

# RECUNOAȘTEREA ȘIRURILOR NUMERICE SCRISE DE MÂNĂ

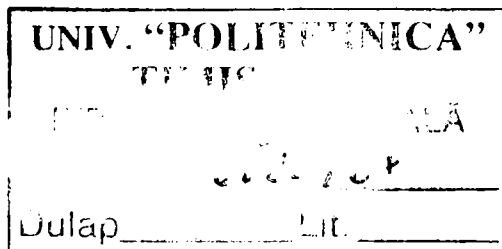
Teză destinată obținerii  
titlului științific de doctor inginer  
la

Universitatea "Politehnica" din Timișoara  
în domeniul ȘTIINȚA CALCULATOARELOR  
de către

**Ing. Dan Cireșan**

Conducător științific: prof.univ.dr.ing. Ștefan Holban  
Referenți științifici: prof.univ.dr.ing. Ștefan-Gheorghe Pentiu  
prof.univ.dr.ing. Viorel Negru  
prof.univ.dr.ing. Marius Crișan

Ziua susținerii tezei: 17.10.2008



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, 2008

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. Ștefan Holban, pentru înțelegerea și îndrumarea pe care mi le-a oferit de-a lungul anilor de doctorat. Fără imboldurile sale în momentele dificile ale activității de doctorat, această teză nu ar mai fi fost finalizată vreodată.

De asemenea, as dori să mulțumesc prof. Augustin Devian, cel care, la vârsta marilor întrebări, mi-a arătat frumusețea matematicii și, implicit, a științei.

Mulțumesc mamei mele pentru că totul a început cu prima carte de popularizare a științei pe care mi-a dăruit-o. Au urmat apoi multe alte cărți, fiecare mai frumoasă decât cealaltă. Pentru aceasta și pentru toate lucrurile minunate pe care le-a făcut pentru mine va rămâne veșnic în amintirea mea.

Timișoara, octombrie 2008

Dan Cireșan

Cireșan, Dan

### **Recunoașterea șirurilor numerice scrise de mână**

Teze de doctorat ale UPT, Seria 10, Nr. 10, Editura Politehnica, 2008, 106 pagini, 42 figuri, 13 tabele.

ISSN: 1842-7707

ISBN: 978-973-625-777-3

Cuvinte cheie: recunoașterea scrisului de mână, rețele-neuronale de convoluție, GPU, optimizare, șiruri numerice, cifre parțial suprapuse, NIST SD 19, transformata Hough.

Rezumat,

Teza abordează problema recunoașterii șirurilor numerice conținând un număr nedefinit de cifre izolate sau parțial suprapuse. Metoda propusă nu utilizează segmentarea. Pentru a rezolva cazurile de cifre alipite sau parțial suprapuse, se folosește pe lângă clasificatorul uzual, la nivel de cifră, un clasificator antrenat cu perechi de cifre parțial suprapuse.

Sistemul creat atinge o rată de recunoaștere de 93.77%, în condițiile în care nu folosește nici segmentare, nici antrenare cu exemple negative. Analiza erorilor arată că o viitoare antrenare cu exemple negative va duce la creșterea ratei de recunoaștere la 97-98%.

Au fost efectuate teste de performanță folosind procesoare grafice (GPU). Rezultatele preliminare sunt promițătoare, obținându-se o accelerare de peste 20 de ori față de varianta scrisă pentru CPU.

pentru Dana



# Cuprins

<b>1</b>	<b>Introducere</b>	<b>1</b>
1.1	Motivație . . . . .	1
1.2	Organizarea tezei . . . . .	3
1.3	Descrierea metodei . . . . .	3
<b>2</b>	<b>Stadiul actual al domeniului</b>	<b>5</b>
2.1	Stadiul actual al recunoașterii cifrelor . . . . .	6
2.2	Stadiul actual al recunoașterii șirurilor numerice . . . . .	7
2.2.1	Recunoașterea holistică . . . . .	7
2.2.2	Recunoașterea cu segmentare . . . . .	8
2.3	Concluzii . . . . .	10
<b>3</b>	<b>Aspecte teoretice</b>	<b>11</b>
3.1	Noțiuni de prelucrare a imaginilor . . . . .	11
3.1.1	Centrul de masă al unei imagini . . . . .	11
3.1.2	Distorsiuni . . . . .	12
3.1.3	Transformata Hough . . . . .	13
3.2	Rețele neuronale artificiale . . . . .	14
3.2.1	Perceptronul . . . . .	15
3.2.2	Rețele neuronale artificiale cu hărți de convoluție (RNC) . . . . .	15
3.2.3	Metode de antrenare a rețelelor neuronale . . . . .	16
<b>4</b>	<b>Baze de date cu cifre și șiruri de cifre</b>	<b>18</b>
4.1	NIST SD 19 . . . . .	19
4.1.1	Structura NIST SD 19 . . . . .	19
4.1.2	Conținutul NIST SD 19 . . . . .	21
4.1.3	Ierarhizarea datelor . . . . .	22
4.1.4	Formatul fișierelor NIST . . . . .	23
4.1.5	Programe utilitare ale NIST SD 19 . . . . .	24
4.2	MNIST . . . . .	24
4.3	LNIST - Liu-NIST . . . . .	25
4.4	TDNS (Two-digit Numeral Strings) . . . . .	26
4.4.1	Setul de antrenare . . . . .	27
4.4.2	Setul de test . . . . .	30
4.5	Seturile de test 3DNS, 6DNS și 10DNS . . . . .	30
4.5.1	Metoda de extragere a câmpurilor din formulare . . . . .	31

4.5.2	Setul de test 3DNS	36
4.5.3	Setul de test 6DNS	37
4.5.4	Setul de test 10DNS	38
4.6	Concluzii	38
<b>5</b>	<b>Experimente</b>	<b>40</b>
5.1	Experimente cu cifre individuale	40
5.1.1	Rețele neuronale de tip perceptron	41
5.1.2	Rețele neuronale cu extractori de trăsături	41
5.2	Experimente cu șiruri compuse din două cifre parțial suprapuse	43
5.2.1	Rețele neuronale de tip perceptron	43
5.2.2	Rețele neuronale cu extractori de trăsături	44
5.2.3	Concluzii	49
5.3	Experimente cu șiruri de cifre din NIST SD 19	49
5.3.1	Sistemul de recunoaștere	49
5.3.2	Analizorul de componente conexe	49
5.3.3	Componenta de clusterizare (agregare)	50
5.3.4	Clasificatorul	52
5.4	Experimente pe setul 6DNS	57
5.5	Experimente pe setul 10DNS	58
5.6	Concluzii	58
<b>6</b>	<b>Optimizarea vitezei de execuție</b>	<b>60</b>
6.1	Variante CPU ale optimizării algoritmului de propagare înainte	61
6.1.1	Varianta open source modificată	61
6.1.2	Varianta C++ fără vectori STL	62
6.1.3	Varianta SSE2	63
6.2	Variante GPU ale optimizării algoritmului de propagare înainte	65
6.2.1	Varianta BROOK+	68
6.2.2	Variantele CUDA	71
6.3	Concluzii	79
<b>7</b>	<b>Concluzii</b>	<b>82</b>
7.1	Contribuții	82
7.2	Sumarul tezei	83
7.3	Continuări posibile	85



# Lista figurilor

2.1	Ratele de recunoaștere pentru șiruri din NIST SD 19. Preluare din [44]. [35] reprezintă [30], [36] reprezintă [9], iar [14] reprezintă [58]. . . . .	10
3.1	Transformata Hough pentru detecția de linii. Fiecărei linii $d$ din spațiul imagine îi corespunde un punct în spațiul parametric Hough. . . . .	13
3.2	Exemplu de matrice de acumulare rezultată în urma aplicării transformatei Hough. Cercurile marchează maximele corespunzătoare celor șase linii detectate. Afișarea grafică folosește negru pentru elemente nule și tonuri de gri cu iluminare proporțională cu valoarea respectivului element pentru valori nenule. Albul se folosește pentru elementul sau elementele cu număr maxim de puncte. . . . .	14
3.3	RNC pentru recunoașterea caracterelor. Preluare din [37]. . . . .	15
3.4	Modul de deplasare a nucleului de convoluție. . . . .	16
3.5	Modul de conectare a neuronilor de pe straturile de convoluție. . . . .	17
4.1	Ierarhia seturilor de test și de antrenare. . . . .	18
4.2	Exemplu de formular din NIST SD 19 (formularul f1816_24). A se vedea și Figura 4.4. . . . .	20
4.3	Arborele aferent ierarhizării formularelor în partiții. . . . .	23
4.4	Arborele aferent ierarhizării imaginilor de caractere după partiții și persoane. . . . .	23
4.5	Comparație între centrarea în centrul de masă (stânga) și centrarea în dreptunghiul delimitator (dreapta). Punctul roșu reprezintă centrul de masă al imaginii. Pe ultimul rând se află suprapunerea celor două imagini de pe coloana corespunzătoare. . . . .	26
4.6	Primul rând - exemple de artefacte în imagini $128 \times 128$ pixeli din NIST SD 19; rândul al doilea - imaginile rezultate în urma procesului de corecție semiautomată. . . . .	27
4.7	Primul rând - imagini defectuos segmentate sau deformate din baza NIST SD 19; rândul al doilea - imagini rezultate în urma corecției manuale. . . . .	28
4.8	Primul rând - imagini $128 \times 128$ pixeli din NIST SD 19; rândul al doilea - imagini $20 \times 20$ pixeli rezultate după tăiere relativă la dreptunghiul delimitator, centrare și redimensionare. . . . .	28
4.9	Interfața programului de corectare semiautomată. . . . .	29
4.10	Cifrele a nouă persoane după etapele de tăiere după dreptunghiul delimitator, centrare și redimensionare la $20 \times 20$ pixeli. . . . .	30

4.11	Procesul de suprapunere parțială: primul rând - imaginile originale; rândul al doilea - imagini deplasate pe verticală; al treilea rând - imaginea finală după deplasarea pe orizontală cu 13 pixeli (distanța minimă dintre imagini) + 3 (distanța de suprapunere) pixeli. . . . .	31
4.12	Limitarea variației parametrului $r$ . . . . .	33
4.13	Imagini obținute în urma procesului de extragere a câmpurilor din formulare. Primul rând - imagini brute, al doilea rând - imagini după faza de curățare automată, rândul al treilea - imagini după faza de curățare manuală. . . . .	36
4.14	Imaginile cu numere de trei cifre care au fost scrise eronat de subiecții din NIST SD 19. . . . .	37
4.15	Imagini cu numere de șase cifre, eronat scrise de subiecții NIST SD 19. . . . .	37
4.16	Imagini cu numere de zece cifre, eronat scrise de subiecții NIST SD 19 (toate ar trebui să conțină șirul 0123456789). . . . .	38
5.1	Exemple de imagini din fișierul de antrenare. . . . .	43
5.2	Arhitectura rețelei neuronale cu hărți de convoluție. . . . .	45
5.3	Suprafețele erorilor rețelelor cu 8 hărți de convoluție pe stratul $L_1$ . . . . .	48
5.4	Arhitectura sistemului de recunoaștere a șirurilor de cifre. . . . .	50
5.5	Componentele conexe înainte (sus) și după (jos) faza de agregare. Fiecare componentă este colorată distinct și este desenat și dreptunghiul delimitator. . . . .	51
5.6	Reconectarea cifrei 5. . . . .	51
5.7	Cazurile de unire defectuoasă de tipul $\overline{XI}$ . . . . .	53
6.1	Codul sursă al algoritmului de propagare înainte pentru varianta C++ cu vectori STL. . . . .	62
6.2	Codul sursă al algoritmului de propagare înainte pentru varianta C++. . . . .	63
6.3	Codul sursă al algoritmului de propagare înainte pentru varianta C++ cu SSE II. . . . .	64
6.4	Evoluția performanțelor procesoarelor și a plăcilor grafice. Preluare după [54]. . . . .	66
6.5	Comparație între evoluțiile lățimii de bandă pentru memoria sistem (RAM) și pentru memoria globală a plăcilor grafice. Adaptare după [54]. . . . .	67
6.6	Comparație între alocarea tranzistoarelor pe suprafața de siliciu pentru GPU și CPU. Preluare după [54]. . . . .	68
6.7	Codul sursă al algoritmului de propagare înainte pentru varianta BROOK+. . . . .	69
6.8	Arhitectura GPU-urilor de la NVIDIA (seriile GeForce 8 și 200). Preluare după [65]. . . . .	71
6.9	Organizarea thread-urilor. Adaptare după [56]. . . . .	73
6.10	Kernel-urile pentru straturile de convoluție. . . . .	75
6.11	Kernel-urile pentru straturile complet conectate. . . . .	76
6.12	Kernel-urile pentru stratul L3. . . . .	78
6.13	Gradul de accelerare a diferitelor variante de kernel pentru stratul L3. Referința este considerată varianta C++ optimizată, fără vectori STL. . . . .	80

# Lista tabelelor

4.1	Conținutul și modul de utilizare ale seturilor de date. . . . .	19
4.2	Versiuni anterioare ale bazei de date NIST cu imagini de caractere scrise de mână. (după [4]) . . . . .	19
4.3	Bazele de date în care au fost publicate partițiile. † - câmpuri completate, dar neprocesate. ∅ - câmpuri necompletate. (după [4]) . . . . .	21
4.4	Exemplu de fișier cu referințe ( <i>hsf_page/truerefs/ref_24.txt</i> ). A se vedea și Figura 4.2. . . . .	22
4.5	Numărul de exemplare din fiecare clasă. . . . .	22
5.1	Rezultate obținute cu perceptronul cu un strat ascuns. . . . .	41
5.2	Rezultate obținute cu rețele neuronale cu hărți de convoluție. . . . .	42
5.3	Rezultate obținute cu perceptronul cu un strat ascuns. . . . .	44
5.4	Rezultate obținute cu rețele neuronale cu hărți de convoluție. . . . .	46
5.5	Influența strategiei de rejecție asupra ratei de recunoaștere (forțat = fără rejecție). . . . .	49
5.6	Rate de recunoaștere pentru setul de test 3DNS (1C - cu deplasarea $RNC_1$ deasupra imaginii, 2C - cu deplasarea $RNC_2$ deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor). . . . .	54
5.7	Rate de recunoaștere pentru setul de test 6DNS (1C - cu deplasarea $RNC_1$ deasupra imaginii, 2C - cu deplasarea $RNC_2$ deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor). . . . .	57
5.8	Rate de recunoaștere pentru setul de test 10DNS (1C - cu deplasarea $RNC_1$ deasupra imaginii, 2C - cu deplasarea $RNC_2$ deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor). . . . .	58



# Capitolul 1

## Introducere

### 1.1 Motivație

De ce recunoașterea scrisului de mână? Deoarece scrisul a reprezentat, de la inventarea sa și până în prezent, modalitatea cea mai facilă și sigură de stocare și de transmitere a informațiilor. Din punct de vedere evolutiv, caracterul social al omului a fost generat de necesitatea creșterii șansei de supraviețuire. La rândul lui, caracterul social a favorizat prin selecție apariția mijloacelor de comunicare din care a emers limbajul în forma orală, pentru ca apoi să apară și scrisul. Acesta reprezintă poate cea mai importantă invenție a omului, dar trebuie să admitem că inventarea sa a fost constrânsă de evoluție.

Scrisul a generat o revoluție în civilizația umană prin cele două avantaje majore pe care le are față de comunicarea orală. Primul este capacitatea de stocare sigură a informației, nealterată de transmiterea acesteia pe cale orală. Al doilea avantaj este reprezentat de faptul că permite transferul cunoștințelor de la o generație la alta, facilitând acumularea acestora și impulsivând progresul. Dacă scrisul nu ar fi existat, având în vedere lipsa memoriei genetice și a duratei limitate de viață a oamenilor, atunci probabil că omul ar fi ramas într-un stadiu incipient de civilizație.

Alți pași importanți în evoluția scrisului au fost reprezentați de inventarea hârtiei și a tiparniței. Aceste două invenții au făcut accesibilă informația unui număr mult mai mare de oameni, atât prin scăderea costului cărților, cât și prin multiplicarea facilă a acestora. Inventarea calculatoarelor electronice în secolul trecut a generat o nouă revoluție în știință, apărând noi modalități de stocare, transmitere și prelucrare a informației. Din acel moment, informația a început să fie scrisă în două forme: cea uzuală, utilizând hârtie, și cea nouă, utilizând suporturi electronice și optice. A devenit necesară conversia informației dintr-un format într-altul. Dacă transformarea din format electronic în format tipărit este banală, conversia inversă s-a dovedit a fi atât complexă, cât și complicată. Conversia unor cazuri particulare de scris, cum ar fi cel tipărit pe hârtie, este rezolvată acceptabil, în schimb conversia scrisului de mână, mai ales în forma cursivă, este într-o faza incipientă, departe de finalizare, rezultatele fiind, în general, mult sub necesitățile impuse de utilizarea practică.

*Problema recunoașterii scrisului de mână este unul din subiectele nerezolvate ale Inteligenței Artificiale. Natura acestei probleme impune o abordare euristică în rezolvarea ei, astfel că multe din părțile procesului de recunoaștere sunt și probabil vor rămâne perfectibile.*

*În acest domeniu nu există soluții definitive.*

Teza abordează o problemă particulară a recunoașterii scrisului de mână, cea a recunoașterii șirurilor numerice alcătuite dintr-un număr nedefinit de cifre. Recunoașterea șirurilor numerice prezintă atât avantaje, cât și dezavantaje față de recunoașterea textelor, propozițiilor etc. Dintre avantaje se pot aminti varietatea mai mică a simbolurilor (nu există variante minuscule și majuscule) precum și numărul mai mic al acestora. În schimb, există un dezavantaj major: la recunoașterea șirurilor numerice nu se pot aplica modele lingvistice [18].

Un exemplu de aplicație practică a recunoașterii șirurilor numerice este recunoașterea codurilor poștale de pe plicuri, a cecurilor bancare. Pe plan mondial există multe firme ce au ca domeniu de activitate producerea de echipamente pentru recunoașterea scrisului. Una dintre cele mai importante, Parascript - Total Recognition Company, are în portofoliul de produse aplicații pentru domeniul financiar (recunoașterea scrisului de pe cecuri, precum și autentificarea de semnături), domeniul poștal (recunoașterea adreselor), procesarea formularelor cu diverse tipuri de date. Faptul că produsele acestei firme sortează zilnic peste un sfert de miliard de plicuri în Statele Unite ale Americii demonstrează utilitatea și necesitatea cercetărilor în acest domeniu.

Recunoașterea șirurilor numerice este mult mai complicată decât cea a recunoașterii cifrelor deoarece apar probleme suplimentare: lungimea necunoscută a șirului (spre deosebire de cifre sau litere [19], unde este unu, și de cuvinte, unde poate fi extrasă din dicționare), segmentarea, suprapuneri parțiale etc.

După cum se va vedea în capitolul următor, Starea actuală a domeniului, rata de recunoaștere actuală este de aproximativ 96%. Având în vedere că precizăm anterior că zilnic peste un sfert de miliard de plicuri sunt sortate automat în SUA, 4% rată de eroare înseamnă extrem de mult, chiar dacă presupunem că 90% din plicuri sunt tipărite, iar în acest caz rata de recunoaștere este de aproape 100%. Ar rămâne aproximativ un milion de plicuri ce trebuie sortate manual, astfel că orice îmbunătățire a ratei de recunoaștere este importantă.

La modul general, îmbunătățirea ratei de recunoaștere este importantă în întreg domeniul recunoașterii formelor și la dezvoltarea inteligenței artificiale.

Într-un domeniu matur din prisma numărului de metode care sunt folosite pentru rezolvarea sa, nu este loc de filozofie. Din acest motiv, această teză va adera cu strictețe la cuvintele lui Sir William Thompson, Lord Kelvin (1824-1907): *"When you can measure what you are speaking about, and express it in numbers, you know something about it; but when you cannot measure it, when you cannot express it in numbers, your knowledge of it is of a meager and unsatisfactory kind; it may be the beginning of knowledge, but you have scarcely, in your thoughts, advanced it to the stage of science."* dedicând cea mai mare parte a ei experimentelor și comparării rezultatelor acestora.

*Teza de doctorat a fost susținută de grantul CNCSIS TD 2006-2007 nr. 58GR / 2006, cu titlul "Recunoașterea textelor scrise de mână în limba română".*

## 1.2 Organizarea tezei

Această teză este organizată după cum urmează:

**Capitolul 2.** Capitolul următor descrie stadiul actual al domeniului, începând cu clasificatorii pentru cifre izolate și finalizând cu clasificatorii pentru șiruri de cifre. În cazul celor din urmă se prezintă cele două abordări existente: holistică și cea care implică segmentarea.

**Capitolul 3.** Acest capitol prezintă succint aspecte teoretice ce vor fi folosite în partea experimentală a tezei. Întrucât nu se dorește încărcarea inutilă a lucrării cu informații care există prezentate în multe alte cărți și articole, aspectele teoretice vor fi doar schițate și apoi se vor oferi referințe pentru detalii suplimentare. Prima parte a capitolului prezintă aspecte legate de procesarea imaginii, iar a doua oferă detalii despre arhitectura rețelelor neuronale utilizate.

**Capitolul 4.** Bazele de date cu cifre și cu șiruri de cifre sunt detaliate în Capitolul 4. Sunt descrise modurile în care acestea sunt construite, relaționarea și alcătuirea lor. Se stabilesc seturile de date pentru antrenare și pentru testare. Seturile sunt astfel construite încât să poată fi efectuată o comparație riguroasă și sistematică cu rezultatele existente deja în domeniul recunoașterii cifrelor și a șirurilor de cifre scrise de mână.

**Capitolul 5.** Folosind seturile de antrenare și de test create, în Capitolul 5 sunt oferite rezultatele experimentelor efectuate. Întâi sunt antrenați și testați clasificatorii pentru cifre izolate și pentru perechi cifre parțial suprapuse, folosind două tipuri de rețele neuronale: perceptronul multistrat și rețele neuronale cu hărți de trăsături. Cei mai performanți clasificatori sunt apoi folosiți la crearea sistemului de recunoaștere pentru șiruri cu număr nedefinit de cifre. Sunt propuse și evaluate trei metode de selecție a clasificatorului corespunzător conținutului imaginii: o cifră sau o pereche de cifre.

**Capitolul 6.** Ca urmare a concluziilor capitolului anterior, penultimul capitol are ca scop studierea accelerării antrenării rețelelor neuronale. Sunt studiate șase variante pentru două arhitecturi distincte. Primele trei sunt concepute pentru arhitectura x86, folosind diverse optimizări ale codului C++, precum și aranjarea cât mai favorabilă a datelor în memorie. Una din variante apelează chiar la limbajul de asamblare și la instrucțiunile SIMD SSE3. Următoarele trei sunt proiectate pentru arhitecturi masiv paralele ce folosesc capacitățile de calcul ale procesoarelor plăcilor grafice (GPU). Prima este dezvoltată în BROOK+ pentru plăci AMD-ATI, iar următoarele două - în CUDA pentru plăci NVIDIA.

**Capitolul 7.** Teza se încheie cu o analiză a rezultatelor și cu evaluarea posibilelor dezvoltări ulterioare.

## 1.3 Descrierea metodei

**Obiectivul tezei este de a oferi o metodă eficientă de recunoaștere a șirurilor numerice scrise de mână, compuse dintr-un număr nedefinit de cifre izolate sau parțial suprapuse.** Întrucât marea majoritate a metodelor au ca pas anterior recunoașterii un proces de

segmentare care poate afecta negativ clasificarea, metoda pe care o propun nu utilizează segmentarea. Pentru a rezolva cazurile de cifre alipite sau parțial suprapuse am folosit pe lângă clasificatorul uzual, la nivel de cifră, un clasificator antrenat cu perechi de cifre parțial suprapuse.

În Figura 5.4 se poate observa calea parcursă de o imagine până când este recunoscut șirul numeric pe care îl conține. Blocul **CCA** extrage componentele conexe ale imaginii, separând cifrele izolate și toate celelalte părți de cifre. Etapa de **clusterizare** se ocupă cu concatenarea părților de cifre extrase în faza precedentă. Am creat un set de reguli care sunt aplicate iterativ asupra componentelor extrase. La finalul etapei se obțin componente care reprezintă cifre și grupuri de cifre unite.

**Clasificatorul** folosește rezultatele obținute de cei doi clasificatori (pentru cifre și pentru perechi de cifre parțial suprapuse) pentru a determina ce se află în imaginea componentei de la intrarea sa. Am implementat trei algoritmi simpli și eficienți care decid rezultatul căruia dintre cei doi clasificatori este corect. Algoritmii au un rol foarte important deoarece nu se cunoaște a priori câte cifre se află în componenta analizată. Alte metode din literatură folosesc pentru această operație antrenarea clasificatorilor cu exemple negative, ceea ce este destul de dificil deoarece exemplele negative trebuie extrase manual sau generate automat. În plus, crește dimensiunea setului de antrenare, ceea ce conduce la mărirea timpului necesar antrenării.

Întregul sistem de recunoaștere a fost testat cu seturi de imagini din NIST SD 19 [4]. Aceste seturile de imagini sunt folosite și de alți cercetători [44], cea mai bună rată de recunoaștere fiind de aproximativ 96%. Sistemul pe care l-am creat atinge o rată de recunoaștere de 93.77%, în condițiile în care nu folosește nici segmentare, nici antrenare cu exemple negative. Analiza erorilor arată că dacă seturile de antrenare ar conține și exemple negative, rata de recunoaștere ar fi la nivelul celor mai bune rezultate existente. Din acest motiv, în viitor setul de antrenare va fi augumentat cu exemple negative. Rezultatele au fost publicate în [17].

Ultima parte a tezei este dedicată analizei posibilităților de a crește viteza de execuție a algoritmului de antrenare a rețelelor neuronale. Acest lucru este necesar deoarece ar permite antrenarea unor rețele neuronale mai mari, obținându-se în acest mod rate de recunoaștere mai bune. Am încercat variante de optimizare pe CPU folosind inclusiv limbajul de asamblare și instrucțiuni SSE3. De asemenea, am efectuat teste de performanță folosind procesoare grafice (GPU) în BROOK+ și CUDA. Rezultatele preliminare sunt promițătoare, obținându-se o accelerare de peste 20 de ori față de varianta scrisă pentru CPU.

---

*Pentru că am evoluat ca specie de la homo habilis la homo sapiens sapiens, aceasta teză va încerca să prezinte cât mai mult din procentul de inspirație și, de multe ori fără succes, cât mai puțin din restul de 99% de muncă istovitoare și de multe ori plictisitoare.*



# Capitolul 2

## Stadiul actual al domeniului

Domeniul recunoașterii scrisului de mână și, în particular, cel al recunoașterii cifrelor și a șirurilor numerice, sunt cercetate intensiv de câteva decade. Din acest motiv, numărul de articole care tratează această problemă este destul de mare.

O căutare în IEEE Xplore returnează 608 articole ce conțin șirul "digit recognition". Numărul acestora crește la 1304 dacă nu impunem prezența alăturată a celor doi termeni. Având în vedere că există articole care sunt indexate în alte baze de date cu articole, numărul lor este și mai mare de atât.

Dacă efectuăm o căutare după "digit", "recognition", "numeral" (termeni care e foarte probabil să apară într-un articol având ca temă recunoașterea șirurilor numerice), obținem "numai" 60 de rezultate. Studiind anii în care au fost publicate aceste articole se observă că majoritatea au fost scrise în ultimii zece ani și niciunul nu e mai vechi de 1982. Aceste observații arată că:

- recunoașterea cifrelor este o problemă unde este foarte dificil de a mai propune metode noi care să fie mai eficiente decât cele deja existente (se va arăta în continuare că rata de recunoaștere a depășit 99.6%)
- problema recunoașterii cifrelor a fost abordată înaintea celei a recunoașterii șirurilor de cifre. Evident, unul din motive a fost și puterea scăzută de calcul a calculatoarelor
- în momentul în care recunoașterea cifrelor a început să ajungă aproape de limita maximă de recunoaștere și odată cu creșterea performanțelor calculatoarelor, atenția a fost îndreptată asupra problemei mult mai complicată a recunoașterii șirurilor numerice

Înainte de a prezenta cele mai importante realizări din domeniu se impune o explicație. Există multe articole care prezintă metode ce sunt verificate experimental pe seturi de test limitate (numărul de exemplare este prea redus pentru a avea însemnătate statistică) sau nedocumentate / private (deci imposibil de reprodus). Unele dintre aceste articole pretind rate de recunoaștere foarte mari, chiar de 100%. Întrucât nu există posibilitatea de a verifica rezultatele pretinse (decât, eventual, prin reimplementarea metodelor și verificarea lor cu alte seturi de test), am decis să consider aceste rezultate ca având o relevanță scăzută. Am considerat ca fiind relevante doar rezultatele articolelor care erau obținute folosind seturi de antrenare și de test accesibile (gratuit sau nu) și care, mai important, erau folosite de cât mai mulți alți autori. Un alt mod de a aprecia însemnătatea rezultatelor unui articol a fost și rangul conferinței

sau al jurnalului unde a fost prezentat / publicat respectivul articol. Sumarizând, rezultatele pe care le-am considerat a fi relevante au fost obținute cu baze de date publice, folosite de cât mai mulți cercetători din domeniu, și au fost publicate la conferințe renumite (ex: NIPS, ICDAR, CVPR, IPCV, AACL, ICPR etc). Bineînțeles, acesta nu este un mod infailibil de a afla care sunt cele mai bune rezultate în domeniu, dar consider că este o soluție cu probabilitate mare de a oferi rezultate relevante fără a examina exhaustiv toate articolele existente, ceea ce este practic imposibil ([http://www.visionbib.com/bibliography/contentschar.html# Numbers, Digits, Zip \(Postal\) Codes](http://www.visionbib.com/bibliography/contentschar.html# Numbers, Digits, Zip (Postal) Codes)).

Nefolosirea acelorași seturi de test de către toate articolele este o problemă ce afectează tot domeniul recunoașterii scrisului de mână. Există câteva baze de date standard și voi considera numai rezultatele obținute folosind seturile de test din aceste baze de date. Două din aceste baze de date sunt disponibile contra cost: NIST SD 19 [4] și CEDAR [22], iar una este oferită gratuit: MNIST [40].

## 2.1 Stadiul actual al recunoașterii cifrelor

Întrucât este foarte puțin probabilă conceperea unui sistem de recunoaștere a șirurilor numerice ce nu încorporează un clasificator pentru cifre, vor fi prezentate cele mai importante realizări în acest domeniu. Experimentele sunt efectuate îndeosebi cu seturi extrase din NIST SD 19 [4] sau cu seturile MNIST [40]. Detalii despre aceste baze de date se găsesc în Capitolul 4.

Jie Zhou [70] propune un model de rețea neuronală cuantică. Este detaliat algoritmul de antrenare și sunt oferite rezultatele experimentelor asupra unui set de 2000 de cifre din baza de date CENPARMI (<http://www.cenparmi.concordia.ca/>). Rata maximă de recunoaștere este de 95.65%.

Nishiwaki [53] propune o metodă de recunoaștere a șirurilor numerice care verifică existența cifrelor care se ating. Sunt propuse șase tipuri de atingere între cifre. Deși această metodă verifică perechi de cifre, ea este dezvoltată în ideea de a fi folosită pentru recunoașterea cifrelor izolate, în scopul rejectării imaginilor care conțin mai mult de un caracter. Nu sunt oferite date referitoare la ratele de recunoaștere.

Yann LeCun a creat baza de date MNIST ca subset al NIST SD 19 pentru a oferi un punct comun tuturor celor care doreau să-și verifice metodele de recunoaștere. Față de NIST SD 19 care este foarte generală și oferă date brute, lasând la latitudinea utilizatorului modul de alcătuire al seturilor de antrenare și de test, MNIST propune două seturi, unul pentru antrenare și unul pentru test, ce conțin 60000, respectiv 10000 de imagini de cifre. Ratele de eroare raportate pentru MNIST variază de la 12% pentru un perceptron fără straturi ascunse [38], deci un clasificator liniar, până la 0.39% pentru rețele neuronale de convoluție [48].

Franke [27] prezintă o metodă de îmbunătățire a clasificatorilor polinomiali prin învățare iterativă. Rezultatul cel mai bun obținut pe setul MNIST este de 1.13% rată de eroare.

Lauer et al. [36] a creat un sistem de recunoaștere ce folosește vectori suport. Rezultatele sunt foarte bune (eroare 0.54%), dar sunt necesare resurse computaționale mari.

Cele mai multe încercări de rezolvare a acestei probleme folosesc rețele neuronale. Simard, Steinkraus și Platt [64] demonstrează că rețelele neuronale oferă cea mai bună rată de recunoaștere pe setul de test MNIST. Rezultatele prezentate de LeCun, Simard și Ranzato sunt la fel de bune, oferind o rată de eroare de 0.39% [60, 12]. Performanțe foarte bune au fost obținute și prin combinarea clasificatorilor [13].

Pe pagina MNIST (<http://yann.lecun.com/exdb/mnist/>) există câteva zeci de articole ale căror metode au fost validate cu seturile de date din MNIST.

**Având în vedere rezultatele foarte bune obținute pentru recunoașterea de cifre cu rețele neuronale de convoluție, am ales drept metodă pentru propriul clasificator la nivel de cifră utilizarea rețelelor neuronale de convoluție (RNC).**

## 2.2 Stadiul actual al recunoașterii șirurilor numerice

Așa cum arătam în Capitolul 1, recunoașterea șirurilor numerice ce conțin cifre care se suprapun este foarte importantă deoarece este utilă într-un număr mare de aplicații practice: recunoașterea codurilor poștale [24], procesarea documentelor complexe [34].

### 2.2.1 Recunoașterea holistică

Recunoașterea holistică încearcă să evite procesul de segmentare prin implementarea de clasificatori care primesc la intrare toată informația. Facând o paralelă cu recunoașterea textelor, un clasificator holistic nu segmentează un cuvânt în litere urmând să recunoască literele, ci recunoaște tot cuvântul dintr-odată. La recunoașterea cuvintelor acest lucru este posibil deoarece numărul de cuvinte, chiar dacă este foarte mare, este limitat. În schimb, la recunoașterea holistică a șirurilor numerice numărul claselor este mult mai mare, fiind imposibilă crearea unei metode de acest fel pentru cazul general. Din acest motiv clasificatorii holistici se rezumă la șiruri numerice de două, maxim trei cifre.

