este relativ nelimitat. Acestea se pot clasifica în două clase:

### 2.5.3.1 Funcții nesimetrice față de origine similare ca alură cu funcția sigmoidă

Din această categorie, alături de funcția sigmoidă, cea mai utilizată este funcția tangentă hiperbolică. Această funcție este definită de relația:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (2.52)$$

Funcția tangentă hiperbolică este reprezentată grafic ca în figura 2.15.

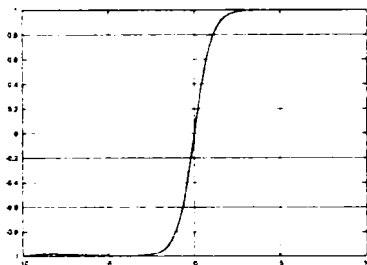


Figura 2.15: Reprezentarea grafică a funcției tangentă hiperbolică

### 2.5.3.2 Funcții simetrice față de origine

Din această categorie fac parte:

**Funcția Gaussiană** — aceasta este o funcție radială, definită de relația:

$$F_{Gauss}(x) = e^{-\frac{x^2}{\nu}} \quad (2.53)$$

în care  $\nu$  este o constantă predefinită, numită *varianță*, cu ajutorul căreia este ajustată amplitudinea funcției. Funcția Gaussiană este redată în figura 2.16.

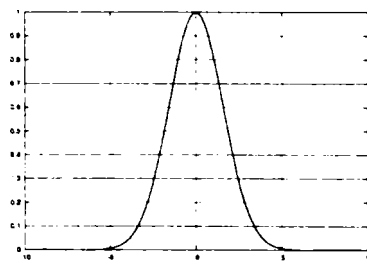


Figura 2.16: Reprezentarea grafică a funcției Gaussiene.

**Funcția secantă hiperbolică** — și această funcție este o funcție radială, fiind definită de relația:

$$\operatorname{sech}(x) = \frac{2}{e^x + e^{-x}} \quad (2.54)$$

Utilizarea unei anumite funcții de activare este o opțiune personală a programatorului, neexistând argumente tehnice privind indici de performanță care să justifice o selecție preferențială.

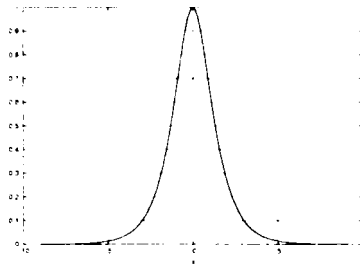


Figura 2.17: Reprezentarea grafică a funcției secantă hiperbolică.

### 2.5.4 Rata de învățare și momentul

Procedura de învățare impune ca actualizarea ponderii să fie proporțională cu  $\frac{\partial F^p}{\partial w}$ . Scăderea reală a gradientului necesită utilizarea unor pași infinitezimali. Constanta de proporționalitate este rata de învățare  $\eta$ . Pentru aplicațiile practice se alege o rată de învățare cât mai mare posibilă, fără a se ajunge însă la *oscilație*. O modalitate de evitare a oscilației la valori mari ale lui  $\eta$  este efectuarea unei actualizări în pondere dependentă de actualizarea anterioară a ponderii prin adăugarea unui termen *moment*,  $\gamma$ :

$$\Delta w_{i,j}(t+1) = \eta \delta_i^p a_j^p + \gamma \Delta w_{i,j}(t) \quad \gamma \in (0,1) \quad (2.55)$$

unde  $t$  indexează numărul epocii, iar  $\gamma$  este o constantă care determină efectul schimbării anterioare a ponderii. Valoarea lui  $\gamma$  poate fi de exemplu 0.1, scăzând la creșterea erorii sau crescând la scăderea ei.

Rolul termenului *moment* este arătat în figura 2.18. Când nu este utilizat termenul *moment*, găsirea minimumului cu o rată mică de învățare necesită o perioadă lungă de timp (1), iar pentru o rată de învățare mare, minimumul nu poate fi găsit din cauza fenomenului de oscilație (2). Adăugarea unui termen *moment* duce la o *găsire mai rapidă a valorii de minim* (3).

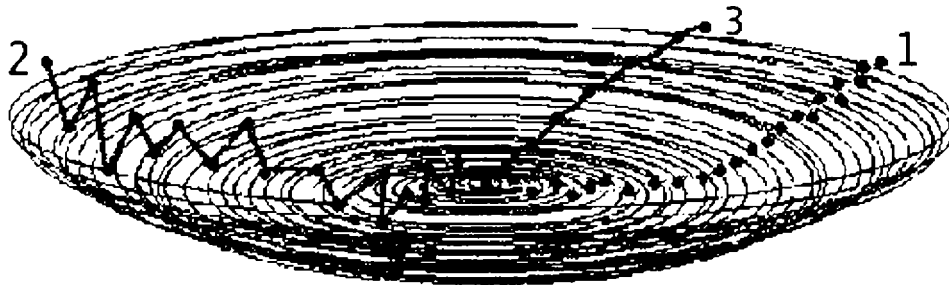


Figura 2.18: Variația erorii rețelei funcție de valoarea momentului.

### 2.5.5 Învățarea unui tipar

Teoretic, algoritmul *backpropagation* determină scăderea gradientului erorii totale a rețelei doar dacă ponderile sunt ajustate pe întregul set de tipare. Se va acorda o atenție mărită ordinii în care setul de tipare este prezentat rețelei deoarece, în cazul în care este utilizată aceeași secvență, rețeaua se poate focaliza pe ultimele tipare

prezentate. Problema se poate depăși dacă se utilizează un algoritm aleator sau o tehnică de permutare la prezentarea tiparelor în etapa de antrenare.

### 2.5.6 Deficiențe ale algoritmului backpropagation

Contrar succesului obținut de algoritmul backpropagation, există câteva aspecte care fac ca utilizarea algoritmului să nu reprezinte garanția succesului universal.

Se vor prezenta mai jos trei dintre problemele cel mai des întâlnite, oferindu-se și câteva din motivele posibile ale manifestării lor.

**O antrenare de lungă durată, cu un gradient al erorii totale scăzut** — se poate datora unei selecții neoptimale a ratei de învățare și a termenului moment. Există algoritmi care pornesc cu o rată a învățării mică, în funcție de variația erorii aceasta modificându-se corespunzător — la o variație decrementală a erorii rata de învățare crescând, în vreme ce ea scade la o variație incrementală.

**Paralizia rețelei** — se poate datora valorilor foarte mari, în valoare absolută, pe care le pot lua ponderile pe parcursul antrenării rețelei, valori care induc valori de activare apropiate de 0 sau 1. Conform relațiilor 2.50 și 2.51, termenul de ajustare a valorii ponderilor, care este direct proporțional cu  $a_j^p(1 - a_j^p)$ , va avea valori nesemnificative, fapt ce duce la un palier staționar de evoluție pentru gradientul de eroare.

**Convergența spre un minim local** — având în vedere că hiper-suprafața de eroare pentru o rețea complexă este accidentată, există posibilitatea ca rețeaua să se afle într-un minim local de valoare nesatisfăcătoare în comparație cu minimul absolut al rețelei. Există metode probabilistice de tipul *algoritmilor genetici*, prin a căror utilizare se pot evita capcanele locale, dar acestea au dezavantajul de a fi lente. O altă posibilitate este de a mări numărul de unități de pe stratul ascuns, metodă care are același dezavantaj — este lentă.

Această problemă este cea mai spinoasă, existând în literatură o serie de articole care pretind rezolvarea parțială a ei, însă *rețelele* oferite nu au o eficiență generală, ci doar evită problema în contextul particular în care este studiată.

Spre exemplificare, în 1994, Klawun și Wilkins publică un articol denumit "Un nou algoritm pentru ieșirea dintr-un minim local" [44], care constă în prezentarea, în cazul sesizării unui minim local (sesizarea se face prin constatarea unei erori maxime pentru unele tipare) pentru câteva epoci, doar a tiparelor cu eroare maximă (metodă cunoscută și sub denumirea de *flashcard-algorithm*). Pentru moment, eroarea totală a rețelei crește dar își urmează cursul descendent, depășind impasul și convergând spre valoarea de minim global.

Această metodă este doar o observație de bun simț foarte des folosită de către utilizatorii de programe personale dedicate, programe care permit monitorizarea în timpul antrenării rețelei a parametrilor de rețea și intervenția în momentele de antrenare ineficace. Este de remarcat însă efortul logistic și financiar alocat pentru elaborarea metodei.

Același rezultat se obține și prin bruiajul rețelei, prin impunerea unei valori anormale unei conexiuni din rețea, bruiajul urmând a se propaga prin întreaga rețea, fapt ce conduce pentru moment la o valoare mare a erorii totale, concomitent cu scoaterea rețelei din minimul local.

## 2.6 Îmbunătățiri ale algoritmului backpropagation

Cercetătorii în domeniu au lansat variațiuni pe tema algoritmului backpropagation care s-au dorit a fi algoritmi îmbunătățiți în ce privește etapa de antrenare a rețelei.

Evaluarea contribuției fiecăruia este prematură, valoarea noilor algoritmi trebuind să treacă testul de durabilitate al aplicării și verificării în condiții în care metodele actuale nu dau rezultate satisfăcătoare. În această secțiune se vor discuta câțiva dintre acești algoritmi.

O îmbunătățire semnificativă în domeniu este înlocuirea metodei de minimizare simplistă pe baza gradientului la fiecare pas, cu o metodă de minimizare direcționată, cuplată cu minimizarea gradientului conjugat. Minimizarea unei funcții  $f$  de-a lungul unei direcții  $u$  aduce funcția  $f$  într-un loc în care gradientul ei este perpendicular pe direcția  $u$ . În caz contrar, procesul de minimizare de-a lungul direcției  $u$  este incomplet. În locul urmăririi gradientului pas cu pas, se construiește un set de  $n$  direcții conjugate între ele, astfel încât minimizarea de-a lungul uneia dintre ele  $u_j$  nu se repercutează în mod distructiv asupra minimizării de-a lungul unei direcții anterioare. Practic minimizările de-a lungul direcțiilor nu interferă. Sistemul de gradul  $n$  va evolua deci, spre un minim global. Metoda este diferită de metoda gradientului descendent, ce presupune minimizarea în direcția cu gradient descendent maxim.

Acceptând aproximarea prin serie Taylor a funcției ce urmează a fi minimizată, se obține:

$$f(x) = f(p) + \sum_i \frac{\partial f}{\partial x_i} \Big|_p x_i + \frac{1}{2} \sum_{i,j} \frac{\partial^2 f}{\partial x_i \partial x_j} \Big|_p x_i x_j + \dots \cong \frac{1}{2} x^T A x - b^T x + c \quad (2.56)$$

unde

- $T$  denotă operațiunea de transpunere
- $c = f(p)$
- $b = -\nabla f|_p$
- $[A]_{i,j} \equiv \frac{\partial^2 f}{\partial x_i \partial x_j} \Big|_p$  este o matrice simetrică de dimensiuni  $n \times n$ , *pozitiv definită* ( $\forall y \neq 0, y^T A y > 0$ ), de fapt, *Hessiana* funcției  $f$  în punctul  $p$ :

$$A_p = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(p) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(p) \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(p) & \frac{\partial^2 f}{\partial x_2^2}(p) \end{bmatrix} \quad (2.57)$$

Gradientul funcției  $f$  este:

$$\nabla f = A x - b \quad (2.58)$$

astfel încât o schimbare în  $x$  induce o schimbare de gradient conform relației:

$$\delta(\nabla f) = A(\delta x) \quad (2.59)$$

Presupunând că funcția  $f$  a fost minimizată de-a lungul direcției  $u_i$  într-un punct în care gradientul  $-g_{i+1}$  al funcției este perpendicular pe direcția  $u_i$ ,

$$u_i^T g_{i+1} = 0 \quad (2.60)$$

se trece la minimizarea pe o nouă direcție  $u_{i+1}$ . Pentru a ne asigura că parcurgerea direcției  $u_{i+1}$  nu afectează minimizarea de-a lungul direcției  $u_i$  este necesar ca gradientul funcției  $f$  să rămână perpendicular pe direcția  $u_i$ , astfel încât:

$$u_i^T g_{i+2} = 0 \quad (2.61)$$

Altfel, va mai trebui parcursă o dată etapa de minimizare în direcția  $u_i$ . Combinând relațiile 2.60 și 2.61 se obține:

$$0 = u_i^T (g_{i+1} - g_{i+2}) = u_i^T \delta(\nabla f) = u_i^T A u_{i+1} \quad (2.62)$$

Când relația 2.59 se aplică pentru doi vectori  $u_i$  și  $u_{i+1}$ , se spune despre aceștia că sunt conjugăți.

Plecând de la același punct  $p_0$ , prima direcție de minimizare  $u_0$  se va lua egală cu  $g_0 = -\nabla f(p_0)$ . Pentru  $i \geq 0$ , direcțiile se calculează cu relația:

$$u_{i+1} = u_i + \lambda_i^u u_i \quad (2.63)$$

unde valoarea coeficientului incremental  $\lambda_i^u$  este aleasă astfel încât să avem  $u_i^T A u_{i-1} = 0$  și gradientii succesivi să fie perpendiculari:

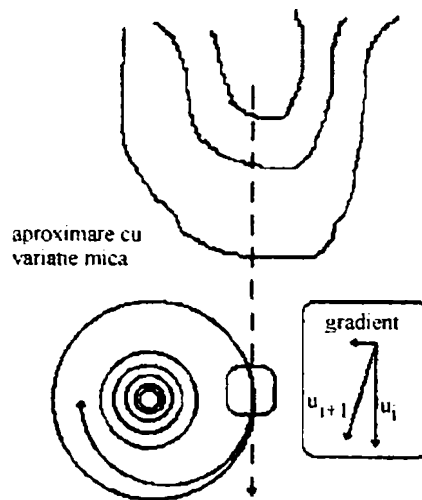
$$\gamma_i = \frac{g_{i+1}^T g_{i+1}}{g_i^T g_i} \quad \text{cu } g(k) = -\nabla f|_{p_k} \quad \text{pentru } k \geq 0 \quad (2.64)$$

Următorul pas este calculul:

$$p_{i+2} = p_{i+1} + \lambda_{i+1}^p \cdot u_{i+1} \quad (2.65)$$

unde coeficientul  $\lambda_{i+1}^p$  se alege astfel încât să minimizeze  $f(p_{i+2})^2$ .

Se poate demonstra că direcțiile  $u$  astfel construite sunt toate interconjugate [86], [69]. Din figura 2.19 se observă că la o nouă selecție de direcție, aceasta se construiește din prima direcție și gradient, rezultând o minimizare în spirală.



**Figura 2.19:** Coborârea cu pantă mică cu gradient conjugat pentru sisteme nepatratice.

Astfel, pentru un sistem pătratic cu  $n$  grade de libertate sunt necesare doar  $n$  iterații. Având în vedere că nu se urmărește minimizarea unui sistem pătratic, este necesară parcurgerea celor  $n$  direcții de mai multe ori.

Powell [69] a introdus câteva îmbunătățiri semnificative pentru a corecta comportarea metodei în sistemele nepătratice.

Van den Boomgaard și Smeulders au arătat în 1989 [94] că pentru o rețea neuronală cu propagare înainte fără unități ascunse, procedura incrementală de determinare a matricii de ponderi optime necesită o ajustare a ponderilor cu:

$$\Delta W(t+1) = \eta(t+1)(d(t+1) - W(t)x(t+1))x(t+1) \quad (2.66)$$

în care  $\eta$  nu este o constantă, ci o variabilă tip matrice  $(N_i + 1) \times (N_i + 1)$  care depinde de vectorul de intrare. Utilizând cunoștințe *a priori* despre semnalul de intrare, necesitățile de stocare pentru valorile termenului moment  $\gamma$  pot fi reduse substanțial.



Silva și Almeida au arătat [82] de asemenea avantajele utilizării unui pas de valoare specifică pentru fiecare pondere din rețea. În algoritmul lor, rata de învățare este adaptată pentru fiecare tipar:

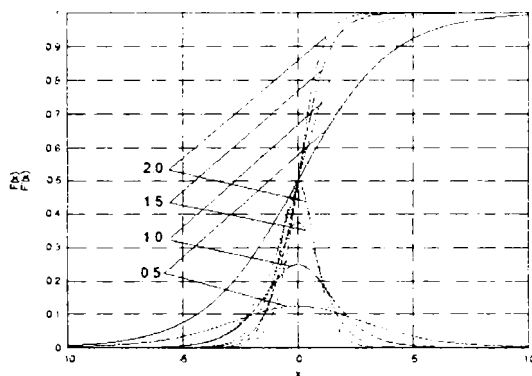
$$\eta_{i,j}(t+1) = \begin{cases} u\eta_{i,j}(t) & \text{dacă } \frac{\partial E_{t+1}}{\partial w_{i,j}} \text{ și } \frac{\partial E_t}{\partial w_{i,j}} \text{ au același semn} \\ d\eta_{i,j}(t) & \text{dacă } \frac{\partial E_{t+1}}{\partial w_{i,j}} \text{ și } \frac{\partial E_t}{\partial w_{i,j}} \text{ au semne diferite} \end{cases} \quad (2.67)$$

unde  $u$  și  $d$  sunt constante pozitive cu valori puțin peste — și respectiv puțin sub unitate. Ideea de bază a acestei metode este de a descrește rata de învățare în cazul oscilației valorilor ponderilor.

### 2.6.1 Parametrii de antrenare a rețelei

În etapa de antrenare a rețelei, conform algoritmului de învățare, există câțiva parametri care definesc această etapă. Enumerarea doar a unora nu conduce la o descriere satisfăcătoare nici a rețelei și nici a capacității de învățare a acesteia.

În figura 2.20 este prezentată variația funcției sigmoide unipolare și a derivatei sale pentru valori ale factorului de temperatură  $\beta$  cuprinse între 0.2 și 2 utilizând un pas de 0.1 unități.



**Figura 2.20:** Variația funcției sigmoide unipolară și a derivatei sale pentru valori ale factorului de temperatură  $\beta$  cuprinse între 0.5–2 cu un pas de 0.5.

Se observă că pentru două valori învecinate ale variabilei, creșterea ambelor funcții este mai substanțială pentru valori ale factorului de temperatură  $\beta$  mai mari. De aici rezultă două consecințe:

1. la normalizarea valorilor utilizând funcția sigmoide, pentru două valori apropiate ale variabilei, valorile funcției sunt progresiv diferite odată cu creșterea factorului de temperatură;
2. având în vedere relația 2.55 de corecție a ponderilor, pentru două valori apropiate ale erorii, corecția este mai drastică la valori ale factorului de temperatură crescute.

Informația despre valoarea absolută a acestui factor nu este relevantă decât dacă este completată cu valorile ratei de învățare  $\eta$  și ale termenului moment  $\gamma$ .

Aceste variabile de rețea intervin doar în corecția ponderilor în vederea reducerii erorii rețelei de la o epocă la alta.

Rata de învățare este de fapt o constantă de proporționalitate cu valori între 0 și 1 cu care se ajustează corecția ponderii. O valoare a ratei de învățare de 0.1 înseamnă

practic că se ține seama doar de 10% din valoarea derivatei funcției sigmoidă din termenul de corecție al ponderii.

Odată caracterizat termenul de corecție a ponderilor prin valoarea factorului de temperatură  $\beta$  și prin valoarea ratei de învățare, în cazul utilizării termenului *moment*, se intervine din nou printr-un proces de ajustare a acestuia conform relației 2.55 în funcție de viteza de convergență a rețelei spre valoarea de minim de eroare.

Importanța termenului *moment* este discutabilă, deoarece acesta face o ajustare a algoritmului de corecție, deci o scădere a vitezei de convergență spre minim chiar și atunci când rețeaua se află pe drumul descendent către minimul global. Ori-cum, valoarea sugerată în literatură [48]  $\gamma = 0.9$  este mult prea mare.

Această valoare impune de fapt o corecție actuală nesemnificativă, accentul fiind pus pe corecția exercitată într-o epocă anterioară. Aplicarea termenului *moment* conduce de fapt la o mediere ponderată și consecutivă a factorului de corecție a ponderilor.

În concluzie, o comparare a funcționării unor rețele în etapa de antrenare este incompletă dacă nu sunt prezentați *toți* factorii discutați anterior.

### 2.6.2 Accelerarea convergenței — rată de învățare adaptivă

În această secțiune se vor descrie și analiza câteva metode euristice descrise de Jacobs ca modalități de accelerare a convergenței învățării prin adaptarea ratei de învățare [39].

**Metoda euristică 1.** — Fiecare parametru ajustabil al rețelei care descrie și funcția de cost a învățării trebuie să fie evaluat în particular la o valoare optimă proprie valorii de moment a funcției de cost.

Această metodă propune o valoare diferită pentru fiecare pondere, argumentând prin faptul că la un anumit moment, o anumită ajustare poate fi benefică pentru o anumită pondere și total necorespunzătoare pentru o altă pondere.

**Metoda euristică 2.** — Valoarea ratei de învățare trebuie să fie diferită de la o epocă la alta.

Această metodă are în vedere faptul că suprafața de eroare diferă de-a lungul regiunii corespunzătoare unei anumite ponderi.

**Metoda euristică 3.** — Când derivata funcției de cost în raport cu o anumită pondere are același semn pentru câteva iterații consecutive ale algoritmului, rata de învățare pentru acea pondere trebuie crescută.

Metoda are în vedere posibilitatea că punctul de operare curent, în spațiul ponderii vizate, poate ajunge pe o porțiune plană a suprafeței de eroare, fapt rezultat din consecvența semnului pe parcursul mai multor epoci. Printr-o creștere a valorii ratei de învățare, numărul de iterații pentru traversarea acestei suprafețe plane va fi mai mic, fapt ce duce la o îmbunătățire a valorii funcției de cost.

**Metoda euristică 4.** — Când semnul algebric al derivatei funcției de cost în raport cu o anumită pondere alternează consecutiv pe parcursul mai multor iterații ale algoritmului, rata de învățare pentru acea pondere trebuie micșorată.

Metoda are în vedere posibilitatea ca punctul de operare curent în spațiul ponderii vizate să ajungă pe o porțiune foarte curbată a suprafeței de eroare, fapt rezultat din inconsecvența semnului pe parcursul mai multor epoci. Astfel, valoarea ponderii oscilează ca semn, rămânând relativ constantă în valoare absolută ajungându-se la o creștere a funcției de cost fără vreun real progres.

Este cert că utilizarea parametrilor de rețea în raport cu fiecare dintre ponderi conform acestor metode euristice duce inevitabil la o modificare fundamentală a algoritmului *backpropagation*, care nu mai corespunde unei căutări în direcția descendentă spre minim.

Deci ajustările aplicate ponderilor se bazează pe:

- derivatele parțiale ale suprafeței de eroare în raport cu ponderile;
- estimarea cuadraturii suprafeței de eroare pe care se situează punctul de operare în spațiul ponderilor.

Astfel, toate cele patru metode euristice satisfac doar necesități locale ale rețelei, fiind benefice pentru anumite puncte și total necorespunzătoare pentru altele.

## 2.7 O implementare secvențială a algoritmului backpropagation

Se prezintă în continuare un exemplu de implementare a unui simulator secvențial pentru algoritmul backpropagation, realizată în limbajul C.

Pentru implementarea structurii rețelei se utilizează două tipuri de structuri de date, una pentru implementarea unui strat al rețelei – LAYER, și una pentru reprezentarea rețelei:

```
typedef struct {          /* Un strat al rețelei:          */
    int    Units;        /* - numărul de unități din strat */
    double* Output;     /* - ieșirile unităților         */
    double* Error;      /* - eroarea în fiecare unitate  */
    double** Weight;    /* - ponderea conexiunilor care */
                    /*   intră în strat             */
    double** WeightSave; /* - ponderi salvate pentru oprirea */
                    /*   învățării                  */
    double** dWeight;   /* - ultimele actualizări de ponderi */
                    /*   pentru termenul moment     */
} LAYER;

typedef struct {          /* O REȚEA:                      */
    LAYER** Layer;      /* - straturile rețelei         */
    LAYER*  InputLayer; /* - stratul de intrare         */
    LAYER*  OutputLayer; /* - stratul de ieșire         */
    double  Alpha;      /* - factorul moment           */
    double  Eta;        /* - rata de învățare          */
    double  Gain;       /* - câștigul funcției sigmoidă */
    double  Error;      /* - eroarea totală            */
} NET;
```

În partea de inițializare trebuie să se stabilească configurația rețelei: număr de straturi, număr de unități pe fiecare strat, valori pentru factorul moment, rata de învățare, funcția sigmoidă – uzual 0.9, 0.25, respectiv 1.

După citirea setului de tipare de antrenare, ponderile conexiunilor între straturi trebuie inițializate la valori mici aleatoare, în intervalul (-0.5, 0.5).

Dacă  $NI$ ,  $NH$ , și  $NO$  reprezintă numărul de unități de pe stratul de intrare, stratul ascuns și respectiv cel de ieșire, următoarele două ecuații dau valoarea activării neuronilor. Pentru un tipar de intrare  $p$ , activarea unității ascunse  $j$ ,  $a_{pj}^H$  este dată de:

$$a_{pj}^H = f(\text{net}_{pj}^H) = f\left(\sum_{i=0}^{NI-1} a_{pi}^I w_{i \rightarrow j}^H\right) \quad (2.68)$$

unde  $w_{i \rightarrow j}^H$  este ponderea conexiunii de la neuronul  $i$  de pe stratul de intrare la neuronul  $j$  de pe stratul ascuns, iar  $a_{pi}^I$  este activarea neuronului  $i$  de pe stratul de intrare;

similar, pentru o unitate  $k$  de pe stratul de ieșire, activarea este dată de:

$$a_{pk}^O = f(\text{net}_{pk}^O) = f\left(\sum_{j=0}^{NH-1} a_{pj}^H w_{j \rightarrow k}^O\right) \quad (2.69)$$

unde  $w_{j \rightarrow k}^O$  este ponderea conexiunii de la neuronul  $j$  de pe stratul ascuns la neuronul  $k$  de pe stratul de ieșire, iar  $a_{pj}^H$  este activarea neuronului  $j$  de pe stratul ascuns.

Funcțiile `PropagateLayer`, respectiv `PropagateNet` implementează pasul înainte conform relațiilor de mai sus:

```
void PropagateLayer(NET* Net, LAYER* Lower, LAYER* Upper)
{
    int i,j;
    double Sum;

    for (i=1; i<=Upper->Units; i++) {
        Sum = 0;
        for (j=0; j<=Lower->Units; j++) {
            Sum += Upper->Weight[i][j] * Lower->Output[j];
        }
        Upper->Output[i] = 1 / (1 + exp(-Net->Gain * Sum));
    }
}

void PropagateNet(NET* Net)
{
    int l;

    for (l=0; l<NUM_LAYERS-1; l++) {
        PropagateLayer(Net, Net->Layer[l], Net->Layer[l+1]);
    }
}
```

Apoi se calculează eroarea rețelei, prin comparare cu valoarea corespunzătoare din tiparul curent:

```
void ComputeOutputError(NET* Net, double* Target)
{
    int i;
    double Out, Err;

    Net->Error = 0;
    for (i=1; i<=Net->OutputLayer->Units; i++) {
        Out = Net->OutputLayer->Output[i];
        Err = Target[i-1]-Out;
        Net->OutputLayer->Error[i] = Net->Gain*Out*(1-Out)*Err;
        Net->Error += 0.5 * sqr(Err);
    }
}
```

După aceea, eroarea este propagată înapoi, conform ecuațiilor 2.50 și 2.51:

```
void BackpropagateLayer(NET* Net, LAYER* Upper, LAYER* Lower)
{
    int i,j;
    double Out, Err;

    for (i=1; i<=Lower->Units; i++) {
```

```

    Out = Lower->Output[i];
    Err = 0;
    for (j=1; j<=Upper->Units; j++) {
        Err += Upper->Weight[j][i] * Upper->Error[j];
    }
    Lower->Error[i] = Net->Gain * Out * (1-Out) * Err;
}
}

void BackpropagateNet(NET* Net)
{
    int l;

    for (l=NUM_LAYERS-1; l>1; l--) {
        BackpropagateLayer(Net, Net->Layer[l], Net->Layer[l-1]);
    }
}

```

În final, se actualizează ponderile conform ecuației 2.55:

```

void UpdateWeights(NET* Net)
{
    int l,i,j;
    double Out, Err, dWeight;

    for (l=1; l<NUM_LAYERS; l++) {
        for (i=1; i<=Net->Layer[l]->Units; i++) {
            for (j=0; j<=Net->Layer[l-1]->Units; j++) {
                Out = Net->Layer[l-1]->Output[j];
                Err = Net->Layer[l]->Error[i];
                dWeight = Net->Layer[l]->dWeight[i][j];
                Net->Layer[l]->Weight[i][j] +=
                    Net->Eta * Err * Out + Net->Alpha * dWeight;
                Net->Layer[l]->dWeight[i][j] = Net->Eta * Err * Out;
            }
        }
    }
}

```

Antrenarea rețelei constă în prezentarea câte unui tipar din setul de antrenare, propagarea lui înainte prin rețea, propagarea înapoi a erorii, urmată de actualizarea ponderilor. O epocă va consta din prezentarea tuturor tiparelor din setul de antrenare, în ordine aleatoare. După terminarea unei epoci, se testează calitatea învățării, prezentând din nou toate tiparele, fără antrenare, doar pentru calculul erorii. Dacă eroarea totală este sub o valoare prestabilită, se poate considera că învățarea a luat sfârșit și se poate salva configurația rețelei.

### 2.7.1 Analiza simulatorului secvențial

Algoritmul backpropagation utilizează intens operații cu virgulă flotantă. În cele ce urmează se va încerca estimarea timpului necesar prelucrării unui tipar, pentru a se putea face o previzionare a timpului necesar pentru o antrenare.

Pentru experimente, au fost utilizate stații IBM RS6000, SunBlade, respectiv PC-uri HP-Compaq, pe care s-au făcut determinări ale timpilor necesari efectuării fiecărei operații.

Din ecuațiile 2.68 și 2.69, rezultă, pentru pasul înainte, că sunt necesare câte o însumare și câte o înmulțire pentru fiecare pondere, atât pentru stratul ascuns cât și pentru stratul de ieșire. În plus, s-a determinat că mai sunt necesare 7.16  $\mu$ secunde

în cazul stațiilor RS 6000, 3.35  $\mu$ secunde pe SunBlade, respectiv 3.86  $\mu$ secunde pe procesoarele Intel utilizate, pentru calculul valorii activării fiecărui neuron de pe stratul intermediar și cel de ieșire, în principal pentru calculul funcției exponențiale.

La pasul înapoi, calculul valorilor delta pentru unitățile de ieșire și cele ascunse, conform ecuațiilor 2.50 și 2.51, vor necesita 2 scăderi și 2 înmulțiri pentru fiecare unitate de ieșire, respectiv 1 scădere și 2 înmulțiri pentru fiecare unitate de pe stratul ascuns, plus câte 1 adunare și 1 înmulțire pentru fiecare conexiune către stratul de ieșire.

La calculul actualizărilor de ponderi, conform ecuației 2.55, algoritmul necesită câte 1 adunare și 3 înmulțiri pentru fiecare pondere, atât pentru stratul ascuns cât și pentru cel de ieșire, iar actualizarea ponderilor mai necesită câte o adunare pentru fiecare pondere, pentru ambele straturi. Tabelul 2.1 sintetizează considerațiile de mai sus, prezentând timpii necesari (în  $\mu$ secunde) pentru unități și ponderi de conexiuni, pentru cele două platforme utilizate (în cazul neuronilor a fost luat în considerație și timpul necesar calculării activării).

Operații	+/-	*	timp RS6000	timp SunBlade	timp Intel
Per unitate ascunsă	1	2	7.502	3.51	4.02
Per unitate de ieșire	2	2	7.566	3.54	4.08
$w_{i \rightarrow j}^H$	3	4	0.748	0.35	0.48
$w_{j \rightarrow k}^O$	4	5	0.951	0.44	0.53

**Tabela 2.1:** Operații în virgulă flotantă necesare.

Pentru exemplificare, pentru o rețea 30-30-30, cu câte 30 de unități pe fiecare strat, sunt necesare 930 de ponderi între stratul de intrare și cel ascuns (900 de ponderi între unitățile de pe cele două straturi plus 30 de ponderi pentru bias către stratul ascuns). Similar, mai sunt necesare 930 de ponderi între stratul ascuns și cel de ieșire, rezultând un total de 1860 de ponderi în rețea.

Rezultă un total de 6600 de adunări și 8490 de înmulțiri. Luând în considerație și timpul necesar calculării funcției de activare, 7.16  $\mu$ secunde pe RS6000, respectiv 3.35  $\mu$ secunde pe SunBlade pentru fiecare dintre cei 60 de neuroni de pe straturile ascuns și de ieșire, rezultă în total un necesar de aproximativ 2.5 milisecunde pe RS6000, respectiv 1.150 milisecunde pe SunBlade, pentru un singur pas înainte și înapoi.

Pentru un număr de 1000 de tipare de intrare și circa 1000 de treceri ale fiecărui tipar (în funcție de rata de învățare, momentul și alte caracteristici ale problemei), vor fi necesare 2500 de secunde pe RS6000, respectiv 1150 de secunde pe SunBlade. În plus, orice implementare mai presupune operații suplimentare – calcule de indici, salvări de ponderi etc – care la rândul lor mai necesită timp suplimentar, care nu a fost luat în considerație în calculele de mai sus.

Există mai multe modalități de compensare a acestor dezavantaje ale rețelelor neuronale.

Prima abordare constă în reducerea dimensiunilor problemei prin pre-procesarea datelor de intrare, rezultând astfel fie o reducere a numărului de iterații necesare antrenării rețelei, sau chiar reducerea dimensiunilor rețelei. Asemenea reduceri sunt aproape întotdeauna specifice problemei și această abordare nu poate fi generalizată în vederea tratării tuturor genurilor de probleme.

O altă posibilitate constă în îmbunătățirea performanțelor algoritmului back-propagation, fie aplicând modificări ad hoc, fie utilizând rezultate din teoria optimizării numerice. Din această categorie face parte metoda gradientului conjugat [42].

O a treia abordare constă în accelerarea algoritmilor existenți prin implementarea lor direct în hardware (utilizând tehnici VLSI [89] sau optice [1]) sau modificarea lor pentru execuție pe o arhitectură paralelă. Acesta este subiectul tratat în această lucrare: paralelizarea algoritmului de învățare back-propagation.

## 2.8 Concluzii

Comparând rețelele neuronale artificiale cu cele biologice, trebuie notat faptul că, dintr-o serie de motive, nivelul curent de detaliu în modelarea neuronilor individuali este destul de simplist. Unul dintre motive îl constituie cunoștințele prezente, întrucâtva limitate, despre fiziologia neuronilor biologici. Un alt motiv important este faptul că nu toate aspectele neuro-fiziologice sunt relevante, dacă principalul scop este atingerea adaptabilității creierelor vii, și nu crearea de sisteme care să modeleze natura cât mai aproape posibil.

Spre deosebire de sistemele de calcul tradiționale, o rețea neuronală artificială este un sistem de procesare a informației neprogramat, adaptiv, care învață din experiență. Pe parcursul fazei de învățare i se prezintă un număr de exemple despre cum trebuie să se comporte la un anumit input. Gradual, rețeaua neuronală se adaptează singură, prin încercări. În loc să i se prezinte un set de instrucțiuni care să precizeze cum trebuie îndeplinită o anumită funcționalitate, rețeaua este capabilă să-și genereze propriile reguli interne care guvernează asocierea dintre intrare și ieșire. Aceste reguli sunt construite și rafinate continuu prin compararea rezultatelor produse de către rețea cu cele găsite în exemple.

O rețea neuronală artificială constă dintr-un set de elemente simple de procesare, numite *unități* (uneori mai sunt denumite și *neuroni*, datorită asocierii cu sistemele nervoase biologice), și un set de *conexiuni* între aceste unități. Activitatea se propagă prin rețea de la o unitate la alta, prin intermediul conexiunilor, fiecare având asociată o *pondere* (sau o *tărie a conexiunii*). Ponderea, care determină mărimea efectului unei unități asupra alteia, este de obicei reprezentată printr-un număr real. În funcție de semnul ponderii, conexiunea va fi fie una *excitatoare*, fie una *inhibitoare*, adică va crește sau va descrește activitatea unității receptoare. Intrarea fiecărei unități din rețea este formată prin combinarea ieșirilor tuturor unităților conectate, cu ponderile conexiunilor corespunzătoare. Activitatea fiecărei unități este apoi determinată prin aplicarea unei *funcții de activare* asupra intrării recepționate și, posibil, asupra activității curente. În final, *funcția de ieșire* transformă activitatea unității într-un semnal de ieșire, care este apoi propagat prin intermediul conexiunilor ca intrare pentru alte unități. De obicei, funcția de ieșire este pur și simplu funcția identitate, astfel încât ieșirea unei unități este egală cu activitatea ei.

Intrarea dinspre mediul înconjurător este furnizată rețelei prin stimularea unor unități speciale, numite *unități de intrare* (sau *senzori*). Tiparele de activitate observate la un anumit set de unități, numite *unități de ieșire*, sunt interpretate ca răspuns al rețelei la intrarea dată.

După cum s-a văzut, răspunsul rețelei la un tipar dat de intrare este determinat de ponderile conexiunilor dintre unități. Funcționalitatea rețelei poate deci să fie modificată prin modificarea acestor ponderi. Ca urmare, tiparul conectivității este ceea ce cunoaște sistemul și determină răspunsul acestuia la o intrare arbitrară. Însă dacă cunoștințele sunt reprezentate de tăria conexiunilor, atunci învățarea trebuie să conste din găsirea unor valori adecvate pentru aceste ponderi. Rezultă că se poate formula o *regulă de învățare* a rețelei, care să precizeze modul în care trebuie modificate ponderile ca răspuns la ieșiri incorecte sau parțial corecte produse de către rețea.

În secțiunea 2.7 a fost prezentată o implementare secvențială a algoritmului backpropagation. Calculele arată că pentru o rețea uzuală timpul de antrenare și necesarul de memorie pot deveni prea mari, putând descuraja investigațiile în domeniu prin experimente cu diferite configurații și parametri de antrenare. Capitolele următoare încearcă să depășească aceste neajunsuri, prin realizarea unor implementări paralele ale algoritmului backpropagation.



### 3 Paralelizarea rețelelor neuronale

#### 3.1 Strategii de paralelizare

Pornind de la paralelismul intrinsec al corespondentului biologic, rețelele neuronale au fost implementate pe o serie de arhitecturi paralele. Pe baza algoritmului backpropagation, au fost identificate mai multe posibilități de paralelizare:

- **Paralelizare la nivelul sesiunilor de antrenare:** Procesul de antrenare se face pornind de la parametri de antrenare inițiali diferiți pe diferite procesoare.
- **Paralelizare la nivelul setului de antrenare:** Tiparele din setul de antrenare sunt grupate pe diferitele procesoare, fiecare procesor executând aceeași rețea neuronală.
- **Pipelining:** Straturile rețelei sunt plasate pe procesoare diferite. Procesorul care posedă neuronii de pe stratul de ieșire poate să calculeze valorile de ieșire și cele de eroare pentru tiparul curent, în vreme ce procesorul care deține stratul intermediar procesează tiparul următor. Această abordare limitează însă numărul de procesoare care pot fi utilizate în paralel.
- **Paralelizare la nivelul nodurilor:** Un procesor deține neuroni de pe mai multe straturi.

Prima paralelizare, cea la nivelul sesiunilor de antrenare, reprezintă doar o valoare teoretică, întrucât ea nu aduce avantaje din punctul de vedere al accelerării antrenării și deci a reducerii timpului de antrenare, decât în faza de determinare a parametrilor pentru care rețeaua prezintă cea mai bună convergență.

Fiecare dintre aceste abordări prezintă un anumit grad maxim de paralelizare posibil, dat de diferiți factori:

- Paralelizarea la nivelul setului de antrenare: numărul de tipare din setul de antrenare.
- Pipelining: Numărul de straturi din rețea.
- Paralelizare la nivelul nodurilor: numărul de neuroni de pe straturi.

Vom prezenta în continuare o serie de realizări din acest domeniu.

##### 3.1.1 Paralelizarea la nivelul setului de antrenare

Această abordare mai poartă și denumirea de paralelizare cu partiționarea datelor, întrucât datele de intrare — tiparele de antrenare — se distribuie pe diferite procesoare.