În [67] este descrisă o metodă de recunoaștere a cifrelor care se ating. Acest articol folosește ca seturi de antrenare și de test imagini de coduri poștale provenite de pe corespondență. Nu se dă niciun alt detaliu asupra datelor folosite, deci experimentele nu pot fi verificate. Oricum, autorii pretind o rată de recunoaștere de 86.8% pe setul limitat pe care l-au folosit (291 de imagini de test). Sunt oferite și rezultate (86.5%) pe un set neprecizat nici măcar ca număr de imagini din CEDAR [22]. Se afirmă că pe același set de test o metodă care implică segmentarea oferă rate de recunoaștere cu până la 6.1% mai mici. Clasificatorul utilizat analizează structura imaginilor, căutând concavități, măsurând gradienti etc, informații cu care apoi construiește un scor ce caracterizează imaginea. O observație interesantă a articolului este că din cele 500 de coduri poștale analizate, 20% conțin cifre unite, 16% fiind perechi de cifre, iar restul de 4% fiind alcătuite din trei sau mai multe cifre unite. În [35] aceiași autori prezintă un sistem multiexpert pentru recunoașterea șirurilor numerice cu cifre care se ating. De această dată setul folosit este mult mai mare, 28464 de exemplare cu cifre unite fiind colectate manual. Și acest set este privat. Cea mai bună rată de recunoaștere este de 91.1%.

Behnke [8] (Capitolul 7.4) a creat un sistem de recunoaștere pentru mărci poștale (timbre

valorice imprimate pe plicuri) care utilizează un clasificator convoluțional pentru a recunoaște perechi de cifre tipărite. Chiar dacă este o aplicație OCR (Optical Character Recognition = recunoașterea caracterelor tipărite) și nu a recunoașterii scrisului de mână, am amintit și acest sistem de recunoaștere deoarece abordează problema în mod holistic.

LeCun et al. [39] propun o metodă de recunoaștere a codurilor poștale folosind rețele neuronale de convoluție. Ideea este dezvoltată mai departe de grupul său în [49, 50]. Pentru a evita segmentarea se construiește o rețea neuronală extinsă pe orizontală în ideea de a încăpea în câmpul de intrare și contextul cifrei ce se dorește a fi recunoscută. Rețeaua este antrenată cu triplete de cifre suprapuse sau nu, generate automat. Acest clasificator este aplicat în fiecare poziție a imaginii pentru a detecta cifrele. Deoarece se fac multe aplicări ale clasificatorului, metoda este destul de lentă (între timp puterea de calcul a crescut de câteva zeci de ori, acest aspect ne mai constituind o problemă). Testele pe care le-au efectuat asupra codurilor poștale de cinci cifre au arătat o rată destul de redusă de recunoaștere: 70% pentru setul de antrenare și 66% pentru cel de test.

## 2.2.2 Recunoașterea cu segmentare

Multe din soluțiile propuse în literatura de specialitate combină procesele de segmentare și recunoaștere. Având în vedere că multe cifre se pot atinge sau pot fi chiar parțial suprapuse, efectuarea unei segmentări corecte care să nu țină cont de conținutul imaginii este foarte dificilă [69]. Dacă se cunoaște conținutul imaginii, înseamnă că deja etapa de clasificare a avut loc.

Casey și Lecolinet [11] prezintă o abordare care integrează segmentarea și recunoașterea folosind o strategie ipoteză-test (recunoaștere bazată pe segmentare). Cordella și Vento [21] analizează în detaliu această metodă.

Shi, Wakabayashi et al. [47] propun o metodă de a crește gradul de recunoaștere când numărul de imagini din setul de antrenare este redus. Algoritmul generează o imagine în oglindă a imaginii unei clase care aparține unei perechi de clase foarte apropiate. Imaginea rezultată este adăugată setului de imagini al celeilalte clase.

Au fost implementate numeroase metode care folosesc segmentarea și o extragere complicată de caracteristici pentru a obține cifre din șirul inițial [45]. Aceste cifre sunt recunoscute folosind un clasificator la nivel de cifră. Procesul de segmentare este aplicat în două moduri distincte [59]. În primul, modulul de segmentare transferă secvențele de cifre izolate la clasificator [63]. Oliveira et al. [58] au rafinat acest concept propunând o verificare efectuată prin suprasedgmentare și prin subsegmentare. Aceiași autori au construit un sistem de recunoaștere folosind o metodă euristică de suprasedgmentare bazată pe clasificatori cu vectori suport [57]. Al doilea mod utilizează ipoteze motivate probabilistic, iar decizia finală este dată de cel mai bun scor segmentare-recunoaștere obținut pentru imaginea de intrare [29]. Acest procedeu este utilizat de Cheong et al. [14] pentru a genera mai multe ipoteze de segmentare bazate pe secvența de tușe ce formează o posibilă cifră.

Autorii din [32] propun o metodă de preprocesare a segmentării cu scopul îmbunătățirii ratei de recunoaștere pentru cifrele care se ating. Experimentele sunt efectuate pe NIST SD 19. Rata de recunoaștere este de 96.4%. Seturile de test sunt similare cu cele pe care le-am

utilizat, dar mai mici, doar 1000 de imagini, în loc de aproape 1500 cât au seturile 3DNS și 6DNS. Oricum, rezultatele pot fi folosite pentru comparații. Aceiași autori au dezvoltat o metodă de recunoaștere [33] care folosește segmentarea simplă și segmentarea combinată cu recunoașterea. Întâi (segmentare fără recunoaștere) sunt separate cifrele izolate de cifrele care se ating. Apoi sunt examinate și segmentate cifrele care se ating. Se folosesc 1500 de imagini din NIST SD 19, dar nu se specifică indicii lor. Rata de recunoaștere este de 91.8%.

**Cel mai important articol recent din domeniul recunoașterii șirurilor numerice scrise de mână este [44].** A apărut în anul 2004 și sumarizează cele mai performante metode apărute anterior. De asemenea, oferă o modalitate clară de comparare a rezultatelor, specificând cu exactitate proveniența și alcătuirea seturilor de test. Am reconstruit aceste seturi de test pentru propriile experimente (Secțiunile 4.5.2 și 4.5.3). Autorii folosesc abordarea clasică a problemei: segmentare în cât mai multe părți posibile, urmată de o concatenare ce are ca scop obținerea de cifre singulare. Etapa de concatenare este combinată cu cea de recunoaștere, analizându-se toate modurile de concatenare. Se obține un graf orientat, ce are două noduri importante între care se alină toate modurile în care pot fi asamblate părțile de imagine. Pentru reducerea numărului de drumuri dintre cele două noduri se impun restricții geometrice, cea mai importantă dintre ele fiind gradul de proximitate a două părți de simboluri. Fiecare astfel de drum este apoi analizat și în urma clasificării primește un scor obținut ca sumă a scorurilor oferite de clasificatorul la nivel de cifră. Drumul cu cel mai mare scor este ales ca fiind cel mai probabil a fi corect. Existând mai multe drumuri, este ușor de implementat un algoritm de rejecție care respinge un anumit rezultat ca probabil incorect în momentul în care diferența de scor dintre cele mai bune două drumuri este sub un anumit prag. În articol sunt analizate mai multe metode de clasificare pentru cifre: trei clasificatori neuronali, doi cu densitate discriminativă și doi care folosesc vectori suport.

Un alt aspect important al articolului este **antrenarea cu exemple negative** (imagini care nu conțin cifre) a clasificatorilor pentru recunoașterea de cifre. Se face o comparație între rezultatele obținute utilizând seturi de antrenare ce conțin sau nu exemple negative. Prezența non-cifrelor în setul de antrenare îmbunătățește rata de recunoaștere cu până la 30% pentru unele metode (rețele neuronale și clasificatori cu vector suport).

La nivel de cifră, rata maximă de recunoaștere este de 99.26%, obținută cu un clasificator cu vectori suport ce folosește funcții radiale. Adăugarea de exemple negative în setul de antrenare produce o scădere a ratei de recunoaștere la 99.24%. La nivel de șir numeric, ratele maxime de recunoaștere se situează la 96.82% pentru șiruri de trei cifre și la 96.74% pentru șiruri de șase cifre (Figura 2.1). Toate rezultatele sunt obținute folosind clasificatori antrenați cu exemple negative. Același autor propune o metodă combinată de recunoaștere [43] folosind atât segmentarea, cât și recunoașterea holistică. Utilizând o metodă adaptată din recunoașterea vorbirii, obține o rată de recunoaștere de 97.42%.

Figura 2.1 sumarizează cele mai bune rezultate obținute în domeniul recunoașterii șirurilor de cifre scrise de mână. Coloana marcată *forced* conține ratele de recunoaștere cu rejecție nulă. RBF înseamnă funcții radiale, PC - clasificator polinomial, DLQDF - clasificator cu funcții discriminante cuadratică, bazat pe învățare discriminativă, SVC-rbf - clasificator cu vectori suport și funcții radiale. Detalii despre implementarea clasificatorilor pot fi găsite în [44].

Method	strlen	#string	Forced	1% err	0.5% err
[35]	3	986	92.7	79.5	70.5
	6	982	90.3	75.5	66.5
[36]	3	956	97.0	N/A	N/A
	6	961	96.7	N/A	N/A
[14]	3	2,385	95.38	90.60	89.81
	6	2,167	93.12	85.68	84.03
RBF	3	1,176	96.00	91.87	89.70
	6	1,471	96.13	91.98	88.72
PC	3	1,176	96.34	93.77	91.33
	6	1,471	96.16	93.20	91.71
DLQDF	3	1,476	96.54	93.36	88.21
	6	1,471	96.19	91.16	87.70
SVC-Rbf	3	1,176	96.82	95.26	91.17
	6	1,471	96.74	94.15	91.64

Figura 2.1: Ratele de recunoaștere pentru șiruri din NIST SD 19. Preluare din [44]. [35] reprezintă [30], [36] reprezintă [9], iar [14] reprezintă [58].

## 2.3 Concluzii

Sistemul propriu de recunoaștere a șirurilor de numere conține părți care există și în alte sisteme de recunoaștere. Mai exact, este vorba de clasificatorul la nivel de cifră. Întrucât ideea tezei nu este de a concepe un nou clasificator (eventual, mai bun) pentru cifre, am decis să folosesc cea mai bună metodă existentă pentru recunoașterea de cifre în propriul sistem de recunoaștere a șirurilor numerice. Astfel, am ales rețelele neuronale de convoluție. De asemenea, pentru acest clasificator am extins setul de antrenare folosind deformări afine și elastice, așa cum au făcut și alți autori (Simard [64] a introdus aceste deformări, aplicându-le pe setul MNIST).

Primul pas în construirea unui clasificator pentru recunoașterea șirurilor numerice cu cifre ce se ating sau se suprapun parțial este reprezentat de rezolvarea cazului cel mai simplu de suprapunere, cel dintre două cifre.

Având în vedere că majoritatea cercetătorilor folosesc metode ce implică segmentarea și că ratele de recunoaștere sunt deja foarte mari (peste 96%), iar recunoașterea holistică este mult mai puțin studiată, am hotărât să aplic o metodă holistică pentru perechile de cifre parțial suprapuse, în ideea evitării segmentării.

Deoarece rețelele neuronale de convoluție s-au dovedit a fi foarte bune pentru recunoașterea de cifre, le-am folosit și pentru construirea unui clasificator la nivel de perechi de cifre parțial suprapuse. Rețelele neuronale s-au dovedit a fi eficiente în foarte multe domenii, fiind și foarte rapide în procesul de clasificare, aspect care este important în aplicații ce trebuie să ruleze în timp real sau cât mai repede posibil (de exemplu un sortator de corespondență).

# Capitolul 3

## Aspecte teoretice

### 3.1 Noțiuni de prelucrarea imaginilor

În această teză, imaginile vor fi abstractizate în felul următor: o colecție de pixeli structurați într-un tabel bidimensional, accesibili prin indicele liniei (coordonata  $y$ ) și cel al coloanei (ordonata  $x$ ). Indicii sunt valori întregi nenegative. Fiecare pixel are o anumită culoare, pentru experimentele efectuate fiind suficient a considera cazul tonurilor de gri cu 256 valori discrete, 0 însemnând alb și 255 - negru. Dacă pixelii pot avea doar două valori, atunci 0 reprezintă alb și 1 - negru. În cazul în care se face excepție de la aceste reguli, se precizează în text. Pentru imagini color se aplică o transformare în tonuri de gri folosind, de exemplu, metoda prezentată în [15] la pagina 9 sau în [16].

#### 3.1.1 Centrul de masă al unei imagini

Din fizică se cunoaște că pentru sistem de  $n$  particule punctiforme, fiecare fiind caracterizată prin masa  $m$  și prin coordonatele  $x_i$  într-un spațiu euclidian  $p$ -dimensional ( $i \in \{1, \dots, p\}$ ), coordonatele centrului de masă pot fi calculate cu formula:

$$x_i^{CM} = \frac{\sum_{j=1}^n m_j \cdot x_i^j}{\sum_{j=1}^n m_j} \quad (3.1)$$

În cazul unei imagini bidimensionale cu  $n \times m$  pixeli ( $n$  linii și  $m$  coloane), intensitatea fiecărui pixel fiind reprezentată prin  $p_{xy}$ , cu  $x \in \{0, \dots, m-1\}$  și  $y \in \{0, \dots, n-1\}$ , putem adapta formula de calcul a coordonatelor centrului de masă în felul următor:

$$x^{CM} = \frac{\sum_{x=0}^{m-1} \sum_{y=0}^{n-1} p_{xy} \cdot x}{\sum_{x=0}^{m-1} \sum_{y=0}^{n-1} p_{xy}} \quad (3.2)$$

$$y^{CM} = \frac{\sum_{x=0}^{m-1} \sum_{y=0}^{n-1} p_{xy} \cdot y}{\sum_{x=0}^{m-1} \sum_{y=0}^{n-1} p_{xy}} \quad (3.3)$$

### 3.1.2 Distorsiuni

În experimentele efectuate cu imagini de cifre am folosit augumentarea setului de antrenare prin aplicarea de distorsiuni asupra imaginilor sale. Transformările efectuate sunt similare cu cele aplicate de Simard [64]. Lauer et al. [36] detaliază aceste transformări. Distorsiunile au constat în mai multe transformări liniare: rotații, scalări și deformări elastice.

#### Transformări afine

Intuitiv, transformările afine într-un spațiu euclidian bidimensional pot fi privite ca transformări geometrice în urma cărora se păstrează coliniaritatea punctelor și proporțiile segmentelor.

Transformările afine aplicate asupra imaginilor de antrenare sunt de două feluri: rotații și scalări. Se calculează câte un câmp de deplasare pentru fiecare direcție: orizontală ( $\Delta X$ ) și verticală ( $\Delta Y$ ). Astfel, pentru punctul de coordonate  $(x, y)$  vom avea deplasările  $\Delta X(x, y)$  și  $\Delta Y(x, y)$ .

**Scalările** pot fi reprezentate matricial în felul următor:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix} \times \begin{pmatrix} x \\ y \end{pmatrix} \quad (3.4)$$

unde  $\alpha$  reprezintă coeficientul de scalare pe orizontală și  $\beta$  reprezintă coeficientul de scalare pe verticală.

**Rotațiile** în formă matricială se calculează cu formula:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} \cos \gamma & \sin \gamma \\ \cos \gamma & -\sin \gamma \end{pmatrix} \times \begin{pmatrix} x \\ y \end{pmatrix} \quad (3.5)$$

unde  $\gamma$  reprezintă unghiul de rotație a imaginii.

#### Transformări elastice

În [64] transformările elastice sunt obținute prin prelucrarea unor câmpuri de deplasare aleatoare. Astfel, se inițializează un câmp de deplasare cu valori reale aleatoare din  $[-1, 1]$ . Distribuția valorilor se consideră a fi uniformă. Asupra acestui câmp se aplică o convoluție cu un nucleu gaussian cu deviația standard  $\sigma$ .

Pentru două dimensiuni, nucleul gaussian normalizat are următoarea formulă:

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} \cdot e^{-\frac{x^2+y^2}{2\pi\sigma^2}} \quad (3.6)$$

Lățimea nucleului este dată de valoarea lui  $\sigma$ .

Rezultatul acestor operații este un câmp de deplasare care aplicat unei imagini produce o deformare elastică în cazul în care  $\sigma$  are valori mici. Dacă  $\sigma$  are valori mari, deformarea devine transformare afină, iar pentru  $\sigma$  foarte mare se ajunge la translație [64]. Se folosește un coeficient  $\alpha$  cu care sunt multiplicare valorile nucleului pentru a controla intensitatea locală a deformării.



### 3.1.3 Transformata Hough

Transformata Hough a fost inventată de Paul Hough în 1962 și patentată de IBM. Utilizarea inițială a constat în detectarea traiectoriilor particulelor încărcate electric din experimentele efectuate în camera cu bule. Richard Duda și Peter Hart ([23]) au creat transformata Hough generalizată în 1972. Dacă metoda inițială detecta numai linii, transformata generalizată poate detecta orice fel de forme [62] (pag. 382-383), în particular segmente de dreaptă, arce de cerc și de elipsă.

Se va prezenta succint transformata Hough pentru detectarea de linii dintr-o imagine. Problema poate fi formalizată astfel: fiind dată o mulțime de puncte prin coordonatele  $x$  și  $y$  dintr-un spațiu euclidian bidimensional, se cere să se determine care din aceste puncte sunt coliniare și care nu. Pentru aceasta, se aplică o transformare a spațiului euclidian la spațiul Hough, utilizând ecuația 3.7.

$$r = r_1 + r_2 = \frac{y}{\cos \theta} + \sin \theta (y - \tan \theta) = x \sin \theta + y \cos \theta \quad (3.7)$$

În ecuația de mai sus,  $\theta$  reprezintă unghiul dreptei cu axa  $Ox$ , iar  $r$  este distanța de la originea sistemului de coordonate la dreaptă (Figura 3.1). Deoarece pentru fiecare punct există o infinitate de perechi în spațiul parametric  $(r, \theta)$ , trebuie găsită o modalitate de reducere a dimensionalității problemei pentru a putea fi aplicată practic. Pentru aceasta, se discretizează problema prin discretizarea valorilor parametrilor. Astfel, dacă imaginea inițială, în forma discretă, avea  $N \times M$  pixeli, putem discretiza spațiul parametrilor alegând, de exemplu,  $\Delta r = 1 \text{ pixel}$  și  $\Delta \theta = 1^\circ$ . Discretizarea se face în funcție de problemă.

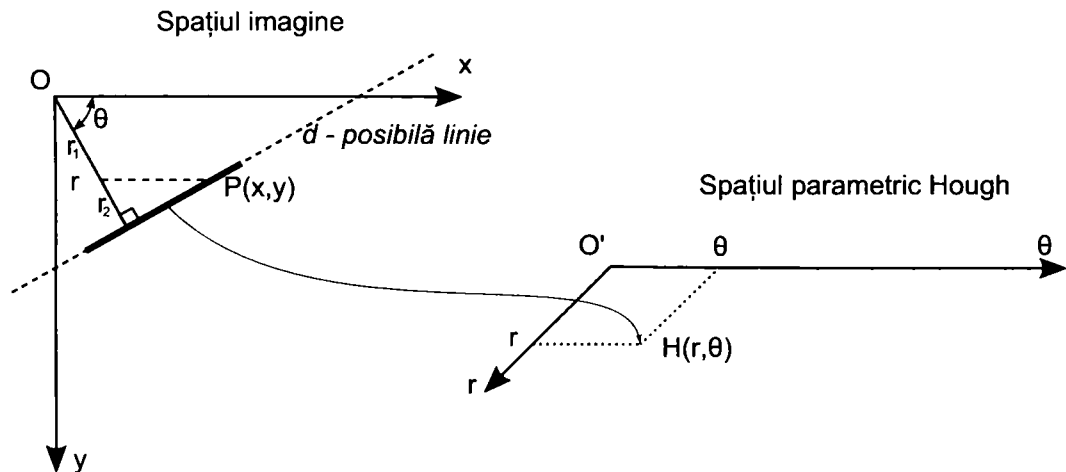


Figura 3.1: Transformata Hough pentru detecția de linii. Fiecărei linii  $d$  din spațiul imagine îi corespunde un punct în spațiul parametric Hough.

Pasul următor rezidă în construirea unei matrici cu  $\frac{x + \sqrt{x^2 + y^2}}{\Delta r} \times \frac{180}{\Delta \theta}$  elemente, unde linia reprezintă distanța, iar coloana reprezintă unghiul. Elementele matricii sunt inițial nule. Pentru toate valorile discrete posibile ale lui  $r$  și  $\theta$ , se construiesc dreptele corespunzătoare și se numără prin câte din punctele imaginii trec. Acest număr reprezintă valoarea celulei respective. Matricea mai poartă numele de acumulator Hough. Un exemplu de astfel de matrice este prezentat în Figura 3.2. Folosind acest algoritm, la sfârșit, în matrice se vor afla numere

întregi nenegative care reprezintă numărul de puncte prin care trece fiecare dreaptă. Cercurile marchează punctele de maxim ale matricei, puncte care reprezintă drepte cu număr maxim de puncte.

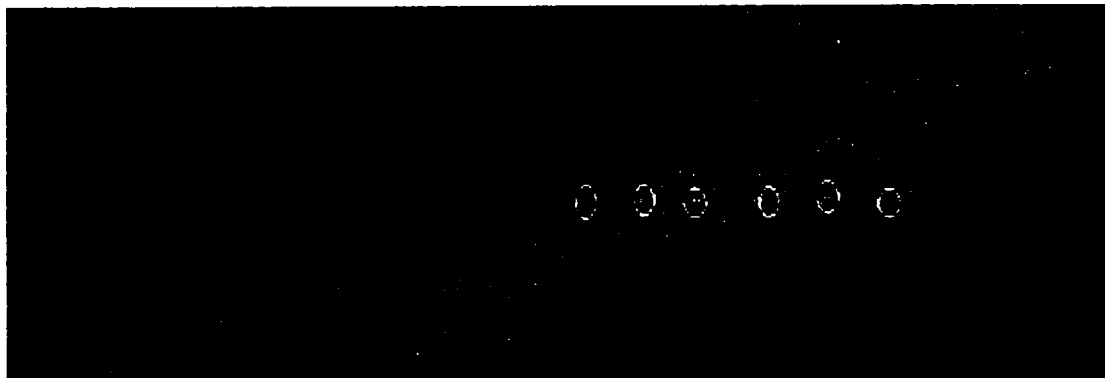


Figura 3.2: Exemplet de matrice de acumulare rezultată în urma aplicării transformatei Hough. Cercurile marchează maximele corespunzătoare celor șase linii detectate. Afișarea grafică folosește negru pentru elementele nule și tonuri de gri cu iluminare proporțională cu valoarea respectivului element pentru valori nenule. Albul se folosește pentru elementul sau elementele cu număr maxim de puncte.

## 3.2 Rețele neuronale artificiale

Rețelele neuronale artificiale (RNA) fac parte din subdomeniul *Machine Learning* (*Învățare automată*) al *Inteligenței Artificiale*. RNA sunt metode statistice de clasificare, cea mai importantă aplicație a lor fiind în recunoașterea formelor. Există multe cărți care tratează acest domeniu, unele din cele mai importante fiind [10] (abordează întreg domeniul recunoașterii formelor) și [9] (dedicată exclusiv rețelelor neuronale artificiale). Pentru o introducere concisă și clară se poate utiliza [51] (Capitolul 6).

RNA sunt alcătuite din unități de procesare numite neuroni. Neuronul artificial este un model matematic simplist al neuronului biologic. Neuronii unei RNA sunt, de obicei, dispuși pe mai multe straturi.

RNA pot fi împărțite în două mari categorii: RNA cu propagare înainte și RNA recurente. Prima categorie este de departe mult mai studiată, RNA cu propagare înainte fiind mult mai bine înțelese și din punct de vedere teoretic. RNA utilizate în această teză sunt de tipul RNA cu propagare înainte.

Într-o RNA cu propagare înainte există conexiuni numai între neuroni de pe straturi consecutive, informația de la intrare propagându-se secvențial prin fiecare strat până la ieșire.

În scopul evitării confuziilor, în această teză am considerat ca strat al rețelei neuronale inclusiv stratul de intrare, strat necomputațional. Astfel, pentru o rețea cu trei straturi, stratul  $L_0$  este stratul de intrare, iar straturile  $L_1$  și  $L_2$  sunt straturi computaționale. Straturile interioare, în acest caz stratul  $L_1$ , se numesc straturi ascunse.

### 3.2.1 Perceptronul

Perceptronul [46] (pagina 458) este un model simplu de rețea neuronală cu propagare înainte. Are două straturi, unul fiind de intrare și unul de ieșire. Stratul de ieșire este unicul strat computațional. Întrucât nu poate rezolva decât probleme liniar separabile, are o arie redusă de aplicabilitate.

Perceptronul multistrat este o variantă mai complexă de perceptron, putând avea unul sau mai multe straturi ascunse. Prezența straturilor ascunse elimină limitarea aplicabilității la probleme liniar separabile.

### 3.2.2 Rețele neuronale artificiale cu hărți de convoluție (RNC)

Având în vedere modul de funcționare al percepției vizuale [8] (Capitolul 2), s-a încercat modelarea acestuia folosind rețele neuronale. Unul din aceste modele folosește hărți de convoluție. RNC sunt tot RNA multistrat cu propagare înainte, unica diferență fiind arhitectura rețelei, mai exact, modul de conectare a neuronilor. Evident, RNC sunt concepute pentru recunoașterea de caracteristici ale imaginilor.

În 1979 Fukushima [28] a propus o rețea neuronală cu hărți de convoluție, numită Neocognitron, care se putea adapta la deplasarea caracterelor în câmpul de intrare. Detalii pot fi găsite la [http://www4.ocn.ne.jp/~fuku\\_k/files/research-e.html#neocog](http://www4.ocn.ne.jp/~fuku_k/files/research-e.html#neocog).

LeCun a folosit în 1995 o rețea de convoluție pentru recunoașterea de caractere. În Figura 3.3 este prezentată arhitectura acestei rețele. În [38] arhitectura este dezvoltată și mai mult, prin adăugarea de noi straturi.

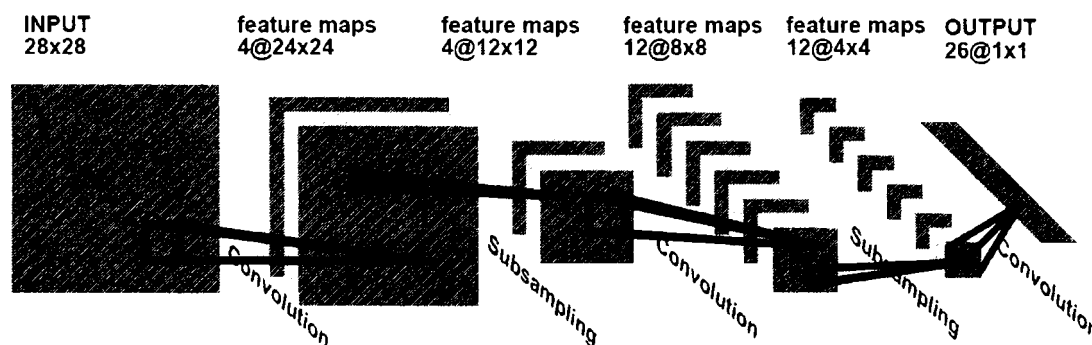


Figura 3.3: RNC pentru recunoașterea caracterelor. Preluare din [37].

Spre deosebire de perceptron, unde fiecare neuron al unui strat era conectat cu toți neuronii stratului anterior, la RNC fiecare neuron este conectat numai la o parte a neuronilor stratului anterior. Pentru a imita funcționarea sistemului vizual uman, neuronii cu care este conectat se află toți în aceeași zonă. O rețea neuronală poate conține atât straturi complet conectate, cât și straturi de convoluție. Straturile de convoluție sunt, uzual, alcătuite din mai multe hărți de convoluție.

În RNC utilizate în Capitolul 5 există două moduri de conectare a neuronilor de pe straturile de convoluție, al doilea mod fiind o extindere a primului [64, 38].

**Primul mod de conectare** (stratul  $L_1$  din Figura 3.5) constă în a conecta fiecare neuron al unei hărți la un nucleu de neuroni ai stratului anterior. Toți neuronii unei hărți împart aceleași

ponderi, astfel numărul de ponderi fiind mult redus față de cazul unui strat complet conectat. Nucleul de convoluție este deplasat asupra stratului anterior din doi în doi neuroni, atât pe orizontală, cât și pe verticală (Figura 3.4). Dimensiunea uzuală [64] pentru nucleu este de  $5 \times 5$  neuroni. Aceasta trebuie să fie impară pentru ca nucleul să poată fi centrat într-un pixel. Simard [64] a arătat că o dimensiune mai mică,  $3 \times 3$  neuroni, este prea redusă pentru a detecta caracteristici primare (colțuri, intersecții, capete de linie) în imagine, iar  $7 \times 7$  neuroni este deja prea mult, depășind dimensiunea caracteristicilor primare.

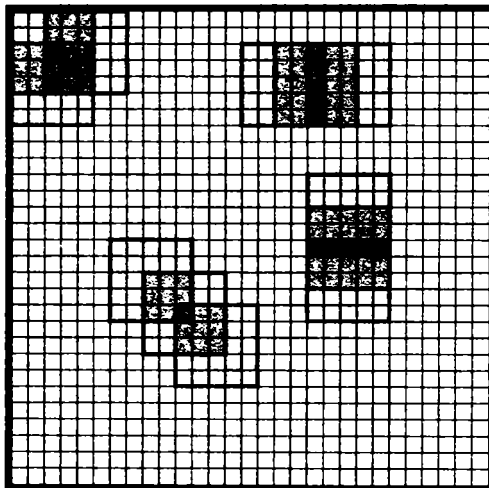


Figura 3.4: Modul de deplasare a nucleului de convoluție.

**Al doilea mod de conectare** poate fi observat între straturile  $L_1$  și  $L_2$  din Figura 3.5. Și în acest caz neuronii dintr-o hartă împart același set de ponderi, dar fiecare neuron de pe stratul  $L_2$  este conectat cu neuroni din nuclee de convoluție din *fiecare* hartă a stratului anterior.

### 3.2.3 Metode de antrenare a rețelelor neuronale

Antrenarea rețelelor cu propagare înainte se face de obicei cu algoritmul de învățare cu propagare înapoi a erorii (backpropagation learning). Acesta este prezentat în toate cărțile ce tratează rețelele neuronale, astfel că nu va mai fi expus în această lucrare. Programul de antrenare a rețelelor neuronale pe care le-am folosit în cadrul experimentelor este conceput după explicațiile din [61](Capitolul 20.5). Am ales această carte deoarece expune într-un mod matematic riguros și detaliat algoritmul de antrenare.

În cazul algoritmilor ce au ca principiu minimizarea erorii medii pătratice, se poate folosi o tehnică de accelerare a antrenării. Algoritmul Levenberg-Marquardt [41] (o descriere a sa pentru RNA este dată de Bishop în [9], Capitolul 7.11) poate fi aplicat pentru minimizarea funcției de eroare.

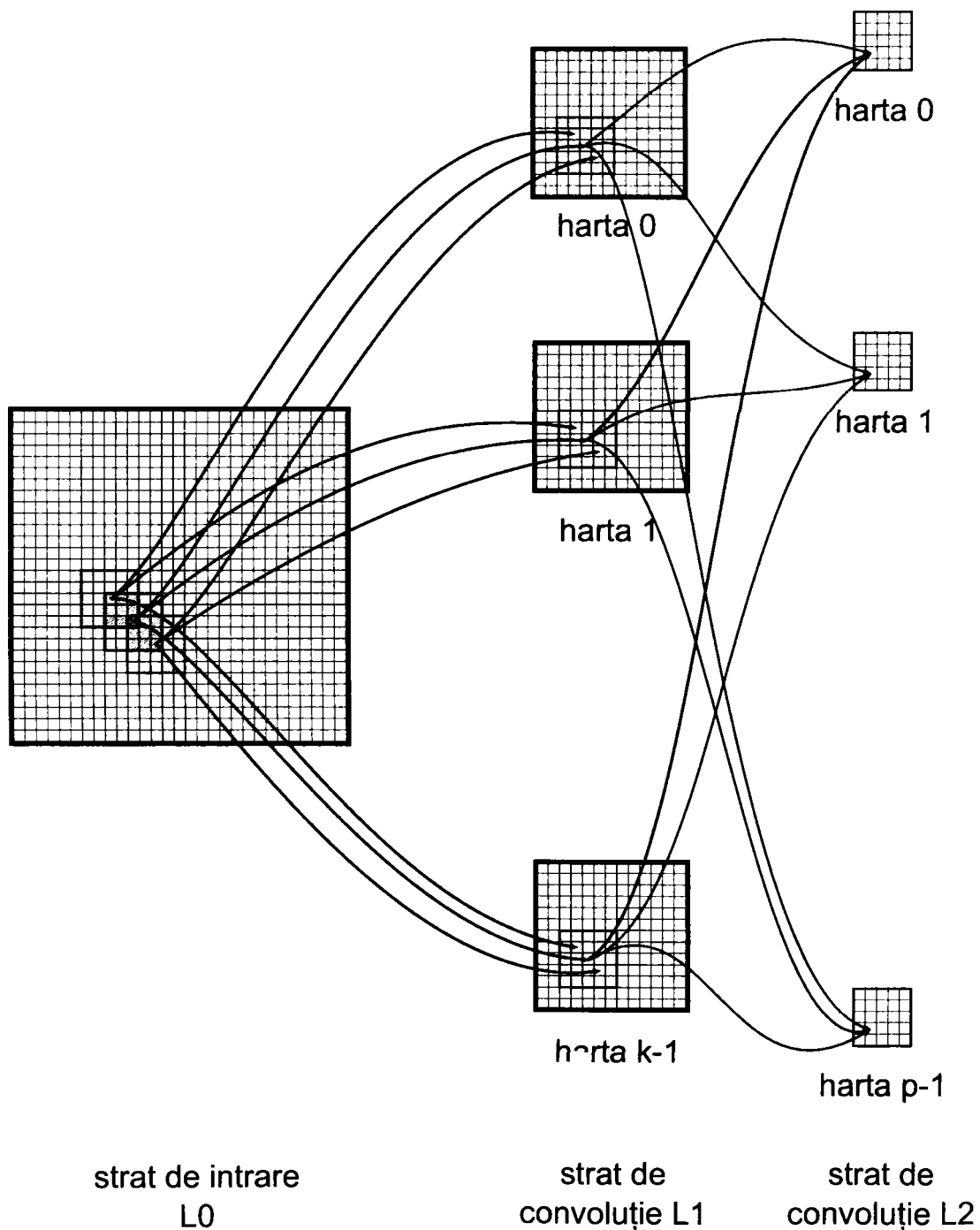
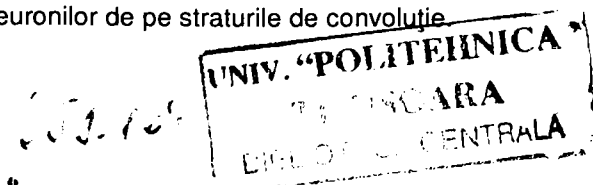


Figura 3.5: Modul de conectare a neuronilor de pe straturile de convoluție



# Capitolul 4

## Baze de date cu cifre și șiruri de cifre

Acest capitol prezintă bazele de date cu cifre și cu șiruri de cifre care au fost folosite pentru obținerea seturilor de antrenare și de test. Cu aceste seturi au fost apoi efectuate experimentele din capitolele ce urmează.

Prima bază de date, NIST SD 19 ([4]), este și cea mai mare dintre toate. Celelalte trei, MNIST ([40]), LNIST și baza de date la nivel de două cifre parțial suprapuse sunt derivate din aceasta. Seturile de test 3DNS, 6DNS și 10DNS sunt obținute tot din NIST SD 19, prin extragerea câmpurilor numerice de 3, 6, respectiv 10 cifre. Figura 4.1 prezintă modul în care relaționează seturile de antrenare și de test, iar Tabelul 4.1 - conținutul și modul de utilizare al seturilor. În Figura 4.1 am evidențiat și clasificatorii la nivel de cifră și de două cifre parțial suprapuse ce vor fi folosiți în sistemul de clasificare general (Figura 5.4).

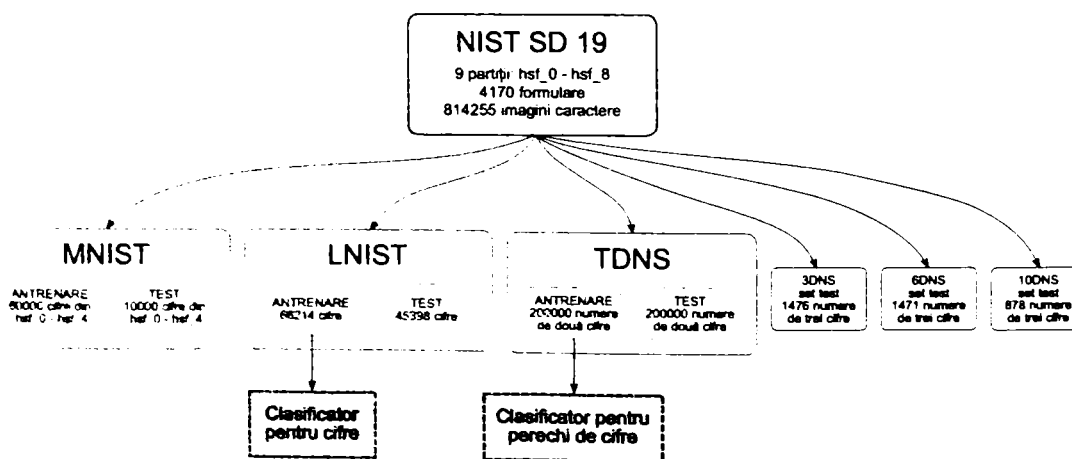


Figura 4.1: Ierarhia seturilor de test și de antrenare.

Tabelul 4.1: Conținutul și modul de utilizare ale seturilor de date.

nume	număr de imagini	tip imagini	utilizare	nr. persoane
NIST SD 19	814255 + 4170	variat	variata	4170
MNIST TRAIN	60000	cifre	antrenare	≈500
MNIST TEST	10000	cifre	test	necunoscut
LNIST TRAIN	66214	cifre	antrenare	600
LNIST TEST	45398	cifre	test	400
TDNS TRAIN	200000	numere a 2 cifre	antrenare	1000
TDNS TEST	200000	numere a 2 cifre	test	1000
3DNS	1476	numere a 3 cifre	test	296
6DNS	1471	numere a 6 cifre	test	296
10DNS	878	numere a 10 cifre	test	296

## 4.1 NIST SD 19

Baza de date NIST SD 19 ([4]) conține imagini de formulare scanate scrise de 4170 de persoane. Un exemplu de astfel de formular este dat în Figura 4.2. Formularele vor fi referite în text prin HSF (Handwriting Sample Form). În continuare vor fi prezentate numai informațiile necesare înțelegerii desfășurării experimentelor.

CD-ul cu baza de date conține gata extrase, eventual segmentate, cifrele și literele din aceste formulare, în total 814255 de imagini. Pe lângă imagini, baza de date conține detalii privind originea, organizarea, data publicării respectivelor imagini, precum și modul recomandat de utilizare a imaginilor. Caracterele segmentate sunt memorate în imagini de  $128 \times 128$  pixeli și sunt etichetate folosind codurile ASCII corespunzătoare caracterelor 0-9, a-z, A-Z. Autorii NIST SD 19 afirmă că etichetarea imaginilor a fost verificată manual, iar eroarea care mai există este de aproximativ 0.1%.

NIST SD 19 conține toate imaginile din versiunile anterioare, SD3 și SD7 (Tabelul 4.2), plus multe altele. Imaginile segmentate sunt organizate în diverse ierarhii pentru a facilita moduri variate de utilizare a lor. Astfel, există referințe pentru ierarhizarea imaginilor în funcție de persoană, clasă, clasă fără informație de specificare a tipului literei (majusculă sau minusculă), câmpul din formularul HSF.

Tabelul 4.2: Versiuni anterioare ale bazei de date NIST cu imagini de caractere scrise de mână. (după [4])

nume	data	versiune	conținut	observații
Special Database 1	mai 1990	HWDB Release 1-1.1	formulare și câmpuri	
Special Database 3	februarie 1992	HWSC Release 4-1.1	caractere	
Special Database 7	aprilie 1992	TST1 Release 6-1.1	caractere	Test Data 1.

### 4.1.1 Structura NIST SD 19

NIST CD 19 conține opt partiții de imagini HSF, numerotate  $hsf_{\{0,1,2,3,4,6,7,8\}}$ . HSF-urile au fost scanate la 300dpi și, așa cum se poate observa și în Figura 4.2, conțin 34 de câmpuri

HANDWRITING SAMPLE FORM

NAME [REDACTED] DATE Sept 16, 1989 CITY SAN ANTONIO TEXAS STATE TEXAS ZIP 78237

This sample of handwriting is being collected for use in testing computer recognition of printed numbers and letters. Please print the following characters in the boxes that appear below.

0123456789	0123456789	0123456789		
<u>0123456789</u>	<u>0123456789</u>	<u>0123456789</u>		
08	395	5268	49171	235969
<u>08</u>	<u>995</u>	<u>5268</u>	<u>49171</u>	<u>235969</u>
111	2956	81207	758298	90
<u>111</u>	<u>2956</u>	<u>81207</u>	<u>758298</u>	<u>90</u>
4671	34560	368704	27	493
<u>4671</u>	<u>34560</u>	<u>368704</u>	<u>27</u>	<u>493</u>
75434	281512	02	564	3000
<u>75434</u>	<u>281512</u>	<u>02</u>	<u>564</u>	<u>3000</u>
335706	48	863	4699	82771
<u>335706</u>	<u>48</u>	<u>863</u>	<u>4699</u>	<u>82771</u>

languoybxqdtpemchrvfisjwz

languoybxqdtpemchrvfisjwz

MOFWTUIXYHSDZKEBQRLPGCVNJA

MOFWTUIXYHSDZKEBQRLPGCVNJA

Please print the following text in the box below:

We, the People of the United States, in order to form a more perfect Union, establish Justice, insure domestic Tranquility, provide for the common Defense, promote the general Welfare, and secure the Blessings of Liberty to ourselves and our posterity, do ordain and establish this CONSTITUTION for the United States of America.

We, the people of the United States, in order to form a more perfect Union, establish Justice, insure domestic Tranquility, provide for the common Defense, promote the general Welfare, and secure the Blessings of Liberty to ourselves and our posterity, do ordain and establish this CONSTITUTION for the United States of America.

Figura 4.2: Exemplu de formular din NIST SD 19 (formularul f1816\_24). A se vedea și Figura 4.4.



Tabelul 4.3: Bazele de date în care au fost publicate partițiile. † - câmpuri completate, dar neprocesate. ∅ - câmpuri necompletate. (după [4])

partiția	persoane	formulare	cifre	majuscule	minuscule	constituția	originea persoanelor
hsf_0	0000-0499	SD1	SD3	SD3	SD3	SD19	Census Field
hsf_1	0500-0999	SD1	SD3	SD3	SD3	SD19	Census Field
hsf_2	1000-1499	SD1	SD3	SD3	SD3	SD19	Census Field
hsf_3	1500-2099	SD1	SD3	SD3	SD3	†	Census Field
hsf_4	2100-2599	SD19	SD7	SD7	SD7	∅	High School
hsf_5	2600-3099	păstrate de NIST pentru teste viitoare				∅	High School
hsf_6	3100-3599	SD19	SD19	SD19	SD19	∅	Census MD
hsf_7	3600-4099	SD19	SD19	SD19	SD19	∅	Census MD
hsf_8	4100-4169	SD19	†	†	†	†	Census MD

cu informații diverse. În această lucrare au fost utilizate doar câmpurile numerice.

În Tabelul 4.3 este dată alcătuirea NIST SD 19, precum și conținutul SD1, SD3 și SD7.

Primele 429 de formulare din *hsf\_5* au fost scrise de elevi de liceu, iar ultimele 71, precum și toate formularele din partițiile *hsf\_{6,7,8}* au fost scrise de angajați ai biroului de recensământ din Suitland, Maryland. Informațiile acestea se vor dovedi utile la construirea seturilor de antrenare și de test, precum și la evaluarea rezultatelor.

## 4.1.2 Conținutul NIST SD 19

Datele fiecărei persoane se află în directoare cu nume de forma *fyyyy\_zz*, unde *yyyy* reprezintă codul persoanei, putând lua valori între 0000 și 4099, iar *zz* este un cod prin care se precizează tipul formularului completat de persoana respectivă. Toate formularele au același număr de câmpuri, câmpurile cu același număr având aceeași dimensiune. Variaza numai ordinea cifrelor, respectiv a caracterelor, din interiorul câmpului. Acest cod poate fi utilizat pentru a accesa formularul corespunzător persoanei respective, în directorul *data/hsf\_pages/truerefs*. Conținutul (Tabelul 4.4) acestor fișiere poate fi utilizat pentru a afla ce ar fi trebuit să scrie persoana în câmpurile formularului. Am folosit aceste informații la construirea seturilor de test cu numere formate din 3, 6 și 10 cifre. În Figura 4.4 este descris conținutul câmpurilor din formularul din Figura 4.2. Formularele există și în variantă necompletată, pentru a putea fi utilizate la colectări suplimentare de date.

Câmpurile de la 3 până la 31 conțin cifre și numere formate din mai multe cifre. Conținutul câmpurilor a fost extras și e gata segmentat la nivel de caracter (cifră sau literă) în baza NIST SD 19. Toate caracterele au fost etichetate automat folosind fișierele cu referințe, și apoi verificate manual ([68]) de două ori de către persoane distincte. În urma segmentării imaginile au fost stocate într-o ierarhie de fișiere. Pentru a afla exact proveniența și conținutul unei imagini se folosește numărul partiției, numărul persoanei și tipul formularului.

Formularele au fost concepute astfel încât pentru fiecare clasă să existe același număr de exemplare. După operațiile de extragere și de segmentare, aceste numere nu mai sunt egale, dar sunt destul de apropiate (Tabelul 4.5).

Din păcate, baza de date nu conține și imaginile cu șirurile de cifre extrase din câmpurile formularelor, ci numai cu cifrele gata segmentate din aceste imagini. Din acest motiv, pentru experimentele asupra șirurilor de cifre a fost necesară dezvoltarea unei metode de extragere a acestor șiruri de cifre.

Tabelul 4.4: Exemplu de fișier cu referințe (*hsf\_page/truerefs/ref\_24.txt*). A se vedea și Figura 4.2.

fld_0		fld_1	
fld_2		fld_3	0123456789
fld_4	0123456789	fld_5	0123456789
fld_6	08	fld_7	395
fld_8	5268	fld_9	49171
fld_10	235969	fld_11	111
fld_12	2956	fld_13	81207
fld_14	758298	fld_15	90
fld_16	4671	fld_17	34560
fld_18	368704	fld_19	27
fld_20	493	fld_21	75434
fld_22	281512	fld_23	02
fld_24	564	fld_25	3000
fld_26	335706	fld_27	48
fld_28	863	fld_29	4699
fld_30	82771	fld_31	languoybxqdtpemchrvisjzkw
fld_32	MOFWTUIXYHSDZKEBQRLPGCVNJA	fld_33	

Tabelul 4.5: Numărul de exemplare din fiecare clasă.

Clasa	0	1	2	3	4	5	6	7	8	9
Exemplare	40363	44704	40072	41112	39154	36606	39937	41893	39579	39533

### 4.1.3 Ierarhizarea datelor

Datele sunt ierarhizate după: formulare și partiții, autor, câmp, numărul clasei în hexa, numărul clasei unite (fără discriminare între literele minuscule și cele majuscule). Din aceste ierarhii prezintă interes pentru lucrarea de față numai cele la nivel de formular și partiție, precum și cea la nivel de autor.

Formatul fișierelor este detaliat în [4]. NIST SD 19 conține patru tipuri de fișiere, tipul lor fiind determinat de extensie:

- .mis* – fișier ce conține (mai multe) imagini de caractere. De obicei este comprimat.
- .pct* – fișier ce conține imaginea unui formular. Există câte un astfel de fișier pentru fiecare din cele 3669 de persoane.
- .cls* – fișier ce conține clasele (etichetele) imaginilor din fișierul *.mis* cu același nume.
- .mit* – fișier ce conține un pointer la un fișier *.mis*. Conține identitatea persoanei, precum și calea și numele fișierului cu imagini segmentate ale acelei persoane. Deoarece există mai multe ierarhii de date, imaginile ar fi trebuit multiplicat pentru fiecare ierarhie, ceea ce ar fi consumat inutil spațiu. Pentru a evita acest lucru, și pentru a permite stocarea întregii baze de date pe un singur CD, imaginile apar o singură dată -în fișiere *.mis*-, iar în rest doar se face referire la ele folosind fișiere *.mit*.

```

hsf_page
|
|
|
hsf_0 ... hsf_8          trurefs          template
|                       |                 |
|                       |                 |
f0000_14.pct ... f0499_10.pct  ref_00.txt ... ref_99.txt  ref_00.ps ... ref_99.ps

```

Figura 4.3: Arborele aferent ierarhizării formularelor în partiții.

### Ierarhizare după partiții și formulare

Imaginile formularelor sunt structurate într-un arbore (Figura 4.3) ce are pe primul nivel partițiile *hsf\_{0,1,2,3,4,6,7,8}* și directorul pentru referințe *trurefs*. De asemenea, mai există și directorul pentru formularele necompletate (*template*).

### Ierarhizare după persoane

Acest arbore (Figura 4.4) conține caracterele organizate după partiție și apoi după persoana care le-a scris. Datele din această ierarhie reprezintă datele prime obținute după procesul de segmentare, celelalte ierarhii fiind doar reorganizări ale datelor din aceasta. În directorul fiecărei persoane se află patru perechi de fișiere *.mis* și *.cls*, câte una pentru cifre, litere minuscule, litere majuscule, caractere din câmpul Constituție. Fișiere *.mit* nu sunt necesare deoarece toate datele se află în directorul curent.

#### 4.1.4 Formatul fișierelor NIST

Imaginile sunt stocate în format binar, alb-negru. Dimensiunea unei imagini ce reprezintă un simbol este de  $128 \times 128$  pixeli. Informația de culoare a pixelilor este memorată la nivel de bit, astfel că o imagine necesită  $128 \times 128 / 8 = 2048$  octeți. Imaginile simbolurilor sunt stocate una după cealaltă, iar blocul rezultat este comprimat cu o variantă a algoritmului CCITT Group 4. Toate fișierele care conțin imagini au un antet de tip IHEAD atașat blocului ce conține imaginile comprimate. Informații detaliate despre modul de codificare pot fi găsite în [4]. Tipul acesta de codificare se folosește la fișierele *.mis*, ele fiind unicele care conțin imagini de caractere. Fișierele care conțin imaginea unui formular întreg sunt de tipul *.pct*.

```

by_write
|
|
|
hsf_0          hsf_7
|
|
f0000_14 f0001_41 ... f0499_10
|
|
d0000_14.mis u0000_14.mis l0000_14.mis c0000_14.mis
d0000_14.cls u0000_14.cls l0000_14.cls c0000_14.cls

```

Figura 4.4: Arborele aferent ierarhizării imaginilor de caractere după partiții și persoane.

Fiecărui caracter segmentat din formulare îi este atribuită o etichetă prin care se specifică din ce clasă face parte. Această eticheta reprezintă chiar codul ASCII al caracterului, în cazul cifrelor având valori între  $30_{(16)}$  și  $39_{16}$ . Etichetele cifrelor dintr-un fișier *.mis* se află memorate într-un fișier *.cls* corespunzător. Acesta este un fișier text ce are pe prima linie numărul de cifre care se află atât în *.mis* cât și în *.cls*. Pe rândurile următoare se află câte o etichetă.

### 4.1.5 Programe utilitare ale NIST SD 19

Pe lângă imaginile din baza de date mai există și utilitare pentru comprimarea și decompimarea imaginilor. Deoarece aceste utilitare sunt scrise în C pentru UNIX, a fost necesară modificarea lor pentru arhitectura x86. Întrucât numărul funcțiilor ce au trebuit adaptate este de peste o sută și se află în câteva zeci de fișiere (peste 300kB de surse C), portarea utilităților pe arhitectura 80x86 nu a fost deloc facilă. Din multele utilitare ale NIST SD 19, am modificat utilitarul *normmis* care normalizează imaginile dintr-un fișier *.mis* la dimensiunea de  $32 \times 32$  pixeli. Pentru a face această normalizare -care era nedorită pentru modul în care doream să folosesc imaginile din *.mis*- *normmis* întâi decompime fișierul *.mis*. Am folosit acest lucru și am modificat programul astfel încât să salveze imaginile în format *.bmp*, câte un fișier pentru fiecare imagine din *.mis*. Am obținut peste o sută de mii de fișiere *.bmp*, fiecare conținând o imagine de  $128 \times 128$  pixeli. La salvarea fiecărui fișier *.bmp* am avut grijă ca numele să conțină toate datele de identificare a imaginii: formular, persoană, etichetă. Pentru crearea ulterioară a seturilor de antrenare și de test am scris un utilitar care pe baza unui fișier de configurare alege doar anumite imagini din cele peste 100000 de bmp-uri.

## 4.2 MNIST

Baza de date MNIST (Modified NIST) [40] a fost creată de Corinna Cortez și apoi modificată de Yann LeCun prin centrarea imaginilor în centrul de masă.

Conține două seturi de imagini etichetate, unul pentru antrenare alcătuit din 60000 de imagini de cifre, și unul pentru testare alcătuit din 10000 de imagini. Așa cum îi spune și numele, este o modificare a NIST SD 19. Imaginile în format alb negru din NIST au fost redimensionate la  $20 \times 20$  pixeli cu păstrarea aspectului. În urma procesului de redimensionare, imaginile sunt reprezentate în tonuri de gri. Imaginea de  $20 \times 20$  pixeli este apoi poziționată într-una de  $28 \times 28$  pixeli, prin suprapunerea centrului de masă al primeia cu centrul imaginii mari.

Baza de date MNIST conține imagini de cifre scrise de mână din NIST SD 3 și din NIST SD 1. În urma multor experimente cu aceste baze de date, Yann LeCun a constatat că SD 3 este mult mai ușor de recunoscut decât SD 1 deoarece conține imagini scrise de personal de la biroul de recensământ al populației, iar SD1 conține imagini cu cifre scrise de elevi de liceu. Deoarece a dorit ca rezultatele sistemului de clasificare să nu fie influențate de persoanele care au scris respectivele cifre, a hotărât să creeze o nouă bază de date prin amestecarea datelor din SD 1 și SD 3. Pentru aceasta, a construit setul de antrenare folosind 30000 de imagini din SD 1 și 30000 de imagini din SD 3. Pentru setul de test a folosit câte 5000 de imagini atât din SD 1, cât și din SD 3. Setul de antrenare conține date de la aproximativ 500 de persoane, iar mulțimile de persoane ale căror imagini de cifre au fost folosite pentru construirea seturilor de antrenare și de test au fost disjuncte.

Întrucât în SD 1 imaginile cifrelor pentru fiecare persoană nu se află la un loc, au fost folosite informațiile din fișierele "cls" (care conțin date despre proveniența fiecărei imagini a unei cifre) pentru rearanjarea cifrelor în funcție de persoana care le-a scris. Detaliile procedurii pot fi urmărite în [40].

Baza de date este structurată în patru fișiere:

- train-images-idx3-ubyte: imaginile setului de antrenare,
- train-labels-idx1-ubyte: etichetele setului de antrenare,
- t10k-images-idx3-ubyte: imaginile setului de test,
- t10k-labels-idx1-ubyte: etichetele setului de test.

În Anexa 7.3 este prezentat modul de codificare a informațiilor din fișiere.

### 4.3 LNIST - Liu-NIST

Aceste seturi de antrenare și de test au fost construite pentru a putea face o comparație cât mai echitabilă între metoda propusă în această lucrare și metodele din [44]. Ele conțin exact aceleași imagini ca seturile folosite în [44]. Setul de test este compus din 66214 cifre extrase din 600 de formulare NIST: 0-399 și 2100-2299. Pentru extragerea acestor cifre am folosit unele din utilitarele menționate în Secțiunea 4.1. După extragerea imaginilor comprimate din fișierele *.mis*, acestea au fost decomprimate și s-au obținut imagini alb-negru cu dimensiunea de  $128 \times 128$  pixeli. Cifrele din aceste imagini au fost încadrate în dreptunghiul ce le delimitează, rezultatul fiind redimensionat cu păstrarea aspectului la  $20 \times 20$  pixeli. În urma redimensionării, imaginile au fost transformate din imagini binare (0/1) în imagini cu tonuri de gri (0 ÷ 255). Apoi s-a calculat centrul de masă al fiecărei imagini și bitmap-ul de  $20 \times 20$  a fost poziționat într-unul de  $29 \times 29$  prin suprapunerea centrului de masă cu centrul imaginii de  $29 \times 29$  (Figura 4.5). Toate aceste prelucrări au fost efectuate cu un program C++ care are ca rezultat final două fișiere de tip *.idx* (a se vedea Anexa 7.3), unul conținând cele 66214 imagini, iar celălalt etichetele corespunzătoare, etichete care au fost obținute din fișierele *.cls* corespunzătoare fișierelor *.mis*.

Setul de test a fost obținut în același mod, dar folosind imaginile cu cifre din 400 de formulare: 500-699 și 2400-2599. În total sunt 45398 de imagini.

Niciuna din cele 111612 imagini nu a fost curățată sau îmbunătățită în niciun fel.

Avantajul folosirii centrării în centrul de masă față de metoda dreptunghiului delimitator rezidă în sensibilitatea scăzută la zgomotele ce ar putea apărea în imagine. La metoda dreptunghiului delimitator e suficient ca un singur pixel de zgomot să se afle la distanță de imaginea cifrei pentru ca imaginea să apară translatată puternic la intrarea clasificatorului. Centrarea în centrul de masă atenuează foarte mult acest efect al zgomotelor exterioare cifrei propriu-zise, păstrând toate imaginile în aproximativ aceeași poziție a câmpului de intrare al clasificatorului. În Figura 4.5 este evidențiat un caz de zgomot care translatează orizontal imaginea în cazul utilizării metodei dreptunghiului delimitator (imaginile din dreapta). Ultimul rând de imagini arată mai clar efectul centrării în centrul de masă, prin suprapunerea imaginilor de pe coloanele corespunzătoare. Se observă că în cazul centrării în centrul de masă imaginea obținută prin suprapunere arată mult mai puțin distorsionată față de imaginea obținută prin suprapunerea imaginilor centrate folosind metoda dreptunghiului delimitator. Cel mai bine este evidențiată diferența pentru golul cifrei nouă, gol care este mult diminuat în cazul din dreapta. Dacă ar fi existat și un exemplar al cifrei nouă cu zgomot în partea opusă

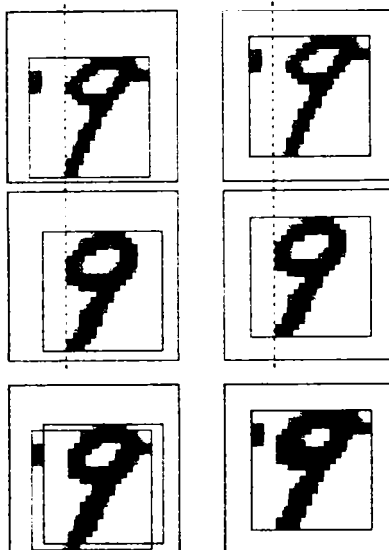


Figura 4.5: Comparație între centrarea în centrul de masă (stânga) și centrarea în dreptunghiul delimitator (dreapta). Punctul roșu reprezintă centrul de masă al imaginii. Pe ultimul rând se află suprapunerea celor două imagini de pe coloana corespunzătoare.

-în dreapta-, iar la un set de antrenare format din zeci de mii de imagini aproape sigur există și astfel de exemplare, atunci golul cifrei nouă până la urmă ar fi dispărut dacă nu s-ar fi folosit metoda centrării în centrul de masă.

#### 4.4 TDNS (Two-digit Numeral Strings)

Această bază de date am creat-o pentru antrenarea unui clasificator care să recunoască șiruri de două cifre parțial suprapuse. NIST SD 19 conține câmpuri cu numere formate din două cifre, dar numărul lor este relativ redus deoarece există numai 4170 de persoane și pentru fiecare persoană sunt câte 4 câmpuri de două cifre. Dar extragerea acestor informații din formulare nu este facilă deoarece trebuia implementat un program care să detecteze și să extragă imaginile din formular. Programul a fost implementat deoarece am avut nevoie și de imaginile cu numere de trei sau mai multe cifre. Însă, multe din imagini cu două cifre nu conțin cifre parțial suprapuse, astfel că numărul final de imagini pentru antrenarea rețelei ar fi insuficient deoarece antrenarea necesită o cantitate considerabilă de date.

Considerând numerele formate din două cifre, există 100 de clase, de la 00 la 99. Pentru fiecare clasă sunt necesare cel puțin câteva sute de exemplare pentru ca antrenarea să reușească. Colectarea manuală, urmată de procesare și etichetare a atât de multor date este aproape imposibilă, oricum, foarte îndelungată și plictisitoare. În consecință, sunt necesare metode alternative. Am prezentat anterior mai multe baze de date care conțin imagini cu cifre etichetate. Metoda alternativă aleasă constă în conceperea unei metode de generare automată a imaginilor cu două cifre parțial suprapuse. Pentru intrare am folosit imagini din NIST SD 19 ([4]) și nu din MNIST, deoarece în MNIST imaginile au fost deja prelucrate în diverse moduri și am dorit să nu influențez în niciun fel rezultatele procesului de recunoaștere.

NIST SD 19 conține nouă partiții, de la  $hsf_0$  la  $hsf_8$ . Am folosit partițiile  $hsf_0 \div hsf_3$  care conțin imagini de la 2100 de persoane. Din aceste date am generat automat atât setul de antrenare, cât și cel de test. În NIST SD 19 se specifică a se folosi partițiile  $hsf_0 \div hsf_3$  pentru antrenare, dar aceste specificații sunt date în ideea procesului de recunoaștere la nivel de cifră, neaplicându-se și pentru șiruri de cifre. Un alt motiv pentru care nu am respectat îndrumările din NIST este acela că partiția  $hsf_4$  care este menționată a se folosi ca set de test conține numai 500 de persoane, și am dorit să creez seturi de test cât mai diverse și cât mai mari. În consecință, am utilizat imaginile primelor 1000 de persoane pentru antrenare, și imaginile următoarelor 1000 de persoane pentru test. Detaliile vor fi prezentate în secțiunile dedicate setului de antrenare, respectiv celui de test.

#### 4.4.1 Setul de antrenare

Acest set conține primele 1000 de persoane, exceptându-le pe cele cu indicii: 48, 52, 741, 785, 802, 825 și 848. Nu s-au utilizat datele de la aceste persoane deoarece erau incomplete în baza de date NIST SD 19. Lipsa unuia din fișierele cu imagini sau cel cu etichete, sau nu existau cifre pentru unele clase. Rezultatul este format din 1000 de grupuri a zece imagini de cifre. Rezoluția în NIST SD 19 este  $128 \times 128$  pixeli, și imaginile sunt codificate alb-negru. Am efectuat inițial experimente direct cu aceste imagini fără a le corecta în vreun fel, dar erorile au fost semnificative, indiferent de complexitatea rețelei neuronale utilizate. Era evident că imaginile au probleme, astfel că le-am studiat mai în amănunt și am observat că multe conțin pixeli disparați sau artefacte generate de preprocesarea NIST. (Figura 4.6)

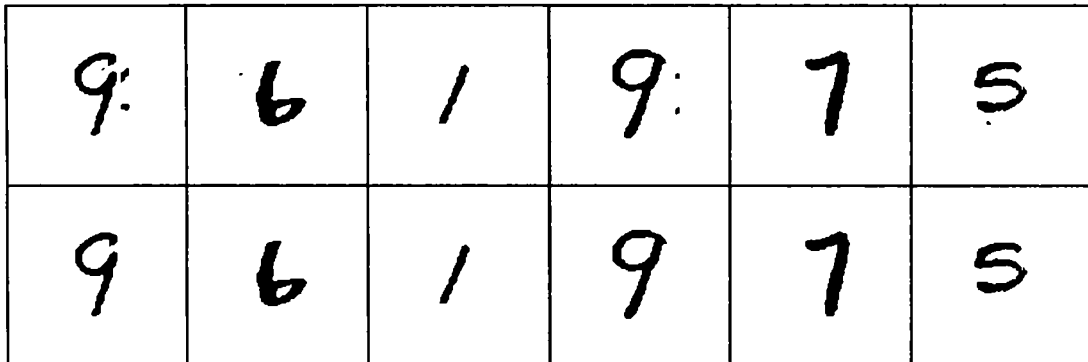


Figura 4.6: Primul rând - exemple de artefacte în imagini  $128 \times 128$  pixeli din NIST SD 19; rândul al doilea - imaginile rezultate în urma procesului de corecție semiautomată.

Acești pixeli și artefacte nu afectează direct capacitatea de recunoaștere a rețelei neuronale. Dacă imaginile sunt utilizate pentru un clasificator la nivel de cifră, nu este necesară nicio corectare a lor. Demonstrația acestui fapt este prezentată în Capitolul 5 la experimentele efectuate cu cifre singulare. Aceste artefacte afectează negativ procesul de unire a cifrelor. Problema majoră rezidă în calcularea poziției de suprapunere, cauzând generarea de imagini nerealistice. Pentru a evita astfel de erori am hotărât să elimin artefactele și pixelii eronați. Am folosit o metodă parțial automată prin scrierea unui program care detectează numărul componentelor conexe din imagine. Dacă acesta este mai mare decât unu, imaginea este prezentată operatorului uman și fiecare componentă conexă este marcată cu o altă culoare

(Figura 4.9). În majoritatea cazurilor eliminarea tuturor componentelor exceptând-o pe cea mai mare a fost suficientă pentru corecție. A existat un număr redus de cazuri de cifre cu linii întrerupte. Soluția aleasă pentru corectarea lor a fost ștergerea manuală a componentelor care ar fi afectat procesul de suprapunere. Într-un număr foarte redus au existat cazuri de imagini foarte deformate rezultate în urma unui proces defectuos de extragere când baza de date NIST SD 19 a fost colectată și procesată. Aceste cazuri au fost corectate manual. Câteva exemple pot fi observate în Figura 4.7.

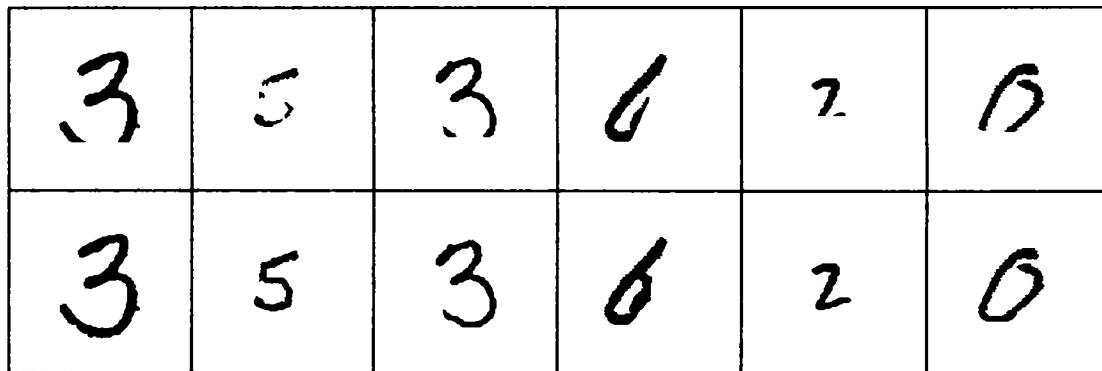


Figura 4.7: Primul rând - imagini defectuos segmentate sau deformate din baza NIST SD 19; rândul al doilea - imagini rezultate în urma corecției manuale.

Fiecare imagine de  $128 \times 128$  pixeli conține un caracter care poate avea dimensiuni foarte diferite. A fost necesară o redimensionare pentru normalizarea cifrelor din setul de imagini. S-a detectat dreptunghiul delimitator al cifrei și apoi conținutul acestuia a fost scalat proporțional până la dimensiunea de  $20 \times 20$  pixeli (4.8). Rezultatul scalării este o imagine cu tonuri de gri afectată de aliasing. Am ales această dimensiune deoarece este suficient de mare pentru a reprezenta complet și în detaliu un caracter, și este mai mică decât marea majoritatea a imaginilor din baza de date NIST. De asemenea, și baza de date MNIST folosește aceeași dimensiune pentru redimensionare și am dorit o tratare unitară pentru toate experimentele următoare.

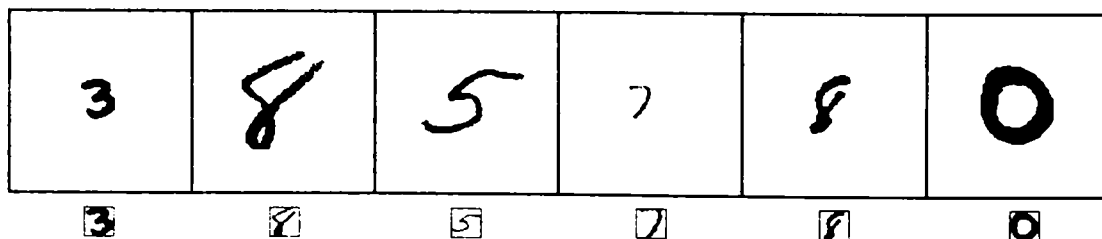


Figura 4.8: Primul rând - imagini  $128 \times 128$  pixeli din NIST SD 19; rândul al doilea - imagini  $20 \times 20$  pixeli rezultate după tăiere relativă la dreptunghiul delimitator, centrare și redimensionare.

Am pornit de la fișierele cu imagini comprimate *.mis* și de la fișierele cu etichete *.cls* utilizate în NIST SD 19 și după ce am extras datele de care am avut nevoie, le-am salvat în formatul *.idx* care are avantajul că este un format necomprimat care poate fi utilizat în mod



direct. De asemenea, poate conține tonuri de gri, nefiind limitat doar la alb sau negru. Acest tip de fișiere e folosit și de Yann LeCun ([40]) și am dorit să tratăm unitar toate fișierele cu imagini și etichete. Tipul fișierelor *.idx* este prezentat în Anexele 1 și 2.