O schemă de distribuire a tiparelor și de colectare a rezultatelor, într-o topologie în inel, este prezentată pentru mașina GF11 în lucrarea lui Witbrock și Zagha "An implementation of backpropagation learning on GF11, a large SIMD parallel computer" [100]. Mașina IBM GF11 este o mașină experimentală cu 566 procesoare, interconectate într-o rețea specifică. Aproape aceeași abordare este utilizată în lucrarea lui Bourrelly "Parallelisation of a neural learning algorithm on a hypercube" [7], însă pe un hipercube.

O abordare alternativă mai eficientă este de asemenea prezentată în [100]. Facilitățile de comunicare ale mașinii GF11 permit o inter-conectivitate mult mai complexă decât în inel. Aici se descrie cum  $P$  procesoare pot colecta și însuma gradientii componenți în  $O(\log_2 P)$  pași, într-o configurație astfel încât orice procesor  $q$  poate comunica cu toate procesoarele  $(q + 2^i) \bmod P$ , pentru  $i = 0, \dots, \log_2 P$ .



### 3.1.2 Paralelizarea la nivelul nodurilor

Această abordare mai este cunoscută și sub denumirea de partiționarea rețelei, întrucât, spre deosebire de metoda anterioară, de data aceasta se distribuie neuronii din rețea.

Strategia de partiționare a rețelei în vederea paralelizării a fost investigată în cadrul câtorva cercetări în domeniu.

Chinn și colegii, în [12] prezintă o asemenea implementare pe mașina MasPar MP-1 — un calculator masiv paralel care poate conține un array de 128x128 de procesoare. Pe aceeași mașină, mai poate fi menționată și lucrarea lui d’Acierno și Vaccaro [16], care combină paralelizarea la nivelul nodurilor cu cea la nivelul setului de antrenare.

Abordări similare se pot întâlni și în lucrările lui Ernoult — “*Performance of backpropagation on a parallel transputer-based machine*”, [22], del Rey Millan și Bofill — “*Learning by backpropagation: Computing in a systolic way*”, [20], și Zhang, McKenna, Mesirov și Waltz — “*The backpropagation algorithm on grid and hypercube architectures*”, [104].

După cum se observă din lucrările enumerate, implementările paralele au fost realizate pe mașini dedicate, cu număr mare de procesoare și o arhitectură deosebită a comunicației. Analizând dotarea laboratoarelor de calculatoare din universități și institute de cercetare, am identificat un potențial ridicat de calcul paralel și distribuit deja existent — un număr de calculatoare secvențiale independente (PCuri) interconectate într-o rețea ethernet. Investigațiile noastre se îndreaptă spre cercetarea posibilităților de utilizare în vederea antrenării cu algoritmul backpropagation, a unei asemenea rețele de calculatoare, fie grupate într-un cluster eterogen, fie grupate într-un cluster Beowulf [85], pe baza instrumentelor software disponibile. Posibilitatea implementării de algoritmi paraleli pe asemenea rețele de calculatoare este prezentată și de Polze și Malek în [67], iar Suresh și colegii, în [88], prezintă o asemenea paralelizare a algoritmului backpropagation pe un set de stații de lucru. Implementarea acestora realizează paralelizarea într-o abordare hibridă, la nivelul setului de antrenare și la nivelul nodurilor. Abordarea lor, însă, prezintă doar posibilitatea unei asemenea implementări și nu a permis investigarea performanțelor și a factorilor care le pot afecta. Această teză propune abordări distincte pentru cele două moduri de paralelizare, și realizarea unei serii de experimente în vederea identificării factorilor care ar putea influența performanțele fiecărei metode utilizate.

## 3.2 Direcții de acțiune

La transformarea unui algoritm secvențial într-unul paralel trebuie luate în considerare cel puțin două strategii. În contextul rețelelor neuronale, cele două strategii sunt cele numite în mod obișnuit *partiționarea datelor*, respectiv *partiționarea rețelei*.

### 3.2.1 Partiționarea datelor

Dacă același algoritm se aplică de un număr mare de ori pe diferite seturi de date, atunci este foarte eficient să se execute aceste aplicații concurrent, pe diferite procesoare, cu condiția ca task-urile să fie complet independente, adică execuția unui task să nu depindă de rezultatele altor task-uri. Această strategie de paralelizare este uneori denumită paralelism *la nivel de job* [25] sau *partiționarea datelor* [68].

Dacă actualizarea ponderilor dintr-o rețea neuronală se face doar după prezentarea mai multor tipare, fiecare dintre aceste prezentări de tipare este un task independent, care se poate executa concurrent. Ca urmare, dacă la antrenarea unei rețele neuronale actualizarea ponderilor se face în mod *batch* sau după o epocă, atunci abordarea prin partiționarea datelor se aplică foarte ușor: datele de antrenare sunt distribuite uniform procesoarelor disponibile, fiecare simulând întreaga rețea, dar pe

sub-seturi diferite de date de antrenare. Apoi rezultatele prezentării diferitelor grupuri de tipare sunt combinate și se face actualizarea ponderilor.

Aplicarea strategiei de partiționare a datelor la rețele neuronale mai este denumită și *paralelizarea antrenării* [20], din motive evidente.

Trebuie subliniat aici faptul că partiționarea datelor nu este posibilă dacă ponderile sunt actualizate după prezentarea fiecărui tipar. În acest caz, starea rețelei se modifică în urma fiecărei prezentări de tipar. Ca urmare, rezultatul prezentării unui tipar depinde de rezultatele *tuturor* tiparelor prezentate anterior, ceea ce înseamnă că task-urile prezentării tiparelor individuale nu mai pot fi considerate independente.

În capitolul 4 se vor descrie și analiza proprietățile unei paralelizări care exploatează strategia de partiționare a datelor.

### 3.2.2 Partiționarea rețelei

Un alt mod de paralelizare este utilizarea paralelismului așa-numit *geometric* [25] sau *spațial* [20]. În această abordare se paralelizează execuția algoritmului pe *un set de date*. Aceasta se realizează prin distribuirea datelor către procesoare astfel încât toate datele necesare unui procesor sunt stocate pe acel procesor sau sunt ușor accesibile la nevoie de la unul din procesoarele vecine.

Aplicarea acestei strategii la rețele neuronale mai este denumită și *partiționarea rețelei* [68], întrucât procesarea unui tipar poate fi paralelizată prin divizarea rețelei, fiecare procesor având de prelucrat o mică porțiune din rețea.

În capitolul 5 se prezintă o discuție a diferitelor moduri în care rețeaua poate fi divizată, precum și o descriere detaliată a unei implementări a strategiei de partiționare a rețelei pentru paralelizarea rețelelor neuronale.

Cele două strategii conduc la aplicații între care se pot identifica o serie de diferențe semnificative:

#### 1. Varianta cu partiționarea datelor:

- prezintă o oarecare limitare în privința dimensiunii rețelei neuronale care poate fi simulată, limitare generată de volumul de memorie disponibilă.
- permite doar actualizarea ponderilor la epocă, după prezentarea unui număr de tipare.

#### 2. Varianta cu partiționarea rețelei:

- permite simularea unor rețele de dimensiuni foarte mari, mai mari decât cele care pot fi procesate într-un algoritm secvențial.
- are o granularitate mai fină în privința actualizării ponderilor, care se poate face după fiecare tipar. Pentru unele aplicații de rețele neuronale această proprietate poate fi esențială pentru calitatea antrenării (vezi aplicația NET-Talk, secțiunea 6.1)

## 3.3 Analiza algoritmilor paraleli

Vor fi introduse aici o serie de concepte, utile în analiza performanțelor algoritmilor paraleli pentru rețele neuronale. Principalul obiectiv constă în cât mai buna exploatare a resurselor procesoarelor. Gradul de exploatare este denumit *eficiență*. La analiza unui algoritm paralel ne interesează modul în care este influențată eficiența algoritmului prin modificarea unui număr de parametri. Asemenea parametri sunt cei specifici rețelei neuronale, cum ar fi numărul de unități de pe fiecare strat, numărul de ponderi, frecvența de actualizare a acestor ponderi etc. Alți parametri țin însăși de paralelizare, cel mai important dintre aceștia fiind numărul de procesoare utilizate și modul de comunicare între ele.

Pentru a putea formula o definiție a eficienței trebuie cunoscută o limită în privința a cât de rapid ar putea să fie un algoritm paralel care utilizează un număr de  $P$  procesoare. Pentru aceasta, va fi introdus un alt concept uzual, denumit *accelerare* ("speed-up") [26]. Întrucât obiectivul principal al oricărei paralelizări este reducerea timpului de execuție, un mod natural de măsurare a performanței unui algoritm paralel este compararea directă a timpului de execuție pentru o problemă dată, cu timpul de execuție al algoritmului secvențial corespondent, pentru a determina cu cât este mai rapid algoritmul paralel.

Formal, se poate spune că, pentru un algoritm paralel, pentru o problemă dată, *accelerarea* este raportul dintre timpul de execuție  $T_{seq}$  al algoritmului secvențial și timpul de execuție  $T_{par}$  al algoritmului paralel, când ambii algoritmi sunt aplicați aceleiași probleme:

$$S(P) =_{def} \frac{T_{seq}}{T_{par}(P)} \quad (3.1)$$

unde  $P$  este numărul de procesoare utilizate la execuția algoritmului paralel. În general, algoritmul secvențial utilizat va fi *cel mai rapid cunoscut*. Totuși, pentru a putea utiliza această măsură a *accelerării* în vederea evaluării succesului în paralelizarea unui algoritm dat, se vor mai impune o serie de restricții asupra modului în care se obține  $T_{seq}$ . Ca urmare, algoritmul secvențial trebuie să fie implementat în același limbaj de programare ca și versiunea paralelă, iar procesorul pe care se execută algoritmul secvențial ar trebui să fie identic cu procesoarele utilizate la execuția algoritmului paralel. Întrucât în studiul de față s-a propus să se utilizeze pentru varianta paralelă o rețea de calculatoare, în care pot coexista mai multe platforme și implicit procesoare diferite, se va lua în considerație pentru  $T_{seq}$  valoarea cea mai mică, adică cea obținută pe mașina cea mai performantă.

Pe de altă parte, ambii algoritmi, și cel secvențial, și cel paralel, trebuie executați pe exact același set de date, având aceiași parametri specifici rețelei neuronale, inclusiv frecvența actualizărilor.

Orice paralelizare a algoritmului backpropagation va fi obligatoriu să conțină toate calculele din algoritmul secvențial. Ca urmare, dacă singurul aspect al unui algoritm paralel pentru rețele neuronale este faptul că acum calculele se pot efectua concurrent, atunci este evident că *accelerarea*  $S(P)$  al unui algoritm paralel pe  $P$  procesoare este limitată de valoarea lui  $P$ . Având de  $P$  ori mai multe resurse de calcul, în cel mai bun caz este de așteptat o reducere a timpului de execuție de  $P$  ori.

Totuși, această valoare nu va putea fi atinsă în cazul de față, când execuția se face într-o rețea cu platforme diferite și deci cu performanțe diferite.

Se poate beneficia însă de reduceri ale timpului de execuție datorită altor factori, cum ar fi parametrii specifici hardware-ului. Un asemenea exemplu este cazul memoriei *cache* existentă pe fiecare procesor, și care este foarte rapidă. Dacă în cazul algoritmului secvențial aceasta nu este suficient de mare pentru a stoca toate ponderile dintr-o rețea, în cazul algoritmului paralel s-ar putea realiza o asemenea partiționare încât ea să conțină toate ponderile necesare pentru acea porțiune din rețea pe care o are de prelucrat un procesor. Având  $P$  procesoare, s-ar putea deci obține un *speed-up* chiar mai mare decât  $P$  (!?) În studiul de față nu este luat în calcul acest fenomen.

Pe baza definițiilor de mai sus, se poate acum exprima eficiența ca raportul dintre *speed-up*-ul obținut experimental și cel optim. Dacă se presupune că *speed-up*-ul se datorează numai efectelor calculelor concurente, adică  $S(P)$  are limita  $P$ , se propune următoarea definiție formală a eficienței:

$$E(P) =_{def} \frac{S(P)}{P} = \frac{T_{seq}}{T_{par} \cdot P} \quad (3.2)$$

După cum se observă, cu presupunerile făcute,  $E(P)$  are valori cuprinse în intervalul  $(0,1]$ .

### 3.4 Principalele obiective ale paralelizării

Cel mai important motiv pentru paralelizarea unui algoritm este de obicei necesitatea de a obține o reducere a timpului de execuție pentru acel algoritm. O asemenea reducere a timpului în procesarea unui set de date este absolut necesară, întrucât nu numai că permite prelucrarea mai multor task-uri de aceeași dimensiune, dar și face posibilă prelucrarea unor task-uri mai mari din punct de vedere computațional, după cum arată și Fox, Johnson, Lyzenga, Otto, Salmon și Walker în *Solving Problems On Concurrent Processors* [26].

De asemenea, pot fi prelucrate task-uri mai mari în privința necesarului de memorie (de exemplu, probleme cu seturi de date mai mari), întrucât paralelizarea permite distribuirea datelor pe care lucrează algoritmul, reducând astfel necesarul de memorie al procesoarelor individuale. Ca urmare, necesarul total de memorie pentru un algoritm paralel este de preferat să nu fie mai mare decât cel al algoritmului secvențial.

Dacă se poate proiecta o paralelizare eficientă a unui algoritm, acesta este un mod simplu și economic de accelerare a execuției unui algoritm. Presupunând că se păstrează eficiența algoritmului paralel chiar și atunci când se utilizează un număr mare de procesoare, un sistem de procesoare care lucrează în paralel va putea să se comporte asemenea unui singur calculator, mai puternic. Și mai mult, un sistem paralel este ușor de extins, astfel încât se poate câștiga în viteză adăugând procesoare.

Totuși, în general, paralelizarea unui algoritm nu constituie o sarcină ușoară. Vor trebui luate în considerație o serie de aspecte, inclusiv proprietățile hardware-ului disponibil. În cazul abordat în lucrarea de față, când fiecare procesor are memoria proprie și este fără acces direct la memoria altor procesoare, este important de cunoscut cât de rapidă este comunicația între procesoare în comparație cu performanțele computaționale ale fiecărui procesor. Interesează în acest caz, când transferul de date și calculul se pot efectua concurrent, să se obțină timpi de comunicație mult mai reduși decât cei pentru calcule.

Pe de altă parte, se va căuta să se asigure o încărcare cât mai uniformă a procesoarelor, întrucât aceasta va asigura timpul de execuție cel mai scurt.

### 3.5 Surse de ineficiență într-un algoritm paralel

Există o serie de motive pentru care un algoritm paralel poate să nu utilizeze resursele computaționale la fel de eficient ca și algoritmul secvențial corespunzător. Atunci când se observă o scădere a eficienței într-un anumit algoritm, cel mai probabil, aceasta va fi rezultatul câtorva dintre cauzele enumerate mai jos. În cadrul acestei lucrări se va analiza în ce măsură aceste cauze se aplică fiecărui algoritm în parte, în secțiunea corespunzătoare.

#### 3.5.1 Software adițional

În vederea procesării datelor provenind de la alte procesoare, poate fi necesar să se introducă o serie de calcule suplimentare, de exemplu pentru indexare, sau chiar mai complicate. De asemenea, uneori este necesară modificarea secvenței calculelor, din diferite motive (vezi mai jos), astfel încât este posibil ca unele rezultate temporare să nu mai fie disponibile și deci să necesite recalcularea. Orice asemenea sarcini suplimentare vor reduce eficiența algoritmului paralel.

Totuși, dacă software-ul suplimentar constituie o fracțiune constantă din sarcinile fiecărui procesor, indiferent de numărul de procesoare care sunt utilizate, atunci cantitatea totală de efort computațional care corespunde software-ului adițional este independentă de numărul de procesoare. Întrucât această sarcină trebuie efectuată în paralel, eficiența algoritmului se reduce întotdeauna cu același factor constant,

indiferent de numărul de procesoare utilizate în execuția algoritmului:

$$E(P) = \frac{T_{seq}}{T_{par}(P) \cdot P} = \frac{T_{seq}}{\frac{T_{seq} + T_{soft}}{P} \cdot P} = \frac{T_{seq}}{T_{seq} + T_{soft}} \quad (3.3)$$

În formula de mai sus s-a presupus că singura cauză a ineficienței algoritmului paralel este software-ul adițional, independent de numărul procesoarelor.  $T_{soft}$  este timpul necesar unui procesor pentru execuția calculelor suplimentare.

Întrucât în formula finală din relația 3.3 nu mai apare numărul de procesoare utilizate, eficiența este redusă printr-un factor *constant*,  $T_{soft}$ . Acest tip de software adițional nu poate să impună o limită a accelerării care poate fi obținută pentru o problemă de rețele neuronale dată. Prezența acestui software duce la o mai slabă utilizare a fiecărui procesor individual implicat în simularea unei rețele neuronale date.

Pe de altă parte, este posibil ca acest software suplimentar să ducă la creșterea volumului de muncă — pentru fiecare procesor și pentru sistem în ansamblu — odată cu creșterea numărului de procesoare. Evident că în acest caz și eficiența se va deteriora.

### 3.5.2 Echilibrarea încărcării

O rețea de mașini lucrând concurent nu a terminat un task decât atunci când *toate* procesoarele și-au terminat sarcinile. Ca urmare, este important ca încărcarea să fie distribuită uniform pe toate procesoarele, astfel încât fiecare mașină să efectueze aceeași cantitate de muncă. Mai mult, încărcarea trebuie echilibrată, uniform, *pe tot parcursul* execuției algoritmului, altfel procesoarele vor alterna stările de lucru, respectiv de așteptare, într-un asemenea mod încât chiar dacă au de efectuat aceeași cantitate de muncă, ele nu o vor termina în paralel.

Este posibil ca problemele de echilibrare a încărcării să devină mai pronunțate odată cu creșterea numărului de procesoare, funcție de configurația utilizată și de natura problemei de rețele neuronale de tratat.

### 3.5.3 Comunicația suplimentară

Orice cantitate de timp utilizată pentru comunicație va reduce eficiența algoritmului, întrucât este vorba despre operații suplimentare, față de algoritmul secvențial respectiv. Ca urmare, se va încerca minimizarea comunicației între procesoare, cel puțin dacă acest transfer de date nu poate avea loc în paralel cu alte calcule.

Un aspect important care trebuie luat în considerare în construcția algoritmului paralel este rearanjarea secvenței de calcule astfel încât să se asigure cel mai înalt grad posibil de paralelism între comunicație și calcule. Uneori, totuși, este posibil să nu mai fie de efectuat calcule care să nu depindă de datele comunicate până în acel moment. De asemenea, trebuie remarcat faptul că deși costul comunicației poate fi redus semnificativ prin efectuarea în paralel cu calculele, acest cost nu va ajunge niciodată neglijabil.

În câteva aplicații, cantitatea de date care se transmit de către fiecare procesor nu depinde de numărul de procesoare implicate. Aceasta înseamnă că timpul de transmisie pentru fiecare procesor nu se reduce la creșterea numărului de procesoare. Din moment ce fracțiunea alocată fiecărui procesor dintr-o problemă de rețele neuronale devine din ce în ce mai mică, odată cu creșterea numărului de procesoare utilizate pentru simulare, aceasta determină ca o fracțiune tot mai mare din timp să fie alocată comunicării, ducând astfel la scăderea eficienței.

Pentru a demonstra aceasta, se presupune că pentru un algoritm paralel oarecare singura cauză pentru care eficiența nu este cea optimă este transmisia, adică totul, în afara comunicației, este perfect paralelizat. Mai mult, presupunând că timpul utilizat în fiecare procesor pentru transmisie depinde numai de dimensiunea problemei



de rețele neuronale, și deci este independent de numărul de procesoare, se va obține

$$S(P) = \frac{T_{seq}}{T_{par}(P)} = \frac{T_{seq}}{\frac{T_{seq}}{P} + T_{comm}} \quad (3.4)$$

De remarcat că dacă  $T_{comm}$  — timpul consumat în fiecare procesor pentru comunicație — este nul, atunci *speed-up*-ul este cel optim. Pe de altă parte, dacă  $T_{comm}$  nu este nul, transmisia va impune o limită superioară pentru *speed-up*,  $\frac{T_{seq}}{T_{comm}}$ .

În general, aceleași considerații legate de efectele software-ului adițional se aplică și la transmisie. Dacă timpul de transmisie nu se reduce la creșterea numărului de procesoare, atunci aceasta va duce la deteriorarea eficienței la mărirea numărului de procesoare.

### 3.5.4 Porțiuni inerent secvențiale

Pot exista porțiuni din algoritm care sunt inerent secvențiale, adică pur și simplu nu pot fi paralelizate, de exemplu datorită faptului că execuția câte unui pas în algoritm depinde de finalizarea pasului anterior. În contextul rețelelor neuronale, asemenea porțiuni inerent secvențiale pot fi adesea găsite în distribuția inițială a ponderilor sau a tiparelor de antrenare, ca și în colectarea finală a rezultatelor parțiale de la procesoarele individuale. Un alt exemplu ar putea fi calculul produsului scalar global (ca în cazul algoritmului gradientului conjugat [42]). Trebuie menționat faptul că o porțiune inerent secvențială a algoritmului uneori poate fi totuși paralelizată, în sensul că toate procesoarele efectuează aceleași calcule, fiecare procesor calculându-și singur rezultatele. Aceasta însă nu constituie o reală paralelizare, din moment ce timpul de rulare necesar obținerii rezultatului nu se reduce în comparație cu algoritmul secvențial. Totuși, ar putea fi o soluție mai elegantă (și mai eficientă) decât varianta în care un singur procesor efectuează calculele, urmând apoi un *broadcast* pentru a le transmite celorlalte procesoare.

Dacă un algoritm conține porțiuni inerent secvențiale, timpul necesar executării acestor porțiuni va determina o limitare a *speed-up*-ului posibil de obținut, independent de numărul de procesoare utilizate. Acest fenomen este cunoscut ca *legea lui Amdahl* [26]. Aceasta spune că dacă o porțiune inerent secvențială a algoritmului necesită  $\frac{1}{\alpha}$  din timpul de execuție atunci când algoritmul se execută pe un procesor, atunci nu este posibil să obținem un *speed-up* mai mare decât  $\alpha$ . Pe măsură ce adăugăm procesoare, timpul necesar execuției porțiunii secvențiale va reprezenta o fracțiune tot mai mare din timpul total de execuție, reducând astfel eficiența.

### 3.5.5 Limitări specifice problemei

Atâta vreme cât doar o singură strategie de paralelizare este utilizată, întotdeauna există o limită teoretică a *speed-up*-ului care poate fi atins pentru o problemă de dimensiune dată.

Orice paralelizare constă în divizarea problemei într-un număr de sub-probleme, care pot fi distribuite pe procesoarele disponibile. Diferitele strategii de paralelizare divid problema în moduri diferite — se poate spune de dimensiuni diferite. Pentru fiecare dimensiune/strategie de paralelizare, există o limită a numărului de subprobleme în care poate fi descompusă problema originală. Acest număr de sub-probleme determină o limită superioară a numărului de procesoare care pot fi utilizate eficient în algoritmul paralel, întrucât fiecare procesor va trata cel puțin o singură sub-problemă.

Ca exemplu, se poate lua în considerație o implementare simplă a strategiei de partiționare a datelor pentru paralelizarea rețelelor neuronale. În contextul unui asemenea algoritm, prezentarea unui tipar din setul curent de tipare de antrenare reprezintă o sub-problemă, care poate fi rezolvată de către un singur procesor. Totuși, în acest caz problema inițială nu poate fi divizată în mai multe asemenea sub-probleme decât numărul de tipare din setul de antrenare, întrucât fiecare procesor trebuie să

trateze cel puțin un asemenea tipar (sau nu va face nimic altceva decât să "deranjeze" celelalte procesoare care "au de lucru"). Cu alte cuvinte, în abordarea partiționării datelor nu se poate obține un *speed-up* mai mare decât dimensiunea setului de antrenare utilizat.

Similar, se observă că într-o implementare a partiționării rețelei nu se pot utiliza mai multe procesoare decât numărul de unități din rețea, limitându-se astfel *speed-up*-ul care poate fi atins la numărul de unități din rețea.

De notat totuși că limitările de mai sus sunt mai mult teoretice, întrucât de obicei *speed-up*-ul este mult mai limitat de efectele unui număr de alte cauze, mai ales software-ul suplimentar și transmisia de date.

### 3.6 Considerații specifice paralelizării rețelelor neuronale

În general, la paralelizarea unui algoritm, se consideră că singura diferență față de algoritmul secvențial va consta în viteza cu care se execută cei doi algoritmi. Cu alte cuvinte, cei doi algoritmi vor fi identici din punct de vedere funcțional, adică vor produce același *output*.

În contextul rețelelor neuronale, aceasta înseamnă că prin paralelizarea algoritmului se va păstra calitatea învățării, în privința puterii de generalizare și a vitezei de învățare (având ca măsură numărul de prezentări de tipare de învățare). Totuși, deși ar putea suna ciudat, la algoritmii de învățare paralelizați nu se întâmplă întotdeauna așa.

Problema apare datorită modificării frecvenței cu care sunt actualizate ponderile conexiunilor din rețeaua neuronală. Calitatea învățării este puternic influențată de numărul de tipare prezentate rețelei între două actualizări (adică de dimensiunea unui set de antrenare). Desigur, algoritmul *backpropagation* secvențial poate utiliza orice frecvență de actualizare a ponderilor, potrivită problemei tratate. Totuși, în cazul paralelizării (în special în cazul partiționării datelor), adesea ponderile nu sunt actualizate cu frecvența cea mai adecvată problemei date. Pentru o paralelizare eficientă, se alege un alt număr de prezentări de tipare între două actualizări de ponderi. Ca rezultat, dimensiunea unui set de învățare depinde de obicei de numărul de procesoare implicate în execuția algoritmului paralel.

Algoritmii de paralelizare pentru rețele neuronale sunt paralelizați din dorința de a mări viteza (măsurată în timp absolut) cu care rețeaua poate fi antrenată pentru a efectua o sarcină dată, cu un grad dat de perfecțiune. Ca urmare, ca o consecință a efectului descris, de modificare a dimensiunii unui set de învățare, modul corect de comparare a diferiților algoritmi paraleli este *nu* prin măsurarea numărului de tipare care pot fi prezentate rețelei pe secundă, întrucât aceasta nu ia în considerație nici o diferență în efectul asupra învățării a fiecărei prezentări de tipar. Ideal ar fi să se compare algoritmii paraleli prin măsurarea timpului în care fiecare algoritm poate să antreneze o rețea neuronală dată pentru a efectua o anumită sarcină, cu un grad dat de perfecțiune. Oricum, efectul dimensiunii setului de antrenare asupra calității învățării diferă de la aplicație la aplicație și deci dacă doi algoritmi paraleli utilizează seturi de învățare de dimensiuni diferite, evaluarea lor va depinde de problema luată ca etalon.

Concluzia la cele de mai sus este că, pentru a realiza într-adevăr o comparație corectă a doi algoritmi paraleli — rulați sau nu pe același sistem de procesoare — ei trebuie aplicați aceleiași probleme, având toți parametrii identici, mai ales dimensiunea setului de învățare.

### 3.7 Experimente în vederea analizei

## performanțelor

În general, în cadrul tezei de față s-a urmărit efectul tuturor tipurilor de parametri asupra eficienței algoritmului paralel, fie că este vorba despre parametri specifici rețelei neuronale, cum ar fi dimensiunea setului de antrenare, fie că este vorba despre parametri specifici hardware-ului, cum ar fi structura configurației de procesoare sau numărul de procesoare utilizate în execuția algoritmului paralel.

După cum s-a menționat mai devreme, prin paralelizarea unui algoritm se urmăresc două obiective principale: reducerea timpului necesar aplicării algoritmului pentru o problemă dată, respectiv creșterea dimensiunii problemelor care pot fi tratate.

### 3.7.1 Probleme de dimensiune fixă

Pentru un număr de probleme, se vor prezenta efectele creșterii numărului de procesoare utilizate în execuția algoritmului paralel. Dacă fracțiunea de timp necesară de fiecare procesor este proporțională cu numărul de sub-probleme tratate de către acel procesor, atunci graficele *speed-up*-ului vor fi liniare, ceea ce indică faptul că *speed-up* este liniar cu numărul de procesoare (cât timp avem mai puține procesoare decât numărul de sub-probleme), iar graficul eficienței va fi o dreaptă orizontală, ceea ce ne arată că eficiența este independentă de numărul de procesoare utilizate.

Dacă însă fracțiunea de timp necesară fiecărui procesor crește odată cu numărul de procesoare, eficiența se va deteriora pe măsură ce adăugăm procesoare în sistem.

Pentru fiecare experiment, se va încerca să se interpreteze graficele, în scopul de a determina care din circumstanțele menționate anterior joacă cel mai important rol în reducerea eficienței algoritmului în cauză.

### 3.7.2 Probleme de dimensiune variabilă

De asemenea, se va examina relația dintre dimensionalitatea problemei de rețele neuronale și numărul de procesoare, pentru a determina care este dimensiunea maximă a problemelor care pot fi tratate (cu același grad de eficiență) când se utilizează mai multe procesoare. Ideal ar fi să se utilizeze de două ori mai multe procesoare atunci când se dublează dimensiunea problemei, însă, după cum se va vedea, este dificil să se măsoare dimensiunea unei probleme de rețele neuronale.

Mai precis, *dimensiunea* unei probleme este dată de cantitatea de efort computațional asociat rezolvării problemei (utilizând algoritmul secvențial). În domeniul rețelelor neuronale, *rezolvarea problemei* înseamnă antrenarea rețelei până la atingerea unui grad dat de perfecțiune.

De aceea, de obicei este foarte dificil să se compare dimensiunile pentru două probleme de rețele neuronale, întrucât adăugarea de unități în rețea sau mărirea setului de antrenare pot sau nu să modifice numărul de cicluri de antrenare necesare atingerii unui nivel dat de perfecțiune, în funcție de natura antrenării dorite.

Pentru un tip dat de antrenare, nu se poate doar să se presupună că o rețea conținând un număr dublu de unități ascunse reprezintă o problemă (aproximativ) de două ori mai mare<sup>1</sup>, întrucât numărul de cicluri de antrenare necesare atingerii unei performanțe date poate fi mare atât în cazul în care avem prea puține unități ascunse, cât și atunci când ele sunt prea multe.

Similar, creșterea dimensiunii setului de antrenare poate conduce la scăderea performanței rețelei în privința numărului de cicluri de antrenare necesare atingerii unui nivel dat de perfecțiune. Dublarea dimensiunii setului de antrenare poate mări dimensiunea computațională a problemei cu un factor mult mai mare de 2. Sau, în cazul altor aplicații, este posibil să reducă numărul de cicluri de antrenare necesare producerii unei soluții, reducând chiar dimensiunea problemei.

<sup>1</sup>La dublarea numărului de unități ascunse, se dublează și numărul de ponderi din rețea. Totuși, întrucât numărul de unități de ieșire este același, efortul de calcul la fiecare ciclu de antrenare va fi doar *aproximativ* dublat.



Datorită dificultăților enumerate mai sus, nu se va încerca să se măsoare efectele scalării dimensiunii problemei cu numărul de procesoare asupra eficienței. În schimb, se va examina influența asupra eficienței a variației altor parametri, ca dimensiunea rețelei și dimensiunea setului de antrenare.

### 3.7.2.1 Număr variabil de sub-probleme

În loc să se studieze eficiența atunci când *dimensiunea* problemei este scalată cu numărul de procesoare, se va examina gradul în care se păstrează eficiența atunci când numărul de *sub-probleme* este scalat cu numărul de procesoare. Cu alte cuvinte, se va verifica dacă se pot utiliza mai multe procesoare, fără a reduce din eficiență, atât timp cât fiecare procesor prelucrează același număr de sub-probleme.

După cum s-a precizat, chiar dacă o problemă nouă poate fi divizată într-un număr dublu de sub-probleme, aceasta nu implică neapărat dublarea dimensiunii problemei. Într-o paralelizare bazată pe partiționarea rețelei, numărul de sub-probleme este dublat prin dublarea numărului de unități din fiecare strat al rețelei. Nu se poate preciza efectul asupra dimensiunii problemei, întrucât acesta depinde de aplicația de rețele neuronale. Însă ceea ce se cunoaște este faptul că dublarea numărului de unități înseamnă mărirea numărului de ponderi de 4 ori, conducând la mărirea de aproape 4 ori a efortului computațional la prezentarea unui tipar.

În cazul paralelizării prin partiționarea datelor, dublarea numărului de sub-probleme se obține prin dublarea dimensiunii setului de antrenare. Din nou, influența asupra dimensiunii problemei depinde de aplicația de rețele neuronale tratată. Se poate afirma, totuși, că, întrucât se dublează numărul de prezentări de tipare dintre două actualizări de ponderi, volumul operațiilor asociate unui singur ciclu de antrenare se va dubla la rândul său<sup>2</sup>.

Ca urmare, se poate afirma că, dintr-un punct de vedere, având o problemă nouă, care conține un număr mai mare de sub-probleme, nu este vorba atât de o problemă mai mare, cât de pur și simplu o problemă *diferită*, cu proprietăți diferite.

Pentru câteva seturi de asemenea numere fixe de sub-probleme per procesor, se va studia efectul mării numărului de procesoare. La fel ca în cazul experimentelor cu probleme de dimensiune fixă, se va încerca să se interpreteze graficele *speed-up*-ului și eficienței rezultate. Dacă graficele eficienței nu sunt drepte orizontale, se poate concluziona că activitatea suplimentară din fiecare procesor nu este independentă de numărul de procesoare utilizate la execuția algoritmului.

### 3.7.2.2 Număr fix de sub-probleme

Se vor examina efectele modificării de parametri specifici rețelelor neuronale care nu influențează numărul de sub-probleme, întrucât cunoașterea acestor efecte este necesară cercetătorului în rețele neuronale, cu scopul alegerii celei mai potrivite paralelizări pentru aplicația dată. Ca urmare, se va examina, de asemenea, efectul modificării dimensiunii rețelei în cazul paralelizării prin partiționarea datelor, precum și influențele asupra eficienței a modificării dimensiunii setului de antrenare în cazul abordării cu partiționarea rețelei.

## 3.8 Concluzii

Acest capitol realizează o investigație asupra diverselor abordări în paralelizarea algoritmului backpropagation. Au fost identificate două strategii majore: cea cu partiționarea datelor și cea cu partiționarea rețelei. În urma acestor investigații se evidențiază o serie de aspecte care ar putea afecta performanțele algoritmilor care urmează a fi prezentați în capitolele 4 și 5. Utilizarea unei rețele de calculatoare impune o atenție

<sup>2</sup>Mai precis, va fi *aproape* dublu, întrucât efortul necesar actualizării ponderilor nu se modifică.

sporită în ceea ce privește comunicația între mașini, aceasta fiind una dintre cauzele pentru care accelerarea obținută nu va putea fi cea ideală.

Capitolul trasează și o serie de experimente care vor trebui efectuate în vederea măsurării performanțelor algoritmilor propuși. În principal, se va urmări comportarea acestor algoritmi atât pentru probleme de dimensiune fixă — cum afectează mărirea numărului de mașini pe care rulează simulatorul — precum și pentru probleme de dimensiune variabilă — respectiv relația dintre dimensiunea problemei și numărul de procesoare utilizate.

Capitolul de față prezintă algoritmi care au stat la baza construcției simulatorului DataParSim – simulator paralel cu partiționarea datelor, implementat și testat pe trei rețele de calculatoare, una conținând 12 stații IBM RS6000, una cu 16 stații SunBlade150 și o rețea de 16 PC-uri cu procesor Pentium la 2,8 GHz. Secțiunile 4.1, 4.2 descriu algoritmi utilizați, în secțiunea 4.3 este prezentat simulatorul realizat pe baza acestor algoritmi, urmând ca în 4.4 să fie prezentate și evaluate rezultatele experimentelor efectuate.

După cum s-a precizat în capitolul anterior, un mod simplu de paralelizare a unei rețele neuronale, cu actualizarea ponderilor după epocă sau în mod *batch*, este așa-numita strategie de *partiționare a datelor*, în care datele de antrenare sunt distribuite uniform pe procesoare. Toate procesoarele simulează întreaga rețea, dar pentru sub-seturi diferite de antrenare. Pe parcursul fiecărui ciclu de învățare, fiecare procesor prezintă tiparele din "felia" sa din setul curent<sup>1</sup>. Gradienții calculați în acest mod de către fiecare dintre procesoare îi vom numi *gradienți component*, întrucât fiecare dintre ei este rezultatul prezentării către rețea doar a unei porțiuni din setul de antrenare. Odată ce toate tiparele din fiecare set au fost prezentate de către diferitele procesoare, gradienții componentii rezultați sunt combinați (însumate) într-un gradient global, care este apoi utilizat pentru calculul modificării fiecărei ponderi. Apoi, ponderile sunt actualizate și noile valori sunt transmise tuturor procesoarelor, ca fiind noul set de ponderi care vor fi utilizate în prezentările de tipare din următorul ciclu de învățare.

Aceasta este cea mai simplă formă de partiționare a datelor pentru paralelizarea rețelelor neuronale. Se va discuta implementarea acestei forme în secțiunea următoare. În secțiunea 4.2 se va prezenta o variațiune a partiționării datelor, în care fiecare procesor nu tratează o întreagă copie a rețelei.

#### 4.1 O implementare simplă a strategiei de partiționare a datelor

Procesul de implementare a acestei forme de partiționare a datelor este foarte simplu și evident, întrucât se dispune deja de o implementare secvențială a algoritmului de backpropagation. Tot ce trebuie făcut este să se pună împreună un set de procesoare, fiecare dintre ele rulând algoritmul secvențial clasic, în mod independent. Având  $P$  procesoare, vor fi deci  $P$  copii identice ale întregii rețele. Pe parcursul fiecărui ciclu de învățare, fiecare procesor  $q$  prezintă rețelei tiparele din porțiunea sa,  $B_q$ , din setul curent  $B$  de tipare de antrenare. Ca rezultat, în procesorul  $q$  se calculează gradientul component  $\bar{g}^{<q>}$ :

$$\bar{g}^{<q>} = \sum_{p \in B_q} \frac{\partial E_p}{\partial \bar{w}} \quad (4.1)$$

unde  $\bar{w}$  denotă setul tuturor ponderilor.

Pe baza ecuațiilor

$$\Delta w_{j \rightarrow k}^o = -\eta \sum_p \frac{\partial E_p}{\partial w_{j \rightarrow k}^o} = \eta \sum_p \delta_{pk}^o a_{pj}^H \quad (4.2)$$

și

$$\Delta w_{j \rightarrow k}^o(n+1) = \eta \sum_p \delta_{pk}^o a_{pj}^H + \alpha \Delta w_{j \rightarrow k}^o(n) \quad (4.3)$$

<sup>1</sup> În această secțiune nu se face distincție între actualizarea la epocă și cea în mod *batch*, întrucât din punctul de vedere al paralelizării unui algoritm, actualizarea la epocă poate fi privită ca un caz special de mod *batch*, în care setul de antrenare este constituit din întregul set de tipare.

modificarea ponderilor conform tiparelor din setul  $B$  poate fi exprimată acum ca sumă a gradientilor componenți  $\bar{g}^{<q>}$ :

$$\Delta_B \bar{w}(n+1) = -\eta \sum_{p \in B} \frac{\partial E_p}{\partial \bar{w}} + \alpha \Delta_B \bar{w}(n) = -\eta \sum_{q=0}^{P-1} \bar{g}^{<q>} + \alpha \Delta_B \bar{w}(n) \quad (4.4)$$

Ca urmare, trebuie efectuată o singură modificare față de algoritmul secvențial: o schemă de colectare și însumare a tuturor gradientilor componenți  $\bar{g}^{<q>}$  și distribuirea ponderilor actualizate (*broadcast*) în modul cel mai eficient posibil.

#### 4.1.1 Configurație inel. Configurație arbore

În funcție de configurația de procesoare, trebuie utilizate scheme de comunicație diferite, în funcție de modul de conectare.

Având o configurație în inel de  $P$  procesoare, se poate implementa însumarea gradientilor componenți astfel încât fiecare procesor, după  $P-1$  pași, să dețină gradientul total. Inițial, fiecare procesor comunică gradientul component propriu către predecesorul său din inel (figura 4.1). Apoi toate procesoarele adună gradientul component recepționat la o sumă temporară. În același timp, gradientul component este transmis predecesorului, concurrent cu recepționarea unui nou gradient component.