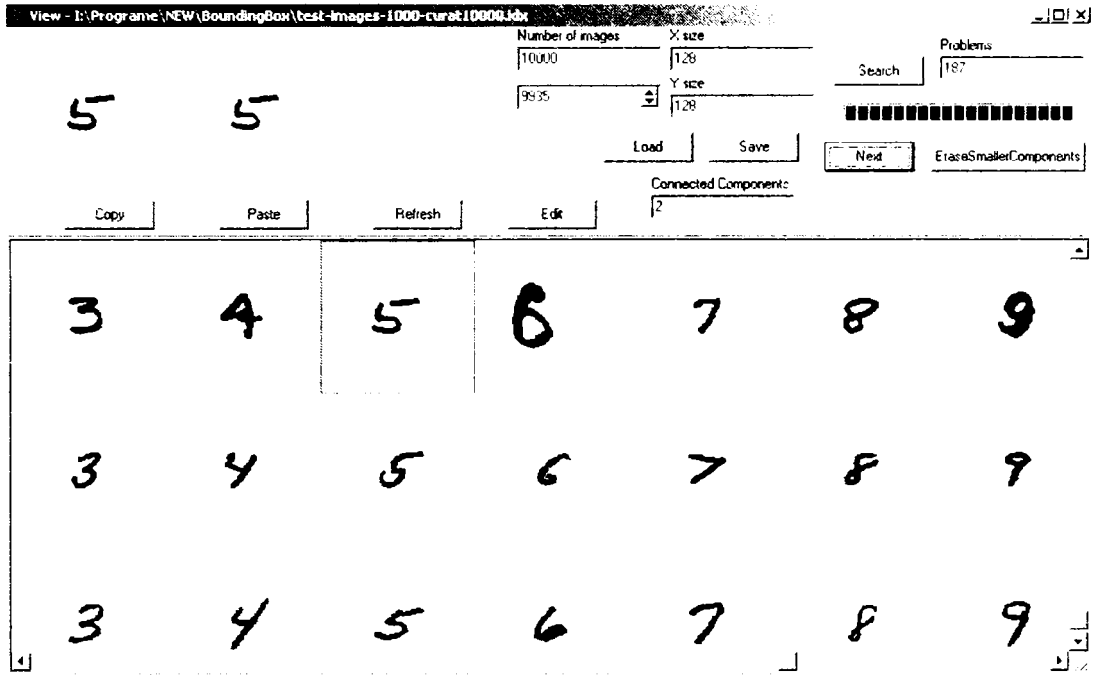


Figura 4.9: Interfața programului de corectare semiautomată.

După prelucrările anterioare, rezultatul este un fișier cu imaginile a 10000 de cifre provenite de la 1000 de persoane diferite. Fiecare rând conține imagini de la o singură persoană. În Figura 4.10 este arătat modul de codificare. Cu aceste zece cifre ale fiecărei persoane construim 200 de imagini cu perechi de cifre parțial suprapuse. Imaginile cu indici pari conțin cifre suprapuse orizontal cu un pixel, iar cele cu indici impari conțin cifre suprapuse orizontal cu 3 pixeli. Am folosit suprapuneri orizontale doar pe 1 și 3 pixeli, deoarece mai puțin de 1 înseamnă doar atingere - caz ușor rezolvat de un clasificator la nivel de cifra care e aplicat pe o fereastră ce se deplasează deasupra imaginii -, iar mai mult de trei (3 reprezentând 15% într-o imagine de  $20 \times 20$  pixeli) e puțin probabil să întâlnim într-un text scris de mână. Cazul suprapunerii pe o distanță de doi pixeli a fost evitat deoarece e prea asemănător cu celelalte două și doar ar fi contribuit la mărirea inutilă cu 50% a setului de antrenare și, implicit, a timpului necesar antrenării.

În direcție verticală ambele imagini de cifre au fost deplasate aleator cu zero sau un pixel, în sus sau în jos. În acest mod s-a dorit simularea scrierii nealiniată orizontal din scrisul de mână. Pixelii obținuți în urma suprapunerii au ca valoare a nivelului de gri minimul dintre 255 și suma celor două nivele de gri ale pixelilor ce au fost suprapuși. Înainte de procesul de suprapunere parțială fiecare imagine de  $20 \times 20$  pixeli a fost centrată într-o imagine de  $28 \times 28$  pixeli pentru a permite deplasările orizontale și verticale. Imaginea rezultată în urma unirii a două cifre are dimensiunea de  $40 \times 28$  pixeli (Figura 4.11). Acestea sunt datele brute care vor fi redimensionate ulterior pentru antrenarea și testarea rețelelor neuronale. Am ales

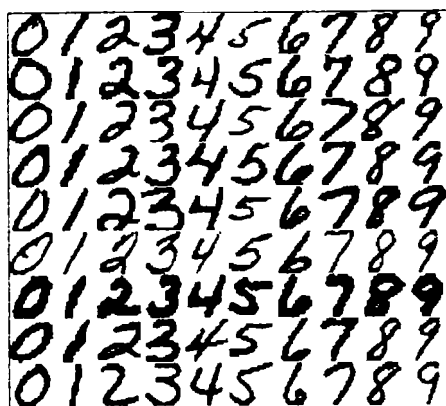


Figura 4.10: Cifrele a nouă persoane după etapele de tăiere după dreptunghiul delimitator, centrare și redimensionare la  $20 \times 20$  pixeli.

lățimea de 40 de pixeli deoarece imaginile originale aveau 20 de pixeli, din acest motiv în urma procesului de suprapunere parțială nu putea rezulta o imagine mai lată de 40 de pixeli. Înălțimea stabilită la 28 de pixeli permite deplasarea pe verticală a imaginilor cu patru pixeli, atât în sus, cât și în jos. Setul de antrenare obținut conține 200000 de imagini, câte 2000 pentru fiecare din cele 100 de clase.

#### 4.4.2 Setul de test

Acest set este compus din persoane începând cu poziția 1010 în partițiile  $hsf_0 \div hsf_3$ . Din motive similare cu cele de la setul de antrenare, formularele următoarelor persoane au fost ignorate: 1102, 1107, 1110, 1209, 1432, 1571, 1572, 1723, 1725, 1726, 1727, 1756, 1767, 1965, 1969, 1974, 1984, 1990, 1992 și 1995. Am aplicat procesul semiautomat de îndepărtare a pixelilor care ar fi afectat nedorit procesul de unire a cifrelor. Pe aceste seturi am îmbunătățit doar aproximativ 10 imagini care au fost prea distorsionate pentru a fi interpretate de un om. Nu am făcut alte îmbunătățiri ale calității imaginilor din acest set deoarece am dorit să se poată observa capacitatea de generalizare a rețelei. Un alt motiv a fost acela că s-a dorit testarea rețelei cu imagini cât mai apropiate de cele întâlnite în realitate. Restul pașilor de procesare sunt identici cu cei de la setul de antrenare:

- tăierea imaginii la limitele dreptunghiului delimitator al cifrei,
- centrarea imaginii,
- redimensionarea la  $20 \times 20$  pixeli,
- deplasări orizontală și verticală,
- unirea a două imagini într-una de  $40 \times 28$  pixeli.

### 4.5 Seturile de test 3DNS, 6DNS și 10DNS

Întrucât lucrarea de față tratează subiectul recunoașterii șirurilor de cifre scrise de mână, sunt necesare seturi de test cu astfel de șiruri. În literatură se folosesc șirurile din NIST SD 19 [4]

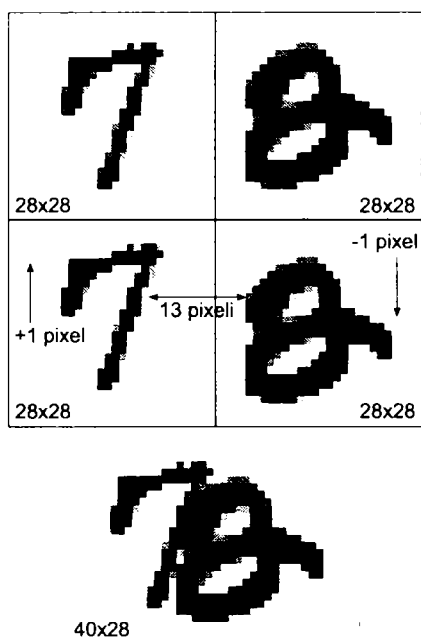


Figura 4.11: Procesul de suprapunere parțială: primul rând - imaginile originale; rândul al doilea - imagini deplasate pe verticală; al treilea rând - imaginea finală după deplasarea pe orizontală cu 13 pixeli (distanța minimă dintre imagini) + 3 (distanța de suprapunere) pixeli.

și din CEDAR [22]. Deoarece am avut la dispoziție doar NIST SD 19, testele au fost efectuate pe această bază de date.

NIST SD 19 conține gata extrase cifrele și literele din cele câteva mii de formulare. Nu conține, însă, șirurile de cifre, astfel că a trebuit scris un program care să le extragă din imaginile cu formulare (Figura 4.2).

Cu metoda descrisă în secțiunea următoare am extras câmpurile de 3, 6 și 10 cifre din formularele 1800-2099. Din cele 300 de pagini cu formulare, autorii din [44] au eliminat patru pagini cu imagini foarte șterse: 1965, 1977, 2027 și 2067. Au fost eliminate și în cadrul acestor experimente, rămânând 296 de imagini de formulare.

#### 4.5.1 Metoda de extragere a câmpurilor din formulare

Studiind formularele de tip HSF (un exemplu este dat în Figura 4.2), se observă că sunt alcătuite din câmpuri dreptunghiulare în care persoana care le completează trebuie să scrie informația tipărită deasupra câmpului respectiv. Se cunoaște că toate formularele arată aproximativ la fel: au același număr de câmpuri de aceeași dimensiune, poziționate în aceleași locuri. Diferă numai conținutul câmpurilor, iar informația cu cifrele sau literele ce ar trebui să se regăsească în acele câmpuri se află în fișierul cu referințe corespunzător formularului. Cunoscând aceste detalii, extragerea imaginilor cu șirurile de cifre pare inițial o operație destul de simplă. Doar că formularele au fost întâi tipărite, apoi completate mai mult sau mai puțin corect și apoi scanate. În urma acestor prelucrări dreptunghiurile nu se mai află în aceleași poziții; liniile au grosimi diferite și pot fi întrerupte; distanțele dintre câmpuri nu

mai sunt egale între ele; toată imaginea poate fi rotită cu un unghi de câteva grade față de orizontală, mai mult, acest unghi poate varia în interiorul documentului; scrisul poate depăși dreptunghiurile ce delimitează câmpurile.

Deoarece câmpurile din formulare sunt dispuse pe rânduri, iar dreptunghiurile aferente sunt aliniate, având aceeași înălțime, un mod de a începe extragerea informației ar putea fi reprezentat de detectarea acestor linii cvasiorizontale. În formular există 20 de astfel de linii (Figura 4.2), delimitând 10 rânduri cu informație. Doar rândurile 1, 2, 3, 4, 5 și 6 conțin șiruri de cifre, restul neprezentând interes. Un mod simplu de a detecta liniile dintr-o imagine este cu ajutorul histogramelor. Dezavantajul acestora este acela că funcționează doar pentru un unghi care trebuie cunoscut a priori. Deoarece imaginile pot fi rotite cu orice unghi -în urma experimentelor am constatat că unghiul de rotire al formularelor este de maxim  $\pm 1.5^\circ$ -, această metodă nu poate fi aplicată. Deși unghiul pare mic, luând în considerare că o linie are lungimea de aproximativ 2000 de pixeli,  $1.5^\circ$  înseamnă o abatere pe verticală de aproximativ  $2000 \cdot \tan 1.5^\circ \approx 52$  de pixeli. Cum grosimea liniei scanate la 300dpi este de aproximativ 8-10 pixeli, devine evident de ce metoda histogramelor orizontale nu funcționează. Soluția ar reprezenta-o mai multe histograme cvasiorizontale cu abaterea față de orizontală de maxim  $\pm 1.5^\circ$ , însă există o altă metodă -transformata Hough- care rezolvă aceste probleme. În secțiunea 3.1.3 au fost prezentate aspecte teoretice generale referitoare la această transformată.

Având în vedere că sunt foarte multe formulare de prelucrat, fiecare având dimensiunile de  $3300 \times 2560$  pixeli, și că transformata Hough este foarte mare consumatoare de resurse de calcul, aplicarea ei pare a fi neoportună. Se știe că transformata Hough va lua fiecare din cei 8448000 de pixeli ai imaginii și va testa să vadă dacă se află pe una din foarte multele drepte posibile. În mod normal parametrii transformatei Hough (a se vedea Secțiunea 3.1.3) ar lua următoarele valori:  $\theta \in [-90^\circ, +90^\circ]$ , iar  $r \in [-3300, +4177]$ . Dacă am alege pașii de discretizare  $dr = 1 \text{ pixel}$  și  $d\theta = 0.2^\circ$  (e necesar un astfel de pas mic pentru  $\theta$  deoarece liniile orizontale ale documentului sunt foarte lungi -de aproximativ 2000 de pixeli- și un pas mai mare ar conduce la abateri prea mari care ar implica eșecul detectării liniilor din imagine), observăm că trebuie să verificăm apartenența fiecăruia din cele 8448000 de puncte la fiecare din cele  $\frac{(3300+4177) \text{ pixeli}}{1 \text{ pixel}} \times \left( \frac{180^\circ}{0.2^\circ} + 1 \right) = 6736777$  drepte. Ajungem la  $8448000 \times 901 = 7611648000$  aplicări ale Formulei 3.7 pentru un singur formular, mult prea mult chiar și pentru procesoarele actuale.

Evident, sunt necesare îmbunătățiri ale algoritmului. Pasul de discretizare  $dr$  nu poate fi mărit deoarece am începe să detectăm puncte din grosimea liniilor, ceea ce ar duce la alte complicații. Mai rămân doi parametri, intervalele de valori pentru  $r$  și pentru  $\theta$ . Am putea căuta linii orizontale numai în anumite zone ale documentului, dar o matrice de acumulare cu indici ce variază necontinuu iarăși complică mult algoritmul. Pe lângă acest motiv, e destul de "periculos" să alegem doar anumite benzi orizontale din imagine, deoarece imaginile pot fi translatate pe verticală destul de mult, chiar și cu câteva sute de pixeli. Rămâne parametrul  $\theta$ , care, din fericire, poate fi limitat la intervalul  $[-1.5^\circ, +1.5^\circ]$ , deoarece cunoaștem ca imaginile formularelor nu pot fi rotite mai mult de atât. Această limitare a domeniului de valori pentru  $\theta$  mărește viteza algoritmului de  $\frac{180/0.2+1}{3/0.2+1} \approx 56$  ori, aducând timpul de procesare la nivelul a câteva secunde pentru un formular (folosind un singur fir de execuție).

O optimizare a fost adusă și parametrului  $r$ , prin limitarea variației sale. Analizând Figura 4.12 se poate observa că pentru variații de  $\pm 1.5^\circ$  ale lui  $\theta$  putem limita variația lui  $r$  doar la domeniul pozitiv dacă nu mai considerăm posibilele drepte ce ar trece numai prin puncte ale

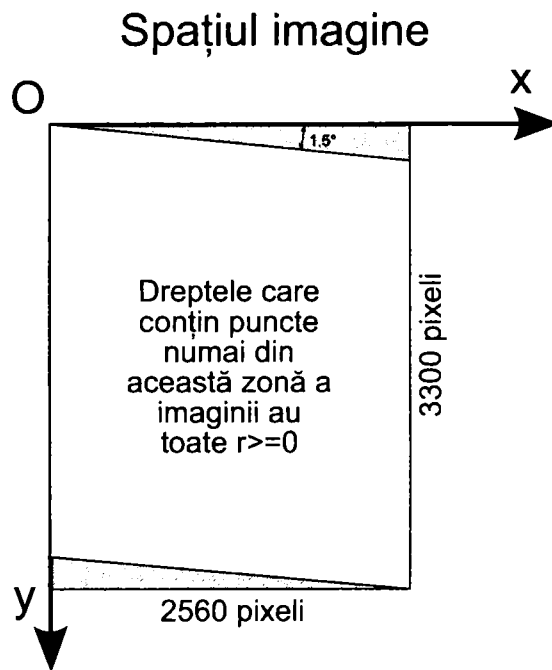


Figura 4.12: Limitarea variației parametrului  $r$ .

zonei hașurate din Figura 4.12. Cum informațiile din formulare se află cu câteva sute de pixeli mai jos, nu afectăm deloc corectitudinea algoritmului. În schimb, va crește viteza de rulare și codul va fi puțin mai simplu din cauză că nu mai există și zone negative pentru indici.

Algoritmul 1 prezintă modul de obținere a matricei de acumulare. Fiecare element al matricei are ca valoare numărul de pixeli negri ai imaginii prin care trece respectiva linie. Folosind informațiile prezentate anterior se poate calcula dimensiunea matricei de acumulare:  $4177 \times 16 = 66832$ . Foarte multe din valorile matricei sunt nenule, astfel că trebuie să găsim o metodă prin care să selectăm doar acele elemente care corespund liniilor căutate. O metodă simplă este aceea de a aplica un filtru care să decidă dacă fiecare valoare a matricei reprezintă o dreaptă utilă sau nu. Arătam mai sus că abaterea maximă pentru o linie de 2000 de pixeli la o înclinație de  $1.5^\circ$  este de 52 pixeli. Putem aplica o fereastră de forma unui pătrat cu latura de  $52 \times 2 + 1 = 105$  pixeli centrată în fiecare pixel al matricei de acumulare. Dacă în această fereastră există cel puțin o valoare mai mare decât cea din centrul ferestrei, valoarea elementului central este setată pe 0. Deși e o metodă foarte primitivă, este eficientă și este suficient de rapidă pentru a nu necesita implementarea unor algoritmi mai complicați. Timpul necesar procesării întregii matricei este de aproximativ o secundă. Pe lângă procesarea decrisă am mai adăugat condiția ca fiecare linie să conțină cel puțin 800 de pixeli. Acest prag a fost obținut studiind suma lungimilor segmentelor coliniare din formulare.

Rezultatul operațiilor de până acum este reprezentat de câteva zeci de linii candidate. Analizându-le am constatat că pe lângă liniile ce se doreau a fi detectate apăreau și câteva linii "fantomă". Acestea erau detectate în textul tipărit al paragrafului din Constituție. Deoarece fontul de acolo era foarte gros și adunat, într-adevăr existau câteva linii care se puteau forma cu pixeli din caracterele tipărite. De obicei se detecta o linie pentru fiecare rând al paragra-

```

Input: Imagine
Output: matricea de acumulare  $H$ 
inițializează toate elementele din  $H$  cu 0
foreach  $0 \leq y < \text{ÎnălțimeaImagine}$  do
  foreach  $0 \leq x < \text{LățimeaImagine}$  do
    if  $\text{Imagine}[y, x] \neq 0$  then
      foreach  $0 \leq t < t_{max}$  do
         $\theta = (0.2t - 1.5)\pi/180$ 
         $r = x \sin \theta + y \cos \theta$ 
        if  $0 \leq r < 4177$  then
          | incrementează  $H[t, r]$ 
        end
      end
    end
  end
end

```

**Algorithm 1:** Transformata Hough.

fului. Transformata Hough doar **numără** punctele care se află pe o linie, nu verifică **continuitatea** ei, neavând nicio informație referitoare la distribuția punctelor pe dreaptă. Din acest motiv gradul de fragmentare al liniilor detectate în paragraful Constituției nu reprezintă un impediment în detectarea lor. Se impune, însă, eliminarea acestor linii false din prisma scopului metodei de extragere a imaginilor cu șiruri de cifre. Pentru aceasta este suficient să măsurăm gradul de fragmentare al liniilor și să le eliminăm pe cele care depășesc un anumit prag maxim admis.

Gradul de fragmentare al unei linii poate fi măsurat calculând numărul de discontinuități al respectivei linii, folosind Algoritmul 2. Algoritmul reface linia, prin căutarea în spațiul imaginii originale a pixelilor de pe dreapta cu parametrii  $r$  și  $\theta$  în spațiul Hough. Sunt numărate trecerile  $0 \rightarrow 1$  și  $1 \rightarrow 0$ , adică trecerile de la pixeli albi la negri și invers. Acest număr e apoi împărțit cu 2 pentru a obține numărul de discontinuități.

```

Input: Imagine, dreapta  $(r, \theta)$ 
Output: numărul de discontinuități ale drepte curente
treceri=0
foreach  $0 \leq x < \text{LățimeaImagine}-1$  do
   $y_1 = r - \frac{x \sin \theta \frac{\pi}{180}}{\cos \theta \frac{\pi}{180}}$ 
   $y_2 = r - \frac{(x+1) \sin \theta \frac{\pi}{180}}{\cos \theta \frac{\pi}{180}}$ 
  if  $\text{Imagine}[y_1, x] \neq \text{Imagine}[y_2, x+1]$  then
    | incrementează treceri
  end
end
return  $\frac{\text{treceri}}{2}$ 

```

**Algorithm 2:** Determinarea numărului de discontinuități ale unei drepte.

După experimentarea algoritmului pe câteva formulare am ajuns la un număr maxim admisibil de întreruperi egal cu 60 pentru liniile candidate la a delimita câmpuri ale formularului. Toate liniile cu mai mult de 60 de discontinuități sunt ignorate, ele apărând în paragraful tiparit al Constituției. În acest moment numărul de linii încă este mai mare decât numărul de linii orizontale care delimitează rândurile cu câmpurile formularului. Liniile detectate sunt mai multe deoarece unele dintre ele sunt multiplu detectate, adică există câteva linii cu exact același număr de puncte și cu parametri foarte apropiați. De obicei diferă numai prin parametrul  $r$ , existând linii multiplicat de 2-5 ori. Dintre aceste linii ne este suficientă una singură. Folosind Algoritmul 3 păstrăm doar liniile care delimitează rânduri cu câmpuri din formular. Dacă există mai multe linii consecutive detectate, atunci o păstrăm pe prima, dacă suntem la începutul unui rând, respectiv doar pe ultima, dacă suntem la sfârșitul unui rând. Fiecare linie conține informații despre  $r$ ,  $\theta$  și numărul de puncte. Variabila *start* are valoarea *true* dacă se caută o linie pentru începutul unui rând, și *false* dacă se caută o linie pentru sfârșitul unui rând.

**Input:**  $L[]$ , vector cu liniile detectate

**Output:**  $D[]$ , vector cu liniile păstrate

*start*=true //început de rând

**foreach** linie  $L[i]$  din  $L[]$  **do**

    află lista de linii  $L[i+j]$  cu proprietatea că  $|L[i].r - L[i+j].r| \leq 10$

**if**  $j=0$  **then**

        | adaugă  $L[i]$  la  $D[]$  //nu există linii alăturate

**end**

**else**

**if** *start* **then**

            | adăugăm  $L[i+j]$  la  $D[]$  //o copiem pe ultima

**end**

**else**

            | adăugăm  $L[i]$  la  $D[]$  //o copiem pe prima

**end**

**end**

*start* = *not*(*start*) //schimbăm starea

**end**

**Algorithm 3:** Selectarea liniilor ce delimitează rânduri cu câmpuri din formular.

În  $D[]$  vom avea 20 de linii care delimitează cele 10 rânduri de câmpuri. Întrucât ne interesează numai câmpurile numerice, păstrăm doar liniile cu indici de la 2 la 13, linii care ar trebui să delimiteze cele cinci rânduri de câmpuri numerice.

Fiecare pereche de linii din cele 10 linii rămase, este trimisă unui modul de extragere a imaginilor câmpurilor numerice. Se verifică dacă numărul de puncte de pe cele două linii nu diferă cu mai mult de 200 de puncte și dacă  $||r_1 - r_2| - 90| > 20$ , deoarece înălțimea medie a unui rând este de 90 de pixeli. Dacă vreunul din aceste teste nu este trecut, procesarea continuă cu următoarea pereche de linii. Urmează delimitarea câmpurilor pe orizontală. Pentru aceasta se reformează cele două linii în spațiul imagine și apoi se suprapun. Avem o singură linie care ar trebui să prezinte discontinuități în spațiile dintre câmpuri. Câmpurile din formulare sunt separate pe orizontală de cel puțin 55 de pixeli albi. Din acest motiv am ales ca prag 40 de pixeli pentru discontinuitatea maximă permisă. Linia anterioară este segmen-

tată în fiecare poziție unde începe o discontinuitate mai mare de 40 de pixeli. Bineînțeles, se tratează și cazurile particulare de început și de sfârșit de linie. Primul rând cu câmpuri numerice conține 3 câmpuri, iar următoarele 5 conțin câte 5 câmpuri. Numărul de câmpuri detectate este comparat cu aceste constante pentru a se detecta eventuale erori de segmentare. Dacă valorile se potrivesc, atunci se apelează un alt modul care extrage din formular imaginea rectangulară cu dimensiunile și poziția detectate anterior. Se construiește și un nume sugestiv pentru imaginea câmpului, de forma *fxxx\_yy\_fldzz\_eticheta.bmp*, unde *zz* reprezintă numărul câmpului, iar *eticheta* - eticheta respectivului câmp. Eticheta este citită din fișierul cu etichete, aferent fiecărui formular. Imaginea încă este în formă brută, conținând în unele cazuri părți din dreptunghiul ce delimita câmpul. Primul rând din Figura 4.13 prezintă câteva astfel de imagini înainte de procesul de curățare.

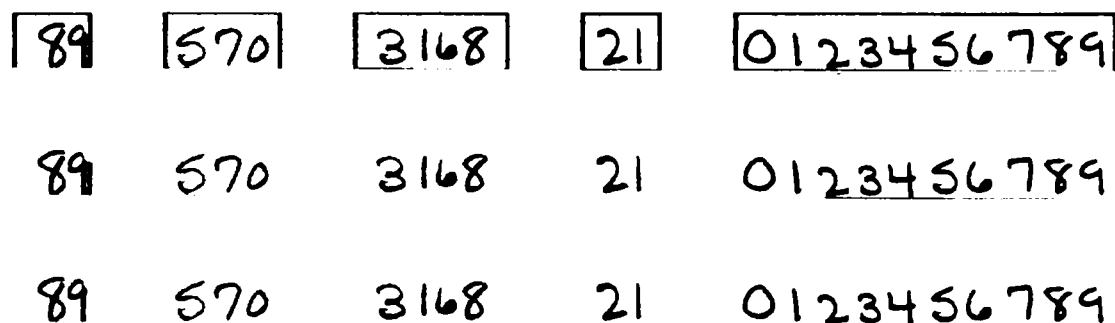


Figura 4.13: Imagini obținute în urma procesului de extragere a câmpurilor din formulare. Primul rând - imagini brute, al doilea rând - imagini după faza de curățare automată, rândul al treilea - imagini după faza de curățare manuală.

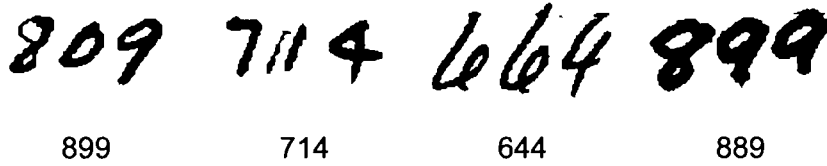
Grosimea liniilor care formează dreptunghiurile de delimitare a câmpurilor este de aproximativ 10 pixeli. Am luat o marjă de 50% și am considerat un dreptunghi interior imaginii, depărtat la 15 pixeli de extremități. Ideea constă în ștergerea condiționată a pixelilor din acest cadru dreptunghiular. Se ia fiecare pixel din această cadru. Dacă este negru, atunci se caută componenta conexă din care face parte. Dacă această componentă conexă nu conține pixeli din afara cadrului dreptunghiular, atunci pixelul împreună cu toată componenta conexă pot fi șterși. Determinarea componentei conexe se face în mod uzual prin căutare în opt direcții, dar condiționat de a nu se depăși zona cadrului dreptunghiular. Această metodă simplă rezolvă peste 80% din cazurile de imagini brute cu resturi ale dreptunghiurilor ce delimitează câmpurile (Figura 4.13, rândul din mijloc). Restul de imagini sunt curățate manual deoarece textul din interior atinge sau suprascris dreptunghiul delimitator (Figura 4.13, ultimul rând).

#### 4.5.2 Setul de test 3DNS

Setul de test 3DNS este folosit pentru a efectua o comparație exactă între metoda propusă în această lucrare și cele prezentate în [44]. Din acest motiv am extras exact aceleași imagini ca în [44]. Formularele conțin câte 5 câmpuri de trei cifre, în total existând  $296 \times 5 = 1480$  imagini de numere formate din trei cifre. În [44] se afirmă că au fost eliminate patru imagini care nu corespundeau etichetelor, dar nu se specifică care sunt aceste imagini. Analizând vizual imaginile detectate eronat de clasificator am identificat patru imagini care, într-adevar, nu corespundeau cu etichetele lor (Figura 4.14). Imaginile greșit scrise sunt: formularul 1902



- câmpul 7, formularul 2035 - câmpul 28, formularul 2056 - câmpul 11, formularul 2076 - câmpul 11. Deci, la final, setul de test 3DNS conține 1476 de imagini cu numere de trei cifre. Toate aceste imagini au fost tăiate la dimensiunea dreptunghiului delimitator.

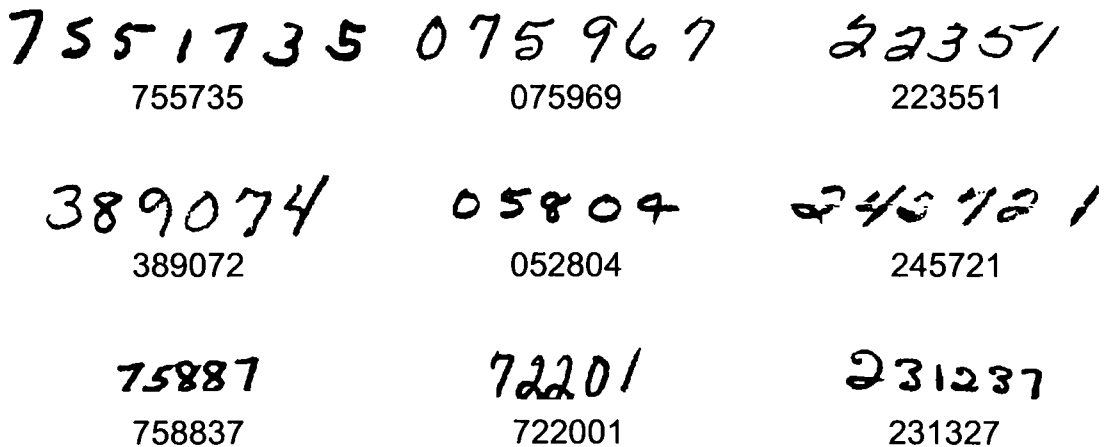


809                      714                      644                      889

Figura 4.14: Imaginile cu numere de trei cifre care au fost scrise eronat de subiecții din NIST SD 19.

### 4.5.3 Setul de test 6DNS

Și setul de test 6DNS este folosit pentru a compara metoda propusă cu cele din [44], astfel că l-am compus din aceleași imagini. Mai exact, din cele 296 de formulare (1800-2099, minus cele patru scrise foarte șters) am extras toate câmpurile de 6 cifre. Fiind câte cinci câmpuri de șase cifre pe formular am obținut 1480 de imagini de numere a șase cifre. Ca în [44], am eliminat imaginile ce nu conțineau textul scris deasupra câmpului. Am identificat opt exemple clare de astfel de erori, al nouălea -rândul doi, coloana trei- a fost ales deoarece autorii din [44] afirmă că există nouă astfel de imagini eronate (Figura 4.15). Am ales cea mai defectuos scrisă imagine, deoarece în [44] nu se specifică indicii imaginilor. Întrucât detectarea acestor imagini este laborioasă și plictisitoare, pentru referințe următoare, voi menționa indicii imaginilor eronate: formularul 1910 câmpul 26, formularul 1913 câmpul 26, formularul 1937 câmpul 18, formularul 2003 câmpul 22, formularul 2035 câmpul 22, formularul 2043 câmpul 14, formularul 2070 câmpul 18, formularul 2080 câmpul 22 și formularul 2094 câmpul 18.



7551735                      075967                      22351  
755735                      075969                      223551

389074                      05804                      245721  
389072                      052804                      245721

75887                      72201                      231237  
758837                      722001                      231327

Figura 4.15: Imagini cu numere de șase cifre, eronat scrise de subiecții NIST SD 19.

La final setul de test 6DNS este compus din 1471 de imagini cu numere formate din șase cifre. Toate imaginile au fost tăiate la dimensiunea dreptunghiului delimitator.

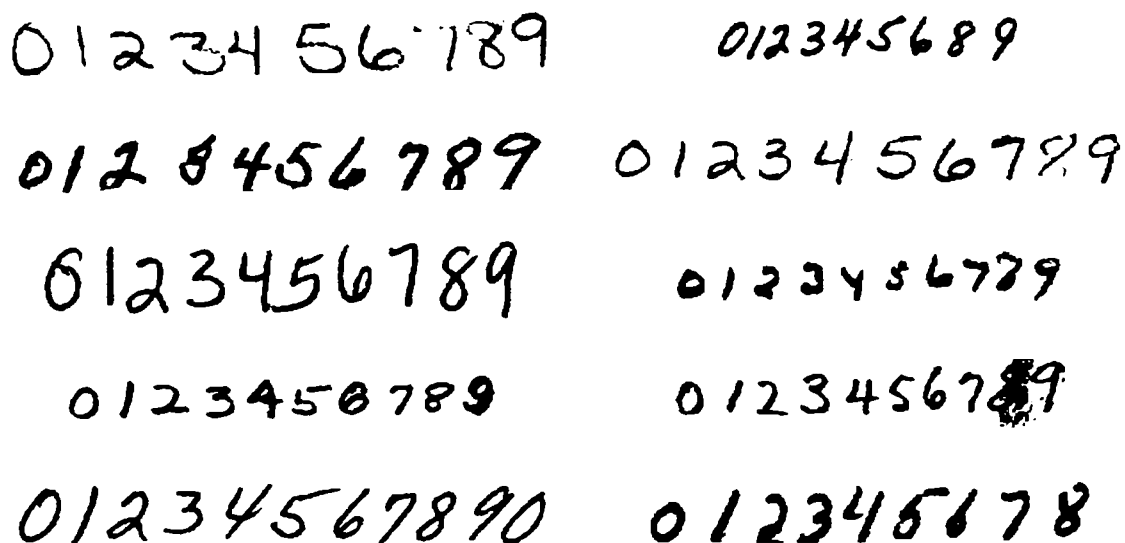


Figura 4.16: Imagini cu numere de zece cifre, eronat scrise de subiecții NIST SD 19 (toate ar trebui să conțină șirul 0123456789).