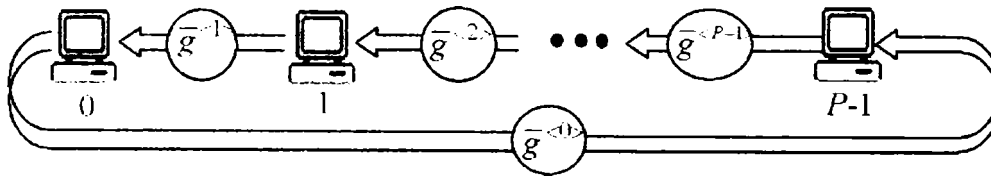


Figura 4.1: Însumarea gradientilor componenți.

După  $P-1$  pași, fiecare procesor a recepționat gradientii componenți de la toate celelalte procesoare. În acest moment, el va deține suma tuturor gradientilor componenți și va putea face actualizarea ponderilor.

O schemă similară este prezentată pentru mașina GF11 în lucrarea lui Witbrock și Zaha "An implementation of backpropagation learning on GF11, a large SIMD parallel computer" [100], și aproape aceeași abordare este utilizată în lucrarea lui Bourrely "Parallelisation of a neural learning algorithm on a hypercube" [7], însă pe un hipercube.

O abordare alternativă mai eficientă este de asemenea prezentată în [100]. Facilitățile de comunicare ale mașinii GF11 permit o inter-conectivitate mult mai complexă decât în inel. Aici se descrie cum  $P$  procesoare pot colecta și însuma gradientii componenți în  $O(\log_2 P)$  pași, într-o configurație astfel încât procesorul  $q$  poate comunica cu toate procesoarele  $(q + 2^i) \bmod P$ , pentru  $i = 0, \dots, \log_2 P$ .

Acest mod de interconectare nu poate fi însă aplicat, datorită traficului încărcat care ar rezulta. Totuși, este posibilă însumarea gradientilor componenți în  $O(\log_2 P)$  pași, configurând procesoarele sub forma unui arbore binar<sup>2</sup>.

După prezentarea tuturor tiparelor din setul curent, toate procesoarele fără succesori în structura de arbore transmit gradientii componenți acumulați către procesorul părinte (Figura 4.2a). Toate celelalte procesoare colectează gradientii componenți de la succesorii lor imediați, îl însumează pe cel propriu și transmit rezultatul în sus în ierarhie (Figura 4.2b). După  $\lceil \log_2 P \rceil$  pași, procesorul rădăcină va deține gradientul

<sup>2</sup>S-ar putea utiliza și arbori de nivel superior, reducând astfel numărul de pași și înălțimea arborelui, însă ar crește necesarul de memorie pentru comunicarea gradientilor componenți, iar traficul ar fi același.

total, care poate fi utilizat în calculul modificării ponderilor. Utilizând structura de arbore din nou, procesorul rădăcină poate transmite ponderile actualizate tuturor celorlalte procesoare în  $\lfloor \log_2 P \rfloor$  pași.

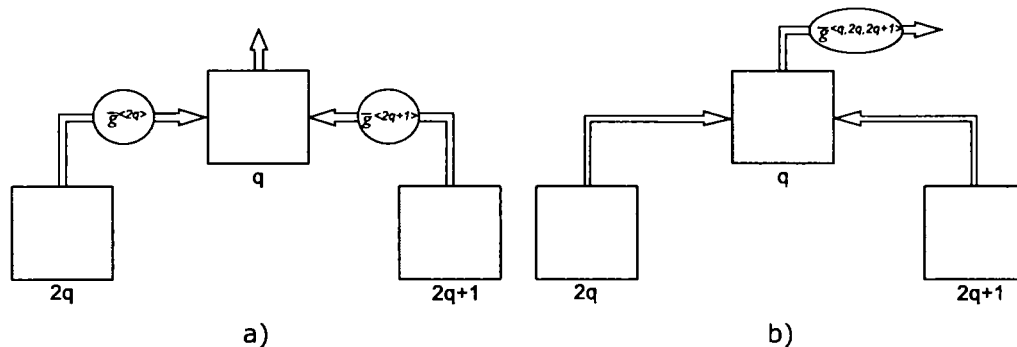


Figura 4.2: Însurarea gradientilor componente într-un arbore binar.

$P$  procesoare în inel pot actualiza ponderile în  $P - 1$  pași, în timp ce sunt necesari  $\lfloor \log_2 P \rfloor$  pași într-o configurație în arbore binar. Aceasta înseamnă că dacă se utilizează mai mult de 5 procesoare, configurația în arbore necesită mai puțini pași. Este deci evident că se va prefera configurația de arbore atunci când este disponibil un număr mare de procesoare. În schimb, în configurația în inel se poate efectua însurarea gradientilor componente concurrent cu transmisia și recepția de (alți) gradienti componente. De remarcat însă faptul că efectuarea în paralel a comunicației și a calculelor economisește cel mult timpul asociat task-ului cel mai scurt<sup>3</sup>.

#### 4.1.1.1 Necesarul de memorie

Necesarul de memorie este aproximativ același pentru ambele configurații. În ambele cazuri, fiecare procesor trebuie să rezerve spațiu în plus pentru câte două valori reale pentru fiecare pondere din rețea, necesare la transmisia gradientilor componente.

Singura diferență semnificativă în privința necesarului de memorie dintre cei doi algoritmi este următoarea: în configurația în inel fiecare procesor trebuie să păstreze modificările anterioare ale ponderilor, pe lângă valorile calculate curent și valorile ponderilor. Aceasta nu mai este necesar în configurația de arbore, întrucât noul set de ponderi se calculează în procesorul rădăcină. Ca urmare, doar procesorul rădăcină trebuie să memoreze modificările anterioare. Toate celelalte procesoare necesită spațiu numai pentru ponderi, gradientii componente calculați și câte doi vectori de ponderi pentru comunicație.

Întrucât fiecare procesor păstrează de asemenea o copie a tuturor unităților din rețea, dacă  $NI$ ,  $NH$  și  $NO$  reprezintă numărul de neuroni de pe stratul de intrare, cel ascuns, respectiv stratul de ieșire, se obține următorul necesar de memorie<sup>4</sup> (măsurat în valori reale) pentru fiecare procesor (cu excepția celui rădăcină):

$$2(NI + NH + NO) + 4((NI + 1)NH + (NH + 1)NO) \quad (4.5)$$

Acest necesar enorm de memorie impune o limită mult prea strictă asupra dimensiunii rețelelor neuronale care pot fi paralelizate utilizând implementarea simplă a strategiei de partiționare a datelor. Există însă diferite moduri de evitare a unui asemenea necesar de memorie. Unul ar fi să se încerce reducerea redundanței cât

<sup>3</sup>Cel puțin în principiu; trebuie ținut însă cont de faptul că transmisia se face printr-o rețea Ethernet.

<sup>4</sup>Au fost incluse doar variabile influențate de dimensiunea rețelei.

mai mult posibil; altă soluție ar fi să nu se transmită un întreg gradient component dintr-odată. Ambele strategii vor fi utilizate într-un algoritm prezentat în secțiunea 4.2.

## 4.2 O implementare avansată a strategiei de partiționare a datelor

Pornind de la o scurtă descriere a unui algoritm similar construit de Pomerleau și colegii [68] pentru o mașină Warp, această teză propune o soluție originală, construită cu scopul minimizării necesarului de memorie și al eliminării datelor redundante.

Se prezintă în această secțiune o implementare mai puțin costisitoare din punctul de vedere al memoriei, a strategiei de partiționare a datelor. Procesoarele individuale nu mai păstrează propria copie a întregii rețele. Ponderile sunt memorate într-un singur procesor (numit *administrator*) și circulă, una câte una, între procesoare, în funcție de necesități.

Administratorul nu procesează tipare de antrenare și nu le prezintă el însuși rețelei neuronale. Această sarcină este efectuată exclusiv de către restul procesoarelor (numite *slave*), cărora li se distribuie în mod uniform tiparele de antrenare. Administratorul controlează circulația ponderilor și efectuează actualizarea lor<sup>5</sup>.

Întrucât fiecare procesor *slave* nu păstrează permanent nici o pondere, este esențial ca fiecare pondere recepționată să fie utilizată în cât mai multe calcule posibil, pentru a evita astfel transmisiile inutile de ponderi. Pe parcursul pasului înainte, aceasta înseamnă că la transmisia de ponderi care conectează un strat cu următorul, aceste ponderi trebuie utilizate în propagarea activității corespunzătoare pentru cât mai multe tipare posibil, adică procesorul  $q$  va utiliza ponderile pentru propagarea activității tuturor tiparelor din porțiunea sa  $B_q$  din setul curent de tipare. Pentru aceasta, fiecare procesor *slave*  $q$  trebuie să păstreze atâtea copii ale tuturor unităților din rețea câte tipare sunt în  $B_q$ .

Ca rezultat al considerațiilor de mai sus, toate tiparele din  $B_q$  sunt propagate înainte de către procesorul  $q$  înainte ca valorile de eroare generate de către aceste tipare să fie propagate înapoi în rețea. Aceasta contravine algoritmului descris în secțiunea 4.1, în care un tipar era propagat înainte și înapoi, și se calcula un singur gradient component, înainte de prezentarea următorului tipar.

### 4.2.1 Topologia procesoarelor

Toate procesoarele slave trebuie să recepționeze fiecare pondere din rețea o singură dată pe parcursul pasului înainte, începând cu ponderile dintre unitățile de pe stratul de intrare și cel ascuns. La prima vedere, o configurație de inel poate să nu pară ideală, din moment ce numărul de pași de comunicare necesari pentru transmiterea unei ponderi de la administrator la câțiva slave poate fi destul de mare. Totuși, toate procesoarele slave necesită aceleași ponderi, însă nu este neapărat necesar să utilizeze aceleași ponderi în același timp. Ca urmare, ponderile pot fi transmise una câte una prin inel, în stil *pipe*. În acest fel, nici un procesor nu va avea de așteptat pentru nici o pondere, cu excepția fazelor de pornire și de oprire, care constituie o mică parte din timpul total de execuție, întrucât numărul de ponderi este de obicei mult mai mare decât numărul de procesoare.

### 4.2.2 Comunicația

După cum s-a menționat în secțiunea anterioară, ponderile rețelei neuronale sunt circulate una câte una prin inelul de procesoare pe parcursul activității de propagare.

<sup>5</sup>În realitate, se va vedea că și administratorul va procesa câteva tipare, mai puține decât procesoarele *slave*. Totuși, pentru moment se va presupune că administratorul nu prezintă tipare rețelei.

Aceasta înseamnă că pe parcursul pasului înainte fiecare procesor trebuie să recepționeze și să transmită fiecare pondere din întreaga rețea, ceea ce înseamnă că fiecare procesor trebuie să participe într-un număr dublu de comunicații față de numărul de ponderi din rețea. La fiecare pondere comunicată avem asociat un anumit volum de calcule. Utilizând un asemenea volum mare de comunicații, este important să construim algoritmul astfel încât să asigure efectuarea concurentă a calculelor și a transmisiei de ponderi. Mai mult, dacă timpul necesar calculelor este mult mai mare decât timpul de comunicație, transmisia ponderilor poate avea loc fără încetinirea semnificativă a calculelor efectuate de *slave*.

De aceea, se recomandă o transmisie asincronă a acestor ponderi.

### 4.2.3 Pasul înainte

Pe parcursul propagării activității ponderile sunt circulate una câte una prin inelul de procesoare de către administrator, așa cum este arătat în figura 4.3.

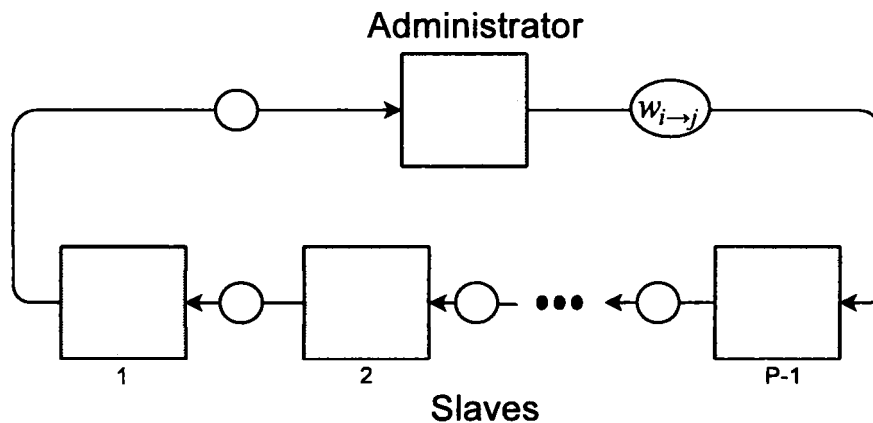


Figura 4.3: Circulația ponderilor în pasul înainte.

Inițial, sunt circulate ponderile dintre stratul de intrare și cel ascuns, permițând astfel propagarea activității de la unitățile de intrare către cele ascunse. Apoi se transmit ponderile dintre unitățile ascunse și cele de ieșire și se calculează tiparul de activitate pe unitățile de ieșire. Acești doi pași de propagare înainte a activității se efectuează în același mod, singura diferență constând în ponderile care se transmit și unitățile utilizate în calcule. În continuare vom descrie în detaliu numai propagarea activității de la unitățile de intrare către cele ascunse.

De fiecare dată când un procesor  $q$  recepționează o pondere  $w_{i \rightarrow j}^H$ , care conectează unitatea  $i$  din stratul de intrare cu unitatea  $j$  din stratul ascuns, transmite în continuare ponderea către următorul procesor din inel (procesorul  $q - 1$ ). Concurrent, această pondere este utilizată pentru calculul contribuției unității  $i$  la intrarea totală a unității  $j$  pentru toate tiparele din porțiunea  $B_q$  a acestui procesor din întregul set  $B$ . Cu alte cuvinte, pentru fiecare tipar  $p$  din  $B_q$ , la recepționarea lui  $w_{i \rightarrow j}^H$  se calculează a  $i$ -a parte din suma:

$$\sum_{i=0}^{NI-1} a_{pi}^I w_{i \rightarrow j}^H = net_{pj}^H, \quad \forall p \in B_q \quad (4.6)$$

După ce toate ponderile care conectează cele două straturi au trecut pe la toate procesoarele, fiecare procesor  $q$  conține pentru toate unitățile din stratul ascuns cele



$|B_q|$  intrări generate de tiparele din  $B_q$ . Apoi fiecare procesor  $q$  calculează activarea fiecărei unități, pentru fiecare tipar  $p$  din  $B_q$ :

$$a_{pj}^H = f(\text{net}_{pj}^H), \quad \forall p \in B_q, \quad \forall j \in H \quad (4.7)$$

Odată calculată activitatea unităților ascunse, este posibil să se determine tiparele activării pentru unitățile de ieșire. Aceasta se efectuează exact în același mod ca mai sus, cu diferența că de data aceasta administratorul transmite ponderile dintre unitățile ascunse și cele de ieșire.

#### 4.2.4 Pasul înapoi

Propagarea înapoi a valorilor delta și calculul gradientilor componenți este realizată în trei pași.

##### 4.2.4.1 Pasul întâi

Inițial, trebuie calculate valorile delta pentru toate unitățile de ieșire, pentru fiecare din tiparele din setul corespunzător fiecărui procesor. Pentru toate tiparele  $p$  din  $B_q$ , procesorul  $q$  calculează valoarea delta pentru fiecare unitate de ieșire:

$$\delta_{pk}^O = a_{pk}^O(1 - a_{pk}^O)(t_{pk} - a_{pk}^O), \quad \forall p \in B_q, \quad \forall k \in O \quad (4.8)$$

După cum se observă, acest pas nu necesită circularea de ponderi și poate fi efectuat independent, de către fiecare dintre procesoare.

##### 4.2.4.2 Pasul al doilea

Următorul pas constă în propagarea înapoi a valorilor delta către unitățile ascunse. Acest pas necesită circulația din nou a ponderilor dintre unitățile ascunse și cele de ieșire, pentru a calcula contribuția fiecărei unități de ieșire la eroarea din fiecare unitate ascunsă. Aceasta se efectuează într-un mod similar celui descris la pasul înainte.

Atunci când procesorul  $q$  recepționează o pondere  $w_{j \rightarrow k}^O$  conectează unitatea  $j$  din stratul ascuns cu unitatea  $k$  din stratul de ieșire, se transmite această pondere către următorul procesor din inel (procesorul  $q-1$ ). Concurrent, ponderea este utilizată pentru a calcula, pentru toate tiparele din  $B_q$ , contribuția unității de ieșire  $k$  la eroarea din unitatea ascunsă  $j$ . La recepționarea lui  $w_{j \rightarrow k}^O$ , pentru fiecare tipar  $p$  din  $B_q$  se calculează a  $k$ -a parte din următoarea sumă:

$$\sum_{k=0}^{NO-1} \delta_{pk}^O w_{j \rightarrow k}^O = e_{pj}^H, \quad \forall p \in B_q \quad (4.9)$$

După ce toate ponderile conexiunilor dintre unitățile ascunse și cele de ieșire au trecut pe la toate procesoarele, procesorul  $q$  conține, pentru toate unitățile ascunse, cele  $|B_q|$  valori de eroare generate de tiparele din  $B_q$ . Fiecare procesor calculează apoi valorile delta corespunzătoare pentru toate tiparele din  $B_q$ :

$$\delta_{pj}^H = a_{pj}^H(1 - a_{pj}^H)e_{pj}^H, \quad \forall p \in B_q, \quad \forall j \in H \quad (4.10)$$

##### 4.2.4.3 Pasul al treilea

Al treilea pas calculează gradientul rezultat în urma prezentării tuturor tiparelor din  $B$ . Calculul gradientilor componenți parțiali asociați ponderilor conexiunilor dintre unitățile de intrare și cele ascunse se efectuează în modul următor. Pe baza ecuației (prezentată în secțiunea 2.1)

$$\Delta_B \bar{w}(n+1) = -\eta \sum_{p \in B} \frac{\partial E_p}{\partial \bar{w}} + \alpha \Delta_B \bar{w}(n) = -\eta \sum_{q=0}^{P-1} \bar{g}^{<q>} + \alpha \Delta_B \bar{w}(n) \quad (4.11)$$



putem obține modificarea  $\Delta_B w_{i \rightarrow j}^H$  a ponderii  $w_{i \rightarrow j}^H$  corespunzător tiparelor din setul  $B$ :

$$\Delta_B w_{i \rightarrow j}^H(n+1) = -\eta \sum_{p \in B} \frac{\partial E_p}{\partial w_{i \rightarrow j}^H} + \alpha \Delta_B w_{i \rightarrow j}^H(n) = -\eta \sum_{q=0}^{P-1} g_{i \rightarrow j}^{H \langle q \rangle} + \alpha \Delta_B w_{i \rightarrow j}^H(n) \quad (4.12)$$

Cu alte cuvinte, pentru actualizarea ponderii  $w_{i \rightarrow j}^H$  administratorul trebuie să cunoască suma gradientelor componente *parțiali*  $g_{i \rightarrow j}^{H \langle q \rangle}$  de la toate procesoarele *slave*. Un asemenea gradient component parțial poate fi calculat cu:

$$g_{i \rightarrow j}^{H \langle q \rangle} = \sum_{p \in B_q} \frac{\partial E_p}{\partial w_{i \rightarrow j}^H} = \sum_{p \in B_q} \delta_{pj}^H a_{pi}^H \quad (4.13)$$

Calculul fiecărui asemenea gradient component parțial poate fi efectuat de către procesorul respectiv, întrucât deja au fost calculate valorile delta și cele de activare pentru toate tiparele  $p \in B_q$ , pentru toate unitățile din rețea. Însă în loc să se transmită ponderi, se transmit gradientii componente parțiali calculați, după cum se vede în Figura 4.4<sup>6</sup>.

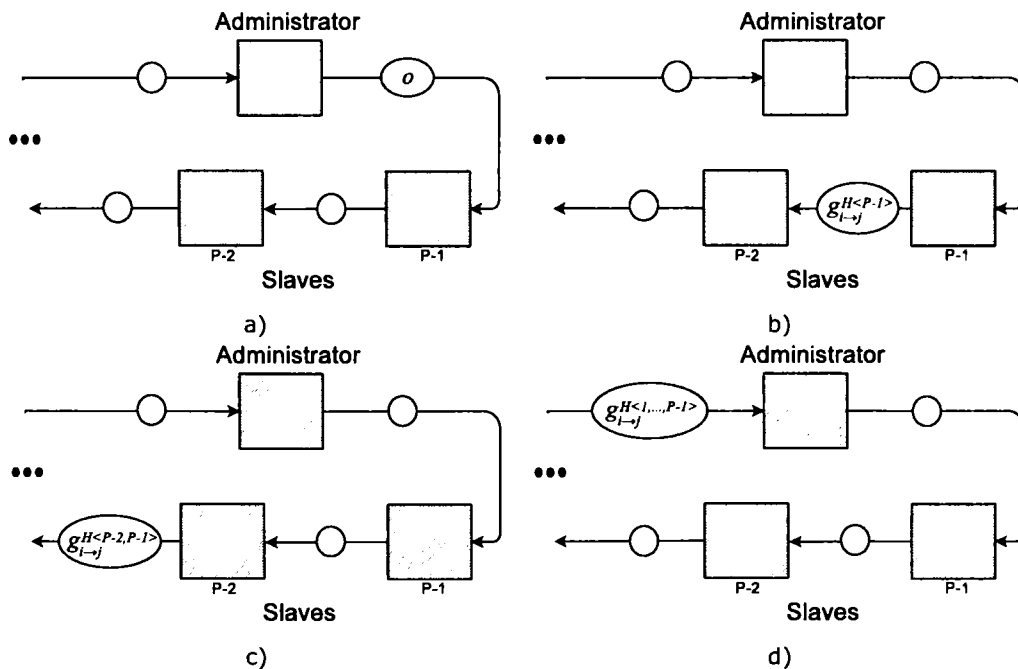


Figura 4.4: Calculul gradientelor parțiali la pasul înapoi.

Inițial, pentru fiecare gradient parțial  $g_{i \rightarrow j}^H$  de calculat, administratorul transmite valoarea 0 procesorului  $P - 1$  (Figura 4.4a). Concurrent, procesorul  $P - 1$  calculează contribuția sa  $g_{i \rightarrow j}^{H \langle P-1 \rangle}$ . Apoi gradientul component parțial al procesorului  $P - 1$  este adunat la 0, iar rezultatul este transmis mai departe procesorului  $P - 2$  (Figura 4.4b).

<sup>6</sup>De fapt, se transmit *negatele* gradientelor componente parțiali. În acest fel se evită negarea gradientului parțial total la calculul modificării ponderii.

În acest timp, procesorul  $P - 2$  calculează  $g_{i \rightarrow j}^{H < P-2 >}$  și, la recepționarea lui  $g_{i \rightarrow j}^{H < P-1 >}$  de la procesorul  $P - 1$ , le însumează și transmite rezultatul  $g_{i \rightarrow j}^{H < P-2, P-1 >}$  către procesorul  $P - 3$  (Figura 4.4c).

În general, procesorul  $q$  calculează  $g_{i \rightarrow j}^{H < q >}$ , în timp ce recepționează gradientul component parțial  $g_{i \rightarrow j}^{H < q+1, \dots, P-1 >}$  de la procesorul  $q + 1$ . Gradientul component parțial calculat de procesorul  $q$  este adunat gradientului component parțial sosit, iar rezultatul  $g_{i \rightarrow j}^{H < q, \dots, P-1 >}$  este transmis procesorului  $q - 1$  concurrent cu calculul următorului gradient component parțial (corespunzător ponderii care urmează după  $w_{i \rightarrow j}^H$ ).

În acest mod, gradientii componente parțiale sunt însumați de către toate procesoarele *slave*, astfel încât la sfârșit (după  $P$  pași) administratorul recepționează  $g_{i \rightarrow j}^{H < 1, \dots, P-1 >} = g_{i \rightarrow j}^H$ , gradientul parțial rezultat în urma prezentării tuturor tiparelor din  $B$  (Figura 4.4d). Administratorul va utiliza acest gradient parțial pentru calculul actualizării  $\Delta_B w_{i \rightarrow j}^H$ , pentru ponderea  $w_{i \rightarrow j}^H$ .

Gradientii componente parțiale  $g_{i \rightarrow j}^O$  asociați ponderilor conexiunilor dintre unitățile ascunse și cele de ieșire se calculează în același mod.

După ce au fost recepționați toți gradientii parțiale și au fost actualizate toate ponderile, administratorul poate transmite prima pondere dintre unitățile de intrare și cele ascunse, iar prezentarea unui alt set de antrenare poate să înceapă.

#### 4.2.5 O îmbunătățire a eficienței

Administratorul efectuează circulația ponderilor, colectarea gradientilor parțiale și actualizarea ponderilor. Numai ultima dintre aceste sarcini presupune un efort computațional. Ca urmare, pe parcursul pasului înainte administratorul nu va fi ocupat majoritatea timpului, având doar de așteptat transmisia sau recepția următoarei ponderi. Pe parcursul pasului înapoi, administratorul va fi ceva mai ocupat, având de actualizat ponderile pe măsură ce recepționează gradientii parțiale. Totuși, avem de-a face cu o diferență importantă între efortul administratorului și cel al procesoarelor *slave*: efortul computațional în slave este proporțional cu numărul de tipare din porțiunea sa din setul curent de antrenare, în timp ce efortul computațional asociat actualizării ponderilor de către administrator este constant, independent de dimensiunea setului de antrenare.

Considerațiile de mai sus conduc la concluzia că probabil se pierde timpul în administrator atunci când porțiunea din setul de antrenare aferentă unui *slave* este suficient de mare. Soluția este simplă: și administratorul va primi de procesat o "felie" din setul de antrenare.

Dificultatea apare în determinarea dimensiunii acestei porțiuni pe care o va procesa administratorul. Evident, va trebui să proceseze mai puține tipare decât procesoarele slave, altfel nu va mai rămâne timp pentru celelalte activități. Însă procesarea de către administrator a unor tipare va crește eficiența numai dacă aceasta reduce încărcarea procesoarelor *slave*, ceea ce va crește viteza procesoarelor *slave*. Nu ajută la nimic însă dacă se reduce încărcarea pe unele dintre procesoarele *slave*, întrucât în intervalul dintre transmisii este determinat de procesorul care transmite cel mai rar. Aceasta înseamnă că administratorul trebuie să reducă încărcarea pe *slave* într-un asemenea mod încât încărcarea să fie echilibrată uniform pe toate procesoarele *slave*. Numărul de tipare de procesat în administrator trebuie să fie astfel încât procesoarele slave să trateze "felii" la fel de mari din setul de antrenare, pentru ca să nu avem procesoare care să fie în așteptarea altuia să fie gata pentru comunicație.

Pe de altă parte, dacă administratorul are de prelucrat același număr de tipare ca și oricare *slave*, nu va mai avea timp pentru actualizarea ponderilor. Mai rău, va încetini activitatea tuturor procesoarelor *slave*, întrucât acestea vor avea de așteptat până ce administratorul termină de actualizat ponderile. Este obligatoriu deci ca administratorul să prelucreze ceva mai puține tipare decât procesoarele *slave*. Va trebui investigat care este valoarea optimă care trebuie aleasă pentru numărul de tipare de

procesat în administrator. Empiric, s-a ajuns la o formulă de calcul a dimensiunilor seturilor de tipare prelucrate de către administrator, respectiv de către procesele *slave*, formulă prezentată în 4.3.2.

#### 4.2.6 Necesarul de memorie

Unul dintre motivele proiectării acestei variante avansate l-a constituit necesarul excesiv de memorie în cazul algoritmului prezentat mai devreme în secțiunea 4.1. Acolo, fiecare procesor păstra o copie a întregului set de ponderi. Proiectarea acestui algoritm s-a bazat pe încercarea de a reduce necesarul de memorie prin evitarea redundanțelor introduse de necesitatea de a memora  $P$  copii identice ale tuturor ponderilor. În acest algoritm numai administratorul păstrează întregul set de ponderi, în vreme ce toate celelalte procesoare au nevoie doar de capacitatea de a memora două sau trei ponderi. Totuși, fiecare procesor *slave* trebuie să păstreze câteva copii ale tuturor *unităților* din rețea, întrucât prezentarea completă a unui tipar (incluzând propagarea înainte a activității și propagarea înapoi a erorii) nu se poate realiza într-un singur pas. Din fericire, acest necesar suplimentar de memorie nu este nici pe departe atât de mare în comparație cu memoria economisită prin eliminarea setului redundant de ponderi în fiecare procesor.

În cele ce urmează, interesul se va concentra asupra variabilelor care variază fie cu dimensiunea rețelei, fie cu cea a setului de antrenare, și vor fi excluse toate variabilele constante în toate execuțiile. Numărul maxim de tipare din porțiunea unui procesor *slave* este  $\lceil B/(P-1) \rceil$ , întrucât ar putea fi posibil să nu se permită administratorului să participe la procesarea de tipare. Ca urmare, se obține următorul necesar de memorie pentru fiecare procesor *slave* (vezi 4.1.1.1):

$$\lceil B/(P-1) \rceil \cdot 2(NI + NH + NO) \quad (4.14)$$

Se poate spune că necesarul de memorie pentru procesoarele *slave* este independent de numărul de ponderi din rețea, contrar cazului primului algoritm de partiționare a datelor. Acesta este un avantaj, întrucât numărul de ponderi dintr-o rețea crește quadratic cu numărul de unități. După cum era de așteptat, algoritmul descris va fi mai puțin consumator de memorie la execuția pe rețele mari.

### 4.3 Simulatorul DataParSim

A fost realizat un simulator — DataParSim constând din două programe: un program *administrator* și un program *slave*. Algoritmul utilizat este cel descris în secțiunea 4.2 – implementarea avansată, în care administratorul transmite în rețea câte o pondere, eliminându-se astfel necesitatea memorării tuturor ponderilor local, la fiecare procesor. S-a încercat îmbunătățirea eficienței simulatorului prin procesarea unui număr de tipare și pe mașina *administrator*.

Informațiile necesare celor două programe se transmit prin intermediul unui *fișier de configurare*, care precizează parametrii rețelei neuronale și cei ai rețelei de calculatoare în care va evolua simulatorul. Acest fișier, precum și modul în care lucrează *administratorul* și programele *slave* sunt descrise în continuare.

#### 4.3.1 Fișierul de configurare

Cele două programe, *administratorul* și procesele *slave*, primesc informațiile necesare prin intermediul unui fișier de configurare, identificat prin intermediul variabilei de mediu `DATAPARSIM_CFG`, fișier care conține informații despre rețeaua neuronală de simulat și despre stațiile care alcătuiesc rețeaua de calculatoare. O informație suplimentară se transmite programelor *slave* în linia de comandă – numărul procesorului în configurația de inel.

Se dă mai jos un exemplu de fișier de configurație. Valorile sunt precizate câte una pe linie, sub forma

<variabilă>=<valoare>

Liniile goale se ignoră, iar textul de la semnul '#' până la sfârșitul liniei se ignoră (comentariu).

```
# Reteaua neuronală:
eta = 0.1           # rata de învățare
alpha = 0.9        # momentul

NrInput = 50       # neuroni pe stratul de intrare
NrHidden = 50      # neuroni pe stratul intermediar
NrOutput = 50     # neuroni pe stratul de ieșire

NrPatt = 800       # număr total de tipare
PattFile = /work/tipare.dat # fișier tipare de învățare

BatchSize = 800   # dimensiune batch

NrIter = 1000     # număr de iterații
Err = 0.01        # eroare acceptată

# Reteaua de calculatoare:
NrProcs = 12      # număr procesoare în rețea
HostsFile = /work/hosts.dat # lista numelor stațiilor

BatDelta = 103   #
```

Parametrul `NrIter` precizează numărul total de iterații care se efectuează. Dacă el este 0, atunci se ia ca și criteriu de oprire atingerea unei erori totale mai mică decât `Err`.

Parametrul `BatDelta` este un coeficient care exprimă raportul dintre mărirea setului de tipare de antrenare pentru un *slave* și cel pentru *administrator*. Valoarea lui a fost determinată empiric, astfel încât setul de tipare pentru *administrator* să fie cât mai mare posibil, însă nu atât de mare încât procesele *slave* să rămână în așteptare. Pentru exemplul dat, s-a constatat o evoluție convenabilă pentru un număr de tipare pentru *administrator* cu circa 3% mai mic decât pentru *slave*.

### 4.3.2 Programul administrator

După cum este descris în secțiunea 4.2, *administratorul* este singurul proces care păstrează întreaga rețea neuronală, transmitând câte o pondere la procesele *slave*. Programul administrator execută mai mulți pași:

- extrage din fișierul de configurare parametrii necesari rețelei neuronale – număr de neuroni pe fiecare strat, rata de învățare și momentul și criteriile de oprire (număr de iterații, eroarea totală), împreună cu structura rețelei de calculatoare și setul de tipare de antrenare;
- determină dimensiunile seturilor de tipare pentru *administrator* și pentru *slave* și mulțimea de tipare corespunzătoare fiecărui proces:

```
NrSlaves = NrProcs - 1;

AdminBatchMax = (100 * BatchSize) /
                ((BatDelta * NrSlaves) + 100);
SlaveBatch = ((BatchSize - AdminBatchMax) + (NrSlaves - 1)) /
              NrSlaves;
AdminBatch = BatchSize - (NrSlaves * SlaveBatch);
```

```

if (AdminBatch < 0) {
    AdminBatch = 0;
    SlaveBatch = (BatchSize + (NrSlaves - 1)) / NrSlaves;
}

AdminPatts = (NrPatt * AdminBatch) / BatchSize;

```

După același algoritm, se vor calcula acești parametri și în procesele *slave*;

- inițializează ponderile conexiunilor din rețeaua neuronală cu valori mici aleatoare;
- lansează câte un proces *slave* pe fiecare stație;
- deschide canalele de comunicație între procesele *slave* și le transmite aceste informații, astfel încât ele să formeze o configurație în inel;
- lansează procesul de simulare — un ciclu, conținând *NrIter* iterații, sau până se atinge eroarea totală dorită. La fiecare iterație, procesul administrator este mai întâi *transmițător* – trimite câte o pondere la primul *slave* din inel, atât pentru pasul înainte, cât și pentru pasul înapoi. Aceeași pondere este utilizată pentru propagarea câte unui tipar din setul propriu de antrenare.

```

/* Sender: */
/* pas inainte */
for (j = 0; j < NrHidden; j++) {
    for (p = 0; p < BatCnt; p++) {
        hidden_net[p][j] = 0;
    }
    for (i = 0; i < NrInput; i++) {
        /* calcul intrare in unitati hidden */
        weight_output = hidden_weight[i][j];
        Send (weight_output);
        for (p = 0; p < BatCnt; p++) {
            pattern = (p + batch_start) % PattCnt;
            hidden_net [p][j] += input_pattern[pattern][i] *
                                hidden_weight[i][j];
        }
    }
}

for (j = 0; j < NrHidden; j++) {
    /* calcul intrare de la bias si activare in
       unitati hidden */
    weight_output = hidden_weight[NrInput][j];
    Send (weight_output);
    for (p = 0; p < BatCnt; p++) {
        pattern = (p + batch_start) % PattCnt;
        hidden_net[p][j] += hidden_weight[NrInput][j];
        hidden_act[p][j] = calculate_activity (hidden_net[p][j]);
    }
}

for (j = 0; j < NrOutput; j++) {
    for (p = 0; p < BatCnt; p++) {
        output_net[p][j] = 0;
    }
    for (i = 0; i < NrHidden; i++) {
        /* calcul intrari in unitati output */
        weight_output = output_weight[i][j];
        Send (weight_output);
        for (p = 0; p < BatCnt; p++) {

```

```

        output_net [p][j] += hidden_act[p][i] *
                           output_weight[i][j];
    }
}

for (j = 0; j < NrOutput; j++) {
    /* calcul intrare de la bias si activare
       unitati output */
    weight_output = output_weight[NrHidden][j];
    Send (weight_output);
    for (p = 0; p < BatCnt; p++) {
        output_net[p][j] += output_weight[NrHidden][j];
        output_act[p][j] = calculate_activity (output_net[p][j]);
    }
}

Send (zero_value); /* dummy */

/* pas inapoi */
/* transmite ponderi hidden-output */
change_work = 0.0; /* dummy val */
for (j = 0; j < NrOutput; j++) {
    /* calcul eroare la output si actualizare ponderi bias */
    change_output = change_work;
    change_work = 0.0;
    Send (change_output);
    for (p = 0; p < BatCnt; p++) {
        pattern = (p + batch_start) % PattCnt;
        oua = output_act[p][j];
        output_err[p][j] = target_pattern[pattern][j] - oua;
        output_dlt[p][j] = output_err[p][j] * oua * (1.0 - oua);
        change_work += output_dlt[p][j];
    }
}

Send (change_work);

change_work = 0.0; /* dummy val */
for (i = 0; i < NrHidden; i++) {
    for (p = 0; p < BatCnt; p++) {
        hidden_err[p][i] = 0.0;
    }
    for (j = 0; j < NrOutput; j++) {
        /* calcul actualizare ponderi si eroare in hidden */
        weight_output = output_weight[i][j];
        change_output = change_work;
        change_work = 0.0;
        Send (weight_output); Send (change_output);
        for (p = 0; p < BatCnt; p++) {
            hidden_err[p][i] += output_dlt[p][j] *
                               output_weight[i][j];
            change_work += output_dlt[p][j] *
                           hidden_act[p][j];
        }
    }
}

Send (zero_value); Send (change_work);

```

```

/* transmite actualizare initiala pentru ponderi
catrta unitati hidden: */
change_work = 0.0;
for (j = 0; j < NrHidden; j++) {
  /* calcul delta in unitati hidden */
  change_output = change_work;
  change_work = 0.0;
  Send (change_output);
  for (p = 0; p < BatCnt; p++) {
    hua = hidden_act[p][j];
    hidden_dlt[p][j] = hidden_err[p][j] * hua * (1.0 - hua);
    change_work += hidden_dlt[p][j];
  }
  for (i = 0; i < NrInput; i++) {
    /* calcul actualizari */
    change_output = change_work;
    change_work = 0.0;
    Send (change_output);
    for (p = 0; p < BatCnt; p++) {
      pattern = (p + batch_start) % PattCnt;
      change_work += hidden_dlt[p][j] *
        input_pattern[pattern][i];
    }
  }
}
Send (change_work);

```

După cum se observă, *administratorul* calculează gradientii corespunzători setului propriu de tipare, după care îi transmite în rețea, la primul proces din inel. Urmează apoi rolul de *receptor* – primește de la ultimul proces din rețea gradientii acumulați și calculează noile ponderi în rețeaua neuronală:

```

/* Receiver: */
/* pas inainte: */
for (j = 0; j < NrHidden; j++) {
  for (i = 0; i < NrInput; i++) {
    Receive (weight);
  }
}
for (j = 0; j < NrHidden; j++) { /* bias */
  Receive (weight);
}
for (j = 0; j < NrOutput; j++) {
  for (i = 0; i < NrHidden; i++) {
    Receive (weight);
  }
}
for (j = 0; j < NrOutput; j++) { /* bias */
  Receive (weight);
}
Receive (weight); /* dummy */

/* pas inapoi: */
/* calcul noile ponderi hidden-output */
Receive (change); /* dummy */
for (j = 0; j < NrOutput; j++) {
  Receive (change);
  change = (alpha * prev_output_change[NrHidden][j]) +
    (eta * change);
}

```

```

    output_weight[NrHidden][j] += change;
    prev_output_change[NrHidden][j] = change;
}
Receive (weight); Receive (change); /* dummy */
for (i = 0; i < NrHidden; i++) {
    for (j = 0; j < NrOutput; j++) {
        Receive (weight); Receive (change);
        change = (alpha * prev_output_change[i][j]) +
            (eta * change);
        output_weight[i][j] += change;
        prev_output_change[i][j] = change;
    }
}

/* calcul noile ponderi input-hidden */
Receive (change); /* dummy */
for (j = 0; j < NrHidden; j++) {
    Receive (change);
    change = (alpha * prev_hidden_change[NrInput][j]) +
        (eta * change);
    hidden_weight[NrInput][j] += change;
    prev_hidden_change[NrInput][j] = change;
    for (i = 0; i < NrInput; i++) {
        Receive (change);
        change = (alpha * prev_hidden_change[i][j]) +
            (eta * change);
        hidden_weight[i][j] += change;
        prev_hidden_change[i][j] = change;
    }
}
}

```

- măsoară timpul necesar unei iterații și întregului proces;
- salvează noua configurație a rețelei neuronale pentru utilizări ulterioare.