#### 4.5.4 Setul de test 10DNS

Din cele 296 de formulare am extras  $296 \times 3 = 888$  de imagini cu numere formate din 10 cifre, fiecare formular conținând trei câmpuri a zece cifre. Din cele 888 de imagini, zece nu corespund cu informația scrisă deasupra respectivului câmp și au fost eliminate (Figura 4.16). Imaginile eronate au următorii indici: formularul 1858 câmpul 4, formularul 1896 câmpul 4, formularul 1909 câmpul 3, formularul 1957 câmpul 4, formularul 1975 câmpul 4, formularul 2021 câmpul 3, formularul 2023 câmpul 3, formularul 2038 câmpul 3, formularul 2039 câmpul 3 și formularul 2049 câmpul 5.

După excluderea imaginilor eronate au rămas 878 de imagini în setul de test 10DNS. La fel ca pentru seturile 3DNS și 6DNS, imaginile au fost tăiate la dimensiunea dreptunghiului delimitator.

## 4.6 Concluzii

Acest capitol a avut drept scop crearea de seturi de imagini pentru antrenarea și testarea clasificatorilor. De asemenea, au fost prezentate bazele de date NIST SD 19 și MNIST.

Pentru antrenarea clasificatorului la nivel de cifră au fost create seturile din LNIST (66214 imagini pentru setul de antrenare și 45398 de imagini pentru setul de test) prin extragerea de imagini din NIST SD 19. Acestea sunt construite identic cu cele folosite de autorii din [44, 45] pentru a permite o comparație cât mai riguroasă a rezultatelor.

Pentru antrenarea clasificatorului la nivel de două cifre parțial suprapuse (2RNC) au fost create seturile de antrenare și de test TDNS. Din NIST SD 19 au fost extrase 20000 de imagini de cifre (10000 pentru antrenare și 10000 pentru test) din care s-au creat câte 20000 de imagini cu perechi de cifre parțial suprapuse atât pentru setul de antrenare, cât și pentru setul de test. Pentru formarea imaginilor cu perechi de cifre parțial suprapuse a fost dezvoltat

un algoritm care unește cifrele în funcție de morfologia acestora.

Dacă imaginile setului LNIST au fost extrase pentru a fi folosite de către clasificatorul la nivel de cifră, pentru cazul general, cel al șirurilor numerice formate din mai multe cifre, au fost construite seturile de test 3DNS, 6DNS și 10DNS. Și în construirea acestor seturi s-a urmărit obținerea unor seturi cât mai apropiate de cele folosite de alți cercetători pentru evaluarea performanțelor metodelor de recunoaștere. Cele trei seturi sunt identice cu cele folosite în domeniu, în limita informațiilor furnizate de autorii care le-au utilizat anterior. Există posibilitatea ca una, maxim două imagini din cele aproape 1500 să fie diferite.

Pentru extragerea din NIST SD 19 a imaginilor necesare construirii seturilor 3DNS, 6DNS și 10DNS a fost elaborat un algoritm care folosește o variantă a transformatei Hough pentru detectarea liniilor din formularele NIST. După detecția acestor linii, sunt detectate dreptunghiurile ce încadrează câmpurile dorite. Apoi are loc extragerea propriu-zisă a imaginilor și se aplică un algoritm de înlăturare automată a posibilelor părți din dreptunghiul delimitator care au fost extrase în același timp cu imaginile șirurilor numerice. La final a fost aplicat un proces de inspectare și curățare manuală a imaginilor. Algoritmul de extragere automată a imaginilor din câmpurile formularelor a diminuat considerabil timpul necesar construirii pe cale manuală a celor trei seturi de test. De asemenea, algoritmul este destul de general, putând fi adaptat facil pentru alte tipuri de formulare.

# Capitolul 5

## Experimente

În acest capitol vor fi evaluate performanțele rețelelor neuronale de tip perceptron și de convoluție prezentate în Capitolul 3. Întâi, rețelele vor fi antrenate și apoi vor fi evaluate folosind seturile de date din Capitolul 4. Se va începe cu evaluarea rețelelor pentru recunoașterea cifrelor individuale, iar apoi se vor antrena rețele pentru recunoașterea de numere formate din două cifre parțial suprapuse. După ce ambii clasificatori - pentru cifre singulare și pentru numere formate din două cifre parțial suprapuse - vor fi antrenați și testați, se va implementa o metodă pentru recunoașterea numerelor formate din 3 sau mai multe cifre. Această metodă va fi testată cu seturile 3DNS, 6DNS și 10DNS construite în Secțiunile 4.5.2, 4.5.3 și 4.6.

Rețelele neuronale și toate programele adiționale pentru procesarea imaginilor au fost implementate în C++. Am folosit și cod open source [3] pe care l-am modificat pentru a putea fi folosit la experimentele pe care le-am efectuat. Pentru toți neuronii rețelei am folosit ca funcție de activare tangenta hiperbolică  $\frac{e^{2x}-1}{e^{2x}+1}$ . Am obținut o creștere semnificativă a vitezei de antrenare prin paralelizarea algoritmului pentru a folosi suportul nativ existent în procesoarele multinucleu existente. Am testat algoritmi pe procesoare cu două și cu patru nuclee.

### 5.1 Experimente cu cifre individuale

Aceste experimente au o dublă motivație. Prima este dată de evaluarea corectitudinii implementării algoritmilor prin compararea rezultatelor cu cele existente deja, iar a doua de obținerea unui clasificator de calitate (care are rată mare de recunoaștere) la nivel de cifre singulare. Acest clasificator va fi utilizat apoi în sistemul de recunoaștere (a se vedea Figura 5.4) care recunoaște șiruri de una, două sau mai multe cifre distincte sau parțial suprapuse.

Ca seturi de antrenare au fost folosite seturile MNIST și LNIST. La prima vedere, poate părea inutilă folosirea a două seturi de antrenare. Scopul acestei lucrări este de a propune o nouă metodă de recunoaștere a șirurilor scrise de mână. Ca orice metodă nouă, trebuie validată prin efectuarea de experimente ale căror rezultate trebuie comparate cu cele existente deja în literatura de specialitate. Arătăm în introducere că articolul [44] prezintă cele mai bune rezultate obținute până în acest moment. Autorii din [44] folosesc ca set de antrenare pentru clasificatorul la nivel de cifră setul LNIST, astfel că pentru o comparație corectă a metodelor am antrenat și clasificatorul din metoda proprie cu acest set. În plus, setul LNIST

are proprietatea că este disjunct cu seturile de test 3DNS, 6DNS și 10DNS. Setul MNIST a fost folosit pentru a verifica corectitudinea implementării algoritmilor de antrenare a rețelelor neuronale prin compararea rezultatelor cu cele din [40]. Din păcate, setul de antrenare din MNIST conține imagini cu cifre din toate partițiile  $hsf_{\{0-4\}}$ , deci nu este disjunct cu seturile de test 3DNS, 6DNS și 10DNS. Din acest motiv nu poate fi folosit la experimentele ce implică seturile 3DNS, 6DNS și 10DNS.

### 5.1.1 Rețele neuronale de tip perceptron

Întrucât perceptronul fără straturi ascunse nu poate fi folosit la clasificarea elementelor mulțimilor în clase neliniar separabile, am folosit perceptronul cu un strat ascuns. Am variat numărul de neuroni de pe stratul ascuns între 60 și 220 (Tabelul 5.1). Antrenarea a fost efectuată cu setul de antrenare LNIST (Secțiunea 4.3), iar testarea cu cele 10000 de imagini din setul de test MNIST (Secțiunea 4.2). Puteam folosi și setul de test LNIST, dar, după cum se va vedea în secțiunea următoare, rezultatele sunt foarte asemănătoare. Nu am mărit numărul de neuroni de pe stratul ascuns mai mult de 220 pentru că timpul ajunsese la aproximativ 600 de secunde pentru o epocă, iar eroarea nu scădea decât foarte încet.

Tabelul 5.1: Rezultate obținute cu perceptronul cu un strat ascuns.

# neuroni pe stratul ascuns	eroarea pe setul de antrenare LNIST(%)	eroarea pe setul de test MNIST(%)
60	2.08	4.83
80	1.59	4.02
100	1.39	3.68
120	1.31	3.54
140	1.40	3.58
160	1.26	3.70
180	1.32	4.11
200	1.23	3.63
220	1.26	3.46

Toate rețelele au fost antrenate fără activarea distorsiunilor pe setul de antrenare. Am variat rata de învățare,  $\eta$ , de la 0.0001 până la 0.000005, înmulțind-o la fiecare epocă cu factorul 0.85. După atingerea valorii minime de 0.000005,  $\eta$  a rămas constant. Când numărul de imagini eronat clasificate din setul de antrenare s-a stabilizat, am oprit antrenarea.

Pentru a verifica dacă aplicarea distorsiunilor ar influența mult rata de recunoaștere pentru setul de test, am antrenat cea mai mare rețea -cea cu 220 de neuroni pe stratul ascuns- și cu activarea distorsiunilor pe setul de antrenare. Eroare a scăzut până la 2.88%, dar această eroare este prea mare față de cele mai bune rezultate raportate în [40]. În plus, a crescut timpul de antrenare cu aproximativ 100 de secunde pe epocă.

### 5.1.2 Rețele neuronale cu extractori de trăsături

Cum perceptronul nu a oferit rate de recunoaștere suficient de bune, am trecut la experimentarea unor rețele de convoluție despre care există date [64, 40] că oferă rate de recunoaștere mult mai bune. Detalii despre arhitectura rețelelor de convoluție se află în Capitolul

3. Am variat cei trei parametri ai rețelei în jurul valorilor despre care Simard scria în [64] că au oferit cea mai bună rată de recunoaștere. Am refăcut și experimentul lui Simard pe setul MNIST și am obținut eroarea minimă pentru setul de test de 0.66%. Simard obținuse 0.4% ceea ce este foarte apropiat de 0.66%. Dacă rețeaua ar fi fost lăsată să se antreneze mai mult, probabil ajungeam la eroarea lui Simard. Însă, scopul era doar verificarea corectitudinii algoritmilor de învățare, nu obținerea unui minim cât mai mic, deoarece, oricum, nu puteam folosi rețeaua antrenată cu setul MNIST, acesta nefiind disjunct cu seturile de test multicifre ce urmau a fi folosite la clasificatorul general.

Am antrenat rețelele cu setul de antrenare LNIST, iar rezultatele se află în Tabelul 5.2. Am testat fiecare rețea și cu setul de test din MNIST pentru a vedea dacă rezultatele sunt asemănătoare cu cele obținute pe setul de test LNIST. Rezultatele testării MNIST sunt puțin mai bune -cu maxim 0.22%-, ceea ce se explică prin faptul ca setul de test MNIST are imagini și din setul de antrenare LNIST. Cel mai bun rezultat este obținut cu o rețea cu 6 hărți de convoluție pe stratul  $L_1$ , 50 de hărți de convoluție pe stratul  $L_2$  și 60 de neuroni pe stratul  $L_3$ .

De această dată am urmat recomandările lui Simard [64] și am antrenat rețelele aplicând distorsiuni asupra setului de antrenare. Valorile maxime ale acestora sunt: 15% pentru factorul de scară, 15% pentru factorul de rotație, 8 pentru  $\sigma$  pentru distorsiuni elastice, și 0.5 pentru factorul de scară pentru distorsiuni elastice. Distorsiunile au fost aplicate fiecărei imagini înainte de a fi trimisă modulului de antrenare. Rata de învățare,  $\eta$ , a pornit de la 0.0001 și am scăzut-o la fiecare epocă înmulțind-o cu factorul 0.85. Când numărul de imagini eronat clasificate din setul de antrenare s-a stabilizat, am dezactivat distorsiunile setului de antrenare și am mai antrenat rețeaua încă 10 epoci cu  $\eta=0.000005$ . Acest pas final a scăzut mult eroarea pe setul de antrenare deoarece acesta nu a mai fost distorsionat, dar a scăzut-o puțin și pe setul de test deoarece a atenuat efectul distorsiunilor care erau prea mari. Alegând o valoare atât de mică pentru  $\eta$ , nu a existat pericolul ca rețeaua să-și piardă capacitatea de generalizare.

Tabelul 5.2: Rezultate obținute cu rețele neuronale cu hărți de convoluție.

$L_1$ (hărți)	$L_2$ (hărți)	$L_3$ (neuroni)	eroare antrenare LNIST (%)	eroare test LNIST (%)	eroare test MNIST (%)
4	20	40	1.37	1.21	1.04
4	40	80	1.18	1.11	0.89
6	50	60	0.63	0.88	0.72
6	50	100	0.80	<b>0.86</b>	0.70
6	60	80	0.54	0.99	0.77
8	60	100	0.92	0.89	0.80
8	60	120	0.96	0.98	0.74

Timpul maxim pentru o epocă de antrenare a fost de 620 de secunde, iar cel minim de 210 secunde. Acești timpi au fost obținuți pe procesoare Core2 Duo la 2.4Ghz. Folosind procesoare Core2 Quad similare ca frecvență, dar cu 4 nuclee, timpii de antrenare s-au redus la 60-70%. Rețeaua neuronală ce a oferit cea mai bună rată de recunoaștere va fi folosită drept clasificator la nivel de cifră în sistemul de clasificare general (Figura 5.4).

## 5.2 Experimente cu șiruri compuse din două cifre parțial suprapuse

În această secțiune vor fi prezentate rezultatele antrenării unui clasificator pentru șiruri de două cifre parțial suprapuse. Acest clasificator va fi utilizat apoi în sistemul de recunoaștere (a se vedea Figura 5.4) care recunoaște șiruri de una, două sau mai multe cifre distincte sau parțial suprapuse.

Toate experimentele au fost efectuate pe imaginile din TDNS micșorate de la dimensiunea originală de  $40 \times 28$  pixeli la  $21 \times 13$  pixeli. Utilizarea imaginilor de aceste dimensiuni pentru imaginile de intrare a permis antrenarea rețelelor într-un timp rezonabil, de maxim 24 de ore pe un calculator cu procesor cu două nuclee rulând la 2.4GHz. Imaginile sunt mici, iar multe dintre ele sunt dificil de recunoscut chiar de om. Câteva exemple de imagini provenite de la două persoane diferite sunt prezentate în Figura 5.1. Dimensiunile exacte ale imaginilor din seturile de antrenare și de test au fost determinate de constrângerile impuse la proiectarea rețelei de convoluție. Redimensionarea de la  $40 \times 28$  pixeli la  $21 \times 13$  pixeli a fost efectuată cu păstrarea aspectului imaginilor. Timpul necesar antrenării unei rețele a variat între o oră pentru cea mai simplă rețea (număr redus de ponderi) și câteva zile pentru rețelele complexe care aveau un număr mare de ponderi. Am folosit 15 calculatoare cu procesoare cu două nuclee pentru mai mult de o săptămână pentru a antrena rețelele. Au existat câteva cazuri de rețele pentru care antrenarea a trebuit reluată deoarece s-a împotmolit într-un minim local.



Figura 5.1: Exemple de imagini din fișierul de antrenare.

### 5.2.1 Rețele neuronale de tip perceptron

Acest tip de rețea neuronală foarte generală a fost ales pentru a stabili o bază de comparație pentru arhitecturile specializate prezentate în secțiunile următoare. Tipul de perceptron ales este unul cu trei straturi, din care unul este ascuns, celelalte două reprezentând straturile de intrare, respectiv de ieșire. Am variat doar numărul de neuroni din stratul intermediar, parametrii celorlalte straturi fiind impuși de dimensiunile imaginilor, respectiv de numărul de clase. Rezultatele din Tabelul 5.3 arată că eroarea descrește cu creșterea numărului de neuroni de pe stratul ascuns. Erorile obținute pe seturile de test și de antrenare sunt foarte asemănătoare, diferența dintre ele fiind de maxim 0.37%. Având în vedere că seturile de antrenare și de test conțin imagini de la persoane distincte, această diferență mică între erorile

pe setul de test și pe cel de antrenare demonstrează capacitatea deosebită de generalizare și de adaptare la nou a rețelei.

Timpul de antrenare necesar fiecărei epoci de antrenare crește cu numărul de neuroni din stratul ascuns, de la 200 de secunde pentru 60 de neuroni până la 710 secunde pentru 220 de neuroni. Eroarea a scăzut lent, astfel că am decis oprirea experimentelor la 220 de neuroni pe stratul ascuns deoarece timpul necesar antrenării creștea prea rapid.

Tabelul 5.3: Rezultate obținute cu perceptronul cu un strat ascuns.

# neuroni pe stratul ascuns	eroarea pe setul de antrenare (%)	eroarea pe setul de test (%)	timpul necesar unei epoci de antrenare (s)
60	26.85	26.74	200
80	24.92	24.65	264
100	23.45	23.47	327
120	22.28	22.21	391
140	22.14	22.29	454
160	21.33	21.54	518
180	21.17	21.35	581
200	20.94	21.31	645
220	20.55	20.74	710

## 5.2.2 Rețele neuronale cu extractorii de trăsături

În [38] Yann LeCun propune o arhitectură sofisticată de rețea neuronală numită LeNet-5. Rețeaua s-a dovedit a fi foarte bună pentru recunoașterea cifrelor singulare, chiar și în prezența zgomotului. Șirurile de cifre distincte pot fi recunoscute facil folosind un astfel de clasificator aplicat pe o fereastră ce se deplasează în lungul imaginii. Patrice Simard ([64]) a folosit rețele neuronale de convoluție pentru recunoașterea cifrelor singulare. Clasificatorul utilizat de mine (Figura 5.2) este derivat din cel al lui Simard, dar este modificat extensiv deoarece am dorit să poată fi folosit pentru recunoașterea perechilor de cifre parțial suprapuse.

În Figura 5.2 stratul de intrare este alcătuit din  $21 \times 13$  neuroni, câte unul pentru fiecare pixel al imaginii. Nivelul de gri ( $0 \div 255$ ) al fiecărui pixel din imaginea de intrare este convertit proporțional la un număr din intervalul  $[-1, 1]$ . Această valoare este transferată neuronului corespunzător de pe stratul de intrare. La fel ca în [64] am conectat stratul de intrare la un strat de convoluție. Acest strat de convoluție este alcătuit din mai multe hărți de convoluție (HC). În experimente am variat acest număr între 2 și 12. Dimensiunea nucleului de convoluție a fost păstrată la  $5 \times 5$  pixeli. Acest nucleu de convoluție a fost aplicat tot din doi în doi pixeli, mai exact, numai pentru pixelii care aveau ambele coordonate pare. Aplicarea nucleului de convoluție în acest mod a generat o reducere a dimensiunilor hărților de convoluție de  $(n - 3)/2$ , unde  $n$  reprezintă înălțimea, respectiv lățimea stratului  $L_0$ . În concluzie, dimensiunea hărților de convoluție din stratul  $L_1$  este de  $9 \times 5$  neuroni. Fiecare din cei 45 de neuroni ai stratului  $L_1$  este conectat cu 25 de neuroni din nucleul de convoluție corespunzător. Pe lângă aceste 25 de conexiuni mai există una pentru bias. Deși există  $9 \times 5 \times (5 \times 5 + 1) = 1170$  conexiuni pentru fiecare hartă, numărul de ponderi ce trebuie antrenate este mult redus,



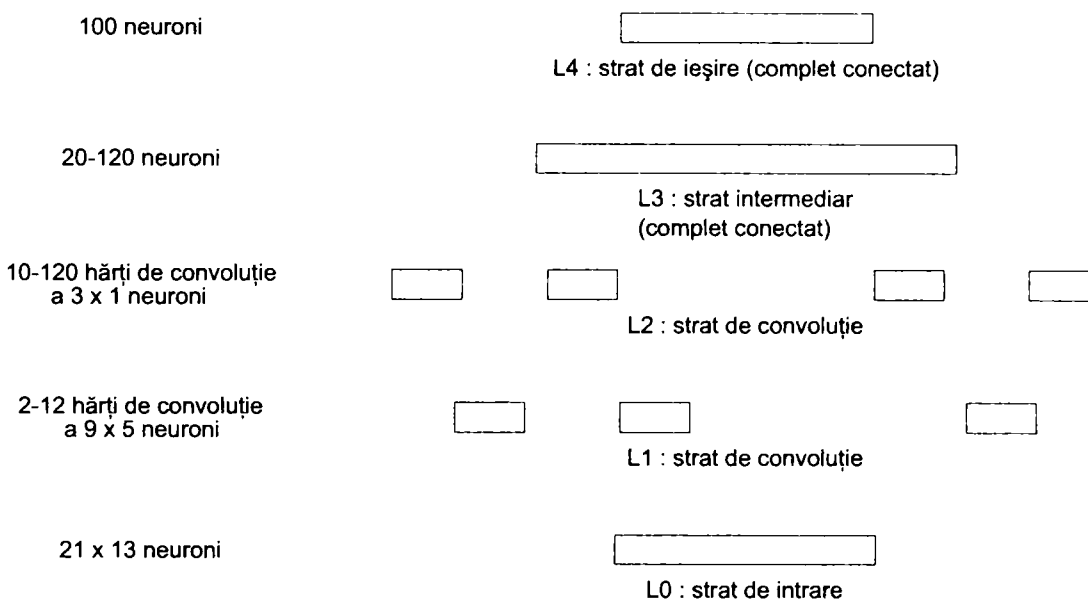


Figura 5.2: Arhitectura rețelei neuronale cu hărți de convoluție.

$5 \times 5 + 1 = 26$ , acestea fiind partajate de toți neuronii din stratul  $L_1$ . Deoarece antrenarea unei rețele neuronale constă în obținerea unor valori apropiate de optim pentru ponderile rețelei, RNC pot fi antrenate mult mai rapid decât rețelele de tip perceptron. Conexiunile și ponderile sunt diferite pentru fiecare hartă de convoluție. În total sunt 11700 conexiuni și 260 de ponderi pentru  $L_1$ , în cazul unei rețele cu 10 hărți pe stratul respectiv.

Pentru stratul  $L_2$  am experimentat cu un număr de 10 până la 120 de hărți de convoluție. Având în vedere că s-a folosit aceeași regulă de aplicare a nucleului de convoluție - doar pentru pozițiile pare - și că dimensiunea nucleului a fost menținută la  $5 \times 5$ , dimensiunea hărților de pe  $L_2$  este de  $3 \times 1$  neuroni. Fiecare neuron din acest strat este conectat cu  $5 \times 5$  neuroni din fiecare hartă a stratului anterior (Figura 3.5) plus o conexiune pentru bias. Deci, o hartă de convoluție din stratul  $L_2$  conține  $3 \times 1 = 3$  neuroni,  $3 \times 1 \times (5 \times 5 \times 10 + 1) = 753$  conexiuni spre neuroni din  $L_1$  (am considerat că  $L_1$  are 10 hărți de convoluție), și  $5 \times 5 \times 10 + 1 = 251$  ponderi. În total, pentru 100 de hărți de convoluție pe stratul  $L_2$ , sunt  $100 \times 3 = 300$  de neuroni,  $753 \times 100 = 75300$  conexiuni și  $251 \times 100 = 25100$  ponderi.

După două nivele de convoluție reducerea este de  $[(n - 3)/2 - 3]/2 = p \Rightarrow n = 4p + 9$ , unde  $p \in 1, 2, 3, 4, 5, \dots$ . În formula anterioară  $n$  reprezintă dimensiunea inițială a imaginii, iar  $p$  - dimensiunea unei hărți de convoluție după două nivele de convoluție. În concluzie,  $n$  ia valori din mulțimea  $13, 17, 21, 25, 29, \dots$ . Pentru a menține dimensiunea imaginilor de intrare la un nivel cât mai mic, dimensiunea stratului de intrare este aleasă ca fiind  $21 \times 13$  neuroni. S-a ținut cont și de aspectul imaginilor cu două cifre concatenate (setul TDNS), acesta fiind  $40/28 \approx 1.43$ . Alegând 13 pentru înălțimea stratului de intrare, lățimea trebuie să fie cel mai mic număr întreg mai mare decât  $40/28 \times 13 = 18.57$ . Valoarea din lista anterioară care corespunde acestei cerințe este 21. Stratul de intrare s-a dorit a fi cât mai mic pentru a obține un timp de antrenare rezonabil pentru rețeaua neuronală, nu mai mare de câteva zile.

Stratul  $L_3$  este total conectat, iar dimensiunea lui variază de la 40 la 120 de neuroni. Ultimul strat,  $L_4$ , are 100 de neuroni, câte unul pentru fiecare clasă, și este, de asemenea,

total conectat cu stratul anterior. Ieșirea fiecărui neuron ia valori în intervalul  $[-1, 1]$ . Pentru o imagine validă, neuronul care reprezintă clasa imaginii va avea o valoare apropiată de 1, iar ceilalți 99 de neuroni vor avea valori apropiate de  $-1$ . Dacă imaginea este ambiguă și seamănă cu mai mult de o clasă, ieșirile neuronilor corespunzător claselor cu care seamănă imaginea vor avea valori corelate cu gradul de asemănare.

Tabelul 5.4: Rezultate obținute cu rețele neuronale cu hărți de convoluție.

$L_1$ (hărți)	$L_2$ (hărți)	$L_3$ (neuroni)	fără deformări		cu deformări		
			Eroare antrenare (%)	Eroare test (%)	Eroare antrenare (cu def.) (%)	Eroare antrenare (%)	Eroare test (%)
2	10	40	35.00	35.43	44.48	39.79	39.88
		60	34.76	34.61	39.24	34.25	34.21
		80	26.25	26.03	37.71	33.42	32.83
	30	40	21.64	21.51	35.82	31.69	31.21
		60	19.98	19.96	32.19	29.95	29.20
		80	18.26	18.15	30.54	25.77	25.31
	50	40	17.80	18.04	28.95	23.62	22.91
		60	15.53	15.68	28.65	22.96	22.47
		80	15.30	15.30	24.67	18.65	18.16
4	10	40	25.63	25.32	30.86	25.76	25.56
		60	22.27	22.08	30.16	27.00	26.60
		80	18.25	18.44	26.08	21.74	21.20
	30	40	15.19	15.25	26.66	23.20	23.06
		60	12.67	13.26	18.88	15.39	15.23
		80	11.80	12.02	18.07	14.29	14.03
	50	40	11.30	12.34	18.78	16.55	16.57
		60	10.67	11.50	15.42	13.04	12.94
		80	10.33	10.64	15.20	11.86	11.72
6	10	40	15.57	15.63	26.35	24.15	23.32
		60	15.50	15.94	23.30	20.40	19.62
		80	15.45	15.67	21.31	17.49	17.34
	30	40	10.72	11.40	25.20	23.07	23.40
		60	10.15	10.81	16.74	13.44	13.57
		80	9.44	10.33	13.95	11.32	11.46
	50	40	9.15	9.83	13.94	11.50	11.82
		60	7.52	8.48	13.76	10.79	11.15
		80	6.73	7.83	12.10	9.63	9.91
8	10	40	15.36	15.78	20.08	17.90	17.86
		60	14.21	14.78	19.91	16.66	16.78
		80	12.97	13.49	17.08	14.38	14.42
	30	40	9.49	10.26	15.78	12.95	13.25
		60	8.20	9.10	13.62	10.01	10.42
		80	7.75	8.72	12.06	9.33	9.74
	50	40	8.06	9.22	15.25	13.16	13.33
		60	7.22	8.31	11.23	8.82	9.51
		80	6.52	7.86	10.70	8.22	8.71
10	60	100	5.63	6.85	10.71	7.89	8.44
	100	100	3.27	5.35	10.15	7.10	7.60
12	120	120	5.23	6.54	-	-	-

Pentru 2, 4, 6 și 8 hărți de convoluție pe stratul  $L_1$ ; 10, 30 și 50 de hărți pe stratul  $L_2$ , respectiv 40, 60 și 80 de neuroni pe stratul  $L_3$  am efectuat toate combinațiile posibile, construind  $4 \times 3 \times 3 = 36$  de rețele neuronale. Fiecare rețea a fost antrenată în două moduri: cu, respectiv fără activarea deformărilor asupra imaginilor din setul de antrenare. Operațiile efectuate asupra imaginilor sunt similare cu cele din [64], dar am ales valori mai mici pentru

rotații, translații și deformări elastice. Motivul descreșterii valorilor acestor parametri este dimensiunea mult mai mică a imaginilor pe care le-am folosit,  $21 \times 13$ , față de  $29 \times 29$  cât a folosit Simard pentru cifre singulare. Explicit, factorul de scalare a fost limitat la maxim 6%, rotațiile la maxim  $\pm 4^\circ$ ,  $\sigma$  pentru distorsiuni elastice a fost ales 10.0, iar scalarea pentru distorsiuni elastice - 0.2.

Antrenarea celor 36 de rețele cu și fără imagini distorsionate a generat 72 de seturi distincte de rezultate (Tabelul 5.4). Pentru rețelele antrenate cu imagini distorsionate există două tipuri de erori pentru setul de antrenare. Prima este eroarea obținută testând rețeaua cu setul de antrenare asupra căruia s-a aplicat o transformare aleatoare (rotație + translație + deformare) la începutul fiecărei epoci. Al doilea tip de eroare este obținut pe setul de antrenare original.

Rezultatele din Tabelul 5.4 arată că eroarea este monoton descrescătoare (cu extrem de mici excepții ce au drept cauză probabilă o subantrenare a rețelei respective sau atingerea unui minim local apropiat de minimul global). Pentru oricare două rețele antrenate similar - ambele cu set de antrenare cu imagini deformate sau ambele cu set de antrenare cu imagini nedeformate - oricare din cele cinci tipuri de erori a variat corespunzător cu modificarea unuia din cei trei parametri ai rețelei. Explicit, dacă doi dintre parametri sunt menținuți constanți (de exemplu, numărul de hărți din  $L_1$  și numărul de neuroni din  $L_3$ ), iar cel de-al treilea (în acest caz, numărul de hărți de pe stratul  $L_2$ ) este mărit, erorile descresc și reciproc. Concluzia evidentă este că erorile pot fi diminuate dacă mărim dimensiunea rețelei. Pentru a verifica această afirmație, am antrenat trei rețele și mai mari (ultimele trei rânduri din Tabelul 5.4). Am obținut cele mai mici erori (5.35% pentru setul de test) cu rețeaua cu 10 hărți pe  $L_1$ , 100 de hărți pe  $L_2$  și 100 de neuroni pe  $L_3$ . Eroarea pe setul de test a fost cu 1.60% mai mică decât precedentul minim obținut, 6.85% (10 hărți pe  $L_1$ , 60 de hărți pe  $L_2$  și 100 de neuroni pe  $L_3$ ). Dezavantajul a fost o creștere accentuată a timpului necesar antrenării, de la 24 de ore la peste 48 de ore.

În concluzie, cele mai bune rezultate, exceptând cazurile speciale cu 10 sau 12 hărți de convoluție pe stratul  $L_1$ , au fost obținute pentru opt hărți pe stratul  $L_1$ . Suprafațele erorilor sunt arătate în Figura 5.3. Cea mai mare rețea antrenată (ultimul rând din Tabelul 5.4) a conținut 12 hărți pe  $L_1$ , 120 de hărți pe  $L_2$  și 120 de neuroni pe  $L_3$ . Chiar și după 96 de ore de antrenare eroarea nu a scăzut sub minimul atins de rețeaua cu 10 hărți pe  $L_1$ , 100 de hărți pe  $L_2$  și 100 de neuroni pe  $L_3$ . Am concluzionat că am detectat un minim suficient de apropiat de minimul global și mărirea în continuare a parametrilor rețelei ar fi inutilă și foarte costisitoare din punct de vedere al timpului de antrenare.

Referitor la timpul necesar unei epoci de antrenare, acesta a variat de la 70 de secunde pentru cea mai simplă rețea până la 670 de secunde pentru cea mai complexă, cea cu 10 hărți pe  $L_1$ , 100 de hărți pe  $L_2$  și 100 de neuroni pe  $L_3$ . Antrenarea a fost mult mai rapidă decât în cazul perceptronului, iar rata de recunoaștere este semnificativ mai bună.

Chiar dacă am scăzut amplitudinea transformărilor aplicate pe imaginile din setul de antrenare, toate rețelele antrenate cu imagini deformate au atins o eroare minimă mai mare, atât pe setul de antrenare, cât și pe setul de test. Adăugarea deformărilor a degradat performanța rețelelor cu cel puțin 0.85% (cazul  $L_1$ : 8 hărți,  $L_2$ : 50 de hărți,  $L_3$ : 80 de neuroni). Se poate concluziona că acest tip de aplicare a deformărilor este eficient doar pentru clasificarea cifrelor singulare, nu și pentru perechi de cifre parțial suprapuse.

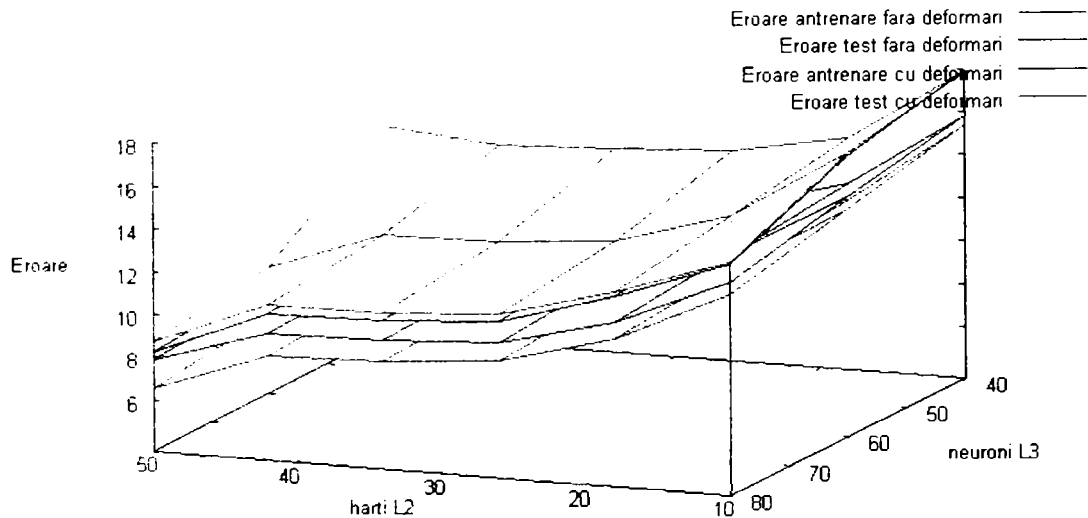


Figura 5.3: Suprafețele erorilor rețelelor cu 8 hărți de convoluție pe stratul  $L_1$ .

### Analiza rezultatelor celei mai bune rețele

Experimentele efectuate au arătat că cea mai bună rețea are zece hărți de convoluție a  $9 \times 5$  neuroni pe stratul  $L_1$ , o sută de hărți de convoluție a  $3 \times 1$  neuroni pe stratul  $L_2$  și o sută de neuroni pe stratul  $L_3$ . Din cele 200000 de imagini din setul de test, 10708 au fost clasificate greșit. Aceasta reprezintă o eroare de 5.35% pe setul de test.