### 4.3.3 Programul *slave*

Procesele *slave* implementează algoritmul backpropagation pentru o singură iterație, conform celor prezentate anterior:

- extrage din fișierul de configurare parametrii necesari rețelei neuronale – număr de neuroni pe fiecare strat, rata de învățare și momentul, precum și structura rețelei de calculatoare și fișierul conținând setul de tipare de antrenare. Din linia de comandă cu care este lansat extrage numărul procesorului în configurația în care este lansat (argumentul `processor_number`);
- determină dimensiunile seturilor de tipare pentru *administrator*, respectiv pentru *slave* și mulțimea de tipare corespunzătoare fiecărui proces, în aceeași manieră ca la *administrator*:

```

NrSlaves = NrProcs - 1;

MaxBatCnt = (BatchSize + (NrSlaves - 1)) / NrSlaves;
MaxPattCnt = (NrPatt + (NrSlaves - 1)) / NrSlaves;

/* calcul batch size pentru acest procesor */
AdminBatchMax = (100 * BatchSize) /
    ((BatDelta * NrSlaves) + 100);
SlaveBatch = ((BatchSize - AdminBatchMax) +
    (NrSlaves - 1)) / NrSlaves;

```



```

AdminBatch = BatchSize - (NrSlaves * SlaveBatch);
if (AdminBatch < 0) {
    AdminBatch = 0;
    SlaveBatch = (BatchSize + (NrSlaves - 1)) / NrSlaves;
}

temp = (BatchSize - AdminBatch) % NrSlaves;
if (temp == 0) {
    BatCnt = SlaveBatch;
}
else if (temp < processor_number) {
    BatCnt = SlaveBatch - 1;
}
else
    BatCnt = SlaveBatch;

/* calcul numar tipare pt acest procesor */
AdminPatts = (NrPatt * AdminBatch) / BatchSize;
SlavePatts = NrPatt - AdminPatts;

PattCnt = (SlavePatts + (NrSlaves - 1)) / NrSlaves;
ptemp = SlavePatts % NrSlaves;
if ((ptemp != 0) && (ptemp < processor_number))
    PattCnt--;

```

- încarcă setul de tipare de antrenare pe care le va procesa;
- se conectează la canalele de comunicație corespunzătoare – este vorba despre procesul precedent și cel următor;
- execută un ciclu infinit, în care propagă prin rețea tiparele din setul propriu și calculează actualizările de ponderi; după cum se va vedea mai jos, aceste două operații presupun recepția de la procesul precedent și transmisia la procesul următor din inel:

```

propagate_activity (batch_start);
calculate_weight_changes (batch_start);
batch_start += BatCnt; batch_start %= PattCnt;

```

Rutina `propagate_activity` primește ca parametru poziția de start în lista de tipare și trece prin rețeaua neuronală `BatCnt` tipare. Calculele se efectuează recepționând câte o pondere de la precedentul proces în inel și trimițând-o mai departe la următorul proces:

```

void propagate_activity (int batch_start)
{
    double weight_input, weight_work, weight_output;
    int pattern;

    Receive (weight_input);

    for (j = 0; j < NrHidden; j++) {
        for (p = 0; p < BatCnt; p++) {
            hidden_net[p][j] = 0.0;
        }
        for (i = 0; i < NrInput; i++) {
            /* calcul intrari in hidden */
            weight_work = weight_input;
            weight_output = weight_input;
            Send (weight_output);
            Receive (weight_input);
        }
    }
}

```

```

    for (p = 0; p < BatCnt; p++) {
        pattern = (p + batch_start) % PattCnt;
        hidden_net[p][j] +=
            input_pattern[pattern][i] * weight_work;
    }
}
}
for (j = 0; j < NrHidden; j++) {
    /* calcul intrare de la bias in hidden si activare */
    weight_work = weight_input;
    weight_output = weight_input;
    Send (weight_output);
    Receive (weight_input);
    for (p = 0; p < BatCnt; p++) {
        hidden_net[p][j] += weight_work;
        hidden_act[p][j] = calculate_activity (hidden_net[p][j]);
    }
}

for (j = 0; j < NrOutput; j++) {
    for (p = 0; p < BatCnt; p++) {
        output_net[p][j] = 0.0;
    }
    for (i = 0; i < NrHidden; i++) {
        /* calcul intrari in output */
        weight_work = weight_input;
        weight_output = weight_input;
        Send (weight_output);
        Receive (weight_input);
        for (p = 0; p < BatCnt; p++) {
            output_net[p][j] += hidden_act[p][i] * weight_work;
        }
    }
}
for (j = 0; j < OUTPT_UNITS; j++) {
    /* calcul intrare de la bias in output si activare */
    weight_work = weight_input;
    weight_output = weight_input;
    Send (weight_output);
    Receive (weight_input);
    for (p = 0; p < BatCnt; p++) {
        output_net[p][j] += weight_work;
        output_act[p][j] = calculate_activity (output_net[p][j]);
    }
}

Send (weight_input); /* dummy */
} /* propagate_activity() */

```

Rutina `calculate_weight_changes` însumează la gradientul ce a fost recepționat de la procesul precedent, propriul gradient, și-l transmite mai departe, la procesul următor în inel:

```

void calculate_weight_changes (int batch_start)
{
    double weight_input, weight_work, weight_output;
    double change_input, change_work, change_output;
    double oua, hua; /* tmp */
    int pattern;

```

```

Receive (change_input); /* dummy */
change_work = 0.0; /* dummy */
for (j = 0; j < NrOutput; j++) {
    /* calcul eroare in output si actualizare ponderi bias */
    change_output = change_work + change_input;
    Send (change_output);
    Receive (change_input);
    change_work = 0.0;
    for (p = 0; p < BatCnt; p++) {
        pattern = (p + batch_start) % PattCnt;
        oua = output_act[p][j];
        output_err[p][j] = target_pattern[pattern][j] - oua;
        output_dlt[p][j] = output_err[p][j] * oua * (1.0 - oua);
        change_work += output_dlt[p][j];
    }
    change_output = change_work + change_input;
    Send (change_output);
    Receive (weight_input); Receive (change_input);
    change_work = 0.0; /* dummy */
    for (i = 0; i < NrHidden; i++) {
        for (p = 0; p < BatCnt; p++) {
            hidden_err[p][i] = 0.0;
        }
        for (j = 0; j < NrOutput; j++) {
            /* calcul actualizari ponderi si eroare in hidden */
            weight_work = weight_input;
            weight_output = weight_input;
            change_output = change_work + change_input;
            Send (weight_output); Send (change_output);
            Receive (weight_input); Receive (change_input);
            change_work = 0.0;
            for (p = 0; p < BatCnt; p++) {
                hidden_err[p][i] += output_dlt[p][j] * weight_work;
                change_work += output_dlt[p][j] * hidden_act[p][i];
            }
        }
    }
    weight_output = weight_input;
    change_output = change_work + change_input;
    Send (weight_output); Send (change_output);
    Receive (change_input); /* dummy */
    change_work = 0.0; /* dummy */
    for (j = 0; j < NrHidden; j++) {
        /* calcul delta in hidden */
        change_output = change_work + change_input;
        Send (change_output);
        Receive (change_input);
        change_work = 0.0;
        for (p = 0; p < BatCnt; p++) {
            hua = hidden_act[p][j];
            hidden_dlt[p][j] = hidden_err[p][j] * hua * (1.0 - hua);
            change_work += hidden_dlt[p][j];
        }
    }

    for (i = 0; i < NrInput; i++) {
        /* calcul actualizari ponderi */
        change_output = change_work + change_input;
        Send (change_output);
    }

```

```

Receive (change_input);
change_work = 0.0;
for (p = 0; p < BatCnt; p++) {
    pattern = (p + batch_start) % PattCnt;
    change_work +=
        hidden_dlt[p][j] * input_pattern[pattern][i];
}
}
}
Send (change_output);
} /* calculate_weight_changes() */

```

În acest fel, la sfârșitul unei iterații, *administratorul* va deține gradientul total, va fi în măsură să calculeze valorile actualizărilor ponderilor din rețeaua neuronală și va putea relua procesul de învățare cu noile valori ale ponderilor.

## 4.4 Rezultate experimentale

S-a exprimat eficiența ca raportul dintre accelerarea obținută și cea ideală. Trebuie subliniat aici faptul că *între* actualizările de ponderi eficiența *trebuie* să fie 100%, întrucât în acest timp toate procesoarele lucrează independent unul față de altul, fiecare procesor evoluând exact în același mod ca și în varianta secvențială. Colectarea tuturor gradientilor componenți și distribuirea ponderilor actualizate constituie sarcini suplimentare în comparație cu algoritmul secvențial și deci va fi nevoie de timp suplimentar. Mai mult, calculul noilor valori ale ponderilor *după* ce toți gradientii componenți au fost colectați nu a fost deloc paralelizat, întrucât doar procesorul rădăcină (0) efectuează acest calcul. Ineficiența generată de acest calcul secvențial al ponderilor crește cu numărul procesoarelor: cu cât mai multe procesoare sunt utilizate în execuția algoritmului, cu atât mai multe procesoare vor fi inactive pe parcursul actualizării. Desigur, aceasta este o proprietate nefericită a oricărui algoritm paralel, mai ales atunci când în realitate calculul noilor valori ale ponderilor nu constituie o fracțiune prea mare din totalul volumului de calcul.

Pe scurt: Întrucât între actualizările de ponderi algoritmul se execută pe fiecare procesor exact în același mod ca și la algoritmul secvențial, se prevede ca reducerea de eficiență care ar fi observată la algoritmul paralel să fie determinată în primul rând de circumstanțele legate de actualizarea ponderilor.

În continuare sunt prezentate cele trei platforme hardware utilizate în acest studiu pentru implementarea algoritmilor distribuiți de antrenare pentru backpropagation. În plus, vor fi discutate condițiile în care a fost efectuate experimentele. În continuare vor fi prezentate rezultatele experimentelor efectuate.

Cele trei platforme hardware pe care s-au efectuat testele prezintă sisteme de operare de tip UNIX, pe care s-a utilizat limbajul de programare C, mediile PVM și MPI, și debugger-ul GNU (pe perioada dezvoltării aplicațiilor). În vederea ridicării rezultatelor experimentale, s-a renunțat la facilitățile de depanare.

Pentru testare, simulatorul a fost executat pe configurații de rețele neuronale generice, cu numărul precizat de neuroni pe fiecare strat, și cu câte un set de tipare de antrenare generate aleator. Nu s-au urmărit performanțele de învățare din punctul de vedere al *calității învățării*, ci doar accelerarea, respectiv eficiența față de varianta corespunzătoare rulată secvențial.

### 4.4.1 Platforme hardware

După cum se poate observa, procesul *administrator* și cele *slave* ar putea să evolueze pe stații de lucru diferite, chiar cu sisteme de operare diferite, atât timp cât se asigură compatibilitatea la nivelul informației transmise. Pentru simplificare, s-a preferat însă omogenitatea rețelei de calculatoare pe care se efectuează simularea.

#### 4.4.1.1 IBM RS/6000

Primul set de stații de lucru pe care s-au efectuat experimentele din acest studiu constă în 12 stații de lucru IBM RS/6000 Model 230, produse de firma IBM, având următoarele specificații:

- procesor POWER Single Chip la 45 MHz
- cache de 128 KB L2
- memorie 64 MB standard PS/2
- disc intern de 1 GB
- adaptor grafic POWER Gt1x
- Ethernet controller 10 BaseT
- sistem de operare AIX 4.1.5

Prin comparație cu celelalte seturi de stații de test, aceste stații IBM RS/6000 prezintă performanțe modeste, atât în privința procesorului, cât și a memoriei și rețelei de interconectare. De aceea, ele au fost utilizate doar în faza de dezvoltare a simulatorului.

#### 4.4.1.2 SunBlade 150

Al doilea set de stații de lucru constă în 16 stații de lucru SunBlade 150, produse de firma Sun, având următoarele specificații:

- procesor UltraSPARC IIi pe 64 de biți, la 650 MHz,
- cache de 512 KB L2 on-die,
- memorie de 1 GB,
- disc intern de 40 GB la 7200 rpm,
- placă de rețea Ethernet/Fast Ethernet, twisted-pair, 100 Base-T,
- sistem de operare Solaris 8.

După cum se observă, este vorba de procesoare mult mai rapide, și, de asemenea, rețeaua în care sunt conectate asigură o viteză mai mare de transmisie, aspect foarte important în algoritmul nostru.

#### 4.4.1.3 Linux

Al treilea set de stații de lucru constă în 16 stații de lucru HP/Compaq, având următoarele specificații:

- procesor Pentium la 2.8 GHz,
- memorie de 1 GB,
- disc intern SATA 160 GB,
- placă de rețea Ethernet/Fast Ethernet, twisted-pair, 100 Base-T,
- sistem de operare Linux.

Și de data aceasta, este vorba de procesoare mai rapide, și, de asemenea, rețeaua în care sunt conectate asigură o viteză mai mare de transmisie, aspect foarte important în algoritmul nostru.

Întrucât performanțele stațiilor IBM RS6000 erau mai slabe decât în celelalte două cazuri, aceste stații au fost utilizate doar în faza de dezvoltare și testare a simulatorului, iar rezultatele experimentale nu au fost incluse în această lucrare, un alt motiv fiind faptul că au fost disponibile doar 12 stații.

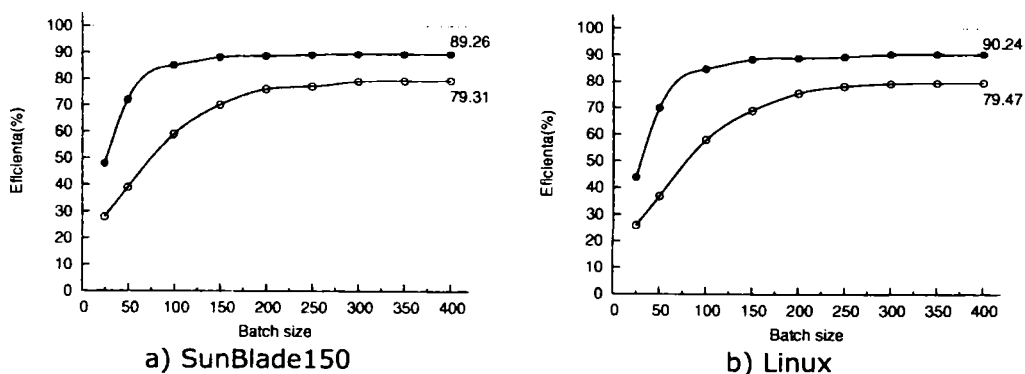
#### 4.4.2 Influența dimensiunii variabile a setului de antrenare

Considerațiile de la începutul acestui capitol sugerează faptul că eficiența va fi scăzută la utilizarea unui set de antrenare de dimensiune mică, și că eficiența va crește dacă se mărește dimensiunea setului de antrenare.

În tabelul 4.1 și în figura 4.5 (a) SunBlade150, b) Linux) sunt prezentate rezultatele pentru o rețea neuronală 50-50-50 (50 de unități pe stratul de intrare, 50 pe cel ascuns, respectiv 50 de unități pe stratul de ieșire), cu varierea dimensiunii setului de antrenare, pe numărul maxim de stații de lucru în fiecare caz (16). Simulatorul de rețea neuronală a fost rulat pentru 10 cicluri de învățare, ceea ce înseamnă că ponderile au fost actualizate de același număr de ori. De remarcat faptul că, întrucât numărul de procesoare și dimensiunea rețelei sunt fixe, timpul necesar efectuării unei actualizări de ponderi este de așteptat să rămână același la toate execuțiile, independent de dimensiunea setului de antrenare utilizat.

Dimensiune batch	SunBlade150		Linux	
	normal	0comm	normal	0comm
25	28.06	48.04	26.02	44.03
50	38.98	72.11	37.06	70.01
100	59.02	85.07	58.00	84.50
150	70.04	88.03	69.11	88.14
200	75.98	88.52	75.51	88.65
250	77.01	89.04	78.03	89.22
300	79.03	89.25	79.13	90.23
350	79.19	89.23	79.34	90.23
400	79.31	89.26	79.47	90.24

**Tabela 4.1:** Eficiența funcție de dimensiunea setului de antrenare



**Figura 4.5:** Eficiența funcție de dimensiunea setului de antrenare

S-a notat cu cerculețe goale rezultatele execuțiilor normale ale algoritmului, adică atunci când procesoarele comunică între ele, cum este normal. Cerculețele pline marchează rezultatele aceluiași execuții, însă fără a avea comunicație. De fiecare dată când un procesor ar dori să comunice, el a executat de fapt o instrucțiune NOP. Observând diferențele dintre cele două grafice, se poate estima măsura în care comunicarea afectează eficiența globală.

După cum era de așteptat, eficiența algoritmului este foarte scăzută pentru dimensiunea cea mai mică a setului de învățare, întrucât fiecare procesor prezintă în acest caz doar un singur tipar, după care se face actualizarea rețelei. Cu alte cuvinte,

ponderile sunt actualizate relativ des. La creșterea dimensiunii setului de învățare timpul de execuție devine tot mai dominat de timpul utilizat de prezentarea tiparelor și calculul gradientilor componenți, întrucât numărul de actualizări de ponderi nu crește. În acest fel, timpul utilizat pentru actualizarea ponderilor devine o porțiune tot mai mică din timpul total de execuție, deci eficiența crește.

Cel de-al doilea grafic, cu ceruculețe pline, confirmă faptul că timpul necesar comunicației pe parcursul fiecărei actualizări de ponderi, este în parte responsabil de eficiența observată, sub cea optimă. Mai mult, timpii totali de execuție arată faptul că timpul necesar comunicării este independent de dimensiunea setului de antrenare, după cum era de așteptat.

Totuși, graficele indică faptul că încă mai pot exista și alte cauze ale ineficienței. Unele dintre acestea ar putea fi generate de instrucțiunile suplimentare pentru determinarea succesorilor în inel, a însumării gradientilor componenți, cât și pentru gestionarea comunicării.

#### 4.4.3 Influența numărului variabil de procesoare

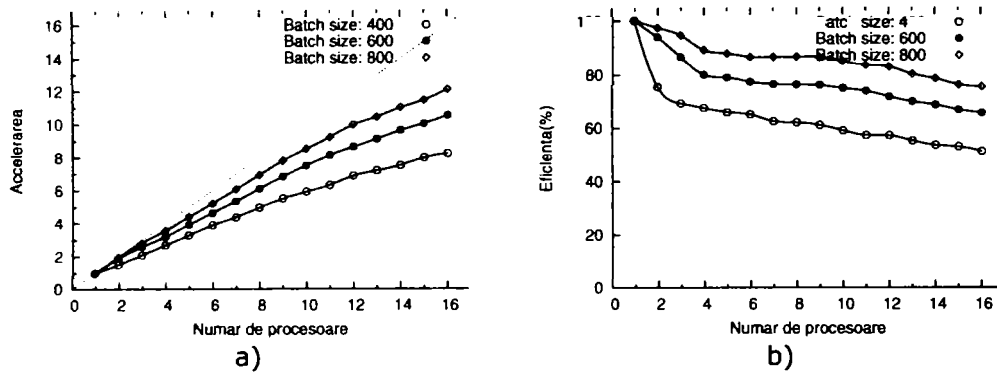
În continuare vor fi prezentate rezultatele execuției pe aceleași stații, când problema este fixă și se mărește numărul de procesoare. Întrucât mărirea setului de antrenare a dus la creșterea eficienței, este de așteptat să poată fi utilizate mai multe procesoare pentru seturi mari de antrenare. Experimentele au fost efectuate utilizând trei seturi de antrenare de dimensiune *fixă*, cu o rețea 50 – 50 – 50.

În tabelul tabelul 4.2 și figura 4.6 (SunBlade150) și respectiv tabelul 4.3 și figura 4.7 (Linux) este prezentată relația dintre accelerare și numărul de procesoare care execută algoritmul paralel. Se observă că pentru toate cele trei dimensiuni ale setului de antrenare (cu excepția, poate, a celui mai mic) accelerarea nu a atins valoarea maximă posibilă chiar și atunci când a fost utilizat numărul maxim de procesoare (16). Aceasta ar însemna, probabil, că s-ar putea obține chiar și o accelerare mai bună având la dispoziție mai multe stații.

Număr procesoare	Accelerarea			Eficiența		
	400	600	800	400	600	800
1	1.00	1.00	1.00	100.00	100.00	100.00
2	1.51	1.88	1.95	75.50	94.00	97.50
3	2.08	2.60	2.84	69.33	86.67	94.67
4	2.70	3.20	3.57	67.65	80.00	89.25
5	3.30	3.95	4.40	66.00	79.00	88.00
6	3.91	4.64	5.20	65.17	77.33	86.67
7	4.38	5.35	6.07	62.57	76.42	86.65
8	4.97	6.10	6.93	62.12	76.25	86.63
9	5.50	6.84	7.79	61.11	76.07	86.60
10	5.91	7.49	8.50	59.10	74.90	85.00
11	6.30	8.12	9.19	57.27	73.82	83.55
12	6.87	8.60	9.95	57.25	71.67	82.92
13	7.18	9.08	10.41	55.23	69.85	80.08
14	7.50	9.61	10.98	53.57	68.64	78.43
15	7.95	10.01	11.42	53.00	66.73	76.13
16	8.20	10.50	12.07	51.25	65.63	75.43

**Tabela 4.2:** Accelerarea – trei dimensiuni ale setului de antrenare: SunBlade

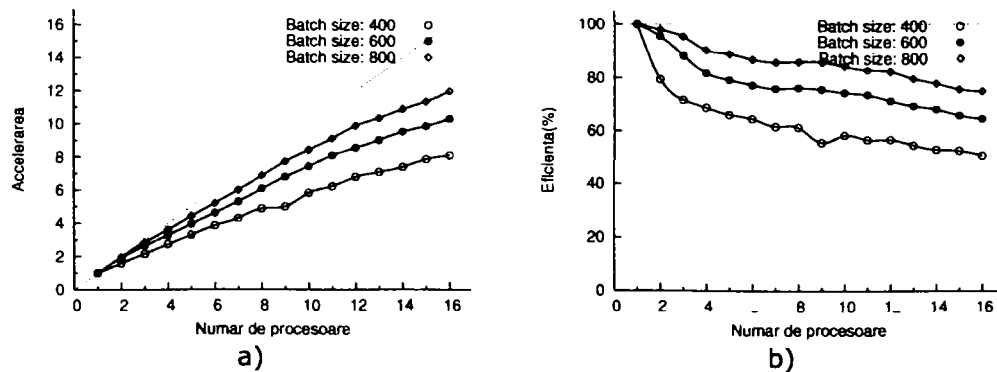
Se mai poate observa, însă, că pe măsură ce se utilizează tot mai multe procesoare, eficiența scade (figura 4.6b, 4.7b). Pe baza rezultatelor obținute în secțiunea anterioară, se poate conchide că există o explicație: pe măsură ce crește numărul de procesoare, fiecare procesor are de prelucrat o porțiune tot mai mică din setul de antrenare. Cu alte cuvinte, setul de antrenare per procesor scade, determinând



**Figura 4.6:** Accelerarea – trei dimensiuni ale setului de antrenare: SunBlade

Număr procesoare	Accelerarea			Eficiența		
	400	600	800	400	600	800
1	1	1	1	100	100	100
2	1.59	1.91	1.96	79.50	95.50	98.00
3	2.15	2.65	2.86	71.66	88.33	95.33
4	2.75	3.27	3.61	68.75	81.75	90.25
5	3.3	3.96	4.44	66.00	79.20	88.80
6	3.87	4.63	5.21	64.50	77.16	86.83
7	4.31	5.31	6.01	61.57	75.85	85.85
8	4.90	6.09	6.88	61.25	76.12	86.00
9	5.0	6.80	7.73	55.55	75.55	85.88
10	5.83	7.43	8.43	58.30	74.30	84.30
11	6.22	8.09	9.11	56.54	73.54	82.81
12	6.8	8.55	9.88	56.66	71.25	82.33
13	7.09	9.02	10.35	54.53	69.38	79.61
14	7.41	9.54	10.90	52.92	68.14	77.85
15	7.88	9.87	11.35	52.53	65.80	75.66
16	8.11	10.32	11.99	50.68	64.50	74.94

**Tabela 4.3:** Accelerarea – trei dimensiuni ale setului de antrenare: Linux



**Figura 4.7:** Accelerarea – trei dimensiuni ale setului de antrenare: Linux



scăderea eficienței după cum s-a demonstrat în secțiunea anterioară. Acest efect poate fi observat și prin simpla comparație a celor trei grafice din figura 4.6b, 4.7b.

Totuși, eficiența depinde nu numai de dimensiunea setului de antrenare per procesor, ci în aceeași măsură și de numărul de procesoare, deci trebuie să mai existe și alte cauze pentru scăderea eficienței, întrucât aceasta depinde nu numai de dimensiunea setului de antrenare per procesor, ci și de numărul de procesoare. Analizând figura 4.6b, 4.7b, se observă că pentru același set de antrenare per procesor, eficiența scade dacă se mărește numărul de procesoare. Aceasta se datorează probabil timpului necesar colectării gradientilor componenți și distribuirii ponderilor actualizate.

#### 4.4.4 Set de antrenare variabil cu numărul de procesoare

A treia clasă de experimente care s-au efectuat examinează efectele scalării dimensiunii setului de antrenare cu numărul de procesoare utilizate la execuția algoritmului. Investigațiile au fost efectuate cu scopul de a determina dacă o creștere a numărului de sub-probleme (tipare în set) va permite să fie mărit corespunzător numărul de procesoare fără a reduce eficiența.

Rețeaua testată a fost una 50 – 50 – 50, deci cu câte 50 de unități pe fiecare strat. Numărul de tipare din seturile de antrenare – generate aleator – a fost ales astfel încât fiecare procesor să proceseze câte 1, 2, 10, respectiv 20 de tipare.

În tabelul tabelul 4.4 și figura 4.8 (SunBlade) și respectiv tabelul 4.5 și figura 4.9 (Linux) dimensiunea setului de antrenare *per procesor* este fixă (în fiecare dintre cele patru serii de experimente). Aceasta înseamnă că, între actualizările de ponderi, fiecare procesor prezintă un număr fix de tipare, independent de numărul de procesoare.

Număr procesoare	Accelerarea				Eficiența			
	1	2	10	20	1	2	10	20
1	1.00	1.00	1.00	1.00	100.00	100.00	100.00	100.00
2	1.19	1.33	1.80	1.88	59.50	66.50	90.05	94.04
3	1.23	1.62	2.39	2.75	41.00	54.00	79.67	91.67
4	1.50	1.80	3.00	3.52	37.50	45.00	75.00	88.00
5	1.60	2.00	3.69	4.34	32.00	40.00	73.80	86.80
6	1.70	2.23	4.35	5.16	28.33	37.17	72.76	86.00
7	1.80	2.48	5.02	5.98	25.71	35.42	71.71	85.43
8	1.90	2.71	5.72	6.68	23.75	33.88	71.58	83.50
9	2.02	2.93	6.43	7.40	22.44	32.56	71.44	82.22
10	2.11	3.05	7.01	8.11	21.10	30.50	70.10	81.10
11	2.19	3.28	7.55	8.72	19.91	29.82	68.64	79.27
12	2.28	3.51	8.08	9.33	19.00	29.25	67.33	77.75
13	2.39	3.62	8.53	9.95	18.38	27.85	65.62	76.54
14	2.51	3.75	8.98	10.49	17.93	26.79	64.14	74.93
15	2.62	3.88	9.41	10.96	17.47	25.87	62.73	73.07
16	2.70	4.00	9.71	11.41	16.88	25.00	60.69	71.31

**Tabela 4.4:** Set de antrenare fix per procesor: SunBlade

În figurile 4.8b, 4.9b se observă că pentru un număr mare de procesoare, dacă se poate mări setul de antrenare, numărul de procesoare afectate problemei poate fi mărit corespunzător, fără scăderea prea mare a eficienței.

În comparație cu figurile 4.6, 4.7, se poate observa că graficele eficienței din figurile 4.8 și 4.9 sunt foarte similare. De fapt, ele sunt aproape identice ca formă. Aceasta confirmă ceea ce s-a enunțat deja, și anume faptul că acest algoritm paralel pierde din eficiență în principal pe parcursul actualizării ponderilor.

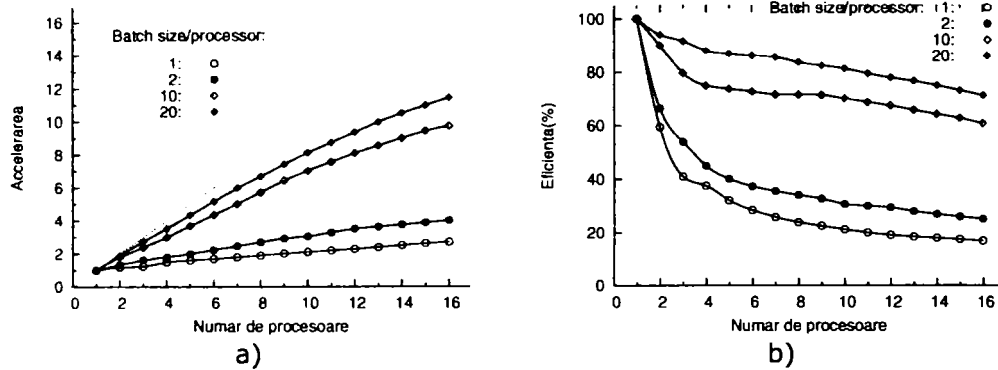


Figura 4.8: Set de antrenare fix per procesor: SunBlade

Număr procesoare	Accelerarea				Eficiența			
	1	2	10	20	1	2	10	20
1	1.00	1.00	1.00	1.00	100.00	100.00	100.00	100.00
2	1.16	1.31	1.78	1.85	58.00	65.50	89.00	92.50
3	1.20	1.59	2.36	2.71	40.00	53.00	78.66	90.33
4	1.47	1.78	2.98	3.49	36.75	44.50	74.50	87.25
5	1.59	1.98	3.67	4.30	31.80	39.60	73.40	86.00
6	1.68	2.20	4.31	5.15	28.00	36.66	71.83	85.83
7	1.78	2.44	4.99	5.94	25.42	34.85	71.28	84.85
8	1.89	2.67	5.68	6.66	23.62	33.37	71.00	83.25
9	2.0	2.90	6.40	7.40	22.22	32.22	71.11	82.22
10	2.10	3.00	6.98	8.08	21.00	30.00	69.80	80.80
11	2.17	3.23	7.49	8.67	19.72	29.36	68.09	78.81
12	2.24	3.48	8.07	9.31	18.66	29.00	67.25	77.58
13	2.34	3.58	8.49	9.90	18.00	27.53	65.30	76.15
14	2.49	3.70	8.95	10.42	17.78	26.42	63.92	74.42
15	2.59	3.85	9.38	10.90	17.26	25.66	62.53	72.66
16	2.68	3.97	9.67	11.37	16.75	24.81	60.43	71.06

Tabela 4.5: Set de antrenare fix per procesor: Linux

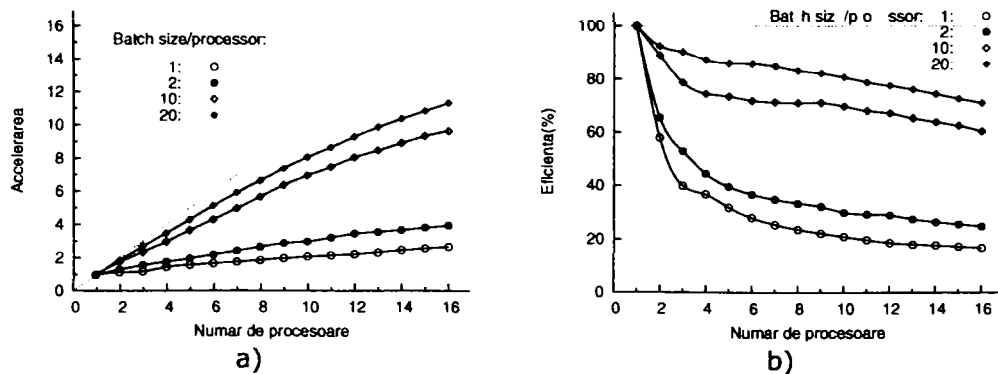


Figura 4.9: Set de antrenare fix per procesor: Linux

Pe baza acestor rezultate se poate decide dacă ar fi util să se mai adauge procesoare în cazul unei probleme date. Cu alte cuvinte, în ce măsură poate fi crescută accelerarea prin adăugarea de procesoare. Crescând numărul de procesoare cu un factor oarecare va determina o creștere a accelerării numai dacă eficiența scade cu un factor mai mic sau cel mult egal, ca rezultat al procesoarelor suplimentare.

Ar fi ideal dacă s-ar putea formula o regulă privitor la cât de multe procesoare pot fi utilizate pentru o problemă specifică. Un exemplu ar fi, să zicem, să se găsească o relație între numărul de procesoare și dimensiunea setului de antrenare. Aceasta ar însemna că atât timp cât fiecare procesor are de prezentat *cel puțin* un anumit număr de tipare, accelerarea va crește la adăugarea de alte procesoare. În general, aceasta ar fi posibil dacă eficiența algoritmului ar depinde numai de dimensiunea setului de antrenare per procesor, și nu ar depinde de însăși numărul de procesoare (în acest caz, graficele eficienței din figurile 4.8b și 4.9b, ar fi fost linii drepte, orizontale). După cum se poate observa, acest algoritm paralel nu prezintă aceste proprietăți și, într-adevăr, se pare că nu se poate formula o regulă generală.

## 4.5 Concluzii

Acest capitol prezintă un simulator pentru antrenarea cu algoritmul backpropagation, care utilizează strategia de paralelizare cu partiționarea datelor. Pornind de la o variantă simplă, în care singura modificare față de varianta secvențială constă în utilizarea și însumarea unor *gradienți parțiali*, a fost descrisă și implementată o variantă avansată, în care s-a urmărit minimizarea necesarului de memorie – oarecum excesiv în varianta simplă – și au fost efectuate o serie de experimente în vederea determinării factorilor care afectează performanțele simulatorului, comparativ cu varianta secvențială.

S-a demonstrat că paralelizarea algoritmului backpropagation este posibilă, într-un mod foarte simplu, utilizând strategia de partiționare a datelor. De asemenea, experimentele au demonstrat că eficiența crește la mărirea dimensiunii setului de tipare de antrenare prezentate rețelei între actualizările de ponderi. Aceasta înseamnă că dacă pot fi utilizate seturi de antrenare foarte mari, este posibil să se obțină o accelerare substanțială utilizând acest tip de paralelizare.

**Contribuțiile** acestui capitol sunt:

1. Dezvoltarea și implementarea unui algoritm de paralelizare cu partiționarea datelor. Un dezavantaj major al acestui algoritm îl constituie necesarul excesiv de memorie, datorat faptului că fiecare procesor trebuie să păstreze o copie a întregii rețele neuronale. O asemenea redundanță ar trebui evitată într-un algoritm paralel, întrucât aceasta înseamnă că paralelizarea nu permite creșterea dimensiunii problemelor care pot fi procesate (din punctul de vedere al necesarului de memorie).
2. Dezvoltarea și implementarea unui algoritm avansat de paralelizare cu partiționarea datelor, cu un proces *administrator* și mai multe procese *slave*. Programul administrator transmite în rețea, către *slave*, pe rând, câte o pondere din rețeaua neuronală, urmând ca procesele *slave* să utilizeze această valoare a ponderii în propagarea tuturor tiparelor asignate. Pentru o îmbunătățire a exploatarei mașinilor, și procesul administrator prelucrează un număr de tipare.
3. Evaluarea necesarului de memorie al algoritmilor prezentați. Algoritmii avansați de paralelizare prezentat, pe lângă faptul că permite scurtarea timpului de antrenare a unei rețele neuronale pentru seturi mari de tipare, mai permite și procesarea unor rețele neuronale de dimensiuni mai mari decât valorile posibile pentru cazul secvențial.
4. Studiul experimental al influenței factorilor care pot afecta performanțele evoluției simulatorului într-o rețea de calculatoare. Concluzia care se poate deduce

este că performanțele simulatorului vor fi cu atât mai bune cu cât numărul de tipare de procesat pe fiecare mașină este mai mare. Din experimentele efectuate, s-au observat valori satisfăcătoare pentru accelerare, respectiv eficiență atunci când au putut fi asigurate cel puțin zece tipare de antrenare per procesor. Investigații pe această temă au fost prezentate în [BaC03] și [BCP07].

În acest capitol se va trata dezvoltarea și analiza unui algoritm care utilizează o strategie de partiționare a rețelei pentru antrenarea unei rețele neuronale într-o rețea de calculatoare.

### 5.1 Construcția algoritmului paralel

În această secțiune vor fi descrise bazele paralelizării, modul în care este divizată rețeaua între procesoare, topologia procesoarelor, iar în final va fi prezentat modul în care se face paralelizarea pasului înainte, respectiv a pasului înapoi.

#### 5.1.1 Divizarea rețelei

Atunci când se dorește să se distribuie unitățile și ponderile pe un număr de procesoare, nu există prea multe alternative la dispoziție. O variantă este divizarea rețelei orizontal, în trei porțiuni, astfel încât un procesor să conțină unitățile unui anumit strat. Aceasta nu e o idee prea bună, întrucât în mod normal o rețea neuronală conține puține straturi și deci nu vor putea fi utilizate prea multe procesoare. Mai mult, dacă se insistă pe utilizarea unei anumite scheme de actualizare a tiparelor, această formă de paralelizare nici nu este posibilă, întrucât pașii înainte și înapoi sunt secvențiali — numai unitățile dintr-un singur strat efectuează calcule la un moment dat.

Rețeaua poate fi divizată vertical, astfel încât unitățile dintr-un strat sunt distribuite uniform pe procesoare. Fiecare procesor va conține deci unități din toate straturile. Sunt posibile două modalități evidente de distribuire a ponderilor:

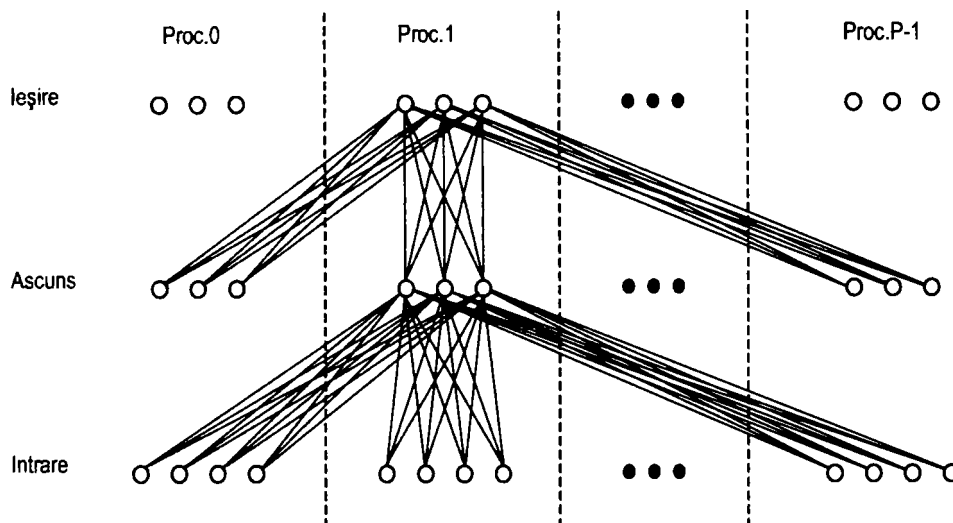
1. Toate ponderile pentru conexiunile care *intră* într-o unitate se memorează în procesorul care simulează acea unitate.
2. Toate ponderile pentru conexiunile care *ies* dintr-o unitate se memorează în procesorul care simulează acea unitate.

În realitate nu există o diferență între cele două strategii. În primul caz pasul înainte va fi simplu de paralelizat, în vreme ce pasul înapoi va fi mai dificil. În al doilea caz — invers. Aceste aspecte vor fi evidente pe parcursul dezvoltării paralelizării. În cadrul lucrării de față a fost aleasă prima strategie.