Din cele 10708 erori doar în 663 (0.33%) de cazuri ambele cifre au fost clasificate defectuos. Această rată de eroare atât de redusă oferă premisele dezvoltării unui sistem de recunoaștere care va utiliza acest clasificator pentru recunoașterea șirurilor de cifre. Cea mai semnificativă cifră - cifra zecilor - a fost clasificată eronat în 5239 de cazuri, iar cea mai puțin semnificativă cifră - cifra unităților - în 6123 de cazuri. Putem concluziona că este mai dificil pentru rețeaua neuronală să identifice cifre suprapuse în partea stângă decât în partea dreaptă.

Distribuția erorilor pe clase arată că clasele 67, 64 și 69 au fost mai dificil de recunoscut, generând 253 (2.36%), 219 (2.04%), respectiv 212 (1.98%) erori de clasificare. În paranteze este trecută procentual contribuția clasei la eroarea totală.

Rata de eroare de 5.35% înseamnă 94.65% rată de recunoaștere. Aceasta reprezintă o rată de recunoaștere forțată, adică fără rejecție. Pentru a îmbunătăți calitatea detecției am calculat ratele de recunoaștere considerând rate de rejecție variabile. O imagine este rejectată dacă  $N(C1) - N(C2) < T$ , unde  $N(C1)$  și  $N(C2)$  reprezintă primele două maxime ale neuronilor de ieșire pentru acea imagine, normalizate la intervalul  $[-1, 1]$ , iar  $T$  reprezintă pragul de rejecție. Atât autorii din [44], cât și din alte publicații folosesc pentru valoarea de prag valorile 0.5% și 1%, astfel că am calculat ratele de eroare pentru aceste rate de rejecție, dar am adăugat și pentru două valori mult mai mari: 2%, respectiv 5%. Diferențele foarte mici dintre ratele de recunoaștere din Tabelul 5.5 arată că rețeaua are un grad mare de siguranță,

rata de recunoaștere cu rejecția maximă de 5% fiind cu numai 0.59% mai mică decât rata de recunoaștere cu fără rejecție.

Tabelul 5.5: Influența strategiei de rejecție asupra ratei de recunoaștere (forțat = fără rejecție).

T (%)	forțat	0.5	1	2	5
rata de recunoaștere (%)	94.65	94.59	94.53	94.42	94.06

### 5.2.3 Concluzii

RNC sunt mult mai performante decât perceptronul cu un strat ascuns, acesta din urmă generând o rată de eroare prea mare pentru a putea fi folosit în clasificatori mai complecși. RNC reprezintă o soluție bună atât pentru recunoașterea de cifre parțial suprapuse, cât și pentru recunoașterea de cifre izolate. Rata de eroare minimă obținută pentru setul de test din TDNS este de 5.35% pentru o rețea cu 10 hărți pe L1, 100 de hărți pe L2 și 100 de neuroni pe L3.

Deformările elastice nu îmbunătățesc performanțele în cazul cifrelor parțial suprapuse, așa cum făceau pentru cifre singulare. Cum voi arăta în secțiunile care urmează, există motive întemeiate pentru care deformările nu s-au dovedit utile.

Rata de eroare de 5.35% este suficient de mică pentru a folosi acest clasificator într-unul general pentru recunoașterea de șiruri numerice formate dintr-un număr indefinit de cifre.

Diferențele extrem de reduse dintre ratele de recunoaștere cu și fără rejecție demonstrează gradul de robustețe al clasificatorului.

Rezultatele au fost publicate în [20] și au fost prezentate la ICCP 2008.

## 5.3 Experimente cu șiruri de cifre din NIST SD 19

După ce am obținut doi clasificatori cu rate bune de recunoaștere - unul pentru cifre izolate și unul pentru perechi de cifre parțial suprapuse - am trecut la construirea unui clasificator care să poată recunoaște numere formate dintr-un număr oarecare de cifre. Acest clasificator a fost apoi testat cu șiruri numerice din NIST SD 19 și comparat cu cele mai bune rezultate existente în domeniu.

### 5.3.1 Sistemul de recunoaștere

Diagrama sistemului de recunoaștere este prezentată în 5.4. Fiecare șir de cifre este întâi tăiat la dimensiunile dreptunghiului delimitator și apoi este prezentat Analizorului de Componente Conex (CCA) care extrage toate componentele conexe (CC) din imagine. În continuare, componentele conexe vor fi clusterizate în imagini ce probabil vor conține cifre sau mulțimi de cifre. Procesul de clasificare este efectuat de Clasificator folosind două Rețele Neuronale de Convoluție (RNC).

### 5.3.2 Analizorul de componente conexe

Toate componentele conexe din imaginea de intrare sunt extrase folosind procedeul recursiv de căutare în doar patru direcții (mai corect: două direcții cu două sensuri fiecare), orizontale

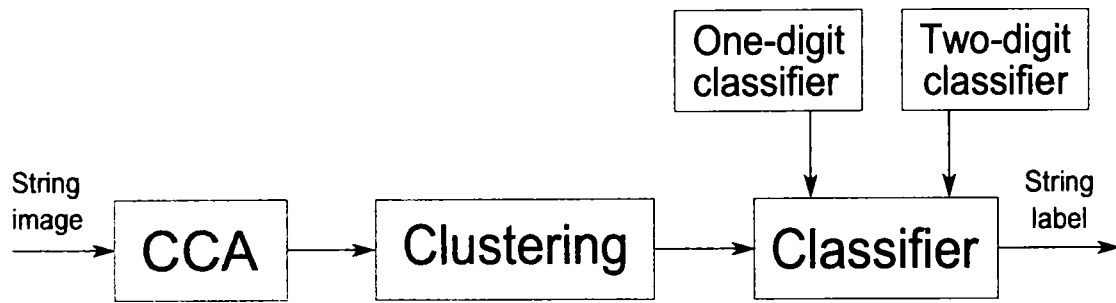


Figura 5.4: Arhitectura sistemului de recunoaștere a șirurilor de cifre.

și verticale. Nu am folosit căutarea în opt direcții deoarece neefectuând segmentare acesta este unicul loc unde imaginea de intrare poate fi separată în cât mai multe componente printr-un simplu proces de extragere. Orice separare a imaginii unei cifre sau a mai multor cifre va fi corectată în faza următoare, cea de clusterizare.

Pentru fiecare componentă se calculează diverși parametri: arie, dreptunghiul circumscris de arie minimă (dreptunghiul delimitator), lățime, înălțime, aspect (lățime / înălțime), distanțe la alte componente. CCA atribuie câte o culoare distinctă fiecărei componente. În contrast cu alte metode ([44, 30, 31]), prezenta metodă poate evita procesul de segmentare prin utilizarea suplimentară a unui clasificator pentru perechi de cifre parțial suprapuse, pe lângă clasicul clasificator la nivel de cifră.

### 5.3.3 Componenta de clusterizare (agregare)

În cazul imaginilor scanate există multe componente conexe care formează o singură cifră sau mai multe cifre alipite. Pentru a putea reconstrui imaginea lor, etapa de agregare efectuează secvențial următoarele operații (constantele din toate testele referitoare la suprafața unei componente - numărul de pixeli - au fost stabilite pentru imagini din setul NIST SD 19. O caracteristică a acestor imagini e aceea că au fost scanate la 300dpi. Constantele pot fi folosite și în cazul altor imagini scanate la alte rezoluții, dar trebuie multiplicare cu pătratul (deoarece suprafața e o arie) factorului de redimensionare de la o rezoluție la cealaltă):

- toate componentele mici conținând maxim 4 pixeli sunt șterse dacă cea mai apropiată componentă este la o distanță mai mare de 5 pixeli. Motiv: sunt prea mici și prea depărtate de alte componente, deci probabil că reprezintă doar zgomot.
- oricare două componente separate prin doar un pixel sunt concatenate dacă respectă cel puțin una din condițiile:
  - ambele au înălțimi mai mari de 60% din înălțimea imaginii, iar distanța dintre ele sau distanța de suprapunere este mai mică de 35% din lățimea componentei celei mai late.
  - ambele au arii mai mari de 200 de pixeli și aspecte mai mari de 0.5, și cel puțin una dintre ele are aspectul mai mare decât 1 și distanța dintre ele sau distanța de suprapunere parțială este mai mică decât 25% din lățimea componentei mai late.
  - ambele au arie mai mare de 400 de pixeli

După acest pas se reia algoritmul pentru componentele separate prin maxim doi pixeli, și, la final, pentru maxim trei pixeli. Aceste operații agreghează cu succes componentele din imaginile scanate defectuos sau care conțin cifre scrise intermitent (Figura 5.5).

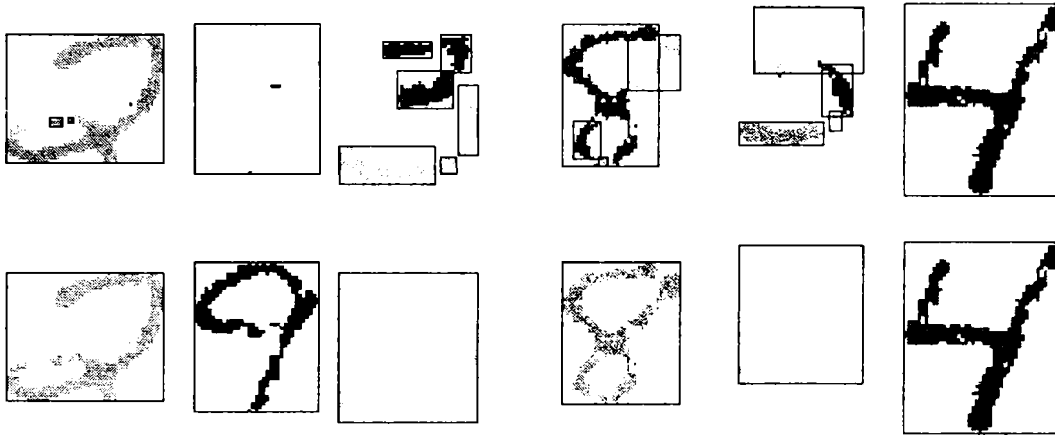


Figura 5.5: Componentele conexe înainte (sus) și după (jos) faza de agregare. Fiecare componentă este colorată distinct și este desenat și dreptunghiul delimitator.

- există multe persoane care scriu cifra 5 din două părți, una pentru corpul cifrei și cealaltă pentru partea superioară. În mod frecvent cele două părți sunt foarte distanțate pe orizontală și etapele anterioare de agregare nu sunt suficiente pentru a reconecta părțile cifrei 5. Acest caz este distinct și apare doar în cazul cifrei 5. Pentru a reforma cifra, se caută întâi partea superioară a lui 5 ca fiind o componentă cu următoarele proprietăți: are cel puțin 10 pixeli, are un aspect alungit pe orizontală (mai mare de 1.4), nu este prea groasă (înălțime mai mică decât 40% din înălțimea imaginii), și este plasată în jumătatea superioară a imaginii. Dacă o componentă de acest tip este detectată, atunci se caută o componentă care să poată reprezenta partea inferioară a cifrei 5. O astfel de componentă e obligatoriu să aibă minim 100 de pixeli.



Figura 5.6: Reconectarea cifrei 5.

Partea superioară a cifrei trebuie să înceapă în partea dreaptă a corpului ei, dar nu mai departe de 50% din lățimea corpului cifrei după aceasta. Dacă nicio componentă nu este detectată, atunci componenta care se presupunea că reprezintă partea superioară a lui 5 rămâne neconcatenată. Dacă ambele părți ale cifrei 5 sunt detectate, atunci

partea superioară este translatată spre stânga până ajunge pe aceeași verticală cu pixelul din dreapta-sus al componentei ce reprezintă corpul cifrei 5 (Figura 5.6).

- acest pas tratează puținele cazuri de cifre scrise din mai multe părți separate vertical. Dacă două componente sunt suprapuse orizontal pe minim 80% din lățimea cel puțin a uneia dintre ele, atunci sunt concatenate.

### 5.3.4 Clasificatorul

Componenta de clasificare (Figura 5.4) folosește două RNC, una pentru cifre și una pentru perechi de cifre parțial suprapuse. Ambele RNC au fost descrise anterior, iar ratele de eroare obținute au fost 0.86% pentru o cifră, respectiv 5.35% pentru două cifre.

Testarea sistemelor de recunoaștere a fost efectuată pe seturile de test 3DNS, 6DNS și 10DNS (detalii sunt prezentate în Secțiunile 4.5.2, 4.5.3 și 4.6), seturi care sunt folosite și de alți autori [44] pentru compararea sistemelor de recunoaștere. Autorii din [44] și-au comparat metodele cu cele mai bune din literatură, prezentând cele mai eficiente metode existente până în anul 2004.

Folosind cele două RNC am implementat mai multe sisteme de recunoaștere. Ambele RNC preiau rând pe rând câte o componentă din imaginea de la intrare și încearcă să determine cărei clase i se potrivește respectiva componentă (să recunoască cifra sau perechea de cifre din imaginea componentei). Înainte de a intra în etapa de clasificare, toate componentele care sunt ori mai mici de 80 de pixeli ori au înălțimea mai mică de 30% din înălțimea imaginii sunt eliminate.

În toți algoritmi ce urmează,  $rez_i$  reprezintă un vector de structuri, fiecare structură conținând două câmpuri: *class* și *score*. Indicele  $i$  este 1 pentru 1RNC și 2 pentru 2RNC. Vectorii  $rez_i$  stochează ieșirea RNC. De exemplu,  $rez_2[3].class$  reprezintă clasa celui de-al patrulea element al rezultatului obținut de 2RNC. Câmpul *class* nu este redundant deoarece este util după ce vectorul este sortat, pentru a ști ce clasă are cel mai bun scor, ce clasă are al doilea scor etc.

#### Sistem de recunoaștere bazat pe scorul maxim

Cel mai simplu sistem de recunoaștere care poate fi implementat ia o imagine și o prezintă la  $RNC_1$  și la  $RNC_2$ . Cele mai mari scoruri ale celor două RNC sunt comparate, iar clasificatorul care a obținut maximum va oferi atât clasa, cât și numărul de cifre prezente în imaginea de la intrare (Algoritm 4).

Din cele 1476 imagini (setul 3DNS) care conțineau 3 cifre, 158 au fost incorect clasificate. Aceasta înseamnă o rată de recunoaștere de 89.29%. Rezultatul pentru perceptronul multistrat din [44] antrenat fără exemple negative este 82.52%. Având în vedere că metoda prezintă nu folosește segmentare deloc, recunoașterea corectă a 89.29% din imagini este promițătoare, dar mai există loc pentru îmbunătățiri. Este foarte important a se observa că rezultatul obținut reprezintă dovada eficienței combinării a doi clasificatori, unul pentru recunoașterea cifrelor și unul pentru recunoașterea perechilor de cifre parțial suprapuse, față de folosirea doar a unui clasificator la nivel de cifră [44], chiar dacă ultima este ajutată de un complicat mecanism de segmentare [44].

După studierea erorilor sistemului de recunoaștere, am ajuns la următoarele concluzii:



```

foreach componentă din imagine do
  rez1[ ] = RNC1(componentă);
  rez2[ ] = RNC2(componentă);
  sortează descrescător rez1 și rez2;
  if rez1[0].score ≥ rez2[0].score then
    | return rez1[0].class;
  end
  else
    | return rez2[0].class;
  end
end

```

**Algorithm 4:** Sistem de recunoaștere bazat pe scorul maxim.

- deoarece am antrenat ambele RNC fără exemple negative, ele nu pot fi utilizate direct pentru a detecta dacă o componentă conține una, două sau mai multe cifre. Prin antrenarea unei RN doar cu exemple pozitive, rețeaua va încerca să clasifice orice imagine de intrare, chiar și pe cele incorecte - care nu conțin reprezentări grafice ale niciuneia dintre clasele învățate -, clasei cu care seamănă cel mai mult. Mai mult, probabil din cauza naturii convoluționale a rețelelor, cele mai confuze cazuri apar la numerele formate din două cifre de tipul  $\overline{XX}$ . Este foarte dificil atât pentru  $RNC_1$ , cât și pentru  $RNC_2$  să distingă dacă o imagine conține o singură cifră sau un număr format din două cifre. De exemplu, dacă imaginea conține doar cifra 4, scorul ambelor RNC poate fi foarte mare și foarte asemănător,  $RNC_2$  detectând numărul 44!
- multe imagini conținând cifrele 0, 1, 3, 4, 7, 8, respectiv 9 sunt interpretate de  $RNC_2$  ca fiind 01, 11, 31, 41, 71, 81, respectiv 91. Problema se datorează metodei de unire ([20]) a cifrelor pentru setul de antrenare al  $RNC_2$ , deoarece dacă cifra 1 este scrisă ca un segment de dreaptă și este suprapusă parțial cu o cifră care are în partea dreaptă pixeli pe toată înălțimea cifrei, rezultatul este foarte confuz (Figura 5.7). Un argument similar se poate oferi și pentru cazurile 10, 16, 18. Aceste cazuri sunt notate cu  $\overline{1X}$  și  $\overline{X1}$ .

0 1 3 4 7 8 9

Figura 5.7: Cazurile de unire defectuoasă de tipul  $\overline{X1}$ .

Evident, antrenarea ambelor RNC cu exemple negative ar fi rezolvat multe probleme, dar colectarea manuală a imaginilor cu exemple negative - imagini care nu conțin cifre sau perechi de cifre - din formularele NIST SD 19 este un proces foarte îndelungat și neplăcut. Numărul de exemple negative trebuie să fie suficient de mare pentru a asigura o antrenare adecvată. În [44], autorii au antrenat clasificatorul la nivel de cifră cu 1000 și apoi cu 4000 de exemple negative. În ambele cazuri rata de recunoaștere a clasificatorului la nivel de cifră s-a redus cu 2-4%, dar acesta a devenit rezistent la imagini ce nu conțin cifre. Metoda prezentată în această teză are doi clasificatori, deci ar trebui colectate exemple negative pentru ambii. Bineînțeles, scrierea unui algoritm de generare automată a exemplilor negative este un alt

mod de a rezolva această problemă, dar ar trebui conceput un astfel de algoritm. Autorii din [44] au ajuns la concluzia că tot mai bune sunt datele colectate manual, astfel că scrierea algoritmului de generare ar putea reprezenta timp irosit. După ce am evaluat toate aceste aspecte, am decis să nu antrenez clasificatorii cu exemple negative și să încerc să dezvolt o metodă care să atenueze problemele prezentate anterior într-o manieră cât mai simplă.

Când am studiat erorile, am observat că în majoritatea cazurilor RNC recunoșteau corect, dar scorul relativ al clasificatorilor genera clasificarea incorectă.

Am antrenat ambele RNC astfel încât să fie cât mai puțin influențate de translații, dar dacă imaginea este perfect centrată în clasificator, atunci scorurile sunt mult mai bune. Considerând acest aspect, am încercat să translatez imaginea în câmpul de intrare al clasificatorilor.  $RNC_1$  a fost antrenată cu baza de date LNIST, iar imaginile din ea sunt centrate în centrul de masă. Din această cauză am aplicat repetat  $RNC_1$  asupra imaginii translate cu  $\pm 1$  pixeli relativ la centrul de masă și am păstrat cel mai mare scor din cele  $3 \times 3 = 9$  aplicări.  $RNC_2$  a fost antrenată cu imagini care au fost centrate prin centrarea dreptunghiului delimitator. Considerând această metodă de generare ([20]) a imaginilor de antrenare, am plasat imaginea de  $18 \times 10$  pixeli în toate pozițiile posibile în câmpul de intrare de  $21 \times 13$  neuroni ai  $RNC_2$ . Există  $(21 - 18 + 1) \times (13 - 10 + 1) = 16$  posibilități. Din nou am reținut scorul maxim pentru fiecare clasă. Din păcate, numărul de clasificări incorecte a crescut (Tabelul 5.6) și timpul de execuție s-a mărit din cauza aplicărilor repetate ale clasificatorilor. Se poate concluda că mutarea clasificatorilor deasupra imaginii (sau, echivalent, a imaginii în câmpul clasificatorilor) nu îmbunătățește rata de recunoaștere a metodei ce folosește clasificarea bazată pe scor maxim.

### Sistem de recunoaștere bazat pe diferența de scoruri

Tabelul 5.6: Rate de recunoaștere pentru setul de test 3DNS (1C - cu deplasarea  $RNC_1$  deasupra imaginii, 2C - cu deplasarea  $RNC_2$  deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor).

	simplu(%)	1C(%)	2C(%)	1+2C(%)
scor maxim	91.46	71.54	69.85	71.54
diferențe	76.49	76.76	92.75	93.36
scor maxim & diferențe	83.40	83.74	93.63	<b>93.77</b>

Arătăm înainte că nu putem să ne bazăm decizia selectării unui anumit clasificator pe scorul maxim oferit de acesta deoarece nu am antrenat clasificatorii cu exemple negative. Mai mult, prin antrenarea rețelelor numai cu exemple pozitive, RNC vor clasifica orice imagine de la intrare într-una din cele mai asemănătoare clase cunoscute. Aceasta înseamnă că dacă un clasificator primește la intrare o imagine care este reprezentarea uneia dintre clasele pe care le cunoaște, atunci ieșirea neuronului corespunzător clasei va avea o valoare apropiată de 1, iar celelalte ieșiri vor fi considerabil mai mici, apropiate de -1. Putem utiliza acest aspect al rețelelor neuronale artificiale pentru a deriva o nouă metodă de a decide dacă o imagine aparține sau nu vreuneia din clasele cu care a fost antrenată rețeaua. Astfel, diferența dintre cele mai mari două scoruri trebuie să fie considerabilă pentru ca imaginea să aparțină unei clase învățate. Se pot imagina cazuri în care această metodă nu va funcționa, de exemplu când imaginea de la intrare nu aparține niciuneia din clase, dar seamănă destul de mult cu

una din ele, iar cu restul nu. Un alt caz ar fi acela când există clase foarte asemănătoare, iar o imagine care seamănă cu una din ele va semăna, evident, și cu celelalte clase, generând din acest motiv scoruri mari și apropiate. Oricum, se poate presupune că aceste cazuri sunt excepții de la regulă, fiind suficient de rare pentru ca metoda să fie eficientă.

Considerând observațiile pe care le-am făcut în urma studierii erorilor generate de Algoritmul 4 pentru cazurile  $\overline{X1}$ , am implementat Algoritmul 5, întărind condiția dacă cel puțin una din cifrele pe care le detectează ca posibil existente  $RNC_2$  este 1.

Contrar cu primul sistem de recunoaștere care a folosit doar scorul maxim al fiecărei RNC, în acest caz mutarea clasificatorilor deasupra imaginilor de intrare a îmbunătățit rata de recunoaștere cu mai bine de 6% până la 92.82%.

```

foreach componentă din imagine do
  rez1[ ] = RNC1(componentă);
  rez2[ ] = RNC2(componentă);
  sortează descrescător rez1[ ] și rez2[ ];
  dif1 = rez1[0].score - rez1[1].score;
  dif2 = rez2[0].score - rez2[1].score;
  if u=1 or z=1 then
    if dif1 ≥ dif2 then
      | return rez1[0].class;
    end
    else
      | return rez2[0].class;
    end
  end
  else
    if dif1 ≥ 2 · dif2 then
      | return rez1[0].class;
    end
    else
      | return rez2[0].class;
    end
  end
end

```

**Algorithm 5:** Sistem de recunoaștere bazat pe diferența de scoruri (z - cifra zecilor, u - cifra unităților).

### Sistem de recunoaștere bazat pe scor maxim și pe diferențe de scoruri

Pentru a vedea dacă se pot obține rezultate mai bune combinând metodele anterioare, am implementat Algoritmul 6. Pentru fiecare clasificator am adunat cel mai bun scor cu diferența dintre cele mai bune două scoruri. Cel mai bun rezultat, **93.77%**, reprezentând 92 de șiruri incorect clasificate, a fost obținut când am aplicat translații asupra clasificatorilor. Acest rezultat este cu 3 până la 18% mai bun decât toate rezultatele din [44] obținute cu numeroase arhitecturi de rețele neuronale, deși acelea au fost antrenate și cu exemple negative. Chiar dacă comparăm metoda prezentată cu cea mai bună din [44], diferența este de doar 3.05%,

și aceasta în contextul în care noua metodă nu folosește nici segmentare, nici antrenare cu exemple negative.

```

foreach componentă din imagine do
  rez1[ ] = RNC1(componentă);
  rez2[ ] = RNC2(componentă);
  sortează descrescător rez1[ ] și rez2[ ];
  dif1 = rez1[0].score - rez1[1].score;
  dif2 = rez2[0].score - rez2[1].score;
  if u=1 or z=1 then
    if rez1.score + dif1 ≥ rez2.score + dif2 then
      | return rez1[0].class;
    end
    else
      | return rez2[0].class;
    end
  end
  else
    if rez1 + dif1 ≥ rez2 + 2 · dif2 then
      | return rez1[0].class;
    end
    else
      | return rez2[0].class;
    end
  end
end

```

**Algorithm 6:** Sistem de recunoaștere bazat pe scor maxim și pe diferențe de scoruri ( $z$  - cifra zecilor,  $u$  - cifra unităților).

Algoritmul poate fi îmbunătățit mai departe prin complicarea regulilor de selecție a rezultatului unuia dintre clasificatori, dar își va pierde simplitatea. Analiza erorilor arată că RNC sunt foarte bune în a recunoaște imaginile de la intrare, dar regulile care selectează rezultatul unuia dintre clasificatori nu oferă totdeauna rezultate corecte. Cele 92 de erori pot fi împărțite în cinci tipuri:

- 4 erori de segmentare, toate fiind cauzate de imagini scanate deficitar,
- 8 erori cu imagini care conțin trei cifre concatenate. Sistemul de recunoaștere nu tratează cazurile de imagini cu mai mult de două cifre parțial suprapuse.
- 23 erori ale  $RNC_1$ . Aproximativ jumătate dintre ele sunt cauzate de o subreprezentare a specimenului din clasa respectivă în setul de antrenare, iar celelalte de imagini foarte confuze.
- 26 erori ale  $RNC_2$ . Majoritatea lor se datorează faptului că cele două cifre parțial suprapuse sunt extrem de diferite în înălțime (chiar cu mai mult de 200%), iar antrenarea  $RNC_2$  a fost efectuată cu perechi de cifre diferite în înălțime prin maximum 10%.
- 63 erori generate de alegerea defectuoasă a clasificatorului. Acestea ar putea fi corectate complicând regulile de selecție sau, preferabil, prin antrenarea și cu cu exemple negative.

Tabelul 5.7: Rate de recunoaștere pentru setul de test 6DNS (1C - cu deplasarea  $RNC_1$  deasupra imaginii, 2C - cu deplasarea  $RNC_2$  deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor).

	simplu(%)	1C(%)	2C(%)	1+2C(%)
scor maxim	88.72	89.73	47.93	50.58
diferențe	56.08	56.83	91.91	92.45
scor maxim & diferențe	66.42	67.51	91.50	<b>92.52</b>

Pentru a verifica dacă cel de-al doilea clasificator - cel la nivel de două cifre - contribuie la îmbunătățirea ratei de recunoaștere, am rulat experimentele dezactivând 2RNC și lasând activ doar clasificatorul la nivel de cifră. Au fost identificate corect 89.16% din imagini în cazul în care nu am deplasat imaginea în câmpul de intrare al clasificatorului, și 89.30% dacă am utilizat deplasările. Se poate observa că ambele scoruri sunt cu peste 4% mai mici decât dacă s-a folosit și 2RNC.

## 5.4 Experimente pe setul 6DNS

Cele mai bune rezultate (Tabelul 5.7) au fost obținute și de această dată aplicând translatarea imaginilor în câmpul de intrare al clasificatorilor. De asemenea, combinația dintre scorul maxim și diferența dintre primele două scoruri s-a dovedit a fi cea mai performantă.

Din cele 1471 de imagini ale setului 6DNS, 110 imagini au fost incorect identificate, reprezentând 7.48% din setul de test 6DNS. Analiza rezultatelor a arătat că și în cazul acestui set de test cele mai multe erori provin din alegerea eronată a clasificatorului. De asemenea, am observat că clasificatorul la nivel de două cifre este mai susceptibil la erori decât clasificatorul la nivel de cifră. Aceasta este normal deoarece complexitatea imaginilor ce trebuie învățate este mult mai mare. Un alt motiv este dimensiunea mai mică a câmpului de intrare a clasificatorului la nivel de două cifre.

Față de rezultatele obținute în [44] pe setul 6DNS, rata de 92.52% este net superioară celei de 63.7% cât au obținut cu perceptronul multistrat antrenat fără exemple negative. Cel mai bun rezultat din [44] (obținut cu Support Vector Classifier with Radial Basis Functions Gaussian kernel), 96.74%, este cu 4.22% mai bun decât cel obținut cu metoda prezentată, dar folosește atât segmentare, cât și antrenare cu exemple negative.

Am verificat și în acest caz comportarea sistemului de recunoaștere în cazul dezactivării 2RNC. Lasând activ doar clasificatorul la nivel de cifră, au fost identificate corect 91.37% din imagini în cazul în care nu am deplasat imaginea în câmpul de intrare al clasificatorului, și 91.91% dacă am utilizat deplasările. Față de setul 3DNS, unde îmbunătățirea adusă de includerea 2RNC a fost de peste 4%, majorarea ratei de recunoaștere s-a diminuat la 0.61% în cazul setului 6DNS. O explicație posibilă ar fi că setul 6DNS conține imagini mai dificil de recunoscut decât setul 3DNS.

Tabelul 5.8: Rate de recunoaștere pentru setul de test 10DNS (1C - cu deplasarea  $RNC_1$  deasupra imaginii, 2C - cu deplasarea  $RNC_2$  deasupra imaginii, 1+2C - cu deplasarea ambelor RNC deasupra imaginii, simplu - fără deplasarea clasificatorilor).

	simplu(%)	1C(%)	2C(%)	1+2C(%)
scor maxim	82.69	85.31	20.27	23.23
diferențe	33.83	34.05	91.80	<b>92.82</b>
scor maxim & diferențe	49.09	49.77	89.64	91.80

## 5.5 Experimente pe setul 10DNS

Rezultatele obținute pe acest set de test pot fi urmărite în Tabelul 5.8. Erorile mari obținute când deplasarea imaginii în câmpul 2RNC a fost oprită arată că dimensiunea de 21x13 pixeli pentru imagini de două cifre este prea mică. Această observație este susținută și de faptul că deplasarea imaginilor în câmpul clasificatorului la nivel de cifră îmbunătățește foarte puțin rata de recunoaștere. Probabil din cauză că în acest caz imaginile sunt mult mai mari, 29x29 pixeli, motiv pentru care clasificatorul este mai robust. Autorii din [44] nu au prezentat rezultatele pentru acest set pe motiv că niciunul din articolele anterioare nu prezenta, de asemenea, rezultatele pe acest set. Din acest motiv nu se poate compara rezultatul obținut în cazul setului 10DNS.

## 5.6 Concluzii

Ratele de eroare obținute au fost de 0.86% pentru o cifră, respectiv de 5.35% pentru două cifre, ceea ce arată că RNC reprezintă o bună alegere atât pentru clasificatorul dedicat cifrelor, cât și pentru clasificatorul dedicat perechilor de cifre parțial suprapuse.

Deși RNC nu au fost antrenate cu exemple negative, rezultatele obținute pentru seturile de test 3DNS și 6DNS sunt mult mai bune decât cele obținute anterior în domeniu, cu rețele neuronale artificiale. Diferențele ajung până la câteva zeci de procente în cazul rețelelor antrenate fără exemple negative.

Utilizarea combinației de clasificatori s-a dovedit a fi eficientă, rata de recunoaștere crescând față de sistemul ce utilizează numai clasificatorul de cifre cu 4.47% pentru setul 3DNS, respectiv cu 0.61% pentru setul 6DNS.

Metoda propusă oferă rate de recunoaștere cu numai câteva procente (3.03% pentru 3DNS și 4.42% pentru 6DNS) mai mici decât cele mai bune rezultate obținute anterior [44]. Având în vedere că rezultatele sunt obținute fără a aplica segmentarea și fără antrenare cu exemple negative, ratele de recunoaștere pot fi considerate foarte bune.

Din analiza erorilor rezultă următoarele două observații: clasificatorul la nivel de perechi de cifre parțial suprapuse ar trebui îmbunătățit și ambii clasificatori ar trebui antrenați cu exemple negative. Pentru a îmbunătăți setul TDNS trebuie detaliată metoda de unire automată a cifrelor, prin crearea de cazuri specifice pentru  $\overline{1X}$  și  $\overline{X1}$ . De asemenea, trebuie mărite imaginile deoarece dimensiunea de  $21 \times 13$  este prea mică pentru reprezentarea suficient de clară a imaginilor. Antrenarea cu exemple negative este foarte importantă pentru ambii clasificatori, nemaexistând apoi alegeri greșite de clasificator. Din analiza manuală a erorilor rezultă că dacă totdeauna s-ar alege clasificatorul potrivit imaginii ( $RNC_1$  pentru cifre,

respectiv  $RNC_2$  pentru perechi de cifre), rata de recunoaștere ar deveni mai bună decât cea mai bună existentă în momentul de față în domeniu.

O mare parte a rezultatele au fost publicate în [17] și au fost prezentate la SYNASC 2008.

# Capitolul 6

## Optimizarea vitezei de execuție

Având în vedere că plăcile grafice folosite pentru accelerarea calculelor în virgulă mobilă au apărut pe piață la sfârșitul lunii iunie 2008, precum și dificultățile legate de scrierea unor programe pe arhitecturi complet noi și extrem de diferite față de arhitectura x86, acest capitol prezintă numai rezultate parțiale, pentru finalizarea portării tuturor aplicațiilor pe plăci grafice fiind necesare cel puțin încă câteva luni. Chiar dacă acest capitol nu aduce nimic nou în privința performanțelor sistemului de clasificare, neavând legătură explicită cu rata de recunoaștere, am considerat că este util a prezenta și aceste rezultate parțiale deoarece tehnologia folosită reprezintă viitorul pentru majoritatea domeniilor din știință care au nevoie de putere foarte mare de calcul, și cu siguranță în domeniul rețelelor neuronale. Mai multe detalii vor fi prezentate în secțiunile dedicate accelerării vitezei de execuție cu ajutorul plăcilor grafice.

Clasificatorul 2RNC a fost antrenat cu o versiune modificată a programelor open source din [3]. Din păcate, acest cod a fost scris într-un mod mai mult "didactic", decât având ca țintă performanța. Oricâte modificări i-am făcut în ideea de a-i îmbunătăți viteza, am ajuns la o limită sub care nu mai puteam scădea timpul de rulare decât dacă rescriam de la zero tot programul. Problema și mai neplăcută era că dacă măream dimensiunea rețelei neuronale peste un anumit nivel, programul începea să funcționeze haotic. Încercarea de a descoperi cauzele acestei comportări s-a soldat cu eșec, chiar după multe săptămâni de depanare (sursele programului au aproximativ 250KB de cod C++). Modul în care se aloca memoria în toată ierarhia de clase, precum și cantitatea imensă care era alocată - mult peste un gigabyte de RAM -, multă din ea chiar pe stivă și nu în heap, generau acest comportament. Un alt factor care a îngreunat mult depanarea a fost arhitectura implementată multi-threading (cu mai multe fire de execuție). Cum scrierea unui program de acest gen este foarte complexă și durează un timp apreciabil, am considerat că utilizarea unor imagini de dimensiuni mai reduse este compromisul optim. Astfel au rezultat toate experimentele și rezultatele din capitolul precedent.