Toate unitățile sunt distribuite uniform pe procesoare, astfel încât toate procesoarele conțin același număr de unități de intrare, ascunse și de ieșire, deși nu este necesar ca numărul de procesoare să dividă numărul de unități de pe un strat. În acest caz, unele procesoare vor simula cu o unitate mai mult decât altele. Pentru simplificare, se va considera însă în continuare că numărul de unități de pe orice strat este divizibil cu numărul de procesoare.

Ponderile conexiunilor din rețea le distribuim în modul următor: ponderile pentru toate conexiunile care intră într-o unitate de ieșire sau ascunsă se memorează în procesorul care simulează acea unitate. Acest mod conduce la o partiționare echilibrată a ponderilor. În Figura 5.1 se prezintă ponderile conținute în procesorul 1.

Strategia de partiționare a rețelei în vederea paralelizării a fost investigată în cadrul câtorva cercetări în domeniu. Pentru abordări similare, se pot consulta lucrările lui Ernoult – "*Performance of backpropagation on a parallel transputer-based machine*", [22], Zhang, McKenna, Mesirov și Waltz – "*The backpropagation algorithm on grid and hypercube architectures*", [104] și del Rey Millan și Bofill – "*Learning by backpropagation: Computing in a systolic way*", [20]. Spre deosebire de lucrările enumerate, care utilizează arhitecturi de transputere, respectiv hypercuburi, arhitecturi cu caracteristici specifice legate de modul de lucru al procesoarelor și modul în care ele sunt interconectate, teza de față propune utilizarea unei rețele de



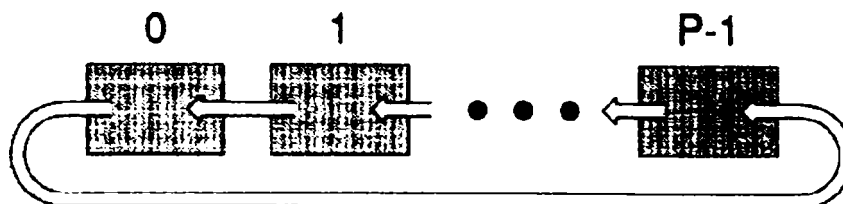
**Figura 5.1:** Distribuirea unităților și a ponderilor.

calculatoare, uzuale într-un laborator universitar. Desigur, interconectarea mașinilor printr-o rețea ethernet a impus diferențe în abordarea diferitelor aspecte, în special în cazul comunicației, în vederea reducerii timpului necesar pentru transmiterea datelor de la un procesor la altul.

### 5.1.2 Topologia procesoarelor

A fost aleasă o conectare în inel. O configurație de inel este simplu de extins pentru a include mai multe elemente și poate să conțină orice număr de procesoare. După cum se va vedea, este o topologie naturală pentru partiționarea rețelei așa cum e prezentată în figura 5.1. Probabil că ar prezenta interes și alte topologii, dar acesta este un subiect pentru dezvoltări ulterioare.

Procesoarele sunt conectate în inel, așa cum este ilustrat în figura 5.2. De notat direcția legăturilor. Un procesor trimite date către predecesorul său și recepționează de la succesor.



**Figura 5.2:** Inel de procesoare.

Chiar dacă fiecare procesor simulează doar o parte din unitățile de intrare și cele ascunse, toate procesoarele trebuie să aibă capacitatea de a stoca activitatea tuturor unităților de intrare, respectiv ascunse, întrucât aceasta este necesară în ambii pași,

cel înainte și cel înapoi, din algoritmul de backpropagation. Memorând activitatea tuturor unităților se evită o serie de comunicații suplimentare.

Contrar paralelizărilor cu partiționarea datelor descrise în capitolul anterior, aici toate procesoarele execută același program. Procesul administrator diferă de cele *s/ave* numai atunci când intervin operații de intrare/ieșire.

### 5.1.3 Notății

Înainte de a începe prezentarea paralelizării, este necesar să fie introduse o serie de notații

Mai întâi, câteva abrevieri convenabile. Mai devreme s-a definit  $P$  ca fiind numărul de procesoare,  $NI$ ,  $NH$  și  $NO$  ca fiind numărul de unități de pe stratul de intrare, ascuns și de ieșire. Se notează cu  $PI$ ,  $PH$ , și  $PO$  — numărul de unități de intrare, ascunse și de ieșire simulate de fiecare procesor. Fie  $O$  mulțimea de indici de unități de ieșire  $(0, 1, \dots, NO - 1)$ . Sub-seturile indicilor de unități de ieșire simulate de către procesoarele individuale vor fi:

$$O = \{ \underbrace{0, 1, \dots, PO - 1}_{O_{[0]}}, \underbrace{PO, \dots, 2 \cdot PO - 1, \dots, (P - 1) \cdot PO}_{O_{[1]}}, \dots, \underbrace{P \cdot PO - 1}_{O_{[P-1]}} \}$$

Indicii pentru unitățile ascunse și cele de intrare se definesc similar.

Se notează cu  $a_j^H$  activarea unității ascunse  $j$ . Vectorul activărilor tuturor unităților ascunse va fi  $\bar{a}^H$ , iar  $\bar{a}_{[p]}^H$  este sub-vectorul activărilor unităților ascunse memorate în procesorul  $p$ . Aceasta se poate scrie ca:

$$\bar{a}^H = ( \underbrace{a_0^H, a_1^H, \dots, a_{P_H-1}^H}_{\bar{a}_{[0]}^H}, \underbrace{a_{P_H}^H, \dots, a_{2 \cdot P_H-1}^H, \dots, a_{(P-1) \cdot P_H}^H}_{\bar{a}_{[1]}^H}, \dots, \underbrace{a_{P \cdot P_H-1}^H}_{\bar{a}_{[P-1]}^H} )$$

Vectorii activărilor unităților de intrare, respectiv de ieșire,  $\bar{a}^I$  și  $\bar{a}^O$  se definesc în mod similar. La fel pentru vectorii valorilor delta pentru unitățile ascunse, respectiv de ieșire,  $\bar{\delta}^I$  și  $\bar{\delta}^O$ . În sfârșit, intrarea de rețea pentru o unitate din stratul ascuns, respectiv de ieșire va fi  $net_j^H$ , respectiv  $net_k^O$ .

De notat faptul că simbolurile pentru intrarea din rețea, activare și delta nu au un indice  $p$  (pentru tipar). Acesta nu este necesar, întrucât algoritmul pe care îl descriem corespunde metodei gradientului stohastic.

Ponderile conexiunilor care intră în stratul de ieșire de la unitatea ascunsă  $j$  la unitatea de ieșire  $k$  vor fi notate cu  $w_{j \rightarrow k}^O$ .  $\bar{w}_{\rightarrow k}^O$  va fi vectorul ponderilor conexiunilor care intră în unitatea de ieșire  $k$  dinspre toate unitățile ascunse, iar  $\bar{w}_{j \rightarrow \rightarrow}$  — vectorul ponderilor conexiunilor spre toate unitățile de ieșire de la unitatea ascunsă  $j$ . La fel ca în cazul activărilor, se notează cu  $\bar{w}_{[p] \rightarrow k}^O$  sub-vectorul de ponderi spre unitatea de ieșire  $k$  dinspre unitățile ascunse simulate de procesorul  $p$ , și cu  $\bar{w}_{j \rightarrow [p]}^O$  sub-vectorul de ponderi care duc spre unitățile de ieșire simulate de către procesorul  $p$  de la unitatea ascunsă  $j$ . Ponderile conexiunilor dinspre stratul de intrare către stratul ascuns se exprimă în mod similar.

### 5.1.4 Paralelizarea

Paralelizarea propusă este formată din șapte pași, doi pași pentru pasul înainte și cinci pași pentru pasul înapoi. Întrucât rețeaua neuronală este distribuită în rețeaua de calculatoare, au fost necesare o serie de modificări în ordinea efectuării calculelor în algoritmul clasic de backpropagation, având ca scop adaptarea la comunicația în rețeaua de calculatoare și minimizarea timpului necesar transmisiei de date. Pentru evidențierea modificărilor efectuate, au fost rescrise formulele de calcul, grupând



termenii pe procesoare la calculul valorilor  $\delta$  pentru stratul ascuns, calcule efectuate pe parcursului pasului al patrulea.

#### 5.1.4.1 Pasul întâi

Primul pas calculează activarea unităților ascunde. De fapt, el constă din doi sub-pași, primul fiind distribuirea activării unităților de intrare, astfel încât toate procesoarele vor avea acces la activarea tuturor unităților de intrare, iar în al doilea sub-pas se efectuează calculul activării unităților ascunde.

Vectorul activării unităților de intrare,  $\vec{a}^I$ , este distribuit astfel încât procesorul  $p$  deține sub-vectorul  $\vec{a}_{[p]}^I$ . Distribuția tuturor sub-vectorilor către toate procesoarele este ușor de efectuat într-un inel: fiecare procesor își transmite propriul sub-vector către predecesorul său în inel. Apoi toate procesoarele recepționează și transmit sub-vectorii calculați de celelalte procesoare, până când dețin întregul vector al activării unităților de intrare. Aceasta necesită  $P(P-1)$  comunicații. Operația este ilustrată în Figura 5.3.

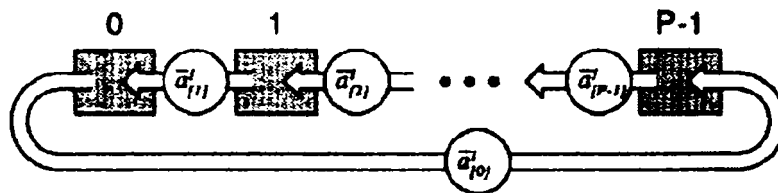


Figura 5.3: Broadcast pentru activarea unităților de intrare.

La terminarea acestei operații, toate procesoarele dețin întregul vector al activării unităților de intrare,  $\vec{a}^I$ . Procesorul  $p$  simulează unitățile ascunde cu indicii  $H_{[p]}$  și ponderile pentru toate conexiunile care intră în aceste unități sunt memorate în procesorul  $p$ . Procesorul  $p$  poate acum să calculeze intrarea de rețea pentru toate unitățile ascunde pe care le simulează:

$$net_j^H = \vec{a}^I * \vec{w}_{i \rightarrow j}^H = \sum_{i=0}^{NI-1} a_i^I w_{i \rightarrow j}^H, \quad \forall j \in H_{[p]} \quad (5.1)$$

unde operatorul  $*$  reprezintă înmulțire vectorială.

Pe baza intrării din rețea, se poate calcula activarea unităților ascunde:

$$a_j^H = f(net_j^H), \quad \forall j \in H_{[p]} \quad (5.2)$$

unde  $f$  este funcția nonliniară de activare. Calculele din cele două ecuații de mai sus pot fi efectuate de către toate procesoarele în paralel, fără comunicație între ele, întrucât toate variabilele care intervin sunt memorate local.

Partea de comunicație din pasul întâi nu poate fi evitată. Totuși, comunicațiile se pot efectua asincron și în paralel cu calculele. Se va exploata această facilitate pentru a combina cei doi sub-pași ai pasului întâi. Dacă se efectuează calculele în timp ce se transmit activările unităților de intrare, iar timpul de calcul este mai mare decât cel de transmisie, volumul mare de comunicații necesar în acest pas nu va avea importanță.

Calculul intrării într-o unitate ascunsă, dat în ecuația 5.1, poate fi rescris în felul următor:

$$net_j^H = \sum_{i=0}^{NI-1} a_i^I w_{i \rightarrow j}^H = \sum_{q=0}^{P-1} \vec{a}_{[q]}^I * \vec{w}_{[q] \rightarrow j}^H, \quad \forall j \in H_{[p]} \quad (5.3)$$



Ultima expresie din ecuație este o sumă pentru toate procesoarele. Fiecare sub-vector  $\bar{a}_{[q]}^I$  din sumă este memorat în câte un procesor.

Combinarea celor doi sub-pași începe prin transmisia de către fiecare procesor a acestui sub-vector, către procesorul precedent din inel, după cum se vede în Figura 5.3. În timpul acestei transmisii, toate procesoarele își calculează contribuția locală la unitățile ascunse pe care le simulează. Pentru toate unitățile ascunse simulate de procesorul  $p$ , se calculează:

$$\bar{a}_{[p]}^I \star \bar{w}_{[p] \rightarrow j}^H, \quad \forall j \in H_{[p]} \quad (5.4)$$

La terminarea acestui calcul, procesorul recepționează sub-vectorul de la succesorul său și îl transmite către procesorul predecesor. Acest sub-vector recepționat de procesorul  $p$  conține activările unităților de intrare ale procesorului  $p + 1$ , iar acum procesorul  $p$  poate să calculeze termenii:

$$\bar{a}_{[p+1]}^I \star \bar{w}_{[p+1] \rightarrow j}^H, \quad \forall j \in H_{[p]} \quad (5.5)$$

care se însumează la termenii calculați în ecuația 5.4.

Această activitate de calcul și transmisie în paralel continuă până la completarea calculului lui  $net_j^H$  pentru toți  $j \in H_{[p]}$ . Se aplică apoi funcția de activare pentru a obține activările unităților ascunse.

Activările unităților de intrare sunt necesare din nou în pasul șase, când se calculează actualizarea ponderilor. Ca urmare, este necesar să se memoreze sub-vectorii activărilor unităților de intrare, pentru utilizare ulterioară.

#### 5.1.4.2 Pasul al doilea

Al doilea pas calculează activarea unităților de ieșire. Acest pas este în mare măsură similar pasului întâi. Calculul intrării într-o unitate de ieșire poate fi rescris după cum urmează:

$$net_k^O = \sum_{j=0}^{NH-1} a_j^H w_{j \rightarrow k}^O = \sum_{q=0}^{P-1} \bar{a}_{[q]}^H \star \bar{w}_{[q] \rightarrow k}^O, \quad \forall k \in O_{[p]} \quad (5.6)$$

Ca și la pasul întâi, fiecare sub-vector  $\bar{a}_{[q]}^H$  din această sumă este memorat în procesoare separate. Este necesar din nou un broadcast, de data aceasta pentru activările unităților ascunse.

#### 5.1.4.3 Pasul al treilea

Al treilea pas calculează valorile delta pentru unitățile de ieșire și aceasta se poate efectua exclusiv pe baza datelor locale. Valoarea delta pentru unitatea de ieșire simulată de procesorul  $p$ , cu indicele  $k$  se calculează pe baza relației 5.7:

$$\delta_k^O = (t_k - a_k^O) a_k^O (1 - a_k^O), \quad \forall k \in O_{[p]} \quad (5.7)$$

#### 5.1.4.4 Pasul al patrulea

În acest pas se calculează valorile delta pentru unitățile ascunse. Pentru calculul valorii delta a unei unități ascunse,  $\delta_j^H$ , în procesorul  $p$  se calculează:

$$\delta_j^H = a_j^H (1 - a_j^H) (\bar{\delta}^O \star \bar{w}_{j \rightarrow}^O) = a_j^H (1 - a_j^H) \sum_{k=0}^{NO-1} \delta_k^O w_{j \rightarrow k}^O, \quad \forall j \in H_{[p]} \quad (5.8)$$

Pasul al patrulea este foarte asemănător cu pasul doi, întrucât în esență constă dintr-un produs vectorial. Totuși, aici ponderile nu sunt memorate la locul potrivit.

Atunci când procesorul  $p$  calculează un produs vectorial el utilizează ponderile  $\bar{w}_{j \rightarrow}$  pentru toți  $j \in H_{[p]}$ . Aceste ponderi, însă, nu sunt memorate în procesorul  $p$ . Ele sunt partiționate în mod egal între toate procesoarele, iar distribuirea ponderilor este mult prea mare consumatoare de timp.

Ecuția 5.8 poate fi rescrisă astfel încât să exprimăm sumele pe procesoare:

$$\delta_j^H = a_j^H (1 - a_j^H) \sum_{k=0}^{NO-1} \delta_k^O w_{j \rightarrow k}^O = a_j^H (1 - a_j^H) \sum_{q=0}^{P-1} \bar{\delta}_{[q]}^O \star w_{j \rightarrow [q]}^O, \quad \forall j \in H_{[p]} \quad (5.9)$$

Ultima expresie este o sumă pe toate procesoarele. Fiecare termen din sumă poate fi calculat numai de către procesorul care memorează ponderile și valorile delta respective. Ca urmare, dacă  $j \in H_{[a]}$  este indicele unei unități ascunse simulată de către procesorul  $a$ , procesorul  $b$  poate calcula termenii  $\bar{\delta}_{[b]}^O \star \bar{w}_{j \rightarrow [b]}^O$  necesari procesorului  $a$  în calculul lui  $\delta_{[a]}^H$ . Se notează cu  $a_{[b]}$  toți acești termeni:

$$a_{[b]} = \left\{ \bar{\delta}_{[b]}^O \star \bar{w}_{j \rightarrow [b]}^O \right\}_{\forall j \in H_{[a]}} \quad (5.10)$$

Ca urmare,  $a_{[b]}$  este un set calculat de procesorul  $b$ , necesar procesorului  $a$ .

Toate procesoarele necesită câte un asemenea set din partea celorlalte procesoare. În total sunt  $P(P-1)$  asemenea seturi. Toate aceste seturi sunt diferite, deci nu există un mod simplu să fie transmise între procesoare, adică nu se poate efectua în exact același mod ca și în pașii 1 și 2.

Există totuși o cale de a efectua transmisia într-un mod similar celui din pasul 2. Pe baza faptului că seturile nu sunt total independente, toate fiind parte a unei însumări, se va demonstra în continuare că este posibil să fie transmise și însumate seturile astfel încât toate transmisiile pot fi efectuate în  $P-1$  pași secvențiali (cu  $P$  transmisiile paralele/asincrone în fiecare pas), la fel ca și în pașii 1 și 2.

Generalizând ecuația 5.10:

$$a_{[b, \dots, c]} = \left\{ \bar{\delta}_{[b]}^O \star \bar{w}_{j \rightarrow [b]}^O + \dots + \bar{\delta}_{[c]}^O \star \bar{w}_{j \rightarrow [c]}^O \right\}_{\forall j \in H_{[a]}} \quad (5.11)$$

În această ecuație se însumează un număr de seturi calculate de procesoarele de la  $b$  la  $c$ , necesare procesorului  $a$ . Suma este o parte a întregii sume utilizată în calculul lui  $\delta_{[a]}^H$ .

Figura 5.4 ilustrează operațiile de efectuat:

1. Procesorul 0 calculează setul său pentru procesorul 1, notat cu  $1_{[0]}$ . Procesorul 1 calculează setul său pentru procesorul 2 —  $2_{[1]}$ , și așa mai departe.
2. Calculul și transmisia se efectuează în paralel. Procesorul 0 calculează  $2_{[0]}$  în timp ce transmite  $1_{[0]}$  către procesorul  $P-1$ . Procesorul 1 calculează  $3_{[1]}$  în timp ce transmite  $2_{[1]}$  către procesorul 0 și așa mai departe.
3. Procesorul 0 recepționează  $2_{[1]}$ , îl adună la setul  $2_{[0]}$  pe care tocmai l-a calculat, obținând  $2_{[0,1]}$ , care este următorul set de transmis. Procesorul 1 recepționează  $3_{[2]}$ , îl adună la setul  $3_{[1]}$  pe care tocmai l-a calculat, obținând  $3_{[1,2]}$ . La fel pentru celelalte procesoare.

În acest moment se repetă pași similari pașilor 2 și 3, până când toate procesoarele dețin întregul set necesar lor, adică procesorul  $p$  deține setul

$$p_{[0, \dots, P-1]} = \sum_{q=0}^{P-1} \bar{\delta}_{[q]}^O \star \bar{w}_{j \rightarrow [q]}^O, \quad \forall j \in H_{[p]} \quad (5.12)$$

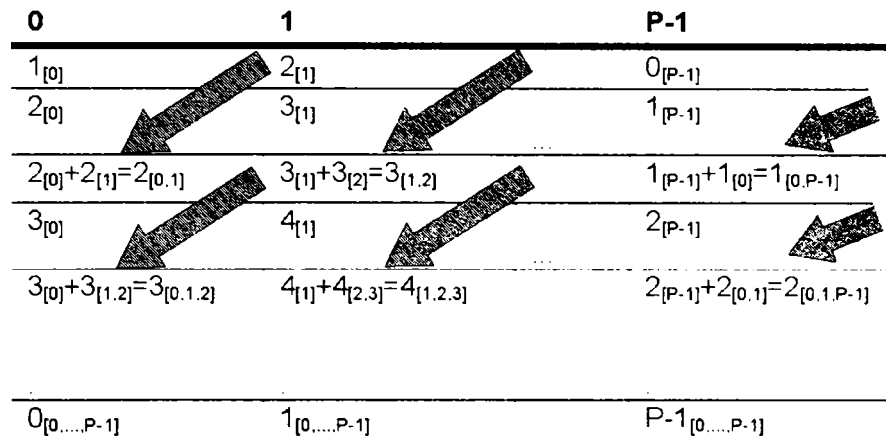


Figura 5.4: Calculul valorilor delta pentru unitățile ascunse.

Aceasta este exact suma din ultima expresie din ecuația 5.9 pentru procesorul  $p$ . Acum este posibil ca procesorul  $p$  să calculeze  $\delta_{[p]}^H$ .

#### 5.1.4.5 Pasul al cincilea

Al cincilea pas calculează actualizările de ponderi pentru conexiunile dintre stratul ascuns și cel de ieșire. Ecuația este:

$$\Delta w_{j \rightarrow k}^O(n+1) = \eta \delta_k^O a_j^H + \alpha \Delta w_{j \rightarrow k}^O(n) \quad (5.13)$$

Întrucât activările unităților ascunse, calculate în pasul doi, au fost memorate, este posibil ca procesoarele să calculeze toate actualizările de ponderi exclusiv pe baza datelor locale.

#### 5.1.4.6 Pasul al șaselea

Acest pas este similar pasului cinci, cu excepția faptului că se calculează actualizările pentru ponderile conexiunilor dintre stratul de intrare și cel ascuns. Ecuația de calcul este:

$$\Delta w_{i \rightarrow j}^H(n+1) = \eta \delta_j^H a_i^I + \alpha \Delta w_{i \rightarrow j}^H(n) \quad (5.14)$$

Din nou, întrucât activările unităților de intrare calculate în pasul întâi au fost memorate, este posibil ca procesoarele să calculeze toate actualizările de ponderi exclusiv pe baza datelor locale.

#### 5.1.4.7 Pasul al șaptelea

Acest pas are drept scop efectuarea actualizărilor de ponderi. Din nou, și aceasta se poate efectua numai pe baza datelor locale, fără nici o comunicație, deci în paralel.

### 5.1.5 Memorarea tiparelor de intrare și de ieșire

În funcție de dimensiunea tiparelor de intrare/ieșire, se pot utiliza strategii diferite de memorare a acestor date în procesoare. Dacă întregul set de tipare de intrare/ieșire este mic, toate perechile de tipare pot fi memorate în întregime în toate procesoarele, astfel încât toate procesoarele să dețină copii identice. În acest caz, pasul întâi al

algoritmului devine foarte simplu. Întrucât tiparele de intrare (activările unităților de intrare) sunt cunoscute tuturor procesoarelor, nu mai este necesar un broadcast.

Dacă setul de tipare de intrare/ieșire este prea mare pentru a-l plasa pe toate procesoarele, ele vor fi distribuite uniform între procesoare. După cum s-a menționat în secțiunea 5.1.4, de fapt se distribuie unitățile de intrare în mod uniform între toate procesoarele, ca și unitățile ascunse și cele de ieșire, iar fiecare procesor va deține din tiparele de intrare exact porțiunea care îi este asociată pentru unitățile de intrare pe care le simulează. Similar pentru tiparele de ieșire.

### 5.1.6 Îmbunătățiri ale algoritmului

Dacă timpul de comunicație din pasul patru constituie o problemă, o îmbunătățire evidentă a acestui algoritm în șapte pași este să se combine pașii patru și cinci, întrucât calculul actualizărilor de ponderi dintre straturile ascuns și cel de ieșire din pasul cinci poate fi efectuat odată cu terminarea pasului trei și deci în paralel cu pasul patru. Actualizarea ponderilor conexiunilor dintre stratul ascuns și cel de ieșire (parte din pasul șapte) poate fi efectuată de asemenea în pasul patru.

## 5.2 Analiza algoritmului

În această secțiune se va face o analiză a factorilor care influențează comportarea algoritmului. Testele care se vor efectua vor fi pe rețele de dimensiuni convenabile; ele nu vor fi rețele din aplicații reale și ca urmare în realitate nu se va efectua nici o antrenare.

Dintre cei șapte pași ai algoritmului, comunicații au loc doar în pașii unu, doi și patru. Cu excepția unor activități suplimentare, pentru comutarea între procese și pentru calculul de indici, probabil că numai comunicația din acești trei pași ar putea împiedica atingerea maximumului de *speed-up* față de versiunea secvențială.

### 5.2.1 Necesarul de memorie

Un avantaj al acestui algoritm îl constituie necesarul scăzut de memorie. Ponderile, activările și valorile delta din rețea sunt distribuite uniform între procesoare. Totuși, valorile de activare pentru toate unitățile de intrare și pentru cele ascunse sunt memorate în toate procesoarele, întrucât este vorba despre valori utilizate de două ori pe parcursul algoritmului. Activările unităților de intrare sunt utilizate în pasul unu și șapte, iar activările unităților ascunse sunt utilizate în pasul doi și pasul cinci.

Necesarul de memorie al fiecărui procesor, conform notațiilor din secțiunea 5.1.3 (măsurat în număr de valori reale pe 64 de biți — *double*) va fi:

$$NI + NH + PH + 2PO + 2((NI + 1) \cdot PH + (NH + 1) \cdot PO) \quad (5.15)$$

Atunci când rețelele simulate sunt de dimensiune mare, singura contribuție semnificativă este cea a ponderilor.

### 5.2.2 Pasul înainte

Cei doi pași care formează pasul înainte sunt în esență similari. Se va analiza în consecință doar unul dintre acești pași. Întrucât s-a discutat pe larg primul pas în secțiunea 5.1.4, aici va fi analizat pe larg acest pas. Se va demonstra că timpul de calcul depinde liniar atât de numărul de unități ascunse simulate de fiecare procesor —  $PH$ , cât și de numărul de unități de intrare simulate de fiecare procesor —  $PI$ . Timpul de comunicație depinde însă liniar numai de  $PI$ . Cu alte cuvinte, mărin  $PH$  este posibil să se mărească timpul de calcul fără să afecteze timpul de comunicație. Astfel, alegând un  $PH$  suficient de mare, timpul de comunicație poate să devină nesemnificativ în comparație cu timpul de calcul.

Procesorul  $p$  calculează intrarea din rețea pentru una din unitățile ascunse pe care le simulează conform ecuației 5.3. Acest calcul se execută în paralel (asincron) cu broadcast-ul activării unităților de intrare, cu observația că transmisia trebuie să se facă cu prioritate maximă, pentru a termina cât mai repede inițializarea și a putea trece cât mai repede la efectuarea calculelor.

La fiecare iterație, calculul termenului  $\bar{a}_{[q]}^I * \bar{w}_{[q]}^H$  conține  $PI$  înmulțiri și  $PI$  adunări. Timpul de calcul pentru fiecare iterație depinde deci liniar de  $PI$ . Întrucât calculul se efectuează pentru toate unitățile ascunse simulate de fiecare procesor, deci pentru toți  $j \in H_{[p]}$ , timpul de calcul va depinde liniar și de  $PH$ .

Situația din pasul doi este exact aceeași cu cea descrisă pentru pasul întâi. Pentru pasul doi timpul de calcul depinde liniar de  $PH$  și  $PO$ , în vreme ce timpul de comunicație depinde liniar numai de  $PH$ .

Concluzia la considerațiile de mai sus este că numărul de unități simulate de fiecare procesor,  $PI$ ,  $PH$  și respectiv  $PO$  determină calitatea paralelizării unei rețele. Dacă numărul de unități per procesor este mare, timpul de calcul este mai mare decât cel de comunicație, caz în care se poate considera că este vorba despre o bună paralelizare. Dacă numărul de unități per procesor este mic, calculele sunt încetinite de comunicație.

### 5.2.3 Pasul înapoi

Atenția se va concentra asupra pasului patru, în care sunt combinate calcule și transmisiile. Seturile transmise sunt vectori de dimensiune  $PH$ . Calculul setului  $a_{[p]}$  corespunde ecuației 5.10. Timpul de calcul pentru un termen din  $\bar{\delta}_{[p]}^O * \bar{w}_{j \rightarrow [p]}^O$  depinde liniar de  $PO$ . Există câte un asemenea termen pentru fiecare unitate ascunsă simulată de un procesor. Ca urmare, timpul de calcul pentru întregul tablou depinde liniar de  $PO$  și  $PH$ .

Ca și mai sus, timpul de calcul depinde liniar de doi factori: numărul de unități ascunse simulate de fiecare procesor și numărul de unități de ieșire simulate de fiecare procesor. Din nou, timpul de comunicație depinde liniar numai de unul dintre acești factori — de data aceasta numărul de unități de ieșire simulate de fiecare procesor. Similar, se poate anticipa o performanță mai bună a algoritmului atunci când numărul de unități simulate de un procesor este mare.

## 5.3 Simulatorul NetParSim

Pe baza considerațiilor teoretice anterioare, a fost realizat un simulator pentru antrenarea rețelelor neuronale cu metoda backpropagation — *NetParSim*.

Antrenarea rețelei neuronale se face distribuit, prin metoda partiționării rețelei. Fiecărui procesor din rețeaua de calculatoare i se atribuie un număr de neuroni de pe fiecare strat al rețelei neuronale. Comunicația între procesoare se realizează conform topologiei în inel, astfel încât este simplu de modificat numărul de mașini care execută simulatorul.

Simulatorul implementează paralelizarea în șapte pași propusă anterior, cu actualizarea ponderilor conexiunilor neuronale după fiecare tipar.

La pasul întâi se face prezentarea unui tipar de antrenare. Pentru a putea calcula activarea fiecărui neuron de pe stratul ascuns, fiecare procesor necesită întregul vector de intrare reprezentat de tiparul curent. Reamintim în acest sens faptul că un neuron de pe un strat are conexiuni de la toți neuronii de pe stratul anterior.

În al doilea pas se calculează activarea unităților de ieșire, similar cu pasul întâi.

Pasul trei calculează valorile *delta* pentru unitățile de ieșire, prin compararea cu valorile de ieșire precizate de tiparul de antrenare.

Pasul al patrulea calculează valorile *delta* pentru unitățile ascunse, conform schemei din figura 5.4.

Pașii șase și șapte calculează și actualizează ponderile conexiunilor care intră în unitățile de ieșire, respectiv ascunse.

Pentru evaluarea performanțelor acestui simulator, prin comparație cu varianta secvențială rulată pe același tip de mașină, se vor efectua o serie de experimente, care vor încerca să determine influența mai multor factori, în principal dimensiunea rețelei neuronale antrenate, respectiv dimensiunea rețelei de calculatoare pe care se rulează simulatorul.

Întrucât *calitatea* antrenării, respectiv numărul de iterații de efectuat necesar atingerii unui anumit grad de antrenare, depinde de parametri de antrenare — rata de învățare, momentul — care oricum sunt aceiași și în cazul variantei secvențiale, vor fi utilizate rețele neuronale generice, cu dimensiunile precizate la fiecare tip de experiment, și seturi de tipare generate aleator. Comparația între varianta secvențială și cea cu partiționarea rețelei se va face prin măsurarea timpului necesar efectuării unui număr de iterații (10). Valorile prezentate pentru accelerare, respectiv eficiența algoritmului, vor reprezenta media unui număr de 10 experimente/rulări din categoria respectivă.

Primul set de experimente vizează efectul modificării dimensiunii rețelei neuronale asupra eficienței algoritmului. Pentru numărul maxim de stații din rețea (16), se va aplica algoritmul pentru rețele neuronale astfel dimensionate încât, pentru fiecare procesor, numărul de neuroni de pe fiecare strat să varieze între 1 și 20, adică rețele neuronale cu 16 neuroni pe stratul de intrare, cel intermediar și cel de ieșire, 32 de neuroni și așa mai departe.

În al doilea set de experimente se va determina influența modificării numărului de procesoare pe care rulează simulatorul. Pentru o rețea neuronală de dimensiune fixă, se va modifica numărul de procesoare, între 2 și 16. Experimentele vor fi efectuate pentru mai multe rețele neuronale, dimensionate astfel încât numărul de neuroni de pe fiecare strat să fie divizibil cu numărul maxim de procesoare utilizate.

În al treilea set de experimente se va determina în ce măsură este posibil să se utilizeze mai multe procesoare pe măsură ce se mărește dimensiunea rețelei neuronale. Utilizând între 2 și 16 procesoare, se vor efectua măsurători pentru rețele neuronale de dimensiuni alese astfel încât pentru  $P$  procesoare să avem pe fiecare strat un număr de neuroni proporțional cu  $P$ . În acest set de experimente, la adăugarea unui procesor în rețeaua de calculatoare, se va adăuga și un număr de neuroni pe fiecare strat al rețelei neuronale.

Al patrulea set de experimente va determina în ce măsură este afectată eficiența algoritmului de modificarea dimensiunii setului de antrenare. Deși algoritmul este astfel proiectat încât actualizarea ponderilor se face după prezentarea *fiecărui* tipar, acest algoritm poate fi ușor modificat pentru însumarea unui număr de gradienti înainte de actualizarea ponderilor.

Simulatorul NetParSim conține un singur program. Prin urmare, procesele rulate pe fiecare stație din rețeaua de calculatoare sunt identice. Diferențierea rolurilor, în *administrator* și *slave* se realizează prin intermediul argumentului transmis în linia de comandă — `processor_number`. În cazul în care acesta lipsește sau este 0, procesul este *administrator* și se va ocupa cu pornirea proceselor *slave*, cu numărul procesului dat de poziția în configurația de inel.

Algoritmul utilizat este cel descris în acest capitol, cu comasarea pașilor al patrulea, al cincilea și, parțial, al șaptelea, după cum s-a sugerat în secțiunea 5.1.6.

### 5.3.1 Fișierul de configurare

Programul primește informațiile necesare prin intermediul unui fișier de configurare, identificat de variabila de mediu `NETPARSIM_CFG`, similar cu cel descris în secțiunea 4.3.1. El conține informații despre rețeaua neuronală de simulat și despre stațiile care alcătuiesc rețeaua de calculatoare. După cum s-a precizat, o informație suplimentară se transmite în linia de comandă — numărul procesorului în configurația de inel.

Se dă mai jos un exemplu de fișier de configurație. Valorile sunt precizate câte una pe linie, sub forma

```
<variabilă>=<valoare>
```

Liniile goale se ignoră, iar textul de la semnul '#' până la sfârșitul liniei se ignoră (comentariu).

```
# Reteaua neuronală:
eta = 0.2           # rata de învățare
alpha = 0.9        # momentul

NrInput = 60       # neuroni pe stratul de intrare
NrHidden = 60      # neuroni pe stratul intermediar
NrOutput = 60      # neuroni pe stratul de ieșire

NrPatt = 10        # număr total de tipare
PattFile = /work/tipare.dat # fișier tipare de învățare

NrIter = 100       # număr de iterații
Err = 0.01        # eroare acceptată

# Reteaua de calculatoare:
NrProcs = 12       # număr procesoare în rețea
HostsFile = /work/hosts.dat # lista numelor stațiilor
```

Parametrul `NrIter` precizează numărul total de iterații care sunt efectuate. Dacă el este 0, atunci se ia ca și criteriu de oprire atingerea unei erori totale mai mică decât `Err`.

### 5.3.2 Codul aplicației

Pentru început, aplicația determină numărul de neuroni de pe fiecare strat pe care îi are de simulat fiecare proces:

```
/* Calcul valori per procesor */
ProcInput = (NrInput + (NrProcs - 1)) / NrProcs;
ProcHidden = (NrHidden + (NrProcs - 1)) / NrProcs;
ProcOutput = (NrOutput + (NrProcs - 1)) / NrProcs;
```

Apoi, se determină indicii corespunzători neuronilor simulați de fiecare proces:

```
/* calcul indici input si lungime date pentru toate procesoare */
/* Strat intrare: */
input_ndx[0] = 0;
input_cnt[0] = ProcInput;
for (i = 1; i < NrProcs; i++) {
    input_ndx[i] = input_ndx[i-1] + input_cnt[i-1];
    tval = NrInput % NrProcs;
    if (tval == 0) {
        input_cnt[i] = ProcInput;
    }
    else if (tval > i) {
        input_cnt[i] = ProcInput;
    }
    else {
        input_cnt[i] = ProcInput - 1;
    }
}
LocalInput = input_cnt[processor_number];
```



```

/* Strat ascuns: */
hidden_ndx[0] = 0;
hidden_cnt[0] = ProcHidden;
for (i = 1; i < NrProcs; i++) {
    hidden_ndx[i] = hidden_ndx[i-1] + hidden_cnt[i-1];
    tval = NrHidden % NrProcs;
    if (tval == 0) {
        hidden_cnt[i] = ProcHidden;
    }
    else if (tval > i) {
        hidden_cnt[i] = ProcHidden;
    }
    else {
        hidden_cnt[i] = ProcHidden - 1;
    }
}
LocalHidden = hidden_cnt[processor_number];

/* Strat ieşire: */
output_ndx[0] = 0;
output_cnt[0] = ProcOutput;
for (i = 1; i < NrProcs; i++) {
    output_ndx[i] = output_ndx[i-1] + output_cnt[i-1];
    tval = NrOutput % NrProcs;
    if (tval == 0) {
        output_cnt[i] = ProcOutput;
    }
    else if (tval > i) {
        output_cnt[i] = ProcOutput;
    }
    else {
        output_cnt[i] = ProcOutput - 1;
    }
}
LocalOutput = output_cnt[processor_number];

```

După cum s-a precizat, singura diferență dintre *administrator* și *slave* este dată de numărul procesului în inel, cel cu numărul 0 fiind *administrator*. Acesta trebuie să pornească celelalte procese *slave* și să deschidă canalele de comunicație dintre ele, pentru configurația de inel. Fiecare proces va recepționa date de la procesul precedent în rețea și va transmite la cel următor.

După pornirea tuturor proceselor, acestea trebuie să aloce spațiu pentru structurile de date utilizate, să inițializeze aceste date și să încarce tiparele de antrenare. Întrucât la inițializarea ponderilor se utilizează valori mici aleatoare, trebuie asigurat ca acestea să fie identice pentru toate procesele – se utilizează aceeași rutină de generare de numere aleatoare, cu *aceeași rădăcină* (*seed*).

Procesul *administrator* pornește măsurarea timpului, după care toate procesele execută un ciclu relativ simplu, de *NrIter* ori, în care se prezintă rețelei neuronale un tipar – selectat aleator – și se actualizează ponderile conexiunilor:

```

/* PatternUpdate () */
for (iteration = 0; iteration < NrIter; iteration++) {
    pattern = (int)random() % NrPatt;
    PropagateActivity (pattern);
    CalculateWeightChanges(pattern);
}

```

Rutina *PropagateActivity* prezintă rețelei neuronale câte un tipar de antrenare. Pentru a putea propaga activarea prin rețeaua neuronală, trebuie să recepționeze



activările neuronilor din stratul de intrare, respectiv cel ascuns, care se află în alte procesoare. Mai întâi se calculează activarea unităților ascunse din procesorul curent:

```

for (j = 0; j < LocalHidden; j++) {
    hidden_net[j] = hidden_weight[TotalInput][j]; /* bias */
}

p = processor_number;
for (k = 0; k < (NrProcs - 1); k++, p++) {
    processor = p % NrProcs;
    Send (input_activity[processor]);
    Receive (input_activity[(processor+1)%NrProcs]);
    StartNdx = input_ndx[processor];
    for (j = 0; j < LocalHidden; j++) {
        for (i = 0; i < input_cnt[processor]; i++) {
            hidden_net[j] += input_activity[processor][i] *
                hidden_weight[StartNdx+i][j];
        }
    }
}

processor = (processor_number + (NrProcs - 1)) % NrProcs;
StartNdx = input_ndx[processor];
for (j = 0; j < LocalHidden; j++) {
    for (i = 0; i < input_cnt[processor]; i++) {
        hidden_net[j] += input_activity[processor][i] *
            hidden_weight[StartNdx+i][j];
    }
    hidden_activity[processor_number][j] =
        CalculateActivity (hidden_net[j]);
}

```

Se transmit activările neuronilor de intrare proprii și cele recepționate de la celelalte procesoare (pentru claritate, în apelurile rutinelor `Send` și `Receive` nu s-a precizat și lungimea tablourilor de valori implicate). Toate aceste valori sunt utilizate în calculul activării unităților ascunse simulate. Apoi se efectuează aceleași operații pentru stratul de ieșire:

```

for (j = 0; j < LocalOutput; j++) {
    output_net[j] = output_weight[TotalHidden][j]; /* bias */
}

p = processor_number;
for (k = 0; k < (NrProcs - 1); k++, p++) {
    processor = p % NrProcs;
    Send (hidden_activity[processor]);
    Receive (hidden_activity[(processor+1)%NrProcs]);
    StartNdx = hidden_ndx[processor];
    for (j = 0; j < LocalOutput; j++) {
        for (i = 0; i < hidden_cnt[processor]; i++) {
            output_net[j] += hidden_activity[processor][i] *
                output_weight[StartNdx+i][j];
        }
    }
}

processor = (processor_number + (NrProcs - 1)) % NrProcs;
StartNdx = hidden_ndx[processor];
for (j = 0; j < LocalOutput; j++) {
    for (i = 0; i < hidden_cnt[processor]; i++) {

```

```

        output_net[j] += hidden_activity[processor][i] *
                        output_weight[StartNdx+i][j];
    }
    output_activity[processor_number][j] =
        CalculateActivity (output_net[j]);
}

```

De data aceasta, se transmit activările neuronilor de pe stratul ascuns simulați de procesul curent și cele primite de la celelalte procesoare.

Pentru calculul actualizărilor ponderilor conexiunilor din rețeaua neuronală (rutina `CalculateWeightChanges`), mai întâi se calculează valorile *delta* în neuronii de ieșire, conform relației 5.7:

```

for (i = 0; i < LocalOutput; i++) {
    tmp = output_activity[i];
    output_delta[i] = (target_pattern[pattern][i] - tmp) *
                      (tmp * (1.0 - tmp));
}

```

Apoi se calculează valorile delta pentru unitățile ascunse și valorile cu care trebuie actualizate ponderile conexiunilor dintre stratul ascuns și cel de ieșire, conform celor descrise în secțiunile 5.1.4.4 și 5.1.4.5:

```

processor = (processor_number + 1) % NrProcs;
StartNdx = hidden_ndx[processor];
for (i = 0; i < hidden_cnt[processor]; i++) {
    out_data[i] = 0.0;
    for (j = 0; j < LocalOutput; j++) {
        out_data[i] += (output_delta[j] *
                       output_weight[StartNdx+i][j]);
        /* calcul actualizare ponderi hidden-output */
        output_change[StartNdx + i][j] =
            (eta * (output_delta[j] *
                  hidden_activity[processor][i])) +
            (alpha * output_change[StartNdx+i][j]);

        /* efectueaza actualizarea ponderi hidden-output */
        output_weight[StartNdx+i][j] += output_change[StartNdx+i][j];
    }
}

p = processor_number + 2;
for (k = 0; k < (NrProcs - 1); k++, p++) {
    Send (out_data);
    Receive (in_data);
    processor = p % NrProcs;
    StartNdx = hidden_uint_index[processor];
    for (i = 0; i < hidden_cnt[processor]; i++) {
        hidden_delta[i] = 0.0;
        for (j = 0; j < LocalOutput; j++) {
            hidden_delta[i] += (output_delta[j] *
                                output_weight[StartNdx+i][j]);
            /* calcul actualizare ponderi hidden-output */
            output_change[StartNdx+i][j] = (eta * (output_delta[j] *
                                                  hidden_activity[processor][i])) +
            (alpha * output_change[StartNdx+i][j]);
            /* efectueaza actualizare ponderi hidden-output */
            output_weight[StartNdx+i][j] +=
                output_change[StartNdx+i][j];
        }
    }
}

```

```

    for (i = 0; i < hidden_cnt[processor]; i++) {
        out_data[i] = in_data[i] + hidden_delta[i];
    }
}

for (i = 0; i < LocalHidden; i++) {
    hidden_delta[i] = out_data[i] *
        (hidden_activity[processor_number][i] *
        (1.0 - hidden_activity[processor_number][i]));
}

for (j = 0; j < LocalOutput; j++) {
    /* calcul actualizare ponderi hidden-output */
    output_change[TotalHidden][j] = (eta * output_delta[j]) +
        (alpha * output_change[TotalHidden][j]);
    /* efectueaza actualizarea ponderi hidden-output */
    output_weight[TotalHidden][j] +=
        output_change[TotalHidden][j];
}

```

Similar se procedează pentru ponderile conexiunilor dintre stratul de intrare și cel ascuns, conform secțiunii 5.1.4.6:

```

/* calcul actualizare ponderi input-hidden */
for (p = 0; p < NrProcs; p++) {
    StartNdx = input_ndx[p];
    for (j = 0; j < LocalHidden; j++) {
        for (i = 0; i < input_cnt[p]; i++) {
            hidden_change[StartNdx+i][j] =
                (eta * (hidden_delta[j] * input_activity[p][i])) +
                (alpha * hidden_change[StartNdx + i][j]);
        }
    }
}

for (j = 0; j < LocalHidden; j++) { /* pentru bias */
    hidden_change[TotalInput][j] = (eta * hidden_delta[j]) +
        (alpha * hidden_change[TotalInput][j]);
}

/* actualizeaza ponderi input-hidden */
for (j = 0; j < LocalHidden; j++) {
    for (i = 0; i < (TotalInput + 1); i++) {
        hidden_weight[i][j] += hidden_change[i][j];
    }
}

```

După cum se poate observa, pasul al șaptelea – actualizarea ponderilor – se efectuează chiar în momentul în care s-a terminat calculul actualizărilor, mai întâi pentru conexiunile dintre stratul ascuns și cel de ieșire, apoi pentru cele dintre stratul de intrare și stratul ascuns.

## 5.4 Rezultate experimentale

Simulatorul NetParSim a fost testat pe aceleași platforme ca și simulatorul DataParSim (vezi secțiunea 4.4.1). Rezultatele prezentate sunt mediile rezultatelor a câte 10 experimente identice din categoria respectivă. Un experiment a constat din câte 100 de cicluri de prezentare, tiparele utilizate fiind generate aleator.

### 5.4.1 Influența dimensiunii variabile a rețelei neuronale

Tabelul 5.1 și figura 5.5 (SunBlade) și respectiv tabelul 5.2 și figura 5.6 (Linux) prezintă eficiența utilizând numărul maxim de stații de lucru în fiecare dintre cele două configurații (16 stații). Atât varianta secvențială cât și cea paralelă au fost executate pe rețele de exact aceleași dimensiuni. În figurile 5.5a) și 5.6a) rețelele sunt astfel dimensionate încât numărul de unități simulate de către fiecare procesor să fie de la 1 la 20, ceea ce înseamnă 16 neuroni de intrare, pe stratul ascuns și la ieșire, 32 de neuroni și așa mai departe.

Unități per procesor (U)	Dimensiune rețea				
	16U 16U 16U n	16U 16U 16U fc	16U 160 160	160 16U 160	160 160 16U
1	32.80	45.69	79.29	72.09	80.71
2	57.71	75.71	80.02	76.02	80.84
3	70.73	82.22	80.50	78.61	81.01
4	75.71	84.09	81.03	79.13	81.16
5	78.69	84.91	81.52	79.59	81.28
6	80.22	85.49	82.03	80.11	81.71
7	81.72	85.74	82.49	80.82	81.83
8	82.39	85.71	83.08	81.29	82.02
9	82.52	85.79	82.51	81.82	82.15
10	82.61	85.91	83.09	82.31	82.29
11	82.73	86.11			
12	82.81	86.09			
13	82.89	86.10			
14	83.02	86.19			
15	83.50	86.20			
16	82.51	84.29			
17	82.49	84.31			
18	82.52	84.29			
19	82.14	84.22			
20	82.81	84.23			

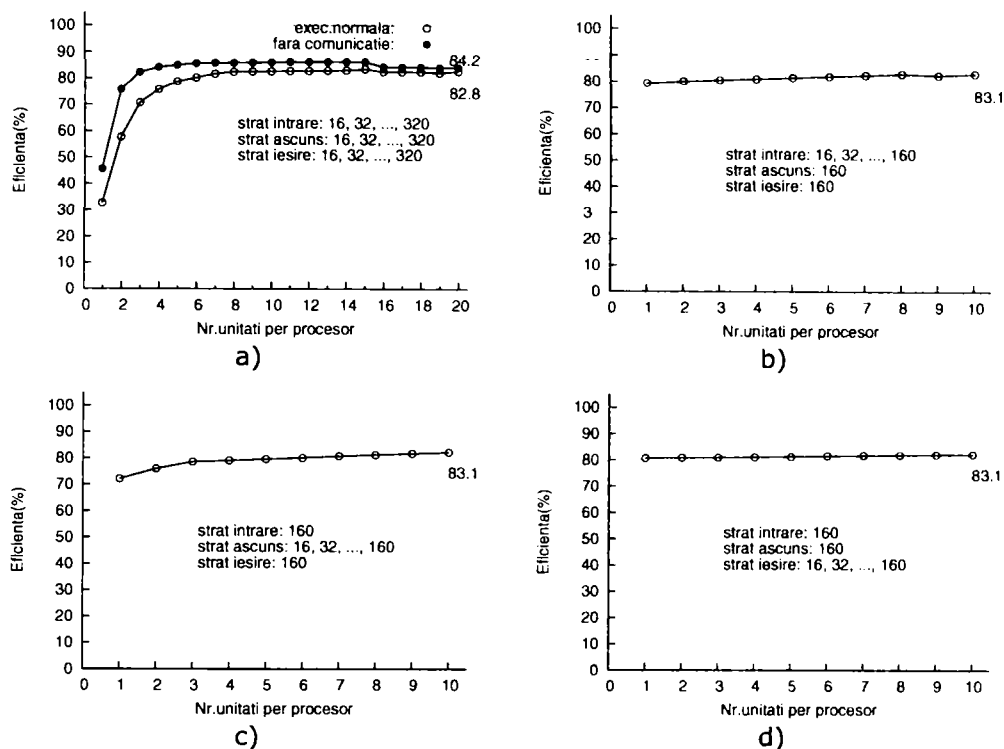
**Tabela 5.1:** Dimensiune variabilă a rețelei: SunBlade.

Se observă că atunci când există doar o unitate per procesor în fiecare dintre cele trei straturi, simulatorul paralel nu prezintă o eficiență satisfăcătoare – doar 35%. Atunci când sunt două unități per procesor, eficiența crește la 59%, iar la circa șapte unități per procesor eficiența se stabilizează la aproximativ 85% – probabil eficiența maximă care se poate obține.

Una dintre cauzele pentru care nu se poate obține o eficiență mai bună s-ar putea datora modului în care algoritmul accesează elementele vectorului  $\vec{a}_{[q]}$ , pe parcursul transmisiei.

O altă cauză este ilustrată în cele două figuri 5.5a și 5.6a. Cerculețele pline reprezintă execuția în absența comunicației (înlocuită cu instrucțiunile NOP). Spațiul dintre cele două curbe reprezintă influența comunicației. Există un spațiu între curba cu cerculețe pline și linia de eficiență 100%, care trebuie să corespundă altor cauze software, întrucât nu există comunicație.

În celelalte trei grafice din figurile 5.5 și 5.6 a fost fixat numărul de neuroni pe două dintre cele trei straturi, și s-a modificat numărul de unități pe al treilea strat. În figurile b) a fost fixat numărul de unități ascunse și de ieșire la 160, și a fost ajustat numărul de unități de intrare între 16 și 160, astfel încât fiecare procesor să simuleze



**Figura 5.5:** Dimensiune variabilă a rețelei – SunBlade.

câte 10 unități de pe straturile ascuns și de ieșire și între 1 și 10 neuroni de pe stratul de intrare. S-a putut astfel evalua efectul unui număr foarte mic de unități de intrare per procesor. Oricum, fără vreun efect notabil. Eficiența este cu foarte puțin mai mică atunci când există doar un neuron de intrare pentru fiecare procesor.

În figurile c) și d) s-a modificat numărul de unități de pe stratul ascuns, respectiv de ieșire. Se observă că atunci când numărul de unități ascunse per procesor este prea mic, eficiența scade puțin, însă nu foarte mult. Numărul de unități de ieșire per procesor nu are efect asupra eficienței câtă vreme sunt simulate suficient de multe unități din celelalte straturi.

Concluzia care se poate deduce este că atunci când numărul de unități per procesor este mic pentru toate straturile eficiența este scăzută, însă faptul că procesoarele dețin puține unități doar într-unul dintre straturi are un efect neglijabil asupra eficienței.

#### 5.4.2 Efectul numărului variabil de procesoare

Tabelul 5.3 și figura 5.7 și respectiv tabelul 5.4 și figura 5.8 reprezintă accelerarea pentru diferite dimensiuni ale rețelei utilizând până la numărul maxim de procesoare (16). În fiecare grafic individual dimensiunea rețelei rămâne fixă și se modifică numărul de procesoare. Dimensiunile rețelei au fost alese astfel încât numărul maxim de procesoare (16) să dividă numărul de unități din fiecare strat. Experimentele au fost realizate pentru rețele: 16-16-16, 32-32-32, 48-48-48, 80-80-80, 160-160-160.

Graficele sugerează, după cum era de așteptat, faptul că pentru un număr mic de neuroni pentru fiecare procesor accelerarea și eficiența sunt departe de valorile

Unități per procesor (U)	Dimensiune rețea				
	16U	16U	16U	16U	16U
	16U	16U	16U	16U	16U
	16U	16U	16U	16U	16U
	n	fc			
1	30.6	43.9	76.9	69.8	77.9
2	55.2	73.1	78.3	73.2	78.1
3	69.1	80.4	78.7	75.1	79.33
4	71.7	82.5	79.8	76.5	79.9
5	76.2	83.1	79.85	76.75	80.1
6	78.9	83.3	80.3	78.1	79.7
7	79.8	83.65	80.7	78.25	79.9
8	80.1	83.5	81.2	79.3	80.1
9	80.4	83.6	80.5	79.8	80.3
10	80.5	83.7	81.7	80.25	80.4
11	80.6	84.0			
12	80.7	84.0			
13	80.8	84.0			
14	81.1	84.1			
15	81.3	84.1			
16	80.9	82.2			
17	80.85	82.2			
18	80.44	82.2			
19	80.1	82.0			
20	80.6	82.1			

**Tabela 5.2:** Dimensiune variabilă a rețelei: Linux.

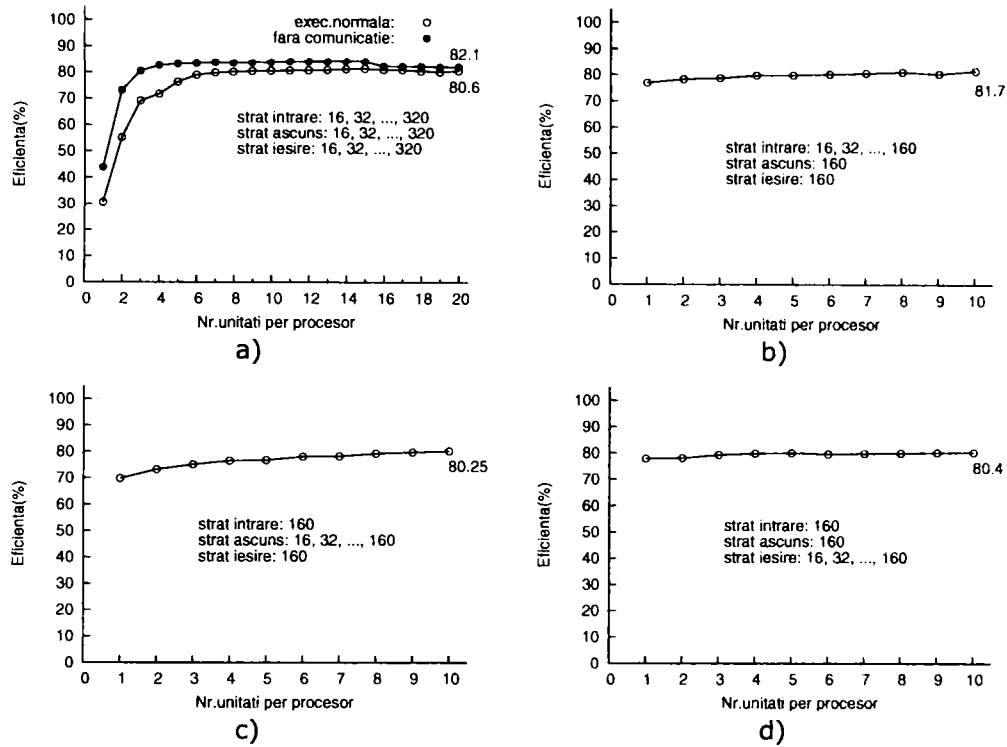
optime.

Pentru rețele neuronale de aceleași dimensiuni, graficele nu diferă semnificativ pentru cele două platforme hardware, SunBlade, respectiv Linux. În privința valorilor obținute pentru accelerare, se observă că în cazul rețelelor de dimensiuni reduse, 16-16-16, respectiv 32-32-32, accelerarea maximă obținută este foarte departe de valoarea ideală, atunci când simulatorul utilizează maximum de procesoare (16), valoarea cea mai bună obținându-se pentru rețeaua neuronală cu dimensiunea cea mai mare (160-160-160). Alura graficelor sugerează faptul că pentru rețelele neuronale cu dimensiunile precizate, performanțele simulatorului (accelerarea, respectiv eficiența) este posibil să se îndepărteze și mai mult de valorile ideale în cazul în care am adăuga mai multe procesoare.

Concluzia care rezultă din acest experiment este că pentru a ne situa cât mai aproape de valorile optime pentru accelerare, respectiv eficiență, este necesar să asigurăm un număr cât mai mare de neuroni pentru fiecare procesor. Pentru cazul experimentelor efectuate, se pare că performanțe acceptabile se pot obține dacă avem cel puțin 10 neuroni de pe fiecare strat asigurați fiecărei mașini din rețeaua de calculatoare.

### 5.4.3 Influența scalării rețelei cu numărul de procesoare

În această secțiune s-a încercat să se determine în ce măsură este posibil să se utilizeze mai multe procesoare pe măsură ce se mărește dimensiunea rețelei. În tabelele 5.5, 5.6 și respectiv în figurile 5.9, 5.10 este reprezentată eficiența când sunt utilizate între 2 și 16 procesoare. Graficul pentru cea mai scăzută eficiență a fost obținut fixând dimensiunea rețelei astfel încât fiecare procesor să simuleze câte o unitate din fiecare strat, ceea ce înseamnă că pentru  $P$  procesoare este simulată o rețea  $P-P-P$ .



**Figura 5.6:** Dimensiune variabilă a rețelei – Linux.

Număr procesoare	Dimensiune rețea				
	16	32	48	80	160
2	1.19	1.20	1.32	1.38	1.60
3	2.01	2.21	2.22	2.29	2.36
4	2.92	3.09	3.09	3.22	3.20
5	2.91	3.41	3.81	4.13	4.14
6	3.93	4.89	4.63	5.11	4.90
7	3.91	5.71	5.11	5.72	5.96
8	4.62	6.31	6.01	6.79	6.89
9	4.61	6.30	7.52	7.03	7.49
10	4.60	6.93	7.51	8.10	8.42
11	4.59	6.92	7.53	8.51	9.47
12	4.57	6.91	8.52	8.50	9.90
13	4.56	6.95	8.51	9.31	11.31
14	4.55	6.94	8.50	9.30	11.62
15	4.54	7.01	8.49	9.29	13.04
16	5.40	7.92	9.51	10.02	13.60

**Tabela 5.3:** Accelerarea – număr variabil de procesoare: SunBlade.

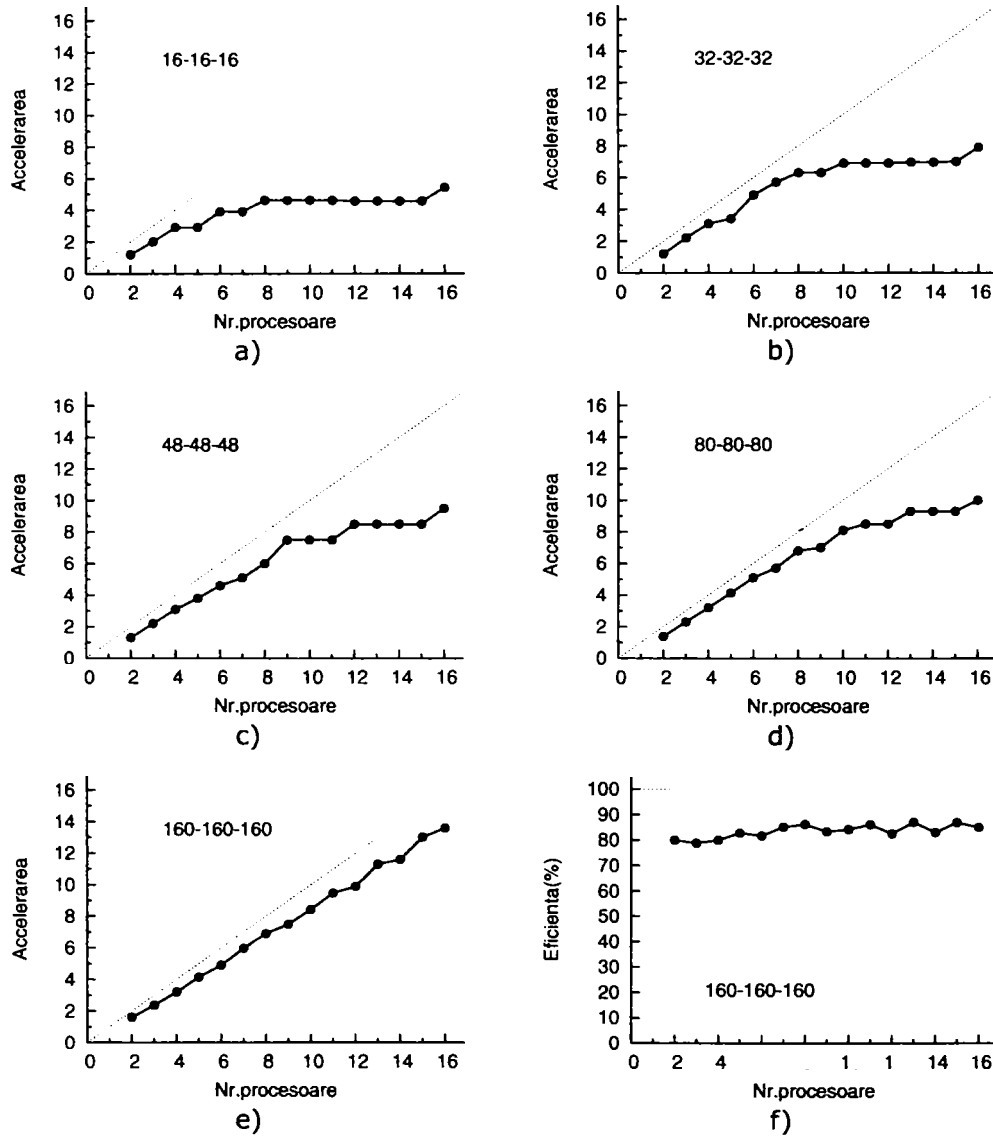


Figura 5.7: Accelerarea – număr variabil de procesoare: SunBlade.



Număr procesoare	Dimensiune rețea				
	16	32	48	80	160
	16	32	48	80	160
2	1.15	1.1	1.15	1.2	1.25
3	1.88	1.9	2.0	2.1	2.1
4	2.7	2.85	2.85	2.95	3.0
5	2.75	3.15	3.5	3.8	3.9
6	3.3	4.4	4.22	4.88	4.6
7	3.5	5.3	4.9	5.3	5.7
8	4.1	5.9	5.8	6.2	6.5
9	4.15	6.0	7.1	6.9	7.1
10	4.2	6.2	7.15	7.8	8.2
11	4.25	6.3	7.2	8.1	9.15
12	4.15	6.35	8.0	8.2	9.6
13	4.2	6.4	8.1	8.9	10.95
14	4.25	6.4	8.15	9.0	11.1
15	4.2	6.85	8.2	9.1	12.0
16	4.9	7.2	9.3	9.8	12.7

**Tabela 5.4:** Accelerarea - număr variabil de procesoare: Linux.

Următoarele grafice sunt pentru câte o rețea  $2P - 2P - 2P$ ,  $3P - 3P - 3P$ ,  $5P - 5P - 5P$  și respectiv  $10P - 10P - 10P$ .

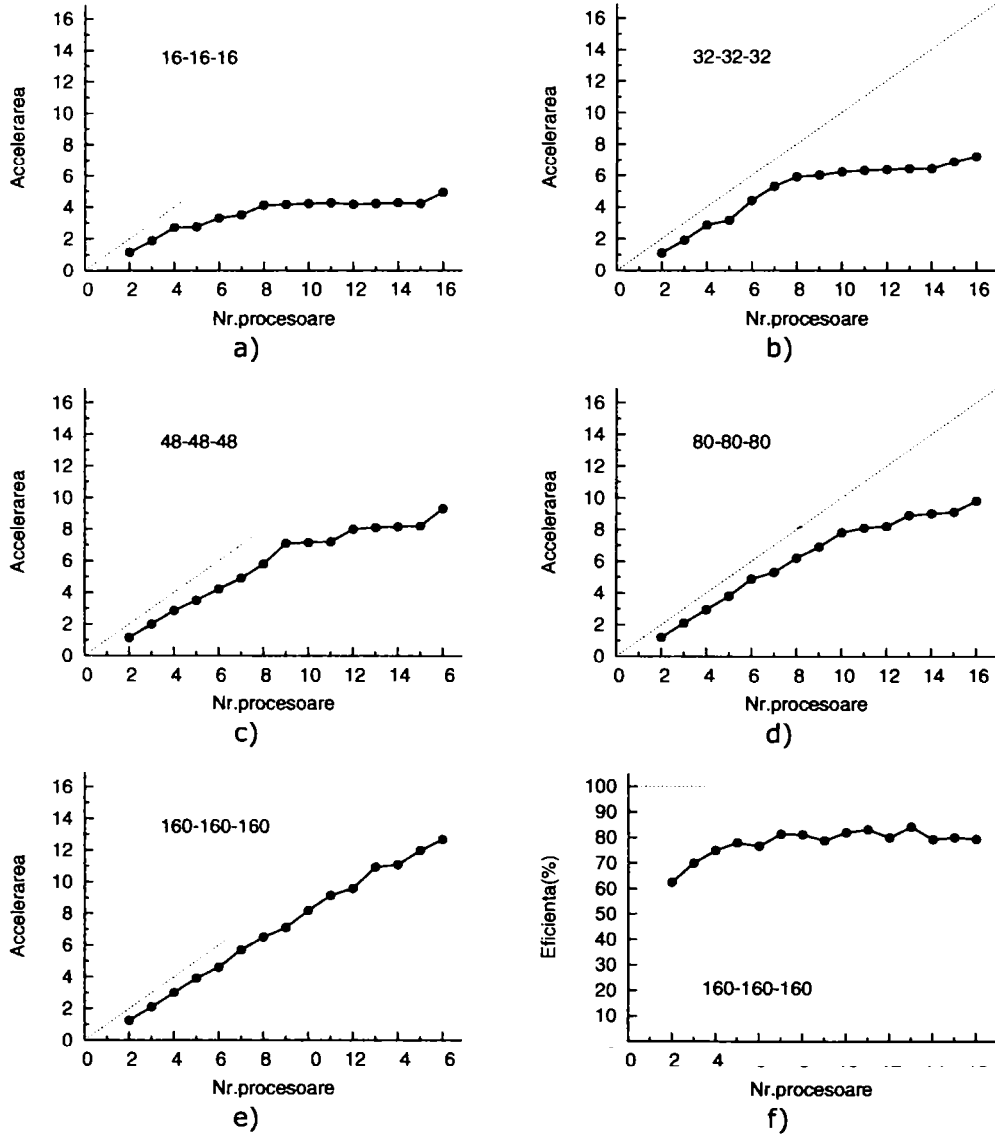
Graficul cu cea mai scăzută eficiență descrește considerabil pe prima porțiune. Aceasta se datorește faptului că pentru un număr mic de procesoare  $P$  care simulează o rețea  $P - P - P$ , evident, rețeaua nu este prea mare. S-a văzut mai devreme că pentru rețele mici, timpul de calcul al funcțiilor de activare reprezintă o porțiune considerabilă a timpului total de execuție. Această parte a algoritmului este desigur bine paralelizată și deci de aici eficiența mai mare la începutul graficelor.

După cum s-a arătat în secțiunea 5.4.1, simularea rețelelor cu număr mic de unități per procesor prezintă o eficiență scăzută. Utilizând cel puțin 5 unități per procesor eficiența este aproape constantă. Ca urmare, este posibil să se mărească numărul de procesoare fără a pierde din eficiență dacă numărul de unități per procesor rămâne constant (deci mărind numărul total de unități proporțional cu numărul de procesoare). De remarcat faptul că numărul de ponderi crește aproape quadratic atunci când numărul de unități de pe straturi crește liniar și deci și volumul de calcul crește aproape quadratic.

Curbele din figurile 5.9, 5.10 sunt linii aproape orizontale. Dacă se presupune că graficele continuă în același mod și pentru un număr mai mare de procesoare, se poate decide dacă este posibil să se utilizeze mai multe procesoare pentru o problemă dată. Dacă într-o simulare procesoarele simulează câte 10 unități din fiecare strat cu o eficiență de 85.6%, este posibil să se utilizeze de două ori mai multe procesoare (5 unități per procesor) cu o eficiență de 80.2%. Este chiar mai preferabil să fie o unitate per procesor în loc de două, chiar dacă procesoarele suplimentare nu sunt utilizate prea bine: având de două ori mai multe procesoare se va obține o eficiență de 31.5% în loc de 58.2%, dar se va constata și o creștere a vitezei de procesare.

#### 5.4.4 Modificarea dimensiunii setului de antrenare

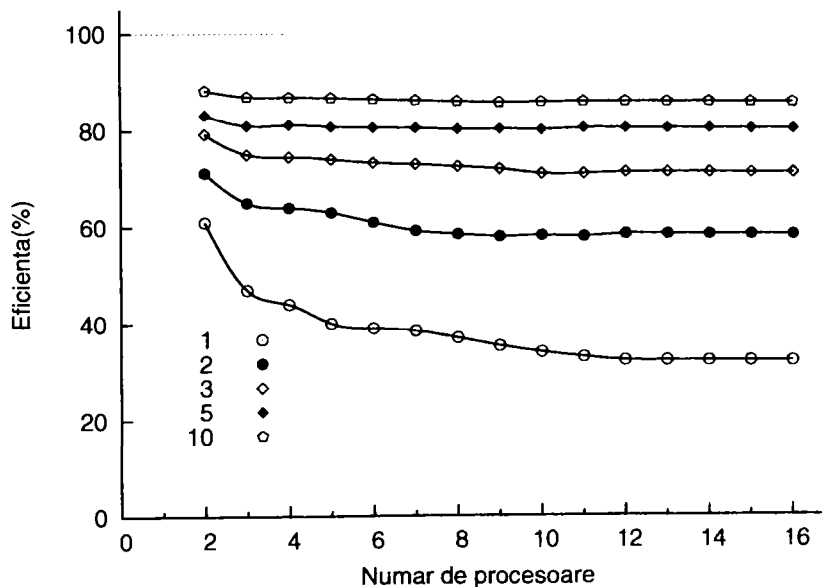
Contrar paralelizării cu partiționarea datelor din capitolele anterioare, schema cu partiționarea rețelei actualizează ponderile *după fiecare tipar*. Desigur, este posibil să se actualizeze ponderile mai puțin frecvent. Modificând foarte puțin programul, a fost posibil să se utilizeze un set de antrenare de dimensiune oarecare. Modificarea



**Figura 5.8:** Accelerarea – număr variabil de procesoare: Linux.

Nr. Procesoare	Dimensiune rețea				
	1P	2P	3P	5P	10P
2	61.01	71.19	79.31	83.21	88.29
3	46.99	64.99	74.98	81.04	87.03
4	44.02	64.03	74.52	81.19	86.87
5	39.98	62.98	74.01	80.79	86.76
6	39.01	61.01	73.27	80.63	86.52
7	38.52	59.19	73.02	80.56	86.24
8	36.97	58.51	72.48	80.24	86.01
9	35.31	58.09	71.99	80.31	85.79
10	34.02	58.21	71.01	80.13	85.88
11	32.98	58.02	70.98	80.52	85.94
12	32.26	58.49	71.31	80.46	85.93
13	32.19	58.46	71.24	80.39	85.91
14	32.13	58.39	71.18	80.34	85.86
15	32.11	58.33	71.16	80.31	85.81
16	32.04	58.31	71.09	80.23	85.72

**Tabela 5.5:** Scalarea rețelei cu numărul de procesoare: SunBlade.



**Figura 5.9:** Eficiența în cazul scalării rețelei cu numărul de procesoare - SunBlade.

Nr. Procesoare	Dimensiune rețea				
	1P	2P	3P	5P	10P
2	58.9	69.1	77.3	81.3	85.8
3	45.8	63.2	73.1	79.2	85.1
4	42.1	61.95	72.7	79.2	84.7
5	38.95	61.15	72	78.7	84.35
6	37.1	59.3	71.15	78.95	84.2
7	36.7	57.3	71.15	79.2	84.15
8	35.2	56.4	70.7	78.1	84.2
9	33.1	56.2	70.1	78.3	83.7
10	32.1	56.3	69.6	78.0	83.8
11	30.9	56.1	69.1	78.3	83.75
12	30.4	56.4	69.33	78.25	83.7
13	30.0	56.35	69.25	78.2	83.8
14	30.2	56.4	69.15	78.15	83.65
15	30.0	56.3	69.0	78.2	83.5
16	30.05	56.25	69.15	78.25	83.5

Tabela 5.6: Scalarea rețelei cu numărul de procesoare: Linux.

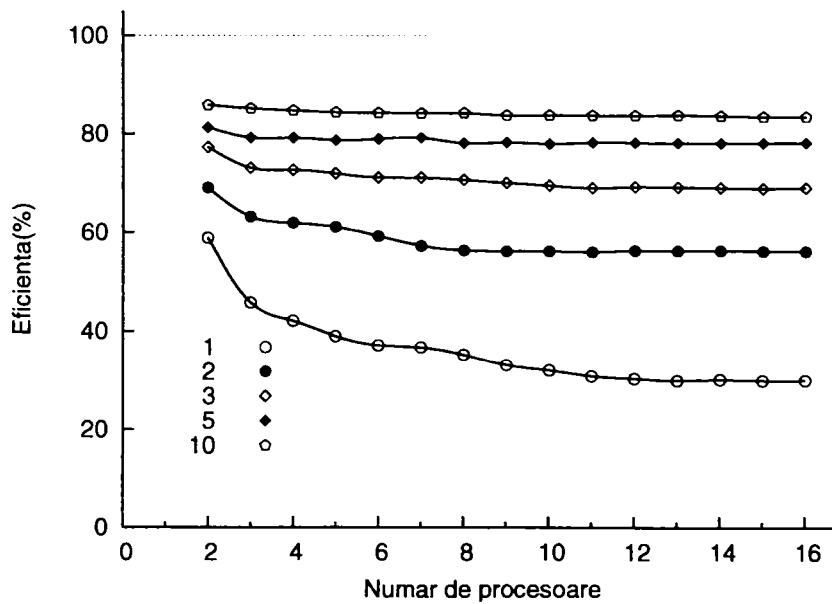


Figura 5.10: Eficiența în cazul scalării rețelei cu numărul de procesoare – Linux.

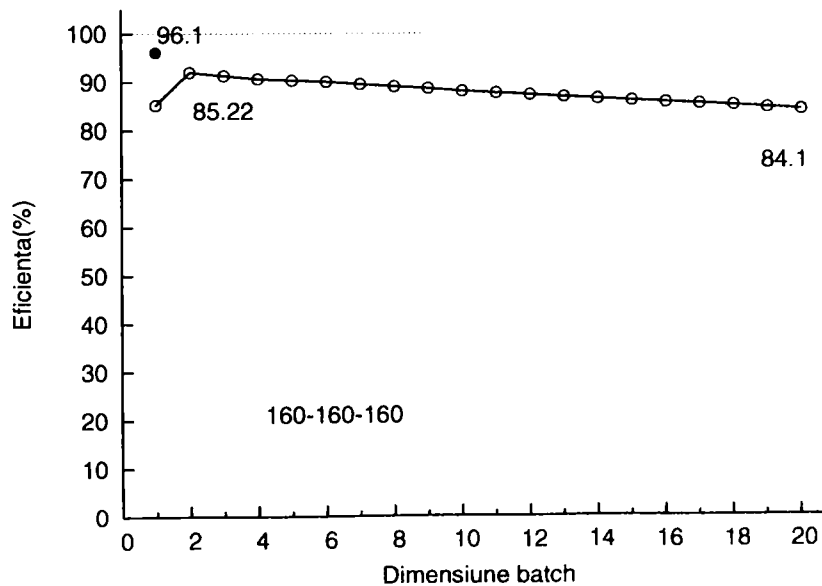
constă în însumarea simplă a gradientilor înainte de actualizarea ponderilor. Această modificare are însă un dezavantaj — nu mai este posibilă combinarea pașilor patru, cinci și șapte din algoritm.

Dimensiune batch	Eficiența	Dimensiune batch	Eficiența
1	85.22	1	83.0
2	92.00	2	90.1
3	91.31	3	89.1
4	90.59	4	88.7
5	90.33	5	88.45
6	90.01	6	88.15
7	89.52	7	87.6
8	88.99	8	87.1
9	88.61	9	86.3
10	88.03	10	86.2
11	87.62	11	85.9
12	87.21	12	85
13	86.84	13	84.7
14	86.41	14	84.45
15	86.02	15	84.15
16	85.64	16	83.7
17	85.33	17	83.25
18	85.00	18	83.1
19	84.51	19	82.8
20	84.09	20	82.0

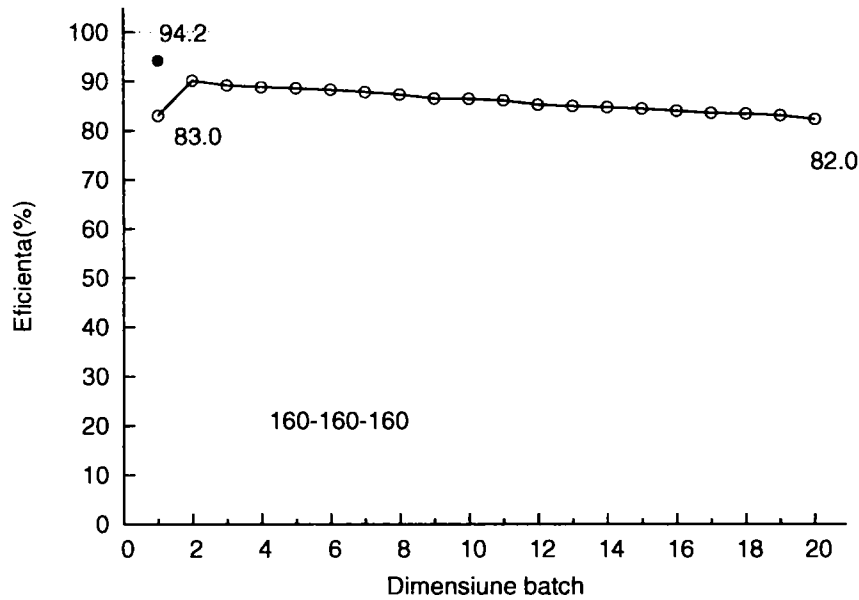
a) SunBlade: rețea 160-160-160

b) Linux: rețea 160-160-160

**Tabela 5.7:** *Eficiența pentru actualizarea la epocă*



**Figura 5.11:** *Eficiența pentru actualizarea la epocă: SunBlade*



**Figura 5.12:** *Eficiența pentru actualizarea la epocă: Linux*

În tabelul 5.7 și figurile 5.11, 5.12 este prezentată eficiența execuțiilor pentru numărul maxim de procesoare (16) cu modificarea dimensiunii setului de antrenare. De remarcă că pentru dimensiunea 1 s-a reprezentat atât varianta modificată (cerculețul negru) cât și implementarea originală (cerculețul alb). Pentru toate celelalte dimensiuni ale setului de antrenare a fost utilizată, desigur, varianta modificată. Implementarea modificată pare să fie mai rapidă decât varianta originală. Totuși, nu este chiar așa, întrucât cele două implementări paralele ar trebui comparate de fapt cu două implementări secvențiale diferite, iar implementarea secvențială a actualizării după fiecare tipar este mai rapidă decât versiunea corespunzătoare pentru dimensiunea 1 a setului de antrenare. În tabelele 5.8, 5.9 sunt prezentați timpii de execuție pentru 10 cicluri de antrenare pentru rețele 160 – 160 – 160.