Una din concluziile din capitolului anterior evidenția faptul că pentru a obține rezultate mai bune este nevoie de un clasificator mai performant pentru perechi de cifre parțial suprapuse. Pe lângă antrenarea cu exemple negative, performanțele clasificatorului pot fi crescute prin mărirea dimensiunii imaginilor cu care este antrenat. O dimensiune adecvată pentru a obține rezultate mai bune pentru 2RNC ar fi  $41 \times 29$  de pixeli. Aceasta înseamnă 1189 de neuroni pentru stratul de intrare al rețelei, și, implicit, numere crescute de neuroni și de hărți de convoluție pe celelalte straturi. Față de imaginile de  $21 \times 13$  pixeli cu care a fost antrenat



2RNC, imagini care necesită numai 293 de neuroni pentru stratul de intrare, imaginile cu dimensiunea de  $41 \times 29$  pixeli ar necesita o rețea de peste patru ori mai mare. Timpul de antrenare ar crește proporțional cu mărimea rețelei. Cum cea mai bună rețea pentru 2RNC a necesitat câteva zile pentru a fi antrenată, timpul pentru noua dimensiune a imaginilor ar ajunge la aproximativ două săptămâni, ceea ce este total nepractic.

Întrucât precizăm că portarea programelor nu este finalizată, pentru a putea face comparații de viteză este necesară alegerea unei porțiuni de cod reprezentativă pentru rețelele neuronale, cod care apoi să fie implementat în diverse moduri. Rețelele neuronale, ca toți clasificatorii, au două părți, una de antrenare și una de testare. Cum algoritmul de propagare înainte (forward propagation) apare atât la antrenare, cât și la testare, unde reprezintă chiar întreg pasul de testare, am decis să aleg o parte a acestui algoritm pentru a fi implementată pe diverse arhitecturi. Am ales cel mai complex strat al rețelei neuronale de convoluție, stratul L3. Acesta este un strat complet conectat, și nu de convoluție, dar conține de departe cele mai multe operații, necesitând cea mai mare parte a timpului de rulare a algoritmului de propagare înainte. Am ales stratul L3 al clasificatorului la nivel de cifră, 1RNC, deoarece este mai mare decât cel de la 2RNC. Este bine ca dimensiunea stratului să fie cât mai mare deoarece scopul final al acestei optimizări este acela de a permite antrenarea unor rețele de dimensiuni mai mari.

În toate experimentele următoare timpul a fost măsurat și raportat pentru 10000 de rulări pentru a putea fi măsurat corect și, în același timp, pentru a putea fi raportat în secunde, nu în microsecunde. De asemenea, setul de test al MNIST are 10000 de imagini, testele fiind efectuate pe întreg setul de test.

Stratul L3 conține 100 de neuroni fiecare din ei având 1251 de conexiuni, una pentru valoarea de prag (bias) și 1250 pentru cei 1250 de neuroni de pe stratul L2. În total există 125100 de conexiuni și tot atâtea ponderi. Calculul acestui strat este format din evaluarea a 100 de produse scalare dintre vectori cu 1251 de componente, urmată de aplicarea funcției de activare asupra rezultatului fiecărui produs scalar.

## 6.1 Variante CPU ale optimizării algoritmului de propagare înainte

În această subcapitol vor fi prezentate rezultatele obținute prin optimizarea codului ce rulează numai pe procesor, fără plăci grafice. Întâi se va prezenta algoritmul inițial, cel care folosește multe clase și vectori STL (Standard Template Library). Timpul său de rulare va fi folosit ca etalon pentru comparațiile următoare. Apoi vor fi prezentate două variante care nu folosesc ierarhii de clase și, mai important, vectori STL, ci memorie alocată dinamic cu `new/malloc`. Ultima variantă este implementată în limbaj de asamblare și folosește instrucțiuni SSE3.

### 6.1.1 Varianta open source modificată

Această variantă este cea cu care am efectuat experimentele, dar care s-a dovedit a fi imposibil de extins la rețele mai mari din cauza timpului mare de rulare și din cauza comportării eronate în cazul creșterii dimensiunii rețelei neuronale. Codul inițial este preluat din [3] și modificat pentru a putea fi utilizat pentru experimentele pe care le-am efectuat. În Figura 6.1 este prezentat codul sursă al algoritmului de propagare înainte. Se poate observa cât de

încărcat este codul și cât de mult face apel la vectori STL. Timpul necesar pentru cele 10000 de imagini ale setului de test a fost de 11.67 secunde. Ca observație, timpul necesar pentru întreg algoritmul de propagare înainte pentru toate cele 10000 de imagini ale setului de test din MNIST a fost de 27 de secunde.

```

void NNLayer::Calculate()
{
    ASSERT( m_pPrevLayer != NULL );

    VectorNeurons::iterator nit;
    VectorConnections::iterator cit;

    double dSum;

    for( nit=m_Neurons.begin(); nit<m_Neurons.end(); nit++ )
    {
        NNNeuron& n = *(*nit); // to ease the terminology

        cit = n.m_Connections.begin();

        ASSERT( (*cit).WeightIndex < m_Weights.size() );

        // weight of the first connection is the bias: neuron is ignored
        dSum = m_Weights[ (*cit).WeightIndex ]->value;

        for ( cit++ ; cit<n.m_Connections.end(); cit++ )
        {
            ASSERT( (*cit).WeightIndex < m_Weights.size() );
            ASSERT( (*cit).NeuronIndex < m_pPrevLayer->m_Neurons.size() );

            dSum += ( m_Weights[ (*cit).WeightIndex ]->value ) *
                ( m_pPrevLayer->m_Neurons[ (*cit).NeuronIndex ]->output );
        }

        n.output = SIGMOID( dSum );
    }
}

```

Figura 6.1: Codul sursă al algoritmului de propagare înainte pentru varianta C++ cu vectori STL.

### 6.1.2 Varianta C++ fără vectori STL

Această variantă (Figura 6.2) reprezintă prima tentativă de optimizare prin rescrierea programului de antrenare a rețelelor neuronale artificiale pe procesoare x86. Scopul ei a fost acela de a verifica gradul de accelerare a execuției față de varianta open source modificată, care se baza excesiv pe clase și pe vectori STL. Datele referitoare la ieșirile neuronilor, precum și la ponderi au fost restructurate astfel încât accesul la memorie să fie cât mai grupate. Pentru creșterea vitezei acceselor la memorie, toți vectorii au fost alocați la adrese de memorie aliniate la 64 octeți [5]. S-a încercat eliminarea indirectărilor inutile prin precalcularea valorilor

```

    for(unsigned l= ;l<nLayers- ;l++)
    {
        double sum = ;
        double *output, *pondere;
        unsigned2 *conex;

        unsigned nr_conexiuni,nr_neuroni;
        unsigned int idx= , indice= ;

10         nr_neuroni=L[l].nrn;
11         output=L[l- ].output;
12         conex = L[l].c;
13         pondere = L[l].p;
14
15         unsigned *pc=(unsigned*)conex;
16         for(unsigned i= ;i<nr_neuroni;i++)
17         {
18             if(pc[ ] ==          ) // adun bara la suma
19                 sum = pondere[pc[ ]];
20             nr_conexiuni=L[l].nrc[i];
21             pc+= ;
22             nr_conexiuni--;
23             for(unsigned j= ;j<nr_conexiuni;j++)
24             {
25                 sum += output[pc[ ] ] * pondere[pc[ ]];
26                 pc+= ;
27             }
28             nr_conexiuni++;
29             L[l].outputSum[i]=sum;
30             L[l].output[i]=SIGMOID(sum);
31             idx+=nr_conexiuni;
32         }
    }

```

Figura 6.2: Codul sursă al algoritmului de propagare înainte pentru varianta C++.

necesare. După toate aceste optimizări timpul de rulare pentru stratul L3 a scăzut la 4.39 secunde. Aceasta înseamnă o creștere a vitezei de 2,65 ori față de varianta originală.

Ca observație, timpul necesar pentru întreg algoritmul de propagare înainte pentru toate cele 10000 de imagini ale setului de test din MNIST a scăzut de la 27 de secunde la 9.97 secunde.

### 6.1.3 Varianta SSE2

Căutând noi metode pentru a crește viteza algoritmului, am recurs la limbajul de asamblare. Mai mult de atât, întrucât procesoarele pe care le foloseam, Intel Core2 Duo, dispuneau și de seturi de instrucțiuni paralele, am scris codul folosind aceste instrucțiuni.

Procesoarele x86 au începând cu Intel Pentium MMX și AMD K6 instrucțiuni SIMD (Single Instruction Multiple Data). Dacă prima extensie adăugată, MMX, a fost pentru numere întregi,

au urmat apoi seturi de instrucțiuni pentru calcule cu numere reale: SSE, SSE2, SSE3, SSE4, SSE4.1, SSE4.2, 3DNow! și variantele sale. Detalii despre aceste seturi de instrucțiuni paralele pot fi găsite în [25].

```

double dotproduct(double *a, double *b, int count)
{
    //xmm0 - registru SSE pe 128 de biti (pot lucra cu 4 float sau 2 double sau...)
    //pentru a stoca rezultatul (pointerul "pr" e necesar in __asm() deoarece pe "r" nu il accepta, daci e pointer si ci)
    double r[ ], *pr=r;
    __asm(
        //definim registrele principale (se-er putea sa fie inutile)
        pusha
        //incaram ecx cu valoarea contorului
        mov     ecx, count
        //incaram adresele vectorilor
        mov     eax, a
        mov     ebx, b
        //xmm0 = 0 (pt ca A sau b = 0), il vom folosi pe post de acumulator al rezultatelor (doi acumulatori!)
        pxor   xmm0, xmm0 //paralel XOR
    repeta:
        //incaram in xmm1 doua valori double din tabloul a[]
        movaps  xmm1, XMMWORD ptr [eax] //transfer din memorie aliniata
        //incaram in xmm2 doua valori double din tabloul b[]
        movaps  xmm2, XMMWORD ptr [ebx]
        //inmultim elementele corespunzatoare din xmm1 si xmm2 (doua inmultiri simultane)
        mulpd  xmm1, xmm2 //inmultire paralela
        //adunam rezultatele inmultirii anterioare
        haddpd  xmm0, xmm1 //adunare orizontala
        //salvăm peste doua valori double in a[]
        add     eax, 8
        //salvăm peste doua valori double in b[]
        add     ebx, 8
        //decrementăm contorul (loopnz il mai decrementează o data!)
        dec     ecx
        //repetăm bucla
        loopnz repeta
        mov     eax, pr
        //salvăm rezultatul xmm0 in vectorul r
        movups  [eax], xmm0 //transfer la adresa realiniata
        //restaurăm registrele principale
        popa
    )
    //returnăm rezultatul final
    return r[ ]+r[ ];
}

```

Figura 6.3: Codul sursă al algoritmului de propagare înainte pentru varianta C++ cu SSE II.

Procesoarele pe care le-am avut la dispoziție, Intel Core2 Duo 6600 și Intel Core2 Quad 6600 suportă maxim SSE3, astfel că a trebuit să limitez optimizarea codului la acest set de instrucțiuni. Modelele noi de procesoare Core2 de la Intel, cele în tehnologie de 45nm, suportă setul de instrucțiuni SSE4 care are o instrucțiune foarte importantă, MAD (Multiply and Accumulate), foarte utilă pentru calculul produsului scalar. Dacă ar fi fost disponibil un astfel de procesor, probabil că viteza ar fi fost și mai mare, dar nu cu mult deoarece așa cum

se poate vedea din Figura 6.3 în bucla de calcul a produsului scalar există și alte instrucțiuni care nu mai pot fi optimizate.

Deoarece a scrie cod în limbaj de asamblare nu e un lucru facil și durează destul de mult, am încercat să evit scrierea în asamblor a întregului algoritm de calcul a stratului L3, scriind o rutină care calculează produsul scalar a doi vectori de mari dimensiuni. Întrucât stratul L3 se bazează pe calculul a 100 de produse scalare de vectori de 1251 de dimensiuni, acest mod de evaluare a performanțelor este adecvat și corect. Figura 6.3 prezintă codul optimizat pentru instrucțiuni SSE3. Pentru documentare am folosit [52]. Codul sursă este comentat în Figura 6.3, astfel că nu va mai fi explicat din nou. Am implementat aceeași rutină pentru produsul scalar și în C++, activând toate optimizările pentru viteză ale compilatorului VS C++ 2008, inclusiv folosirea de instrucțiuni SSE3. În urma măsurării timpilor de execuție am constatat că varianta scrisă în limbaj de asamblare este de 1.3 ori mai rapidă decât cea scrisă în C++. Extrapolând acest rezultat la problema calculului stratului L3, s-ar obține 3.37s, adică o accelerare de 3.46 ori față de algoritmul inițial.

Trebuie menționat că am utilizat valori de tip double, astfel că instrucțiunile SSE au prelucrat câte două valori simultan. Puteam folosi și valori de tip float, caz în care calculele propriu-zise s-ar fi accelerat de două ori deoarece instrucțiunile SSE ar fi prelucrat câte patru valori float simultan. Pe ansamblu, însă, accelerarea nu ar fi fost mult mai mare deoarece algoritmul de calcul al produsului scalar este limitat de accesul la memorie.

Optimizările pe procesor ar putea continua cu scrierea de cod paralel, care să folosească mai multe fire de execuție. Nu am evaluat această posibilitate deoarece numărul de nuclee dintr-un procesor este de maxim 4, iar plăcile grafice au sute de nuclee de procesare, astfel că este imposibil ca orice optimizare pe procesor să fie mai rapidă decât cea pentru plăcile grafice.

## 6.2 Variante GPU ale optimizării algoritmului de propagare înainte

În ultimii ani viteza de procesare a plăcilor grafice a crescut extrem de mult față de viteza procesoarelor de uz general. Inițial plăcile grafice au fost proiectate pentru a efectua foarte multe operații grafice în paralel. Deși aveau o putere de calcul destul de mare, aceasta nu putea fi folosită decât foarte dificil - prin transformarea problemelor respective în operații grafice - și în alte domenii. Jocurile și grafica 3D sunt cele care vând cele mai multe plăci grafice performante, astfel că firmele producătoare au proiectat cipuri capabile de a efectua un număr din ce în ce mai mare de operații. Grafica 3D folosește intensiv operații cu numere reale, de obicei în simplă precizie. Operațiile sunt simple, de obicei adunări și înmulțiri, și mult mai rar complicate: radical, funcții trigonometrice etc. În timp, cipurile grafice s-au transformat în procesoare grafice (GPU - Graphics Processing Unit), devenind mult mai flexibile, putând fi programate pentru a executa și alte operații decât cele dedicate graficii.

Pentru a putea efectua un număr atât de mare de operații există două metode, un procesor unic la o frecvență foarte mare sau multe procesoare la frecvențe moderate. Prima variantă este folosită de procesoarele de uz general, iar tactul acestora nu a depășit niciodată limita de 4GHz. Puterea consumată crește rapid cu creșterea frecvenței, tehnologia actuală nefiind capabilă de a oferi frecvențe mult mai mari la consumuri rezonabile de curent, de exemplu un maxim de 150W. A doua variantă s-a dovedit a fi mult mai practică, zeci, chiar

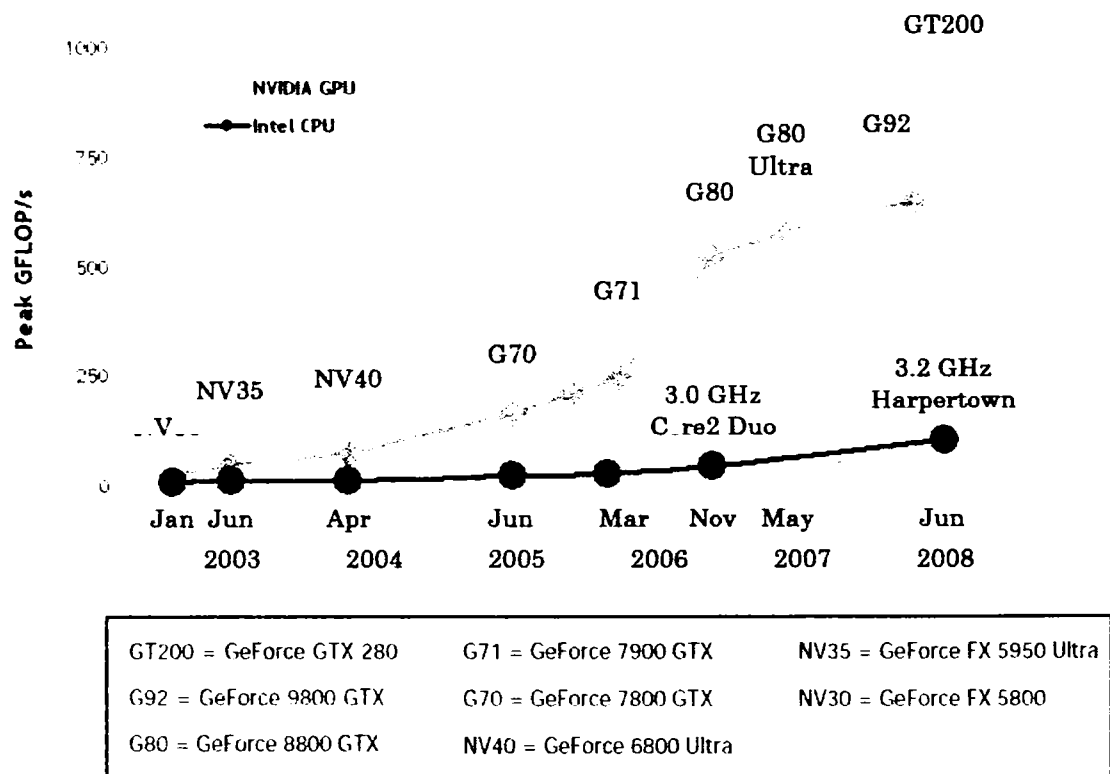


Figura 6.4: Evoluția performanțelor procesoarelor și a plăcilor grafice. Preluare după [54].

sute de procesoare simple ce rulează la frecvență redusă pot efectua împreună un număr impresionant de calcule în virgulă mobilă. În Figura 6.4 sunt prezentate vitezele în GFLOPS ale plăcilor grafice de la NVIDIA, precum și a procesoarelor de la Intel. Se poate observa foarte clar că pentru calcule masive plăcile grafice sunt înaintea procesoarelor cu câteva generații, ceea ce se traduce în 5-10 ani. Cealaltă firmă producătoare de plăci grafice, AMD-ATI, a scos pe piață în ultimele luni plăci grafice chiar și mai performante: ATI Radeon HD 4850 atinge 1 TFLOPS, ATI Radeon HD 4870 - 1.2TFLOPS, și cea mai nouă, ATI Radeon HD 4870X2, - 2.4TFLOPS. Un alt aspect important în care plăcile grafice excelează este lățimea de bandă la memoria plăcii grafice, care este cu aproximativ un ordin de mărime mai mare decât cea a procesoarelor la memoria RAM a sistemului. Trebuie avut în vedere că absolut toate cifrele prezentate în graficele anterioare reprezintă maxime teoretice - atât pentru GPU-uri, cât și pentru CPU-uri -, în practică performanțele obținute depinzând mult de problemă și de abilitatea programatorului. Însă, chiar și ineficient programate, GPU-urile sunt în marea majoritate a cazurilor de câteva ori mai rapide decât CPU-urile.

Avantajele și dezavantajele GPU față de CPU pot fi observate analizând Figura 6.6. Procesoare moderne dedică un spațiu foarte mare memoriei cache (60-90% din numărul de tranzistoare) și relativ puține tranzistoare nucleelor de procesare a datelor care sunt în număr de maxim patru (în viitorul apropiat firma Intel va scoate pe piață procesoare cu șase nuclee). În schimb, GPU-urile folosesc cele mai multe tranzistoare pentru unitățile de calcul și foarte puține pentru memoria cache. Sumarizând, CPU-urile au puține nuclee de execuție, dar care rulează la frecvențe ridicate, au un număr mare de instrucțiuni și sunt foarte flexibile,

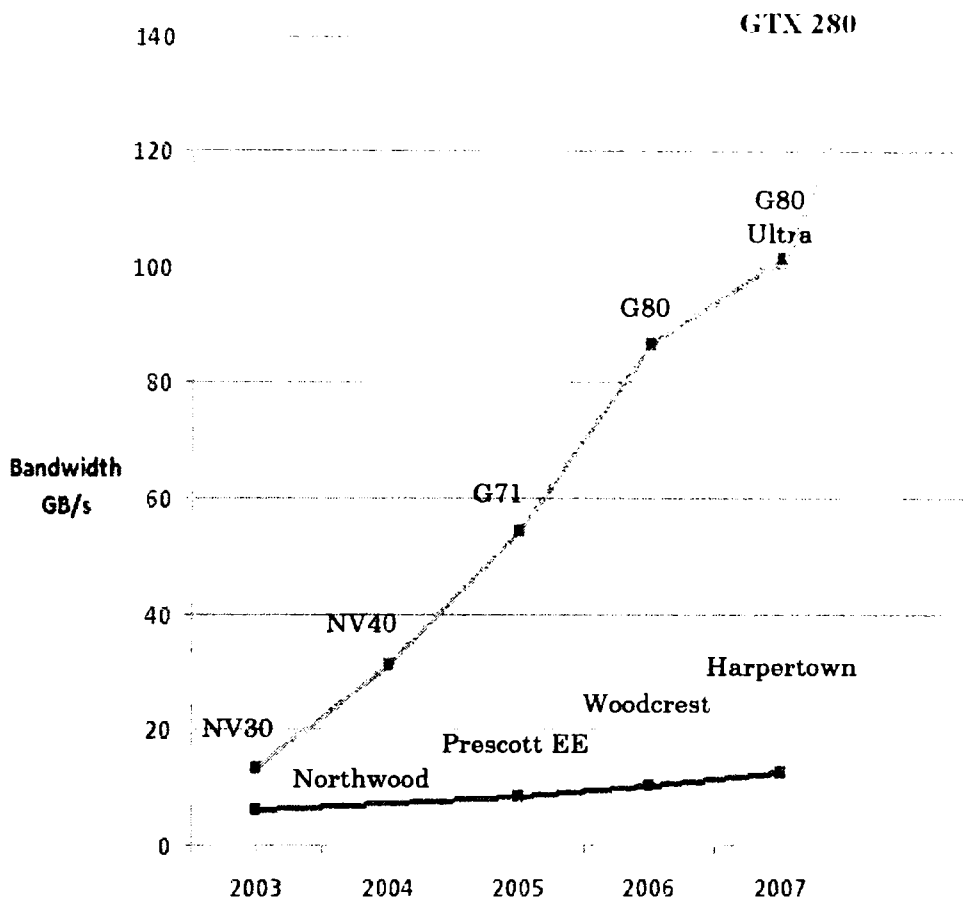


Figura 6.5: Comparație între evoluțiile lățimii de bandă pentru memoria sistem (RAM) și pentru memoria globală a plăcilor grafice. Adaptare după [54].

au o memorie cache considerabilă -ultimele procesoare de la Intel au 12MB cache L2- ce le oferă acces rapid la date aflate la adrese aleatoare de memorie. GPU-urile au foarte multe nuclee de execuție - ATI 4870X2 are  $2 \times 800$  nuclee dispuse în două cipuri - care rulează la frecvențe scăzute, nucleele sunt simple și nu oferă flexibilitatea celor din CPU, memoria cache este foarte puțină, accesul la memoria globală fiind foarte rapid numai dacă se face la adrese consecutive de memorie.

O observație importantă este că atât pentru CPU, cât și pentru GPU, raportul dintre numărul de operații pe secundă și lățimea de bandă este aproximativ același, 7. Acest lucru arată că pentru a folosi la maxim puterea de calcul a unui CPU sau a unui GPU, accesul la memorie trebuie să fie mult mai redus decât numărul de operații asupra datelor. Din acest motiv se spune că problemele sunt limitate de memorie (memory bound) sau limitate de CPU/GPU (CPU/GPU bound).

Odată cu dezvoltarea GPU-urilor, acestea au început să fie denumite GPGPU (General-

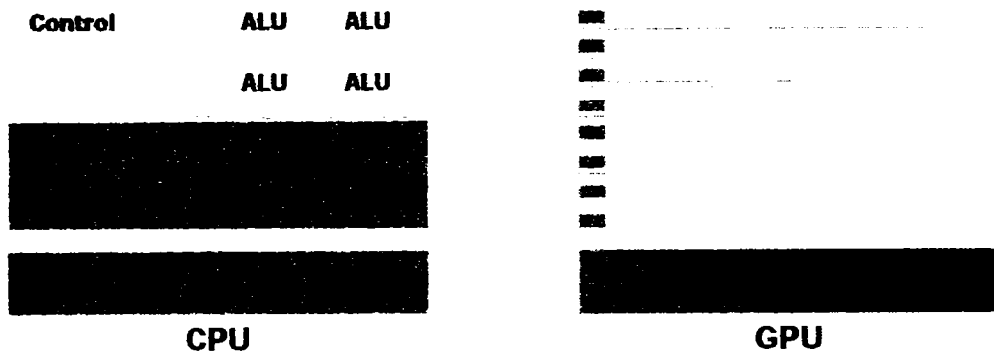


Figura 6.6: Comparație între alocarea tranzistoarelor pe suprafața de siliciu pentru GPU și CPU. Preluare după [54].

Purpose GPU) pentru a evidenția că au devenit mult mai flexibile și pot fi utilizate și la alte aplicații decât cele grafice. Există și o organizație și un site dedicat programării GPGPU-urilor. Primele utilizări majore au fost în supercalculatoare distribuite la nivel planetar și au fost folosite pentru studiul plierii proteinelor - Folding@home - [2], pentru detectarea stelelor neutronice rotitoare - Einstein@Home [1].

Programarea GPU-urilor se face folosind stream computing (programarea fluxurilor de date) și este un mod particular de paralelism care impune anumite limite felului în care sunt create, executate și sincronizate firele de execuție. Toate aceste particularizări fac ideal acest mod de paralelism pentru probleme care au următoarele proprietăți:

- densitate mare de calcule pentru fiecare element care este prelucrat,
- paralelism la nivelul datelor: aceeași rutină (kernel) poate fi aplicată pentru multe date diferite,
- datele sunt consistent localizate: informațiile necesare unei instanțe a kernel-ului se află distribuite în memorie într-un mod similar cu datele celorlalte instanțe ale kernel-ului. Aceasta constrângere a devenit din ce în ce mai laxă odată cu creșterea flexibilității GPU-urilor.

### 6.2.1 Varianta BROOK+

Limbajul BROOK+ [7] este derivat din limbajul BROOK dezvoltat la Stanford [66] pentru programarea fluxurilor de date (stream computing). BROOK+ este specific plăcilor grafice AMD/ATI și reprezintă un nivel de abstractizare peste nivelul limbajului de asamblare al procesoarelor grafice.

În secțiunea anterioară precizăm că performanțele ultimelor plăci de la ATI sunt impresionante, între 1TFLOPS și 2.4 TFLOPS, toate având același cip grafic, RV770 [6], diferind prin frecvența acestuia și prin tipul și frecvența memoriei utilizate (GDDR3 sau GDDR5). Placa pe care am avut-o la dispoziție este ATI Radeon HD 4850, având 800 de procesoare rulând la 625MHz, pentru o viteză maximă de  $625MHz \times 800processors \times 2FLOP/Hz/processor = 1TFLOPS$ . Deși este cea mai puțin performantă dintre plăcile cu GPU RV770 (seria 48xx),



```

kernel void compute_layer(double intr[], double idx_p[][], float idx_s[][],
                          unsigned nrc, out double stari[])
{
    float2 idx = { , };
    int i;
    double sum = ;
    float2 step = { , };

    idx.y = indexof(stari).x;
    sum = idx_p[idx.xy];
    idx.xy += step;
    for(i= ; i<nrc; i++)
    {
        sum += idx_p[idx.xy]*intr[idx_s[idx.xy]];
        idx.xy += step;
    }
    stari = SIGMOID(sum);
}

```

Figura 6.7: Codul sursă al algoritmului de propagare înainte pentru varianta BROOK+.

pentru prețul de 199USD oferă o viteză maximă de calcul care este mult superioară față de orice CPU existent.

Am început optimizarea algoritmului de propagare înainte pe GPU-uri folosind această placă și limbajul BROOK+. Firma AMD-ATI are o mare problemă în partea de software a stream computing. Dacă plăcile grafice la nivel hardware sunt excelente, softul (SDK-ul) lasă foarte mult de dorit. Este documentat foarte puțin, implementarea limbajului BROOK+ nu este finalizată, implicând multe limitări în modul în care trebuie scrise programele. Unicele surse de documentare găsite au fost [7] și [66]. Pe site-ul AMD există un forum dedicat stream computing, dar cei care se ocupă de dezvoltarea limbajului BROOK+ răspund, dacă o fac, cu întârziere destul de mare. Cum nu există alte surse de informații, iar a experimenta mai mult sau mai puțin haotic în încercarea de a afla cum funcționează unele aspecte ale limbajului BROOK+ nu este deloc eficient, timpul necesar învățării acestui limbaj pentru a-l putea aplica cu succes este îndelungat, iar rezultatele pot fi mediocre. La momentul când au apărut plăcile din seria 48xx (iunie-iulie 2008) SDK-ul era la versiunea 1.1. Profiler-ul nu funcționa pe plăcile din noua serie. Cel mai descurajant aspect este lipsa posibilității de depanare automatizată, neexistând un debugger! Astfel că orice depanare trebuie făcută terminând toate firele de execuție pentru că numai din codul scris pentru CPU se pot face afișări.

Un alt aspect negativ al limbajului BROOK+ este nivelul înalt al acestuia, nivel care nu permite accesul la nivel scăzut asupra hardware-ului. Având în vedere complexitatea unui program cu foarte multe fire de execuție, accesul la nivel scăzut este obligatoriu pentru a putea scrie programe performante. Compilatorul nu poate decât în puține cazuri să producă optimizările necesare. Există, într-adevăr, și un limbaj de asamblare pentru GPU, dar este la fel de slab documentat. Se va vedea în secțiunea referitoare la limbajul CUDA ce diferență există între BROOK+ și CUDA în privința documentației, a accesului la nivel scăzut al hardware-ului, ușurinței de a scrie și depana programe etc.

Cu toate aceste limitări am încercat scrierea de cod BROOK+ pentru algoritmul de propagare înainte. După îndelungate încercări nereușite, am ajuns la algoritmul din Figura 6.7.

Nu voi detalia modul de funcționare al limbajului BROOK+, însă câteva explicații asupra

codului din Figura 6.7 sunt necesare. Cuvântul cheie **kernel** specifică faptul că această rutină va fi executată în paralel pentru fiecare element al stream-ului de ieșire *stari*. Parantezele <> specifică natura de stream a parametrului *stari*. Se poate conceptualiza în felul următor: dacă stream-ul (vector-ul) are  $n$  elemente, atunci pentru fiecare element al său se va executa logic, în paralel câte o instanță a kernel-ului, pentru un total de  $n$  fire de execuție. Din considerente fizice, numărul de instanțe executate în paralel este limitat de numărul de procesoare și de alte aspecte ale arhitecturii procesorului grafic. Dacă un stream este precizat folosind parantezele <>, atunci în instanța curentă poate fi utilizat numai elementul corespunzător instanței. Dacă este necesar ca un kernel să folosească date de oriunde dintr-un parametru de tip stream, atunci se folosesc parantezele drepte pentru respectivul stream. Recapitulând, <> înseamnă acces numai asupra elementului curent, iar [] - acces la orice element al stream-ului. Este bine a se evita accesul cu paranteze [], deoarece accesele la locații aleatoare de memorie este foarte lent.

Funcția nu este optimizată pentru un strat anume al rețelei, putând fi utilizată pentru orice tip de rețea cu propagare înainte. Parametrul *intr* este un stream unidimensional care conține ieșirile neuronilor de pe stratul anterior, *idx\_p* este un stream bidimensional ce conține ponderile tuturor conexiunilor neuronilor de pe stratul curent, *idx\_s* este un stream bidimensional ce conține indici ai neuronilor stratului anterior, iar *nrc* este numărul de conexiuni ale fiecărui neuron, considerat în această implementare a fi același pentru toți neuronii unui strat.

Tipul **float2** împachetează două valori de tip float, accesibile cu *.x* și *.y*. Variabila *step* este inițializată cu valoarea {1,0}, deoarece 1 reprezintă pasul necesar trecerii la conexiunea următoare. Este de tip **float2** deoarece tabloul în care se folosește acest pas pentru incrementarea indicelui *idx* este bidimensional. Pe rândul 9 se selectează neuronul corespunzător instanței curente a kernel-ului. Funcția **indexof()** returnează indicele elementului care este prelucrat de instanța curentă a kernel-ului. Pe rândul 10 se încarcă bias-ul și apoi se trece la următoarea conexiune pe rândul următor. Buclele de la rândurile 12-16 calculează restul produsului scalar. La final se aplică funcția de activare asupra rezultatului produsului scalar și se stochează această valoare în elementul din stream-ul de ieșire, element corespunzător instanței kernel-ului.

Ca performanță, această rutină lasă foarte mult de dorit. Timpul necesar pentru cele 10000 de imagini din setul de test al MNIST a fost de peste 15 secunde pentru L3, deci de câteva ori mai mare decât varianta pe procesor. Acest lucru demonstrează cât de slab este compilatorul de BROOK+ pentru anumite probleme. În cazul de față nu este folosită nici măcar 0.5% din puterea procesorului grafic.

În data de 10 septembrie 2008 a apărut noua versiune a SDK-ului de la AMD-ATI. Versiunea 1.2 are o documentație îmbunătățită și limbajul BROOK+ a fost și el puțin dezvoltat. Având în vedere aceste aspecte și faptul că în ultimele luni am început să înțeleg mai bine (scriind programe în CUDA pentru GPU-uri NVIDIA - detalii în secțiunea următoare) modul în care trebuie scrise programele pentru a fi eficiente pe arhitecturi GPU, în viitor voi reevalua posibilitatea scrierii algoritmilor de învățare ai rețelelor neuronale artificiale în limbajul BROOK. Însă, atâta timp cât nu vor fi documentate detaliile arhitecturale ale RV770 și nu vor fi puse la dispoziția programatorului prin extensii BROOK+, este total improbabil că se va ajunge la performanțe apropiate de cele de care este capabil la nivel hard GPU-ul RV770.