Versiunea	16 procesoare paralele	Secvențial
Actualizare la tipar	0.52 sec	4.5 sec
Actualizare batch	0.64 sec	6.2 sec

**Tabela 5.8:** *Timpii de execuție – SunBlade.*

Versiunea	16 procesoare paralele	Secvențial
Actualizare la tipar	0.54 sec	4.65 sec
Actualizare batch	0.68 sec	6.34 sec

**Tabela 5.9:** *Timpii de execuție – Linux.*

În figurile 5.11, 5.12 se observă că eficiența implementării cu actualizarea în mod batch scade puțin atunci când sunt folosite seturi de antrenare de dimensiune

mare. Aceasta se datorează faptului că operațiile de calcul al actualizărilor ponderilor sunt efectuate mai puțin frecvent.

## 5.5 Concluzii

În această secțiune a fost implementat și testat un simulator paralel de backpropagation cu actualizarea la tipar – NetParSim. Algoritmul utilizează o partiționare a rețelei. Mai mult, comunicația a fost redusă la un nivel la care are un efect redus asupra algoritmului.

Algoritmul utilizează schema de actualizare la tipar, ceea ce înseamnă că actualizări foarte frecvente pot fi utilizate fără a influența algoritmul. Totuși, se poate utiliza eficient un număr mare de procesoare doar dacă se simulează rețele neuronale de dimensiune mare.

Datorită acestor caracteristici, această strategie de paralelizare, cu partiționarea rețelei este foarte diferită de varianta cu partiționarea datelor discutată în capitolul anterior. Eficiența acelor algoritmi cu partiționarea datelor era aproape independentă de dimensiunea rețelei simulate, iar dimensiunea setului de antrenare era un factor determinant pentru creșterea eficienței algoritmului de partiționare a datelor.

O observație importantă rezultă din experimentele care au fost prezentate în secțiunea 5.4.1: performanțele simulatorului sunt afectate în principal de numărul de neuroni de pe stratul de intrare și cel intermediar, și mult mai puțin de numărul de neuroni de pe stratul de ieșire.

**Contribuțiile** acestui capitol sunt:

1. Dezvoltarea și implementarea unui algoritm de paralelizare cu partiționarea rețelei. Algoritmul divide rețeaua neuronală, astfel încât fiecare mașină care participă la simulator va deține un număr de neuroni de pe fiecare strat al rețelei neuronale.
2. Reformularea calculelor care se efectuează în algoritmul backpropagation, în vederea distribuirii lor pe o rețea de calculatoare, cu minimizarea timpului necesar transmisiei de date.
3. Evaluarea necesarului de memorie al algoritmului prezentat. Pe baza formulelor prezentate (formula 5.15) se poate face o evaluare a dimensiunii maxime a rețelei neuronale posibil de simulat de către NetParSim.
4. Studiul experimental al influenței factorilor care pot afecta performanțele evoluției simulatorului într-o rețea de calculatoare. Concluzia care se poate deduce este că performanțele simulatorului vor fi cu atât mai bune cu cât numărul de neuroni de pe fiecare strat de procesat pe fiecare mașină este mai mare. Din experimentele efectuate, s-au observat valori satisfăcătoare pentru accelerare, respectiv eficiență atunci când rețeaua neuronală a avut asemenea dimensiuni încât să fie alocați cel puțin zece neuroni de pe fiecare strat, pentru fiecare mașină din rețea. Investigații pe această temă au fost prezentate în [Ba07a] și [BCP07].

## 6 Experimente cu aplicații uzuale

Capitolul de față prezintă rezultatele obținute pe simulatoarele prezentate în capitolele anterioare pentru un set de aplicații uzuale, des utilizate pentru măsurarea performanțelor rețelelor neuronale.

Aplicațiile utilizate în experimente sunt listate în tabelul 6.1.

Aplicația	Dimensiune rețea	Nr. tipare
NETtalk [81]	203 × 240 × 26	5438
Recunoaștere caractere (set de fonturi X11)	40 × 30 × 26	26
Compresie de imagine [13]	64 × 16 × 64	4096
Recunoaștere de imagini/OCR	1024 × 512 × 64	4096
Recunoaștere ținte sonar [29], [28]	60 × 24 × 2	208

**Tabela 6.1:** Aplicații de back-propagation utilizate pentru măsurarea performanțelor antrenării.

### 6.1 NETtalk

În capitolele anterioare am analizat două implementări paralele ale unor simulatoare de rețele neuronale strict din punctul de vedere al eficienței. În cazul utilizării unui set de antrenare de dimensiuni mari, varianta cu paralelizarea datelor poate părea superioară celei cu partiționarea rețelei întrucât se pot obține eficiențe ridicate. Totuși, asemenea considerații nu sunt suficiente în contextul rețelelor neuronale. Modificarea dimensiunii setului de antrenare modifică de asemenea chiar procesul de învățare. Se va evidenția în continuare importanța actualizării frecvente a ponderilor și influența nefastă asupra antrenării în cazul NETtalk când se utilizează un set de antrenare de dimensiune mare. Ca urmare, avantajul utilizării unei scheme cu partiționarea datelor poate să nu mai existe.

În continuare vor fi prezentate rezultatele rulării unei aplicații cunoscută sub numele de NETtalk. Au fost utilizate frecvențe specifice de actualizare a ponderilor, întrucât ele au fost testate și de către alți cercetători pe diferite arhitecturi paralele (vezi [22], [7], [20], [100] și [104]).

#### 6.1.1 Setul de date NETtalk

În această secțiune va fi descris setul de date NETtalk, care a fost construit de Sejnowski și Rosenberg [81] în 1986. Sejnowski și Rosenberg au antrenat o rețea neuronală feed-forward pentru pronunțarea cuvintelor în limba engleză. Setul de date este larg utilizat ca termen de referință în domeniul rețelelor neuronale pentru a testa viteza de antrenare și abilitățile de generalizare. A fost astfel posibilă compararea rezultatelor proprii cu cele obținute de alți cercetători.

Setul complet de date constă din 20008 cuvinte în limba engleză, împreună cu fonemele corespunzătoare și informațiile de accentuare. Modul normal de utilizare a datelor este însă să se facă învățarea doar pe un subset, circa 1000 de cuvinte uzuale. Cuvintele rămase pot fi apoi utilizate pentru a testa abilitățile de generalizare ale rețelei.

Setul de date este astfel construit încât să existe o corespondență unu-la-unu între literele unui cuvânt și foneme și informația de accentuare. De exemplu cuvântul "pronounce" are ca și corespondent un șir de foneme de aceeași lungime (prɒnʌns-) și un șir de simboluri de accentuare (>>>1<<<<). De fapt, simbolurile de accentuare constă atât din informație despre accente, cât și despre silabe. În total sunt 52 de foneme diferite și 5 simboluri de accentuare diferite. [81] conține informații mai detaliate despre aceste simboluri.



Traducerea unei singure litere dintr-un cuvânt în fonemul corect depinde de literele înconjurătoare din cuvânt, de aceea se utilizează o "fereastră" de 7 caractere. Rețeaua neuronală trebuie să genereze fonemul corespunzător literei din centrul ferestrei. Cuvântul defilează prin fereastră astfel încât toate literele din cuvânt sunt plasate, pe rând, în centrul ferestrei de 7 caractere. Se prezintă rețelei neuronale toate literele din fereastră, după cum este ilustrat în figura 6.1. În figură, rețeaua trebuie să genereze fonemul "x", care corespunde literei "o" (din centrul ferestrei). Pentru aceasta, se prezintă rețelei toate caracterele din fereastră, adică "\_pronou".

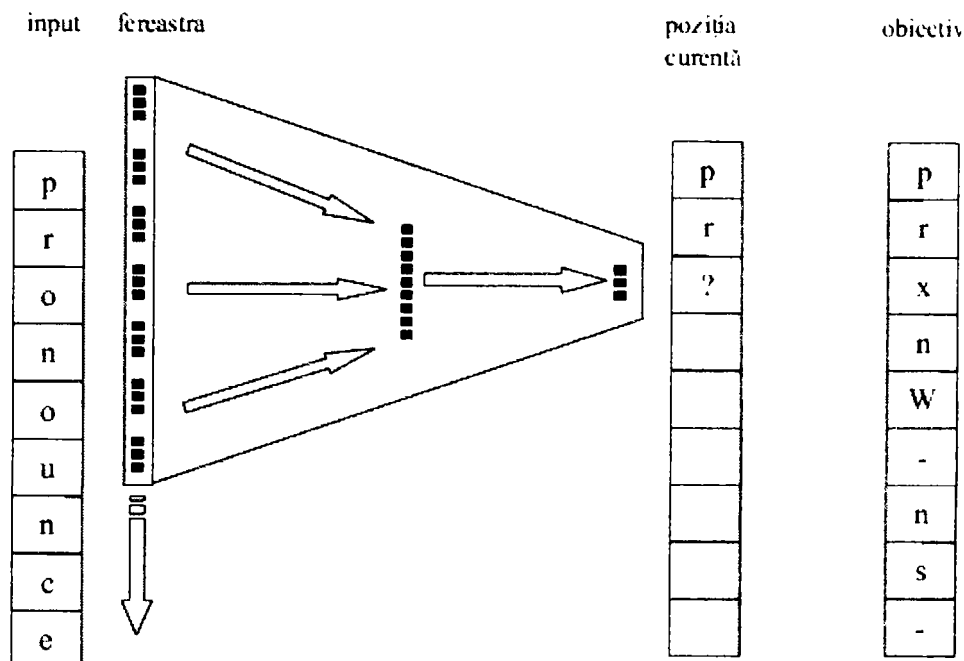


Figura 6.1: Fereastră de prezentare utilizată de NETtalk.

Utilizând o fereastră de 7 caractere, ca și Sejnowski și Rosenberg, este evident că rețeaua neuronală nu va putea pronunța corect toate cuvintele. Un exemplu ar fi cazul cuvintelor *finite* și *infinite*. Prezentând rețelei șirul "finite\_" în fereastră de 7 caractere, pentru a determina fonemul pentru litera din centrul ferestrei, rețeaua nu are cum să știe dacă întregul cuvânt este *finite* sau *infinite*. Pronunțarea literei 'i' (*finite* sau *infinite*) diferă la cele două cuvinte. Au mai fost date ca exemple *though* și *thought*, și altele. Desigur, problema se rezolvă utilizând o fereastră mai mare. Totuși, nu s-a urmărit construirea rețelei neuronale perfecte pentru NETtalk, ci doar să se poată compara rezultatele obținute cu ale altor cercetători, iar majoritatea au utilizat o fereastră de 7 litere.

### 6.1.2 Implementare

A fost utilizată o reprezentare ortogonală a fiecărei litere pentru unitățile de intrare, cu un alfabet cu 26 de litere și o dimensiune a ferestrei de citire de 7, rezultând un total de 182 de unități de intrare. Aceasta înseamnă că pentru fiecare literă din fereastră este activă numai una dintre cele 26 de unități de intrare corespunzătoare. Dacă una dintre litere este caracterul spațiu – nu este activată nici una dintre cele 26 de unități de intrare.

Spre deosebire de unitățile de intrare, Sejnowski și Rosenberg utilizează o reprezentare distribuită pentru unitățile de ieșire. În total sunt 52 de foneme, exprimate prin intermediul a 21 de atribute de articulare, cum ar fi punctul de articulare, înălțimea vocalei și așa mai departe. Cele cinci simboluri de accentuare sunt reprezentate ortogonal, având ca rezultat 26 de unități de ieșire. Pentru a putea compara rezultatele, în general se utilizează 60 de unități pe stratul ascuns. Desigur, ar fi posibil să se obțină performanțe mai bune în privința antrenării dacă s-ar utiliza mai multe unități ascunse, însă, după cum s-a mai precizat, nu s-a căutat să se îmbunătățească performanțele rețelelor neuronale, ci doar să se poată compara rezultatele.

Setul de date NETtalk, compus din cuvinte, foneme și simboluri de accentuare, este foarte mic în comparație cu alte aplicații de rețele neuronale. În total, el necesită circa 15 Kbytes pentru un set de date de 1000 de cuvinte.

Disponând de un set de date atât de mic, a fost posibil să se plaseze câte o copie a întregului set de date în fiecare dintre procesoare, ceea ce prezintă un avantaj pentru paralelizarea cu partiționarea rețelei. În acest fel, a fost nevoie numai de o distribuire inițială a datelor, iar pe parcursul simulării nu a mai fost necesară altă comunicație. Toate procesoarele au acces direct la toate tiparele. Ca urmare, pasul unu din paralelizarea cu partiționarea rețelei este simplu, așa cum s-a menționat în secțiunea 5.1.5.

Întrucât este utilizată o reprezentare binară (și ortogonală) pentru unitățile de intrare, calculul activării unităților ascunse poate fi simplificat. Pentru o singură literă din fereastră, numai una dintre cele 26 de unități de intrare este activă (cu valoarea 1), celelalte 25 fiind inactive (valoarea 0). Ca urmare, în formula:

$$a_{pj}^H = f(\text{net}_{pj}^H) = f\left(\sum_{i=0}^{NI-1} a_{pi}^I w_{i \rightarrow j}^H\right) \quad (6.1)$$

nu este necesar să se efectueze toate înmulțirile. Este suficientă însumarea ponderilor de la unitățile de intrare cu activarea 1, lucru întotdeauna posibil atunci când se utilizează o reprezentare binară.

### 6.1.3 Simulări și rezultate

În această secțiune va fi prezentat efectul alegerii dimensiunii setului de antrenare asupra capabilităților de antrenare. Au fost utilizate câteva dimensiuni specifice, întrucât ele au mai fost aplicate și de alte implementări paralele pe diferite arhitecturi [68], [100].

Rezultatele obținute corespund celor obținute de Sejnowski și Rosenberg. Pentru fiecare execuție s-a precizat atât rata de învățare cât și valoarea momentului, în vederea reproducerii acestor rezultate.

Sejnowski și Rosenberg au utilizat două măsuri ale performanței. Ieșirea este considerată *potrivire perfectă*<sup>1</sup> dacă valoarea fiecărei unități de ieșire diferă cu cel mult 0.1 față de valoarea dorită.

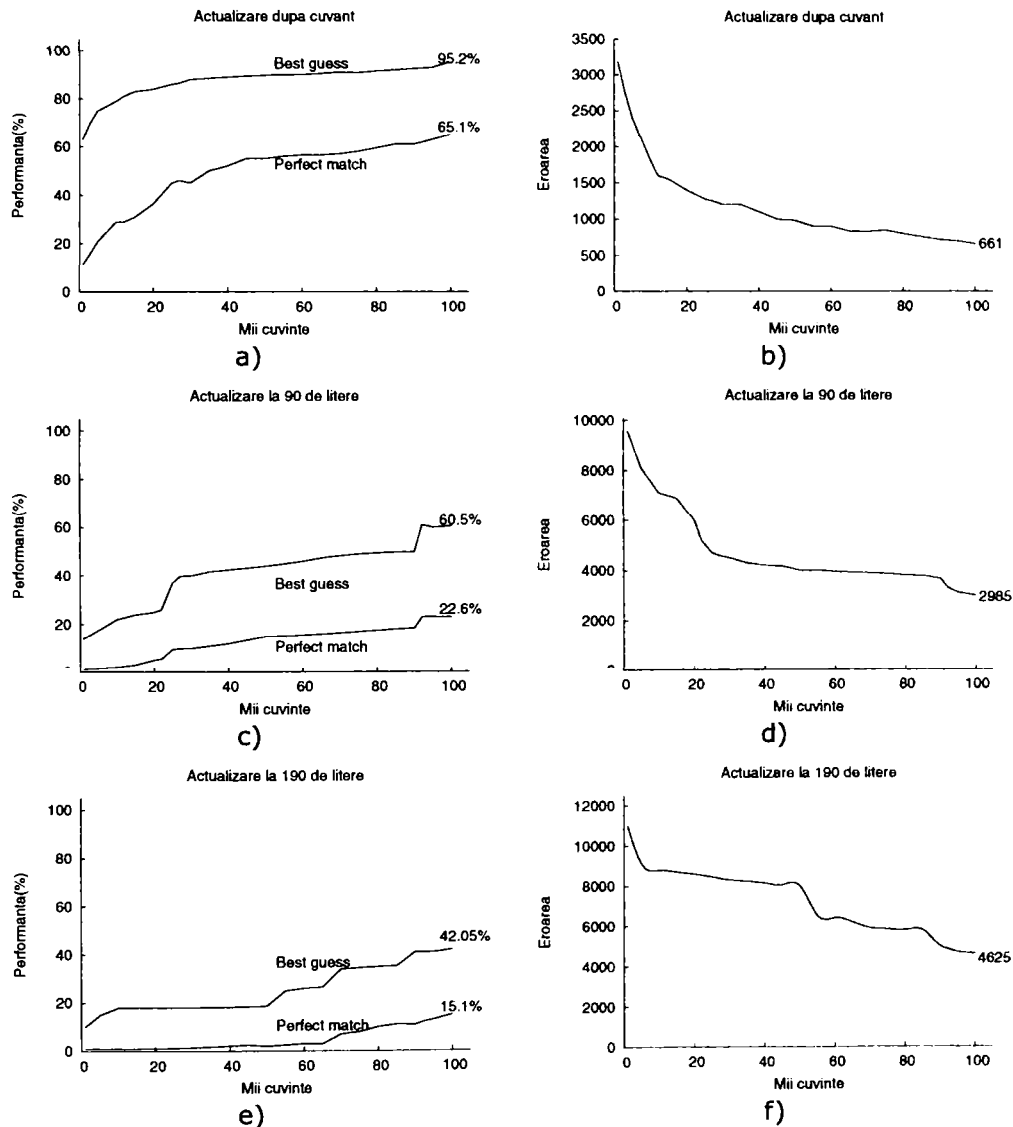
Ieșirea este considerată *cea mai bună alegere*<sup>2</sup> dacă vectorul țintă corect este acel vector dintre cei 52 de vectori țintă posibili care formează cel mai mic unghi cu vectorul actual al activării unităților de ieșire. Aceasta înseamnă că toți cei 52 de vectori țintă posibile (fiecare compus din 26 de valori) se compară cu vectorul generat la ieșire și se determină vectorul țintă cel mai apropiat. Dacă acest vector țintă este vectorul țintă corect, ieșirea reprezintă cea mai bună alegere. Ca urmare, vectorul de ieșire poate să fie foarte diferit de vectorul țintă corect atunci când se utilizează acest criteriu.

<sup>1</sup>"perfect match"

<sup>2</sup>"best guess"

Evident, criteriul celei mai bune alegeri va da cel mai mare procentaj de răspunsuri corecte (o potrivire perfectă este întotdeauna cea mai bună alegere, dar inversa nu este întotdeauna adevărată). Din păcate, criteriul celei mai bune alegeri necesită resurse computaționale mult mai mari.

Pentru antrenare a fost utilizat un set de 1000 de cuvinte mai frecvente din limba engleză. După fiecare prezentare a celor 1000 de cuvinte, rețeaua a fost testată pe aceleași 1000 de cuvinte, pentru a-i măsura performanțele. Graficele din figura 6.2 prezintă performanțele rețelei la antrenarea cu diferite frecvențe de actualizare a ponderilor. În toate sesiunile de antrenare, fiecare cuvânt a fost prezentat de câte 100 de ori, în ordine aleatoare. Valoarea termenului moment a fost cea uzuală de 0.9.



**Figura 6.2:** Performanțele NETtalk după prezentarea a 100000 de cuvinte.

Graficele reprezintă media mai multor rulări, cu puncte de pornire aleatoare în

spațiul ponderilor, însă diferențele au fost ne semnificative, lucru explicabil dacă se iau în considerare numărul relativ mare de ponderi din rețea (12566, incluzând și ponderile de bias).

În figurile 6.2a și 6.2b sunt prezentate rezultatele în cazul actualizării după cuvânt — ponderile sunt actualizate numai după ce au fost prezentate toate literele unui cuvânt. Aceasta este frecvența de actualizare utilizată de Sejnowski și Rosenberg. De remarcat faptul că ponderile sunt actualizate foarte frecvent, însă neregulat, întrucât cuvintele diferă în lungime (între 1 și 14, având o medie de 5.4 pentru cele 1000 de cuvinte. Rata de învățare utilizată a fost de 0.2.

În figura 6.2a este prezentată performanța în procente, măsurată atât ca potrivire perfectă cât și pentru cea mai bună alegere. Se observă că în cazul celei mai bune alegeri performanța crește rapid până la aproape 90% și apoi mult mai încet, atingând 94.3% după 100000 de prezentări de cuvinte. Pentru potrivirea perfectă avem o creștere mai moderată, atingând circa 64%, însă se pare că am putea obține un procentaj și mai bun dacă s-ar continua antrenarea.

În figura 6.2b este prezentată eroarea standard, calculată ca:

$$E = \sum_{p=1}^{5438} \sum_{j=1}^{26} (t_{pj} - a_{pj})^2 \quad (6.2)$$

pentru cele 26 de unități de ieșire și cele 5438 de tipare (numărul de litere din cele 1000 de cuvinte).

În figurile 6.2c și 6.2d sunt prezentate rezultatele în cazul actualizării ponderilor după prezentarea a 90 de litere (dimensiunea setului de antrenare este 90). Pomerleau și alții [68] au utilizat această frecvență de actualizare a ponderilor pe o arhitectură Warp cu 10 procesoare. A fost utilizată o rată de învățare de 0.02, necesară pentru a preveni *oscilațiile*: dimensiunea setului de antrenare este mare, deci se însumează mulți gradienti componenți, iar pentru a păstra modificarea ponderilor la o magnitudine acceptabilă este necesară o asemenea rată mică de învățare. Dacă modificările ponderilor sunt prea mari se obțin oscilații și nu are loc un proces de învățare.

Se observă în figura 6.2c că rețeaua are o performanță scăzută pentru ambele criterii, în comparație cu actualizarea la cuvânt. De asemenea, se remarcă faptul că graficele sunt neregulate. Mai mult, eroarea rămâne mare chiar și după 100000 de prezentări de cuvinte, după cum se observă în figura 6.2d.

În ultimul set de grafice sunt reprezentate rezultatele pentru actualizarea ponderilor după prezentarea a 190 de litere. A fost necesară o rată de învățare de 0.01, pentru a preveni oscilațiile. Se constată o performanță și mai slabă decât în cele două cazuri precedente.

Cele două grafice — pentru actualizarea după 90 de litere, respectiv după 190 de litere, par să indice faptul că s-ar putea crește performanțele dacă s-ar prelungi durata de antrenare. Experimentele efectuate au confirmat această supoziție, în cazul actualizării după 90 de litere și circa 256000 de prezentări de cuvinte, însă pentru actualizarea după 190 de litere rezultatele au fost nesatisfăcătoare și indicau necesitatea unui număr mult mai mare de reluări.

Au fost efectuate experimente și pentru dimensiuni ale setului de antrenare mai mari — 512, respectiv 800 de litere, însă performanțele obținute au fost mult prea slabe.

Concluzia care se poate deduce din rezultatele obținute este că pe măsură ce ponderile se actualizează mai rar, tot mai multe prezentări de tipare sunt necesare pentru a obține aceeași calitate a antrenării. La o concluzie similară au ajuns atât Bourrelly [7], la paralelizarea algoritmului de backpropagation pe un hiper-cub cu 32 de procesoare, cât și Witbrock [100], pe o arhitectură GF11 cu 512 procesoare.

## 6.2 Compresie de imagine

O altă utilizare a rețelelor neuronale multistrat este comprimarea datelor dintr-o imagine prin reducerea redundanței spațiale [83].

### 6.2.1 Suportul teoretic

Structura unei astfel de rețele neuronale este ilustrată în figura 6.3. Rețeaua conține trei straturi: un strat de intrare, unul de ieșire și unul ascuns. Stratul de intrare și stratul de ieșire au același număr de neuroni ( $NI = NO$ ). Compresia de imagine are loc prin alegerea unei valori pentru  $NH$  — numărul de neuroni de pe stratul ascuns — mai mică decât numărul de neuroni de pe stratul de intrare și cel de ieșire. Astfel, ieșirea stratului ascuns este o imagine comprimată, iar stratul de ieșire de-comprimă imaginea. Un număr mic de unități pe stratul ascuns conduce la un raport mare de compresie.

Imaginea de intrare este divizată în blocuri, de exemplu  $8 \times 8$ ,  $4 \times 4$  sau  $16 \times 16$ . Aceste blocuri trebuie convertite în vectori de intrare cu valori normalizate ( $x_i \in [0, 1]$ ) corespunzătoare valorilor pixelilor pentru nivele de gri în  $[0, 255]$ , întrucât s-a demonstrat [6], [41] faptul că rețelele neuronale lucrează mai eficient atunci când atât intrarea cât și ieșirea sunt limitate la domeniul  $[0, 1]$ .

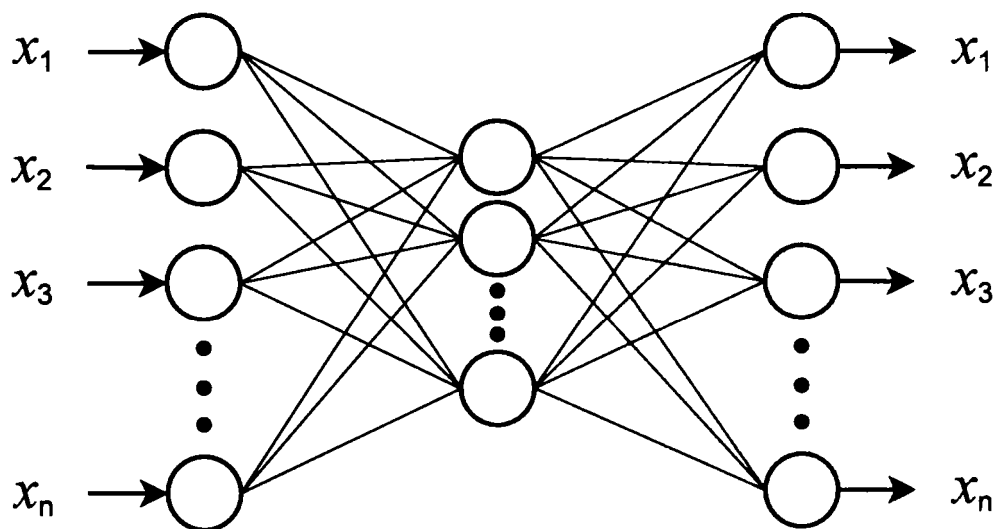


Figura 6.3: Rețea neuronală pentru compresia de imagine.

Dacă un vector de intrare este de dimensiune  $NI$ , egal cu numărul de pixeli din fiecare bloc, atunci se pot nota cu  $w_{ji}$  ponderile conexiunilor dintre stratul de intrare și cel ascuns, cu  $i \in [1, NI]$  și  $j \in [1, NH]$ , care mai pot fi considerate ca o matrice de dimensiune  $NH \times NI$ . Ponderile conexiunilor dintre stratul ascuns și cel de ieșire, de asemenea pot fi notate cu  $w'_{ij}$ ,  $i \in [1, NO]$  și  $j \in [1, NH]$ , care formează de asemenea o matrice, de dimensiune  $NO \times NH$ . Compresia imaginii poate fi descrisă matematic după cum urmează:

$$h_j = \sum_{i=1}^{NI} w_{ji} x_i, \quad 1 \leq j \leq NH \quad (6.3)$$

pentru compresie și

$$\bar{x}_i = \sum_{j=1}^{NH} w'_{ij} h_j, \quad 1 \leq i \leq NO \quad (6.4)$$

pentru decompresie.

Utilizând o asemenea rețea neuronală, compresia de imagine necesită o fază de antrenare. În această fază, trebuie selectat un set de eșantioane pentru antrenarea rețelei cu algoritmul backpropagation, având fiecare vector de intrare ca ieșire dorită. Astfel, se comprimă intrarea printr-un canal mai îngust, reprezentat de stratul ascuns, după care stratul de ieșire reconstruiește imaginea.

Pentru utilizare, după prezentarea unui bloc de imagine, se preia ieșirea stratului ascuns, ceea ce reprezintă imaginea comprimată, care apoi poate fi reconstruită utilizând ponderile conexiunilor dintre stratul ascuns și cel de ieșire, conform relației 6.4.

### 6.2.2 Simulări și rezultate

Compresia de imagine necesită o rețea de dimensiuni reduse [13], însă presupune un volum mare de date, atât pentru antrenare cât și pentru compresie în timp-real. Ca urmare, s-a considerat că varianta cea mai indicată este cea de paralelizare cu partiționarea datelor.

Pentru efectuarea experimentelor au fost necesari mai mulți pași:

**Achiziția de imagine** — imaginile au fost scanate cu ajutorul unui scanner și salvate în format TIFF (Tagged Image Format). Apoi, cu ajutorul programului GIMP, au fost convertite în format BMP, cu nivele de gri (monocrom).

**Normalizare** — întrucât valorile pixelilor rezultați erau întregi, între 0 și 255, a fost necesară preprocesarea lor, pentru a avea valori în intervalul [0, 1].

**Segmentarea imaginii** — imaginea a fost divizată în blocuri mai mici,  $8 \times 8$ , și convertite în vectori. Dacă dimensiunea imaginii nu a fost divizibilă cu dimensiunea blocului dorit, imaginea a trebuit completată cu pixeli cu valoarea zero.

**Pregătirea setului de antrenare** — s-a generat un fișier în care fiecare bloc a fost scris sub forma a doi vectori identici – perechea de antrenare.

**Eliminarea blocurilor similare** — după segmentare este posibil să rezulte blocuri identice, ceea ce reprezintă informație redundantă și care ar putea chiar să deranjeze în faza de antrenare.

**Antrenarea rețelei neuronale** — se furnizează rețelei neuronale tiparele de antrenare și se efectuează procesul de învățare până când se atinge nivelul de eroare acceptabil.

**Realizarea unui program pentru reconstrucția imaginii** — s-au construit două rețele neuronale, cu câte două straturi în loc de trei, având ponderile conexiunilor rezultate la faza de antrenare. Una dintre rețele va efectua compresia imaginii, iar cealaltă – reconstrucția.

În figurile 6.4 și 6.5 sunt reprezentate performanțele unei rețele cu 64 de unități la intrare și 64 de unități la ieșire, pentru diferite rapoarte de compresie – numărul de unități de pe stratul ascuns de 4, 8, respectiv 16. Graficele prezintă accelerarea față de varianta secvențială pentru faza de antrenare, pentru 16 imagini  $256 \times 256$ , măsurând timpul necesar antrenării pentru a ajunge la o eroare sub 0.1, atât pentru varianta secvențială, cât și pentru cea paralelă.

După cum se poate observa, s-a obținut o accelerare mai bună – și o eficiență mai bună – comparativ cu varianta secvențială, în cazul în care rețeaua neuronală a

Număr procesoare	Dimensiune rețea					
	64-16-64		64-8-64		64-4-64	
	acc	ef	acc	ef	acc	ef
2	1.80	90.00	1.79	89.50	1.75	87.50
3	2.52	84.00	2.50	83.33	2.32	77.33
4	3.27	81.75	3.06	76.54	2.83	70.75
5	4.00	80.00	3.72	74.40	3.35	67.00
6	4.68	78.00	4.36	72.67	3.94	65.60
7	5.43	77.50	5.03	71.86	4.55	64.97
8	6.16	77.00	5.67	70.88	5.13	64.10
9	6.88	76.50	6.36	70.72	5.73	63.62
10	7.62	76.25	7.06	70.60	6.34	63.38
11	8.36	76.00	7.76	70.55	6.95	63.18
12	9.09	75.75	8.46	70.48	7.56	63.02
13	9.81	75.50	9.13	70.21	8.19	62.98
14	10.50	75.00	9.75	69.62	8.80	62.90
15	11.18	74.50	10.35	69.01	9.67	62.47
16	11.90	74.38	10.95	68.44	9.98	62.38

Tabela 6.2: Performanțele compresiei de imagini: SunBlade150

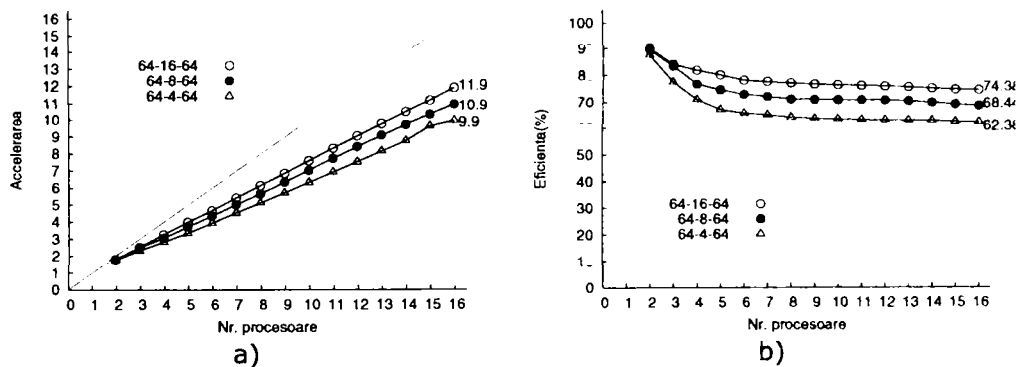


Figura 6.4: Performanțele compresiei de imagini: SunBlade150

avut cele mai multe unități pe stratul ascuns (16), deși aceasta duce la compresia cea mai slabă, însă au fost comparați timpii necesari atingerii unei rate a erorii totale de sub  $10^{-2}$ , fără a lua în considerație performanțele compresiei de imagine.

Nu se observă diferențe semnificative între rezultatele pentru cele două platforme hardware utilizate, SunBlade, respectiv Linux. În ambele cazuri, accelerarea, respectiv eficiența cea mai bună se obține pentru numărul cel mai mare de neuroni de pe stratul ascuns. Numărul mai mare de neuroni pe stratul ascuns duce la un număr mai mare de ponderi în rețea, deci un volum mai mare de calcul, iar rezultatele confirmă concluziile din capitolele anterioare, și anume faptul că eficiența este cu atât mai bună cu cât volumul de calcul *per procesor* este mai mare.

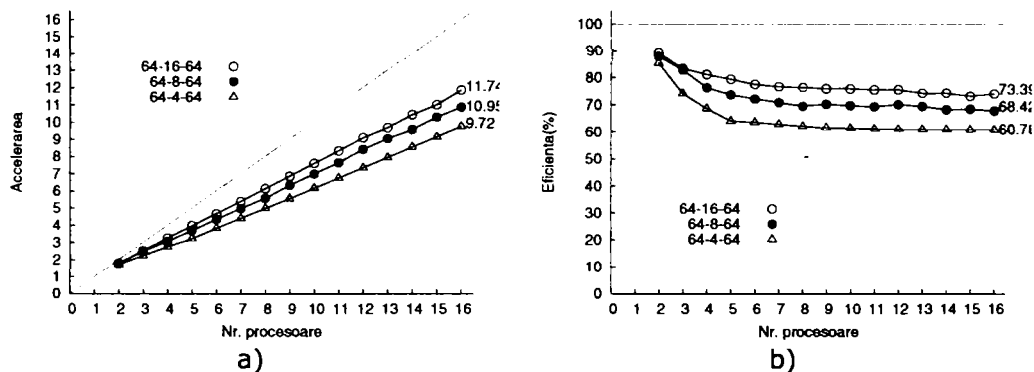
### 6.3 Recunoașterea țintelor de sonar

O altă aplicație utilizată în testarea rețelelor neuronale este cea de recunoaștere a țintelor de sonar. Gorman și Sejnowski [29], [28] au realizat un studiu în care au



Număr procesoare	Dimensiune rețea					
	64-16-64		64-8-64		64-4-64	
	acc	ef	acc	ef	acc	ef
2	1.79	89.29	1.76	88.09	1.71	85.51
3	2.51	83.53	2.49	82.99	2.23	74.25
4	3.25	81.30	3.06	76.41	2.75	68.71
5	3.98	79.55	3.69	73.78	3.21	64.10
6	4.66	77.65	4.34	72.25	3.82	63.60
7	5.38	76.81	4.97	70.94	4.40	62.90
8	6.12	76.51	5.56	69.51	4.97	62.10
9	6.83	75.90	6.31	70.15	5.54	61.59
10	7.59	75.88	6.97	69.65	6.14	61.38
11	8.31	75.51	7.62	69.27	6.73	61.14
12	9.07	75.57	8.40	70.02	7.32	61.02
13	9.65	74.26	9.02	69.40	7.92	60.95
14	10.40	74.28	9.54	68.18	8.53	60.90
15	10.97	73.11	10.26	68.38	9.12	60.81
16	11.83	73.97	10.83	67.72	9.72	60.78

**Tabela 6.3:** Performanțele compresiei de imagini: Linux



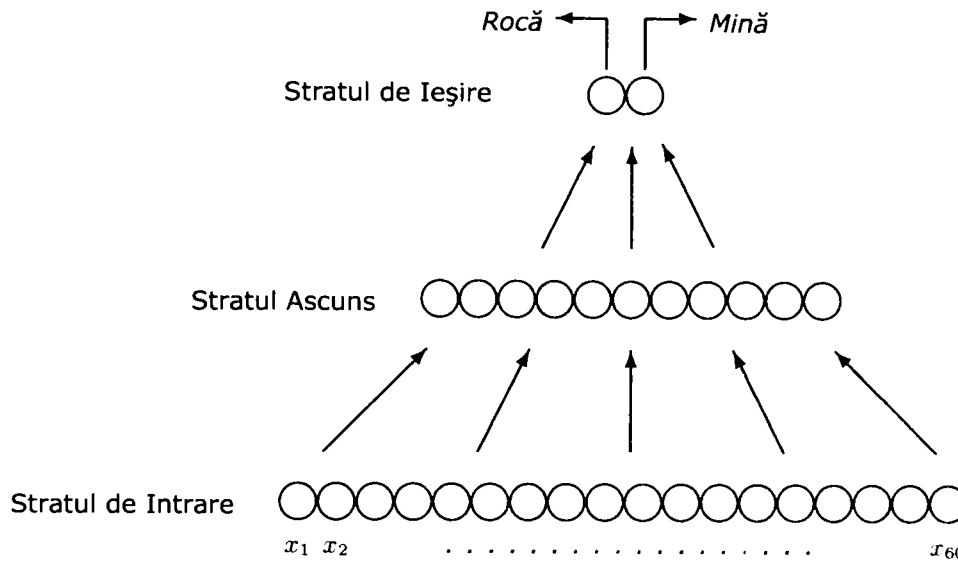
**Figura 6.5:** Performanțele compresiei de imagini: Linux

utilizat o rețea neuronală pentru diferențierea rocilor și a minelor de pe fundul mării cu ajutorul ecoului sonarului. Semnalele de sonar reflectate prezintă diferențe în privința volumului de energie pentru anumite frecvențe, în funcție de obiectul întâlnit. În experimentele lor, Gorman și Sejnowski au cuantizat frecvența răspunsului într-un set de câte 60 de eșantioane în care fiecare valoare reprezintă energia semnalului de răspuns, în cazul unui cilindru metalic, respectiv al unei roci de formă aproximativ cilindrică.

### 6.3.1 Configurația rețelei

Rețeaua neuronală utilizată este prezentată în figura 6.6. Rețeaua este complet conectată, cu 60 de unități de intrare cu valori reale, 24 de unități pe stratul ascuns și două unități pe stratul de ieșire ("rocă" / "mină").





**Figura 6.6:** Configurația rețelei pentru recunoașterea țintelor de sonar

### 6.3.1.1 Setul de Date

Pentru antrenare au fost utilizate două seturi de date: 111 tipare obținute prin reflexia pe un cilindru metalic, la diferite unghiuri, respectiv 97 de tipare provenind de la roci, în condiții similare.

Fiecare tipar constă dintr-un set de 60 de valori între 0.0 și 1.0, fiecare număr reprezentând energia într-o anumită bandă de frecvență, integrată peste un anumit interval de timp.

### 6.3.2 Simulări și rezultate

A fost aleasă o rată de învățare de 2.0 și un moment de 0.0. Ponderile inițiale au fost valori aleatoare uniform distribuite în intervalul  $[-0.3, +0.3]$ . Rețeaua a fost antrenată pe parcursul a 300 de epoci.

Întrucât pe stratul de ieșire sunt numai două unități, au fost efectuate experimente numai cu partiționarea datelor, pe un număr variabil de procesoare. Figurile 6.7, 6.8 prezintă graficele accelerării, respectiv eficienței în cazul celor două tipuri de mașini pe care s-au efectuat experimentele.

## 6.4 Concluzii

Algoritmul de antrenare cu backpropagation a fost paralelizat pe baza celor două abordări – cea cu partiționarea datelor, prezentată în capitolul 6, respectiv partiționarea rețelei, prezentată în capitolul 7. După cum s-a văzut, atunci când dimensiunea setului de antrenare este mare, varianta cu partiționarea datelor poate utiliza în mod eficient mai multe procesoare, independent de dimensiunile rețelei neuronale<sup>3</sup>.

Prin comparație, varianta cu partiționarea rețelei prezintă o serie de avantaje:

- nu depinde de dimensiunea setului de antrenare;

<sup>3</sup>De fapt, aproape independent de dimensiunile rețelei neuronale

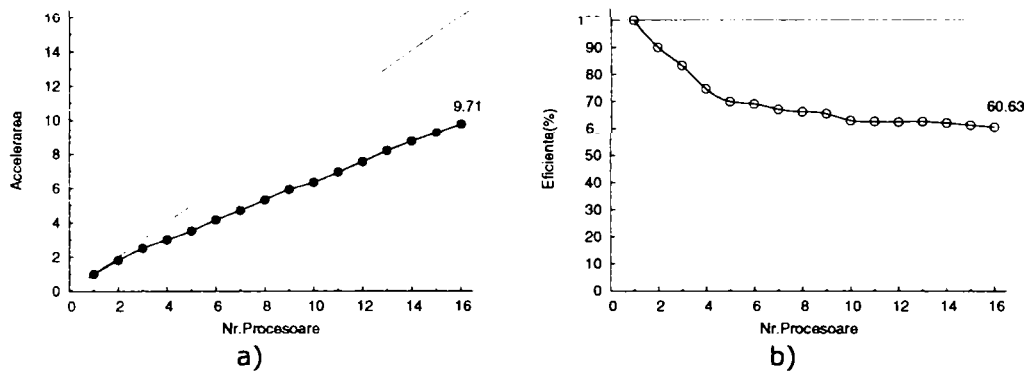


Figura 6.7: Performanțele aplicației Sonar: SunBlade150.

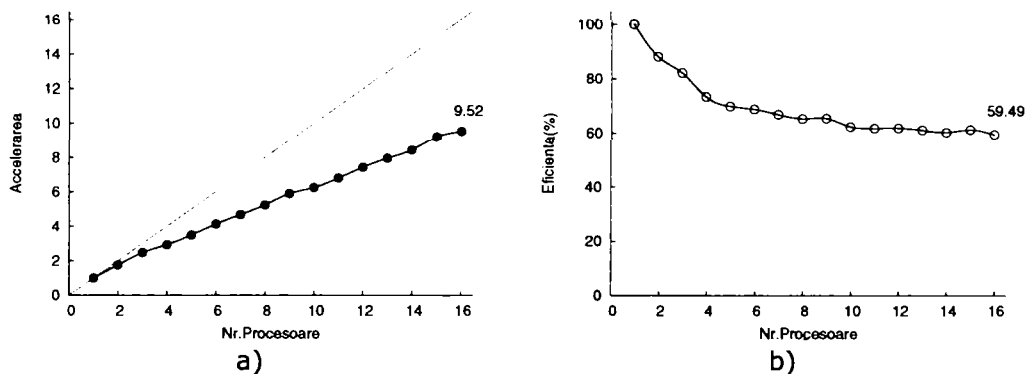


Figura 6.8: Performanțele aplicației Sonar: Linux.

- permite actualizări frecvente ale ponderilor în rețea;

însă necesită dimensiuni mari ale rețelei neuronale pentru a utiliza cât mai eficient mai multe procesoare.

Experimentele efectuate cu aplicația NETtalk au arătat că paralelizarea cu partiționarea datelor este superioară celei cu partiționarea rețelei, dacă luăm ca unic criteriu numărul de operații efectuate, însă este posibil ca rezultatul antrenării să fie nul. Acest lucru a fost sesizat și de către alți cercetători, pe implementările paralele realizate [100], [7]. [68] arată că: *Metoda tradițională de măsurare a performanțelor unui simulator sub forma numărului de conexiuni pe secundă poate conduce la rezultate eronate. Mai precis, este posibil ca implementările de rețele neuronale care simulează un număr mare de tipare în paralel, și care deci necesită simularea mai multor tipare înainte de actualizarea ponderilor, să învețe mai încet și mai puțin robust decât implementările care actualizează ponderile mai frecvent.*

Ca urmare, singura modalitate corectă de comparație între diferitele implementări ale algoritmului backpropagation este să se utilizeze exact aceiași parametri: setul de antrenare, dimensiunea rețelei, frecvența de actualizare a ponderilor.

Alte aplicații au arătat că pentru obținerea unor performanțe bune sunt necesare volume mari de date de antrenare – vezi compresia de imagine – mult mai mari decât numărul de procesoare utilizate.

Implementarea aleasă trebuie deci să fie în funcție de aplicația de rețele ne-

urionale *dată*. Uneori, se obțin rezultate mai bune dacă se utilizează *actualizarea la epocă*, alteleori *actualizarea la tipar* se poate dovedi mai eficientă. Deși eficiența cu care sunt utilizate procesoarele este departe de cea optimă, timpul necesar obținerii unor rezultate concludente se reduce drastic, ceea ce va permite reluarea experimentelor, până la determinarea valorilor optime pentru parametrii antrenare — rată de învățare, momentul, frecvența de actualizare a ponderilor, dimensiunea rețelei neuronale.

### 7.1 Concluzii

În cadrul acestei teze a fost studiată paralelizarea algoritmului de antrenare backpropagation pentru rețele neuronale feed-forward. Principalele obiective ale studiului de față au vizat realizarea unor algoritmi de antrenare bazați pe metoda backpropagation, care, spre deosebire de alte cercetări în domeniu, care au vizat implementări paralele pe arhitecturi paralele și masiv paralele, să evolueze într-o rețea de calculatoare, precum și studiul performanțelor acestor algoritmi:

- Dezvoltarea și implementarea unui algoritm de antrenare bazat pe metoda partiționării datelor – mai multe procesoare rulează același algoritm de antrenare, pe seturi diferite de date; rezultatele antrenării se colectează, se evaluează și se decide dacă să se continue antrenarea, până la atingerea nivelului de învățare prescris.

Experimentele efectuate au demonstrat posibilitatea utilizării unei asemenea abordări în cazul în care pot fi utilizate seturi de tipare de antrenare de dimensiuni foarte mari. S-a constatat o relație directă între dimensiunile setului de tipare de antrenare și eficiența acestui algoritm.

Pe baza acestei abordări, s-au obținut accelerări considerabile ale procesului de învățare, în privința timpului total necesar, avantaj contrabalansat însă de necesarul excesiv de memorie, provenit din faptul că fiecare procesor trebuie să păstreze o copie a întregii rețele neuronale simulate. Acest dezavantaj major duce la o limitare – uneori drastică – a dimensiunii rețelei neuronale, și deci a problemelor abordate.

- Dezvoltarea și implementarea unui algoritm de antrenare bazat pe metoda de partiționare a rețelei – fiecare procesor simulează câte o porțiune din rețeaua neuronală, respectiv neuronii de pe fiecare strat sunt distribuiți cât mai uniform către toate procesoarele. Principala problemă care a trebuit rezolvată în această abordare este reducerea la minimum a comunicației între procesoare, având în vedere că este necesară transmisia de date între neuroni de pe straturi diferite, simulați de procesoare diferite.

Deși, comparat cu metoda partiționării datelor, accelerările obținute – și implicit eficiența – au fost mai mici, această variantă a prezentat o serie de avantaje, dintre care amintim faptul că este posibilă actualizarea ponderilor *la tipar*, modalitate care poate fi importantă în unele aplicații reale.

Un alt avantaj este posibilitatea extinderii limitărilor în dimensiunile rețelei neuronale posibil de simulat, la valori virtual imposibil de simulat în condiții acceptabile pe o singură mașină. Rezultatele experimentelor efectuate în capitolul 5 demonstrează faptul că eficiența acestui algoritm a fost cu atât mai bună și a fost posibilă utilizarea mai multor procesoare cu cât rețeaua neuronală de simulat a fost mai mare, și independent de dimensiunile setului de tipare de antrenare.

Rezultatele experimentelor efectuate pentru cele două simulatoare sugerează o serie de valori minime pentru dimensiunea setului de tipare de antrenare, pentru metoda cu paralelizarea setului de antrenare, respectiv pentru dimensiunea rețelei, în cazul metodei cu paralelizarea nodurilor, valori pentru care eficiența obținută va fi cel puțin satisfăcătoare. Aceste valori sunt precizate sintetic în tabelul 7.1. În tabel au putut fi incluse și recomandări pentru rețele de 32 de stații, întrucât graficele eficienței au fost linii aproape orizontale, după cum s-a precizat și la capitolele respective. Aceleași concluzii sugerează faptul că tabelul ar putea continua și pentru un număr mai mare de mașini care rulează fiecare dintre simulatoare.

În cazul metodei cu partiționarea datelor, s-a constatat că eficiența, respectiv accelerarea, iau valori acceptabile atunci când numărul de tipare din setul de

Nr. procesoare	Partiționarea datelor (Nr. tipare)	Partiționarea rețelei (dimensiune rețea neuronală)
4	40	40 × 40 × 4
8	80	80 × 80 × 8
16	160	160 × 160 × 16
32	320	320 × 320 × 32

**Tabela 7.1:** Dimensiuni minime optime.

antrenare, alocate fiecărei mașini, este de cel puțin zece tipare. Astfel, volumul computațional dintre două actualizări de ponderi, devine mult mai mare decât timpul necesar comunicației, și astfel accelerarea, respectiv eficiența nu vor fi afectate într-o măsură prea mare de timpul necesar transmisiei de date între mașini.

Pentru paralelizarea cu partiționarea rețelei, valori bune ale eficienței și accelerării au fost obținute atunci când fiecărei mașini i-au fost alocate cel puțin 10 unități de pe cel puțin două dintre cele trei straturi. Performanțele nu au fost prea mult afectate de faptul că pe cel de-al treilea strat au fost mai puțini neuroni — dar având cel puțin câte un neuron pentru fiecare procesor.

## 7.2 Rezumat al contribuțiilor

Lucrarea aduce o serie de contribuții în domeniul procesării paralele a fazei de antrenare a rețelelor neuronale:

1. dezvoltarea unui algoritm de distribuire în vederea execuției cu partiționarea datelor, cu scopul minimizării timpului de antrenare a unei rețele neuronale;
2. este propusă, și a fost implementată, o schemă de distribuire a task-urilor și de colectare a rezultatelor, care să respecte tehnica de antrenare cu backpropagation;
3. realizarea simulatorului DataParSim, utilizând o strategie avansată de partiționare a datelor: un program *administrator* transmite câte o pondere a rețelei neuronale către procesele *slave*, care o utilizează pentru propagarea tuturor tiparelor de antrenare;
4. evaluarea necesarului de memorie pentru fiecare dintre cei doi algoritmi cu partiționarea datelor, cu posibilitatea de determinare a dimensiunii maxime a rețelei neuronale care se poate antrena utilizând acești algoritmi;
5. dezvoltarea unui algoritm pentru partiționarea rețelei, care să permită depășirea limitelor impuse asupra problemelor de tratat. Algoritmul se bazează pe o secționare a rețelei neuronale și distribuirea cât mai echilibrată, având în vedere și minimizarea comunicației în rețeaua de calculatoare.
6. o schemă de calcul care reformulează expresiile utilizate în algoritmul clasic de backpropagation, care permite calculul valorilor  $\delta$  în cazul strategiei cu partiționarea rețelei.
7. realizarea simulatorului NetParSim, bazat pe strategia de paralelizare cu partiționarea rețelei.
8. evaluarea performanțelor simulatoarelor DataParSim și NetParSim, atât pentru rețele neuronale generice, cât și pentru aplicații uzuale.
9. enunțarea de recomandări în vederea utilizării cât mai eficiente a resurselor de calcul.

Investigațiile pe aceste teme au fost prezentate în [BaC03], [Ba07a] și [BCP07].

### 7.3 Concluzii finale

Au fost prezentate două strategii de paralelizare – cu partiționarea datelor, respectiv cu partiționarea rețelei – care au fost implementate pe o arhitectură alcătuită din mai multe calculatoare interconectate printr-o rețea locală (LAN). Nu se poate însă preciza care dintre cele două abordări este de preferat în general. Alegerea depinde de problema de rețele neuronale tratată și de frecvența cu care se fac actualizările de ponderi.

Dacă, din rațiuni specifice problemei tratate, se preferă un set mare de tipare de antrenare, probabil că se va utiliza o paralelizare cu partiționarea datelor, care, în cazul utilizării mai multor procesoare va conduce la o accelerare mai bună decât în cazul paralelizării cu partiționarea rețelei, cel puțin în cazul unei rețele neuronale de dimensiuni relativ mici.

Atunci când sunt necesare actualizări frecvente ale ponderilor trebuie utilizată paralelizarea cu partiționarea rețelei, întrucât eficiența acestui algoritm nu depinde de o dimensiune mare a setului de antrenare. În plus, este singurul algoritm care să permită actualizarea după un tipar. Totuși, pentru a obține o eficiență ridicată, fiecare procesor trebuie să prelucereze mai multe unități din fiecare strat al rețelei neuronale. Aceasta conduce la limitarea drastică a numărului de procesoare care pot fi utilizate pentru simularea unor rețele mici. Însă atunci când atât dimensiunile rețelei cât și ale setului de antrenare sunt reduse, nici durata unui ciclu de antrenare nu este mare, nici chiar pentru algoritmul secvențial.

Trebuie subliniat faptul că la alegerea simulatorului paralel nu trebuie să se aleagă o dimensiune mare a setului de tipare după care se face actualizarea doar pentru a obține o aplicare eficientă a uneia dintre metodele de paralelizare, întrucât mărirea dimensiunii setului după care se face actualizarea poate să reducă drastic viteza de învățare și să scadă calitatea învățării în privința posibilităților de generalizare. Acest lucru a fost subliniat în cazul aplicației NETtalk (secțiunea 6.1), când actualizarea la cuvânt – sau la tipar – a condus la performanțe mult mai bune.

De aceea, la compararea diferitelor implementări de simulatoare de rețele neuronale trebuie să se ia ca măsură timpul absolut de execuție necesar pentru a atinge un anumit nivel de perfecțiune pe o problemă specifică de rețele neuronale. Desigur, aceasta nu este o soluție satisfăcătoare, întrucât influența dimensiunii setului de antrenare după care se face actualizarea asupra vitezei de învățare și a posibilităților de generalizare variază de la problemă la problemă. Pot exista probleme în care actualizarea ponderilor după o epocă este superioară actualizărilor mai frecvente, în privința vitezei de învățare și a calității. Cu alte cuvinte, probabil că nu este posibil să se facă o recomandare *generală* în privința schemei de paralelizare de utilizat, ci doar pentru tipuri specifice de probleme de rețele neuronale.

### 7.4 Perspective de cercetare

Studiul de față demonstrează posibilitatea accelerării timpului de antrenare pentru rețele neuronale feed-forward cu algoritmul backpropagation, utilizând resursele disponibile în laboratoarele departamentului. Prin reducerea considerabilă obținută, teza de față oferă o unealtă utilă celor care explorează domeniul rețelelor neuronale, permițând reluarea multiplă a experimentelor, cu parametri diferiți, până la obținerea rezultatelor dorite.

Teza deschide calea unor investigații ulterioare, dintre care au fost identificate câteva:

- determinarea dimensiunilor optime a rețelelor neuronale
- adaptarea tehnicilor utilizate și pentru alte tipuri de rețele – nu doar feed-forward, ci și cele recurente sau cele parțial conectate.

- investigarea posibilităților de paralelizare a altor algoritmi de antrenare, un prim exemplu fiind algoritmul de antrenare *quickprop* [24] și cel al gradientului conjugat [59].
- studiul posibilităților de paralelizare a tehnicilor de preprocesare în vederea reducerii dimensionalității problemelor.
- experimente în vederea determinării utilității mai multor straturi intermediare.

- [1] Yaser S. Abu-Mostafa and Demetri Psaltis. Optical neural computers. *Scientific American*, 256(3):88–95 (Intl. ed. 66–73), March 1987.
- [2] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski. A learning algorithm for Boltzmann machines. *Cognitive Science*, 9:147–169, 1985.
- [3] I. Aleksander and H. Morton. *An Introduction to Neural Computing*. Chapman and Hall, 1990.
- [4] James A. Anderson. A simple neural network generating an interactive memory. *Mathematical Biosciences*, 14:197–220, 1972.
- [5] W. R. Ashby. *Design for a Brain*. Wiley & Sons, New York, 1952.
- [6] Y. Benbenisti, D. Kornreich, H. B. Mitchell, and P. A. Schaefer. New simple three-layer neural network for image compression. *Optical Engineering*, 36:1814–1817, June 1997.
- [7] Jean Bourrely. Paralelisation of a neural learning algorithm on a hypercube. In F. André and J. P. Verjous, editors, *Hypercube and Distributed Computers*, pages 219–229. Elsevier Science Publishers B.V. (North Holland), 1989.
- [8] V. Braitenberg. *Vehicles: Experiments in Synthetic Psychology*. MIT Press, Cambridge, MA, USA, 1984.
- [9] A. E. Bryson and Y.-C. Ho. *Applied Optimal Control*. Blaisdell, New York, 1969.
- [10] S.R. Cajal. *Histology of the nervous system of man and vertebrates*. Oxford University Press, New York, 1995.
- [11] G. F. Carey, Y. Shen, and R. T. McLay. Parallel conjugate gradient performance for least-squares finite elements and transport problems. *International Journal for Numerical Methods in Fluids*, 28:1421–1440, December 1998.
- [12] G. Chinn, K. A. Grajski, C. Chen, C. Kuzmaul, and S. Tomboulian. Systolic array implementations of neural nets on the maspar mp–1 massively parallel processor. In *International Joint Conference on Neural Networks*, volume 2, pages 169–173, San Diego, CA, USA, 1990.
- [13] K. B. Cho, Cheol Hoon Park, and Soo-Young Lee. Image compression using multi-layer perceptron with block classification and SOFM coding. In *Proc. WCNN'94 – World Congress on Neural Networks*, volume III, pages 26–31, Hillsdale, NJ, 1994. INNS, Lawrence Erlbaum.
- [14] M. A. Cohen and Steven Grossberg. Stability of global pattern formation and parallel memory storage by competitive neural networks. *IEEE Transactions on Systems, Man and Cybernetics*, 13:815–826, 1983.
- [15] Y. Le Cun. Une procédure d'apprentissage pour Réseau à seuil assymétrique. *Cognitiva 85: A la Frontière de l'Intelligence Artificielle des Sciences de la Connaissance des Neurosciences*, pages 599–604, 1985.
- [16] A. D'Acierno and R. Vaccaro. A parallel implementation of the back-propagation of errors learning algorithm on a simd parallel computer. In Kappen B. and Gielen S., editors, *Proc. of International Conference on Artificial Neural Networks (ICANN '93)*, pages 1074–1077, Berlin/Heidelberg/New York/Tokyo, 1993. Springer.
- [17] Y.-H. Dai. A family of hybrid conjugate gradient methods for unconstrained optimization. *Mathematics of Computation*, 72:1317–1328, 2003.



- [18] Y.-H. Dai, L.-Z. Liao, and D. Li. On Restart Procedures for the Conjugate Gradient Method. *Numerical Algorithms*, 35:249–260, 2004.
- [19] Y. H. Dai and Y. Yuan. A three-parameter family of nonlinear conjugate gradient methods. *Mathematics of Computation*, 70:1155–1167, 2001.
- [20] J. del R. Millan and P. Bofill. Learning by back-propagation: Computing in a systolic way. In *Parallel Architectures and Language Europe - Vol.2: Parallel Languages (PARLE '89)*, pages 235–252, Berlin - Heidelberg - New York, June 1989. Springer.
- [21] T. DeLillo and T. Hrycak. A stopping rule for the conjugate gradient regularization method for ill-posed problems. *Acoustical Society of America Journal*, 112:2381–2381, November 2002.
- [22] C. Ernault. Performance of backpropagation on a parallel transputer-based machine. In *Neuro-Nimes '88: International Workshop on Neural Networks and their Applications*, pages 311–24, Nimes, France, 1988.
- [23] Federico Faggin. Neural network hardware. In *IEEE International Joint Conference on Neural Networks (6th IJCNN'92)*, volume I, page 153, Baltimore, MD, June 1992. IEEE/INNS.
- [24] S. E. Fahlman. Faster-learning variations on back-propagation: an empirical study. In D.S. Touretzky, G.E. Hinton, and T.J. Sejnowski, editors, *Proc. of the 1988 Connectionist Models Summer School*, pages 38–51, San Mateo, CA, 1989. Morgan Kaufmann.
- [25] B. M. Forrest, D. Roweth, N. Stroud, D. J. Wallace, and G. V. Wilson. Implementing neural network models on parallel computers. *The Computer Journal*, 30(5):413–419, October 1987.
- [26] Geoffrey C. Fox, Mark A. Johnson, Gregory A. Lyzenga, Steve W. Otto, John K. Salmon, and David W. Walker. *Solving Problems On Concurrent Processors, Volume 1, General Techniques and Regular Problems*. Prentice-Hall, Englewood Cliffs, NJ, 1988.
- [27] C. K. Gan, P. D. Haynes, and M. C. Payne. Preconditioned conjugate gradient method for the sparse generalized eigenvalue problem in electronic structure calculations. *Computer Physics Communications*, 134:33–40, February 2001.
- [28] R. P. Gorman and T. J. Sejnowski. Analysis of hidden units in a layered network trained to classify sonar targets. *Neural Networks*, 1:75–89, 1988.
- [29] R. Paul Gorman and T. J. Sejnowski. Learned classification of sonar targets using a massively parallel network. *IEEE Tran. on Acoustics, Speech, and Signal Processing*, 36(7):1135–1140, July 1988.
- [30] E. Hartman, J. D. Keeler, and J. M. Kowalski. Layered neural networks with Gaussian hidden units as universal approximations. *Neural Computation*, 2(2):210–215, 1990.
- [31] Simon Haykin. *Neural Networks*. Macmillan College Publishing Company, New York, 1994.
- [32] D. O. Hebb. *The Organization of Behavior*. Wiley & Sons, New York, 1949.
- [33] A. L. Hodgkin and A. F. Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *Journal of Physiology*, 117:500–544, 1952.
- [34] J. J. Hopfield. Neural networks and physical systems with emergent collective behavior. *PNAS USA*, 79:2554–2558, 1982.

- [35] J. J. Hopfield. Collective processing and neural states. In C. Nicollini, editor, *Modeling and Analysis in Biomedicine*, pages 371–389, 1984.
- [36] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [37] K. Hornik, M. Stinchcombe, and H. White. Universal approximation of an unknown function and its derivatives using multilayer feedforward networks. *Neural Networks*, 3:551–560, 1990.
- [38] R. A. Hyman, J. Tetzlaff, and B. Doporcyk. Simple and unified derivation of conjugate gradient and variable metric minimization. *APS Meeting Abstracts*, pages 1069–+, March 2003.
- [39] Robert A. Jacobs. Increased rates of convergence through learning rate adaptation. *Neural Networks*, 1(4):295–307, 1988.
- [40] H. Jiang and W. Yang. Conjugate-gradient optimization method for orbital-free density functional calculations. *jcp*, 121:2030–2036, August 2004.
- [41] J. Jiang. Image compression with neural networks – a survey. In *Signal Processing: Image Communication 1999*, volume 14, pages 737–760, 1999.
- [42] E. M. Johansson, F. U. Dowlah, and D. M. Goodman. Back-propagation learning for multi-layer feed-forward neural networks using the conjugate gradient method. *International Journal of Neural Systems*, 2:291–302, 1991.
- [43] S. Kirkpatrick, C. D. Gelatt Jr., and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [44] Christoph Klawun and Charles L. Wilkins. A novel algorithm for local minimum escape in back-propagation neural networks: application to the interpretation of matrix isolation infrared spectra. *Journal of Chemical Information and Computer Sciences*, 34(4):984–993, 1994.
- [45] B. Knutson, J. Taylor, M. Kaufman, R. Peterson, and G. Glover. Distributed neural representation of expected value. *Journal of Neuroscience*, 25(19):4806–4812, 2005.
- [46] Randal A. Koene and Yoshio Takane. Discriminant component pruning: Regularization and interpretation of multi-layered back-propagation networks. In *Neural Computation*, volume 11, pages 783–802, 1999.
- [47] Teuvo Kohonen. Correlation matrix memories. *IEEE Transactions on Computers*, C-21:353–359, 1972.
- [48] Ben J. A. Kröse and P. Patrick van der Smagt. *An Introduction to Neural Networks*. University of Amsterdam, Amsterdam, 1993.
- [49] Toshiro Kubota. Massively parallel networks for edge localization and contour integration-adaptable relaxation approach. In *Neural Networks*, volume 17, pages 411–425, 2004.
- [50] Han-Wook Lee and Chan-Ik Park. An efficient parallel block backpropagation learning algorithm in transputer mesh-connected parallel computers. In *IEICE Trans. Inf. & Syst.*, volume E83-D, pages 1622–1630, 2000.
- [51] Tara M. Madhyastha and Daniel A. Reed. Learning to classify parallel input/output access patterns. *IEEE Trans. Parallel Distrib. Syst.*, 13(8):802–813, 2002.

- [52] Q.M. Malluhi, M.A. Bayoumi, and T.R.N. Rao. An efficient mapping of multilayer perceptron with backpropagation anns on hypercubes. In *Proc. of Symposium on Parallel and Distributed Systems*, pages 368–375. IEEE Computer Society Press, 1994.
- [53] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5:115–137, 1943.
- [54] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equations of state calculations by fast computing machine. *Journal of Chemical Physics*, 21:1087–1091, 1953.
- [55] G. Meurant. Estimates of the  $l_2$  norm of the error in the conjugate gradient algorithm. *Numerical Algorithms*, 40:157–169, 2005.
- [56] M. L. Minsky. *Theory of neural-analog reinforcement systems and its application to the brain-model problem*. Ph.D. dissertation, Princeton University, N.J., 1954.
- [57] Marvin Minsky. Steps toward artificial intelligence. pages 406–450, 1995.
- [58] Marvin Minsky and Seymour Papert. *Perceptrons*. MIT Press, Cambridge, MA, 1969.
- [59] Martin Fodslette Moller. A scaled conjugate gradient algorithm for fast supervised learning. *Neural Networks*, 6:525–533, 1993.
- [60] F.C. Morabito. Fast backpropagation using modified sigmoidal functions. In Marinaro M. and Morasso P., editors, *ICANN'94: Proceedings of the International Conference on Artificial Neural Networks*, pages 537–540, Berlin/Heidelberg/New York/Tokyo, 1994. Springer.
- [61] K. Nakano. Associatron - a model of associative memory. *IEEE transactions on Systems, Man, and Cybernetics*, SMC2(3):380–388, 1972.
- [62] Tomas Nordström. *Highly Parallel Computers for Artificial Neural Networks*. PhD thesis, Luleå University of Technology, Luleå, Sweden, 1995.
- [63] D. B. Parker. Learning logic. Technical Report TR-47, Center for Computational Research in Economics and Management Science, MIT, Cambridge, MA, 1985.
- [64] D. B. Parker. Second order back propagation: an optimal adaptive algorithm for any adaptive network. *Proceedings of the IEEE First Annual International Conference on Neural Networks*, 1987.
- [65] Alpha V. Pernia-Espinoza, Joaquin B. Ordieres-Mere, Francisco J. Martinez de Pison, and Ana Gonzalez-Marcos. Tao-robust backpropagation learning algorithm. *Neural Netw.*, 18(2):191–204, 2005.
- [66] M. Pethick, M. Liddle, P. Werstein, and Z. Huang. Parallelization of a backpropagation neural network on a cluster computer. In *Proceedings of the 15th IASTED International Conference on Parallel and Distributed Computing and Systems (PDCS 2003)*, pages 574–582, Marina Del Rey, California, 2003.
- [67] A. Polze and M. Malek. Parallel computing in a world of workstations. In M. H. Hamza, editor, *Proceedings of the Seventh IASTED/ISMM International Conference Parallel and Distributed Computing and Systems*, pages 72–4. IASTED-ACTA Press, Anaheim, CA, USA, 1995.
- [68] D. A. Pomerleau, G. L. Gusciora, D. S. Touretzky, and H. T. Kung. Neural network simulation at warp speed: How we got 17 million connections per second. In *IEEE Second International Conference on Neural Networks (2nd ICNN'88)*, volume II, pages 143–150, San Diego, CA, jul 1988. IEEE.

- [69] J.M.D. Powell. Restart procedures for the conjugate gradient method. *Math. Prog.*, 12:241–254, 1977.
- [70] M. Raus and W. Ameling. A parallel algorithm for a dynamic eta/alpha estimation in backpropagation learning. In Marinaro M. and Morasso P., editors, *ICANN'94: Proceedings of the International Conference on Artificial Neural Networks*, pages 639–642, Berlin/Heidelberg/New York/Tokyo, 1994. Springer.
- [71] N. Rochester, J. H. Holland, L. H. Haibt, and W. L. Duda. Tests on a cell assembly theory of the action of the brain, using a large digital computer. *IRE Transactions on Information Theory*, 2:80–93, 1956.
- [72] F. Rosenblatt. The perceptron: A probabilistic model for information storage in the brain. *Psychological review*, 65:386–408, 1958.
- [73] F. Rosenblatt. Perceptron simulation experiments. *Proceedings of the IRE*, 48:301–309, 1960.
- [74] F. Rosenblatt. *Principles of Neurodynamics*. Spartan, New York, 1962.
- [75] Frank Rosenblatt. On the convergence of reinforcement procedures in simple perceptrons. Report, Cornell Aeronautical Laboratory, Ithaca, New York, 1960.
- [76] D. E. Rumelhart and J. L. McClelland. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Vol. I: Foundations*. MIT Press, Cambridge, MA, 1986.
- [77] David E. Rumelhart, Geoffrey E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In D. E. Rumelhart, J. L. McClelland, and the PDP research group, editors, *Parallel distributed processing: Explorations in the microstructure of cognition, Volume 1: Foundations*. MIT Press, 1986.
- [78] P. Saratchandran, N. Sundararajan, and Shou King Foo. *Parallel Implementations of Backpropagation Neural Networks on Transputers: A Study of Training Set Parallelism*. World Scientific Publishing Co., Inc., River Edge, NJ, USA, 2000.
- [79] J. T. Schwartz. The new connectionism: Developing relationships between neuroscience and artificial intelligence. In S. R. Graubard, editor, *The Artificial Intelligence Debate: False Starts, Real Foundations*, pages 123–141. MIT Press, Cambridge, MA, 1988.
- [80] Udo Seiffert. Artificial neural networks on massively parallel computer hardware. In *ESANN 2002, 10th European Symposium on Artificial Neural Networks*, pages 319–220, Bruges, Belgium, 2002.
- [81] Terrence J. Sejnowski and Charles Rosenberg. NETtalk: A parallel network that learns to read aloud. Technical Report JHU-EECS-86-01, Johns Hopkins University, 1986.
- [82] Fernando M. Silva and Luis B. Almeida. Speeding up backpropagation. In R. Eckmiller, editor, *Advanced Neural Computers*, pages 151–158. Elsevier North Holland, Amsterdam, 1990.
- [83] Noboru Sonehara, Mitsuo Kawato, Sei Miyake, and Kazunari Nakane. Image data compression using a neural network model. In *Proc. of Int. Joint Conference on Neural Networks*, volume 2, pages 35–40, Washington, DC, 1989.
- [84] K. Steinbuch. Die lernmatrix. *Kybernetik (Biological Cybernetics)*, 1:36–45, 1961.

- [85] Thomas L. Sterling, John Salmon, Donald J. Becker, and Daniel F. Savarese. *How to build a Beowulf: a guide to the implementation and application of PC clusters*. MIT Press, Cambridge, MA, USA, 1999.
- [86] Josef Stoer and Roland Bulirsch. *Introduction to Numerical Analysis*. Springer Verlag, 2002.
- [87] N. Sundararajan and P. Saratchandran, editors. *Parallel Architectures for Artificial Neural Networks – Paradigms and Implementations*. IEEE Computer Society,, 1998.
- [88] S. Suresh, S. N. Omkar, and V. Mani. Parallel implementation of back-propagation algorithm in networks of workstations. *IEEE Trans. Parallel Distrib. Syst.*, 16(1):24–34, 2005.
- [89] David W. Tank and John J. Hopfield. Collective computation in neuronlike circuits. *Scientific American*, 257(6):104–114, December 1987.
- [90] W. K. Taylor. Electrical simulation of some nervous system functional activities. In E. C. Cherry, editor, *Information Theory*, pages 314–328, London 1985, 1956. Butterworths, London.
- [91] Jim Torresen, Hiroshi Nakashima, Shinji Tomita, and Olav Landsverk. Exploiting parallel computers to reduce neural network training time of real applications. In *International Symposium on High Performance Computing*, Fukuoka, Japan, 1997.
- [92] Jim Torresen and Shinji Tomita. *Parallel Architectures for Artificial Neural Networks*, chapter Implementation of Backpropagation Training of Neural Networks On Large Parallel Computers. IEEE CS Press, 1998.
- [93] L.A. Trejo and C. Sandoval. Improving back-propagation: Epsilon-back-propagation. In Mira J. and Sandoval F., editors, *From Natural to Artificial Neural Computation, Proc. of International Workshop on Artificial Neural Networks*, pages 427–432, Malaga-Torremolinos, Spain, 1995. Springer.
- [94] Rein van den Boomgaard and Arnold W. M. Smeulders. Self learning image processing using a priori knowledge of spatial patterns. In Takeo Kanade, F. C. A. Groen, and L. O. Hertzberger, editors, *Intelligent Autonomous Systems 2, An International Conference, Amsterdam, The Netherlands, 11-14 December 1989*, pages 305–314. IOS Press, 1989.
- [95] J. Wawrzynek, K. Asanovic, B. Kingsbury, J. Beck, D. Johnson, and N. Morgan. Spert-i1: A vector microprocessor system and its application to large problems in backpropagation training. In et al. Touretzky D.S., editor, *Advances in Neural Information Processing Systems 8 (NIPS'96)*, pages 619–625, Cambridge/Boston/London, 1996. MIT Press.
- [96] P. J. Werbos. *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. PhD thesis, Harvard University, 1974.
- [97] B. Widrow. Generalization and information storage in networks of adaline "neurons". In M. C. Yovits, G. T. Jacobi, and G. D. Goldstein, editors, *Self-Organizing Systems 1962*, pages 435–461, Chicago, Illinois, 1962. Spartan.
- [98] B. Widrow and M. E. Hoff. Adaptive switching circuits. In *1960 IRE WESCON Convention Record*, pages 96–104, New York, 1960.
- [99] S. Winograd and J. D. Cowan. *Reliable Computation in the Presence of Noise*. MIT Press, Cambridge, Massachusetts, 1963.

- [100] M. Witbrock and M. Zagha. An implementation of backpropagation learning on GF11, a large SIMD parallel computer. *Parallel Computing*, 14(3):329–346, August 1990.
- [101] Y.S. Wong, B.H.K. Lee, and T.K.S. Wong. Parameter extraction by parallel neural networks. *Intelligent Data Analysis*, pages 59–72, 2001.
- [102] Kunihiro Yamamori, Toru Abe, and Susumu Horiguchi. Theoretical learning-speed evaluation of parallel back-propagation algorithms. In *Systems Research and Information Systems*, volume 9, pages 121–148, 2000.
- [103] H.H.T. Yeung. Online semantic extraction by backpropagation neural network with various syntactic structure representations. In McGuinness D.L. and Ferguson G., editors, *Proceedings of the 19th National Conference on Artificial Intelligence, 16th Conference on Innovative Applications of Artificial Intelligence (AAAI-04)*, pages 1042–1043. AAAI Press, Menlo Park, CA, 2004.
- [104] Xiru Zhang, M. McKenna, J. P. Mesirov, and D. L. Waltz. The backpropagation algorithm on grid and hypercube architectures. *Parallel Computing*, 14(3):317–327, August 1990.

- [Ba94] Sorin Babii Experiments with a Distributed Ray Tracing Algorithm *Buletinul Științific al UTT, Seria Automatică și Calculatoare*, 39(53):81–87, 1994.
- [BaC03] Sorin Babii, Vladimir Crețu. Parallelizing the Training Phase of BackPropagation in a LAN of Workstations. *Concurrent Information Processing and Computing – NATO Advanced Research Workshop*, pages 29–36, IOS Press, ISBN 1-58603-502-9, ISSN 1387-6694, Sinaia, July 2003.
- [BBD95] Sorin Babii, Laura Barbulescu, Romeo Dumitrescu. A Client–Server Application Model for Distributed Programming in the Netware Environment. *Buletinul Științific al UTT, Seria Automatică și Calculatoare*, 40(54):67–77, 1995.
- [BC03] Sorin Babii, Vladimir Crețu. A Parallel Algorithm for Neural Networks. *Buletinul Științific al Universității "Politehnica" din Timișoara, Seria Automatică și Calculatoare*, 48(62):35–39, Timișoara 2003.
- [BC04] Sorin Babii, Vladimir Crețu. A Distributed Algorithm for Neural Networks Training in a Network of Computers. The 6th International Conference on Technical Informatics, in *Buletinul Științific al Universității "Politehnica" din Timișoara, Seria Automatică și Calculatoare*, 49(63),3:233–236, Timișoara 2004.
- [BDB95] Sorin Babii, Romeo Dumitrescu, Laura Barbulescu. High Level Interface to the IPX Based Netware Communication Services. *Buletinul Științific al UTT, Seria Automatică și Calculatoare*, 40(54):57–66, 1995.
- [BS94] Sorin Babii, Nicolae Szirbik. Intermediate Representation Based on Conceptual Graph Model. *10th Conference on Control Systems and Computer Science – CSCS10*, 3:84–93, București 1995.
- [SB96] Nicolae Szirbik, Sorin Babii. The Basic Ideas for Developing Connectionist–Symbolic Cognitive Systems. *Buletinul Științific al UTT, Seria Automatică și Calculatoare*, 41(55), 1996.
- [SBD93] Nicolae Szirbik, Sorin Babii, Călin Dragomirescu. A Method for Inserting PROLOG Rules in a Neural Network. *Conferința internațională de calculatoare*, București, 2:202–206, mai 1993.
- [Ba06] Sorin Babii. Rețele neuronale – Stadiu actual. Perspective. *Raport de cercetare științifică*, Universitatea "Politehnica" din Timișoara, Departamentul de Ingineria Calculatoarelor și a Programării, 2006.
- [Ba06a] Sorin Babii. DataParSim: Simulator paralel cu partiționarea datelor. *Raport de cercetare științifică*, Universitatea "Politehnica" din Timișoara, Departamentul de Ingineria Calculatoarelor și a Programării, 2006.
- [Ba07] Sorin Babii. NetParSim: Simulator paralel cu partiționarea rețelei. *Raport de cercetare științifică*, Universitatea "Politehnica" din Timișoara, Departamentul de Ingineria Calculatoarelor și a Programării, 2007.
- [Ba07a] Sorin Babii. Performance Evaluation for Training a Distributed Back-Propagation Implementation. *4th International Symposium on Applied Computational Intelligence and Informatics*, 273–276, ISBN: 1-4244-1234-X, Timișoara, 2007.



- [BCP07] Sorin Babii, Vladimir Crețu, Emil Petriu. Performance Evaluation of Two Distributed BackPropagation Implementations. *The 2007 International Joint Conference on Neural Networks (IJCNN)*, Orlando, FL, 2007, ISBN: 1-4244-1380-X, ISSN: 1098-7576, IEEE Catalog Number: 07CH37922C.