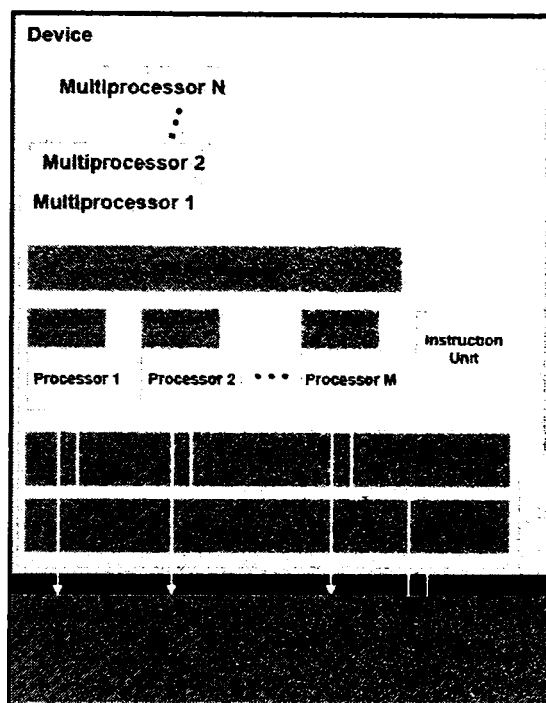


Figura 6.8: Arhitectura GPU-urilor de la NVIDIA (seriile GeForce 8 și 200). Preluare după [65].

## 6.2.2 Variantele CUDA

Celălalt mare producător de plăci grafice este NVIDIA. Este și cel mai mare, având o cotă de piață de peste 60%. În anul 2006 a introdus limbajul CUDA (Compute Unified Device Architecture) [55] special conceput pentru GP-GPU (General Purpose GPU). Acesta funcționează pe toate plăcile grafice din seriile GeForce 8 și GTX 200, în lume existând aproape 100 de milioane de astfel de mini supercomputere, cum mai sunt numite datorită puterii mari de calcul. Spre deosebire de limbajul BROOK+ de la AMD-ATI, CUDA este mult mai dezvoltat și permite mult mai facil accesul la resursele hardware ale GPU-urilor. De asemenea, există mod debug și depanarea se face aproape ca în orice compilator de C. În plus, există și un program de analiză a performanțelor programelor (profiler) scrise în CUDA.

De la NVIDIA am avut la dispoziție modelul de vârf, GTX 280. Acesta a apărut în luna iunie 2008. Este un procesor care conține foarte multe tranzistoare, peste 1.4 miliarde, gravate în tehnologie de 65nm. Frecvența procesorului este 602MHz, iar a ALU-urilor (unități aritmetico-logice) 1297MHz. Conține 30 de multiprocesoare fiecare având 8 procesoare, pentru un total de 240 de procesoare. Viteza maximă de procesare este  $1296MHz \times 240processors \times 3FLOP/Hz/processor = 933GFLOPS$ . Memoria folosită este de tipul GDDR3, rulând la 1107MHz și fiind conectată la GPU printr-un bus foarte lat, de 512 biți (8 canale a 64 de biți). Lățimea de bandă maximă este de  $1107MHz \times 64bytes/transfer \times 2transfers/Hz = 141GB/s$ . În configurație standard, memoria este de 1GB și este împărțită în 8 canale de 128MB fiecare.

În Figura 6.8 este descrisă la modul general arhitectura ultimelor două generații de pro-

cesoare grafice, cele din seriile GeForce 8 și GTX 200. În cazul lui GTX 280  $N=30$ , iar  $M=8$ . Fiecare multiprocesor are 16KBytes de memorie partajată ce poate fi accesată de toate cele 8 unități de calcul pe care le conține multiprocesorul.

Memoria partajată este foarte rapidă, apropiată de viteza regiștrilor, scrierile și citirile durând numai 4 impulsuri de tact. Un aspect important al memoriei partajate este acela că este împărțită în 8 bancuri de câte 1KBytes, accesul la aceste bancuri putându-se efectua simultan. Adresele bancurilor sunt întrepesute la nivel de float, astfel că bancul 0 are adresele 0, 32, 64 ..., bancul 1 - 1, 33, 65 ... etc. Nerespectarea accesării în acest mod a memoriei duce la conflicte de banc, care generează întârzieri în execuție. GTX 280 are o facilitate în plus, aceea de a permite tuturor celor 8 procesoare dintr-un multiprocesor să citească simultan aceeași locație de memorie partajată. În literatură această memorie mai este denumită și memorie cache gestionată de programator (user managed cache). Având în vedere că memoria globală - 1GB GDDR3 - este mult mai lentă decât memoria partajată, unul din modurile în care se pot accelera calculele este acela de a aduce în memoria partajată datele care sunt accesate de mai multe ori.

Fiecare multiprocesor are propriul set de regiștri pe care îi distribuie celor opt procesoare pe care le conține. Numărul maxim de regiștri disponibili unui multiprocesor este 8192.

Oricare din cele 8 procesoare poate citi din memoria globală, sau din memoria pentru texturi sau memoria constantelor. În cazul ultimelor două tipuri de memorie, accesul se face prin intermediul unei memorii cache de 8KBytes.

Accesul la memoria globală nu sunt intermediat de nicio memorie cache. Din acest motiv, accesul la memoria globală sunt eficiente numai dacă se face la adrese consecutive. În plus, e bine ca prima adresă să fie aliniată la nivel de 64 de bytes, cât este dimensiunea magistralei, altfel efectuându-se cel puțin un acces suplimentar. Timpul de acces aleator la memoria globală este de 400-600ns. Accesul consecutiv se face apoi la câte un tact distanță. În cazul cel mai nefavorabil, când toate accesul sunt aleatoare, neexistând localizare spațială sau temporală, lățimea de bandă a memoriei scade de la 141GB/s la aproximativ 8MB/s! Din acest motiv este obligatoriu ca datele să fie accesate astfel încât să existe localizare spațială și temporală. Profiler-ul care există în SDK analizează și aceste aspecte, oferind informații utile în păstrarea coeziunii acceselor la memorie.

În continuare va fi descris pe scurt cum funcționează un program CUDA. Fiecare program este împărțit în foarte multe fire de execuție (thread-uri) care apoi sunt rulate logic în paralel. De fapt, thread-urile sunt organizate într-un grid bidimensional care este format din mai multe blocuri de thread-uri (Figura 6.9). Fiecare bloc este executat de către un singur multiprocesor, toate procesoarele dintr-un multiprocesor executând aceeași instrucțiune la un anumit moment de timp. Pe un multiprocesor pot rula mai multe blocuri logic în paralel, dar fizic iterativ. Pentru detalii poate fi consultată documentația din SDK [54]. Detalii interesante pot fi găsite și în seria de articole din [26]. Există foarte multe articole și prezentări despre modul de funcționare al CUDA, precum și despre cum trebuie scrise programele pentru a rula cât mai optim din prisma vitezei, dar nu vor fi înșirate aici pentru a nu aglomera explicațiile. Întrucât limbajul CUDA este predat și în cadrul universităților, un mod bun de a găsi documentație este de a căuta cursurile afișate pe pagina [www.nvidia.com/object/cuda\\_university\\_courses.html](http://www.nvidia.com/object/cuda_university_courses.html).

Deoarece acest capitol are ca scop optimizarea performanțelor pe diverse arhitecturi, se vor descrie în continuare aspecte ce trebuie avute în vedere pentru a obține programe performante. Informațiile oferite reprezintă o sinteză a articolelor și prezentărilor de pe pagina [www.nvidia.com/object/cuda\\_university\\_courses.html](http://www.nvidia.com/object/cuda_university_courses.html). Există și alte aspecte mai puțin importante și nu atât de eficiente pentru procesul de optimizare, dar nu vor fi expuse deoarece nu

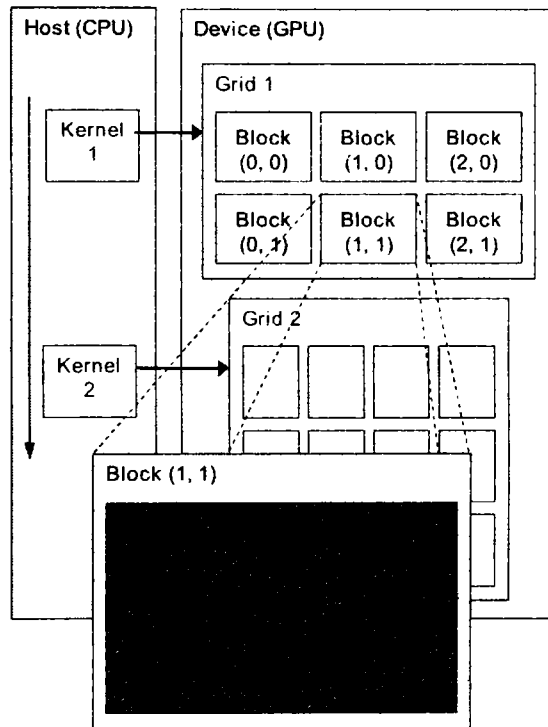


Figura 6.9: Organizarea thread-urilor. Adaptare după [56].

le-am folosit în rutinele ce urmează.

- eliminarea buclelor
  - se face cu directiva de compilare `#pragma unroll n`, unde `n` este o constantă care specifică numărul de iterații care vor fi efectuate
  - utilă doar dacă numărul de iterații este cunoscut aprioric, altfel trebuie efectuate alte artificii de programare
  - am obținut o reducere a timpilor de execuție cu aproximativ 40% pentru bucle mari ce calculau produse scalare (de exemplu în Figura 6.11 rândurile 198-200)
- memoria globală
  - accesul trebuie să fie localizat spațial și temporal
  - este bine a fi accesată cât mai puțin
  - accesul aleatoriu sunt de peste 100 de ori mai lente decât cele la memoria partajată
- memoria partajată
  - trebuie folosită de fiecare dată când se citește repetat aceleași locații de memorie
  - trebuie evitate conflictele de banc

- thread-urile trebuie sincronizate înainte de a citi din memoria partajată
- fiecare bloc este bine să folosească cât mai puțină memorie partajată pentru a putea fi rulate logic în paralel mai multe blocuri de thread-uri pe fiecare multiprocesor
- numărul de blocuri ale gridului
  - trebuie să fie suficient de mare pentru ca multiprocesoarele să aibă ce rula
  - cu cât sunt mai multe blocuri alocate unui multiprocesor, cu atât mai bine sunt ascunse latențele la memoria globală: când un bloc așteaptă după memoria globală, este pus în coada de așteptare și controlul este dat următorului bloc
- numărul de thread-uri din fiecare bloc e bine să fie multiplu de 32 pentru a fi trimise spre execuție cât mai eficient

### Optimizarea kernel-urilor pentru straturile rețelei neuronale

Din multitudinea de variante încercate vor fi prezentate două dintre ele, cele mai performante. Prima variantă este prezentată în figurile 6.10 și 6.11 și implementează întreg algoritmul de propagare înainte. În prima figură sunt implementate straturile de convoluție L1 și L2.

Stratul L1 conține 6 hărți de convoluție de dimensiunea  $13 \times 13$ . Este creat un grid cu 6 blocuri a 169 de thread-uri fiecare. Specificatorul **\_\_global\_\_** arată că funcția (kernel-ul) este apelată de către CPU, dar rulează pe GPU. Întrucât ieșirile stratului anterior vor fi folosite la calculul mai multor produse scalare, au fost încărcate în memoria partajată (declarația este la rândul 97, iar încărcarea la rândurile 102-104). Deoarece memoria partajată are 16 bancuri, pentru a nu genera conflicte de banc doar thread-urile corepunzătoare primilor 16 neuroni se ocupă de încărcarea datelor în memoria partajată. Variabila *btidx* conține indicii thread-ului în blocul curent, iar *idx* - indicii în întregul grid. La încărcarea ieșirilor stratului anterior se folosește *btidx* deoarece fiecare bloc trebuie să încarce datele întrucât memoria partajată este la nivel de bloc, iar ieșirile stratului anterior sunt necesare în toate blocurile.

Specificatorul **\_\_shared\_\_** precizează că respectiva variabilă este alocată în memoria partajată. Pentru acest strat, și numărul ponderilor este redus, existând numai 156 de ponderi, câte 26 pentru fiecare hartă de convoluție. Fiecare bloc trebuie să încarce (rândurile 106-108) numai 26 de ponderi deoarece execută calcule doar pentru una din hărți. Deoarece sunt numai 26 de ponderi, am lăsat un singur thread să se ocupe de încărcare, thread-ul 0 al fiecărui bloc. Apelul *\_\_syncthreads()* asigură că fiecare bloc are toate datele încărcate în memoria partajată. În acest moment poate începe calculul produselor scalare. La rândul 112 fiecare thread inițializează valoarea produsului scalar cu valoarea bias-ului, urmând ca pe rândurile 114-117 să fie adunate produsele dintre ponderi și ieșirile corespunzătoare ale stratului anterior. Cele două for-uri aplică calculul pe un nucleu de  $5 \times 5$  neuroni. La final (rândul 118) este calculată ieșirea neuronului curent prin aplicarea funcției de activare asupra valorii produsului scalar. Rezultatul este memorat în memoria globală.

Stratul L2 conține 50 de hărți de convoluție fiecare a  $5 \times 5$  neuroni. Fiecare neuron al fiecărei hărți este conectat cu câte 25 de neuroni din fiecare hartă de pe L1. Se crează un grid cu 50 de blocuri, fiecare bloc calculând valorile neuronilor dintr-o hartă, având, deci, 25 de thread-uri. Ca și la stratul precedent, ieșirile de pe stratul anterior, precum și ponderile dintre straturi sunt încărcate în memoria partajată. Avantajul straturilor de convoluție este că folosesc un număr restrâns de ponderi pe care le reutilizează. Rândurile 129-135 produc

încărcarea datelor în memoria partajată și asigură sincronizarea thread-urilor la final. Rândurile 137-145 calculează valorile neuronilor de pe stratul L2 într-un mod similar cu algoritmul stratului L1. Al treilea for, cel mai interior, selectează hărțile din stratul L1.

```

__global__ void computeConvolutionLayer1(float *prevLN, float *currLU, float *currLN)
{
    float sum= 0;
    __shared__ float output[ 11 ][ 11 ];
    __shared__ float w[ 5 ][ 5 ];
    unsigned int btdix=threadIdx.y* 11 + threadIdx.x;
    unsigned int idx = blockIdx.x* 11 * 11 + btdix;

    if(btdix< 11*11)
        for(int i=btdix;i< 11*11;i++)
            output[i/ 11][i% 11]=prevLN[i];

    if(threadIdx.x== 6&& threadIdx.y== 6)
        for(int i= 0;i< 11;i++)
            w[i]=currLU[blockIdx.x* 11+i];

    __syncthreads();

    sum = w[ 0 ][ 0 ];
    #pragma unroll 5
    for(int y= 0;y< 5;y++)
        #pragma unroll 5
        for(int x= 0;x< 5;x++)
            sum += output[ *threadIdx.y+y][ *threadIdx.x+x] * w[ +y* 5+x];
    currLN[idx] = SIGMOIDF(sum);
}

__global__ void computeConvolutionLayer2(float *prevLN, float *currLU, float *currLN)
{
    float sum= 0;
    __shared__ float output[ 11 ][ 11 ][ 11 ];
    __shared__ float w[ 5 ][ 5 ][ 5 ];
    unsigned int btdix=threadIdx.y* 11 + threadIdx.x;
    unsigned int idx = blockIdx.x* 11 * 11 + btdix;

    if(btdix< 11*11){
        for(int i=btdix;i< 11*11;i++)
            output[i/ 11][i% 11][i% 11]=prevLN[i];
        for(int i=btdix;i< 11*11;i++)
            w[i]=currLU[blockIdx.x*( 5* 5 + 1)+i];
    }
    __syncthreads();

    sum = w[ 0 ][ 0 ][ 0 ];
    #pragma unroll 5
    for(int y= 0;y< 5;y++)
        #pragma unroll 5
        for(int x= 0;x< 5;x++)
            #pragma unroll 5
            for(int m= 0;m< 5;m++)
                sum += output[ *threadIdx.y+y][ *threadIdx.x+x] * w[ +y* 5* 5 + x* 5 + m];
    currLN[idx] = SIGMOIDF(sum);
}

```

Figura 6.10: Kernel-urile pentru straturile de convoluție.

Stratul L3 (Figura 6.11) este total conectat cu stratul anterior. Fiecare din cei 100 de neuroni ai stratului este conectat cu fiecare din cei 1250 de neuroni de pe L2, deci fiecare neuron are 1251 de conexiuni, una fiind pentru bias. Există  $100 \times 1251 = 125100$  conexiuni și tot atâtea ponderi. La acest caz ponderile nu se refolesc, ci sunt utilizate o singură dată. Oricum, sunt mult prea multe pentru a putea fi încărcate în memoria partajată. Din aceste motive ele vor rămâne în memoria globală și vor fi citite numai când este nevoie de ele. În schimb, cele 1250 de ieșiri ale neuronilor de pe stratul L2 sunt utilizate fiecare de circa 100 de ori. Fiind puține și des accesate este foarte indicat a fi încărcate în memoria partajată. La

execuție este trimis un grid cu un singur bloc a 100 de thread-uri.

Rândurile 185-189 din Figura 6.11 declară variabilele locale necesare și calculează indicele thread-urilor la nivel de grid. Urmează (rândurile 191-194) încărcarea ieșirilor neuronilor de pe stratul L2 în memoria partajată a blocului, finalizată cu sincronizarea thread-urilor. Este obligatorie această sincronizare deoarece numai primele 16 thread-uri ale blocului de 100 de thread-uri participă la încărcarea memoriei partajate. Fără sincronizare thread-urile ce nu efectuează încărcarea ar trece la execuția instrucțiunilor următoare care citesc date din memoria partajată care încă nu ar fi într-o stare coerentă, procesul de încărcare nefiind finalizat.

```

185 __global__ void computeFullConnectedLayer3(float *prevLN, float *currLW, float *currLN)
186 {
187     unsigned int ic= ;
188     float sum= ;
189     __shared__ float output[ ];
190     unsigned int tid = threadIdx.x;
191     unsigned int idx = blockIdx.x*MAX_THREADS + threadIdx.x;
192
193     //loading previous layer outputs in shared memory
194     for(int i=tid;i< ;i+= )
195         output[i]=prevLN[i];
196     __syncthreads();
197
198     //compute output
199     sum = currLW[idx];
200     #pragma unroll 1000
201     for (ic= ; ic<= ; ic++)
202         sum += output[ic- ] * currLW[ic* +idx];
203     currLN[idx] = SIGMOIDF(sum);
204 }
205
206 __global__ void computeFullConnectedLayer4(float *prevLN, float *currLW, float *currLN)
207 {
208     float sum= ;
209     __shared__ float output[ ];
210     unsigned int idx = threadIdx.x;
211
212     //loading previous layer outputs in shared memory
213     for(int i=idx;i< ;i+= )
214         output[i]=prevLN[i];
215     __syncthreads();
216
217     //compute output
218     sum = currLW[idx];
219     #pragma unroll 1000
220     for (int ic= ; ic<= ; ic++)
221         sum += output[ic- ] * currLW[ic* +idx];
222     currLN[idx] = SIGMOIDF(sum);
223 }

```

Figura 6.11: Kernel-urile pentru straturile complet conectate.

În continuare (rândurile 196-201) se calculează produsul scalar corespunzător și se aplică funcția de activare asupra rezultatului. La final are loc scrierea ieșirii neuronului în memoria globală. Sunt necesare câteva precizări pentru operațiile de la rândul 200. Ca și la celelalte straturi, *output* se află în memoria partajată. În schimb, ponderile *currLW* se află în memoria globală care pentru a fi rapidă trebuie accesată începând de la adrese multiplu de 64bytes,



cât este lățimea magistralei, iar adresele locațiilor accesate este bine să prezinte localizare spațială și temporală. Acest lucru înseamnă că citirile pe care le efectuează thread-urile gridului trebuie să fie concentrate, nu răsirate aleator prin memorie, iar aceste citiri (accesări, la modul general) trebuie să fie făcute cvasisimultan. Dacă sunt respectate aceste principii, atunci accesesele vor fi grupate de către GPU și executate în paralel. În mod normal, ponderile erau stocate în memorie în felul următor: cele 1251 de ponderi ale neuronului 0, urmate de cele 1251 de ponderi ale neuronului 1, și tot așa până la neuronul 99. Considerând că thread-urile rulează în paralel, la un moment dat, fiecare thread accesează ponderea cu același indice din lista de ponderi corespunzătoare lui. Aceasta înseamnă accese la adresele:  $1251 \times 4 \times 0$ ,  $1251 \times 4 \times 1$ ,  $1251 \times 4 \times 2$  etc (4 apare din cauză că  $\text{sizeof}(\text{float})=4$ ). Evident, nu este respectată regula localizării spațiale. Din acest motiv vectorul de ponderi a fost reordonat în felul următor: primele 100 de valori reprezintă cele 100 de bias-uri ale neuronilor, următoarele 100 de valori reprezintă ponderile spre neuronul 0 al stratului L2 și tot așa. În acest mod se respectă regula localizării spațiale. Numărul 128 apare în loc de valoarea 100 pe rândul 200 deoarece am dorit ca accesesele să înceapă la adrese multiplu de 64 bytes. În acest mod nu sunt făcute două accese în loc de unul în cazurile în care datele se întind în două blocuri diferite chiar dacă dimensiunea lor cumulată este de cel mult 64 bytes. Bineînțeles, tabelul de ponderi este creat astfel încât ponderile corespunzătoare unui neuron de pe stratul anterior să se aple la multiplii de 128. Aceasta înseamnă că tabelul conține câte 28 de elemente neutilizate pentru fiecare 100 de ponderi. Se "pierd" astfel  $28 \times 1251 \times \text{sizeof}(\text{float}) = 140112$  bytes, ceea ce este nesemnificativ având în vedere că memoria globală este de 1GB.

Pentru GTX 280 alinierea la multiplii de 128 nu aduce nicio creștere de viteză deoarece circuitele hard care accesează memoria globală au fost mult îmbunătățite față de seria GeForce 8 unde acestea erau mult mai simple. La seria GeForce 8 alinierea datelor la multiplii de 128 se manifestă într-o creștere observabilă a vitezei de execuție (10-40%), în funcție de problemă).

Stratul L4 nu va mai fi detaliat fiind foarte asemănător cu stratul L3. Codul din Figura 6.11 este suficient pentru înțelegerea funcționării sale.

Cu această variantă am obținut un timp de rulare pentru kernel-ul stratului L3 pentru cele 1000 de imagini ale setului de test MNIST de 1.60 secunde, de 2.75 ori mai rapid decât algoritmul optimizat pentru procesor. Este o accelerare semnificativă, dar totuși mică având în vedere diferența de performanță dintre placa grafică și CPU. Era evident că se poate mult mai bine din următoarele două motive:

- am testat algoritmul pe două plăci grafice diferite GTX 280 și Quadro NVS 290. Deși ultima are performanțe teoretice de peste 30 de ori mai slabe decât GTX 280, timpii obținuți cu ea au fost de numai 3-4 ori mai mari. Aceasta arată că nu foloseam decât puțin din potențialul plăcii GTX 280,
- întrucât problema este limitată de viteza memoriei (memory bound problem) am calculat pentru stratul L3 lățimea de bandă utilizată. Citirea ieșirilor stratului L2, a ponderilor stratului L3, precum și scrierea ieșirilor stratului L3 totalizează  $1250 + 1251 \times 100 + 100 = 126450$  float-uri, adică 505800 octeți citiți/scriși pentru o imagine (acesta ignorând octeții citiți inutil pentru păstrarea alinierii datelor). Pentru toate cele 10000 de imagini, înseamnă 505800000 octeți, adică 4.71 GB accesați în 1.6 secunde. Lățimea de bandă folosită este de 2.94 GB/s. Având în vedere că GTX 280 are o lățime de bandă de 141GB/s se poate observa cât de ineficient este utilizată. În schimb, 2.94GB/s este aproape de limita maximă a plăcii NVS 290, arătând încă o dată că modul în care am

scris kernel-ul nu este unul scalabil, nefiind capabil de a utiliza resursele plăcii GTX - 280.

```

1 void wrapperComputeFullConnectedLayer3(float *prevLN,
2                                         float *currLW, float *currLN, float *partial)
3 {
4     dim3 dimGrid(1, 1, 1); //30 de blocuri
5     dim3 dimBlock(1, 1, 1); //100 de thread-uri fiecare
6     computeFullConnectedLayer3A<<<dimGrid,dimBlock>>>
7         (prevLN,currLW,currLN,partial);
8
9     dimGrid=dim3(1, 1, 1); //34 blocuri
10    dimBlock=dim3(1, 1, 1); //25 de thread-uri
11    computeFullConnectedLayer3B<<<dimGrid,dimBlock>>>
12        (prevLN,currLW,currLN,partial);
13 }
14
15 __global__ void computeFullConnectedLayer3A(float *prevLN, float *currLW,
16                                             float *currLN, float *partialsum)
17 {
18     #define nrp 40 //numarul de ponderi dintr-o suma partiala
19     float sum=0; //suma partiala
20     __shared__ float output[nrp]; //rezultate stratului anterior
21     unsigned int tid = threadIdx.x; //index thread in bloc
22     unsigned int bid = blockIdx.x; //index bloc in grid
23     //incarcarea rezultatelor stratului anterior in memoria partajata
24     if (tid < nrp)
25         output[tid]=prevLN[bid*nrp+tid];
26     __syncthreads();
27     //calcularea sumelor partiale
28     #pragma unroll nrp
29     for (unsigned int ic=0; ic<nrp; ic++)
30         sum += output[ic] * currLW[nrp*bid+(ic+1)*nrp+tid];
31     partialsum[bid*nrp+tid]=sum; //scrierea sumei partiale
32 }
33
34 __global__ void computeFullConnectedLayer3B(float *prevLN, float *currLW,
35                                             float *currLN, float *partialsum)
36 {
37     unsigned int idx = blockIdx.x*nrp + threadIdx.x; //index thread-ului
38     float sum=currLW[idx]; //incarcarea bias-ului
39     #pragma unroll 30
40     for (unsigned int ic=0; ic<nrp; ic++) //adunarea sumelor partiale
41         sum += partialsum[ic*nrp+tid];
42     currLN[idx] = SIGMOIDF(sum);
43 }

```

Figura 6.12: Kernel-urile pentru stratul L3.

În a doua variantă am continuat optimizarea kernel-ului pentru stratul L3. Varianta anterioară era ineficientă deoarece pentru a nu încărca repetat în memoria partajată ieșirile neuronilor de pe stratul L2, am ales o configurație cu un singur bloc a 100 de thread-uri. Cum un bloc este trimis la execuție la un singur multiprocesor, se poate observa cât de ineficient am utilizat cele 30 de multiprocesoare din GTX 280. Ar trebui să avem cel puțin 30 de blocuri pentru a ocupa GPU-ul, asta dacă ignorăm latențele la memorie pentru ascunderea cărora se precizează că este bine să existe mai multe blocuri disponibile pentru fiecare multiprocesor,

în timp ce unul așteaptă după memoria globală, altul poate deveni activ și executa operații.

Soluția simplă ar fi fost împărțirea blocului de 100 de thread-uri în 34 de blocuri a 3 thread-uri fiecare, ultimul bloc conținând numai un thread. Sunt trei dezavantaje. Primul este numărul mic de thread-uri dintr-un bloc, numai 3, față de cel puțin 32 cât e recomandat. Apoi, fiecare thread are mult de calculat - un produs scalar între doi vectori a 1251 de elemente. La final, memoria partajată a fiecăruia din cele 34 de blocuri ar fi trebuit umplută cu date din memoria globală. Se observă că această soluție are ceva comun cu cea anterioară: numărul de 100 de thread-uri.

Soluția mai complicată rezolvă problema numărului mic de thread-uri prin divizarea calculului produselor scalare în mai multe părți. Din acest motiv avem nevoie de două kernel-uri, unul care calculează sumele parțiale, și unul care formează rezultatul final din aceste sume parțiale (Figura 6.12). Comunicația dintre kernel-uri se face prin memoria globală folosind parametrul *partialsum*. Deoarece dorim să ocupăm toate cele 30 de multiprocesoare ale GTX 280, vom crea 30 de blocuri. Avem 1250 de produse ce trebuie apoi adunate, deoarece pentru bias nu trebuie să efectuăm nicio înmulțire. Numărul 1250 nu se divide cu 30. Ar fi indicat să folosim o repartizare uniformă deoarece fiecare bloc ar avea exact același număr de thread-uri și nu ar trebui să mai punem instrucțiuni condiționale în kernel. Pentru a rezolva acest aspect, completăm cu neuroni cu ieșire nulă stratul L2. Avem nevoie de  $\lceil 1250 \text{ div } 30 \rceil \times 30 = 1260$  neuroni, deci mai adăugăm 10 neuroni fictivi pe stratul L2. Deci, vom avea 40 de blocuri a 100 de thread-uri fiecare (rândurile 3 și 4), iar fiecare thread va calcula suma parțială formată din  $\lceil 1250 \text{ div } 30 \rceil = 42$  de produse. Dimensiunea necesară pentru *partialsum* este de  $30 \times 100 \times \text{sizeof}(\text{float})$ , dar vom aloca  $30 \times 112 \times \text{sizeof}(\text{float})$  pentru a avea datele aliniate la 64 de octeți.

Kernel-ul 3A calculează sumele parțiale. Începe prin a încărca în memoria partajată cele 42 de ieșiri de pe stratul anterior (rândurile 23-25). Numai primele 42 de thread-uri efectuează încărcarea, restul așteaptă la bariera de sincronizare (rândul 25). Urmează calculul sumelor parțiale și scrierea rezultatelor în *partialsum*.

Kernel-ul 3B folosește sumele parțiale calculate de 3A pentru a obține rezultatul final. În urma mai multor teste, am găsit drept configurație optimă 4 blocuri a 25 de thread-uri fiecare (rândurile 8-10). La rândul 37 se inițializează produsul scalar cu valoarea de prag (bias), iar la rândurile 39 și 40 se adună cele 30 de sume parțiale, urmând ca la rândul 41 să se aplice funcția de activare și să se memoreze rezultatul în variabila de ieșire.

Acest mod de a scrie algoritmul de calcul pentru L3 a dus la scăderea timpului de rulare pentru stratul L3 până la 0.2 secunde, ceea ce este de 21.95 ori mai rapid ca varianta C++ optimizată. Folosind calcule similare cu cele de la prima variantă scrisă în CUDA, am obținut pentru lățimea de bandă utilizată de acest strat valoarea de 38 GB/s.

## 6.3 Concluzii

În capitolul dedicat experimentelor am ajuns la concluzia că unul din modurile în care rata de recunoaștere poate fi îmbunătățită este antrenarea clasificatorilor cu imagini mai mari. Acest lucru implică rețele neuronale mai mari, antrenarea cărora va dura mult mai mult timp. Folosind programele existente deja, timpul de antrenare s-ar întinde pe mai multe săptămâni. De aceea devine necesară căutarea unor soluții mai rapide.

Pentru că scrierea tuturor programelor pentru fiecare tip de optimizare ar fi durat multe luni de zile, am ales un mod mai simplu de a testa diferitele moduri de optimizare. Mai exact, am

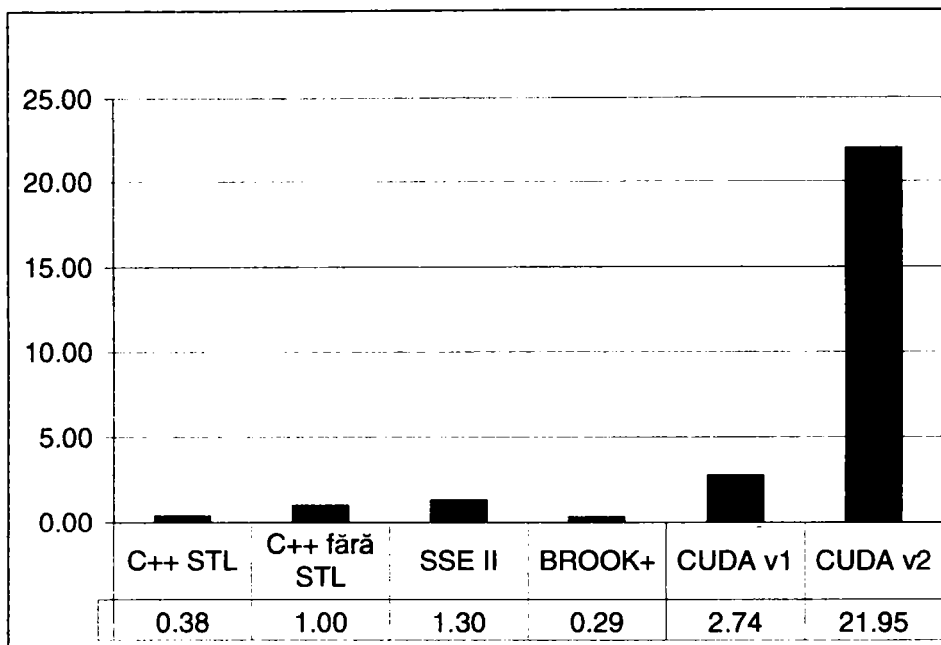


Figura 6.13: Gradul de accelerare a diferitelor variante de kernel pentru stratul L3. Referința este considerată varianta C++ optimizată, fără vectori STL.

ales o parte reprezentativă a algoritmului rețelelor neuronale și am implementat-o în diverse moduri. Pentru experimentele de optimizare am ales stratul L3 al 1RNC (rețeaua neuronală la nivel de cifră). Acest strat este cel mai mare din ambii clasificatori folosiți.

Am efectuat două tipuri de optimizări folosind arhitectura x86 și apoi folosind GPU-urile. Prima variantă optimizată folosind numai procesorul a implicat rescrierea în C++ a stratului L3 și aranjarea datelor de intrare astfel încât să fie cât mai ușor accesibile. S-a obținut o îmbunătățire semnificativă față de algoritmul inițial, de 2.63 ori. Această versiune de optimizare a fost folosită drept referință pentru celelalte tipuri de optimizări. Deși ar fi fost avantajos din prisma mărimii cifrelor, nu am folosit drept referință varianta inițială deoarece este un algoritm open source care este scris într-un mod didactic, nu având ca principiu un timp cât mai mic de rulare. Al doilea motiv pentru care am ales algoritmul C++ optimizat pentru a fi referință este că doream să compar arhitecturi diferite și nu este normal a compara o variantă neoptimizată scrisă pe o arhitectură cu una optimizată scrisă pe cealaltă arhitectură.

Pe arhitectura x86 am mai scris și o variantă în limbaj de asamblare ce folosește instrucțiuni SIMD SSE3. Aceasta s-a dovedit a fi mai rapidă ca varianta C++ optimizată, dar nu am folosit-o ca referință deoarece nu intenționăm să scriu în limbaj de asamblare variantele pentru plăcile grafice și s-ar fi ajuns la comparații nefondate. În Figura 6.13 sunt prezentate

accelerările obținute cu diferitele variante implementate.

Celălalt tip de arhitectură folosit a fost reprezentat de plăci grafice, așa numite GPU-uri. Am efectuat experimente cu GPU-uri atât de la AMD-ATI, cât și de la NVIDIA. Varianta scrisă în BROOK+ pe placa AMD-ATI este foarte lentă deoarece nu este optimizată, documentația oferită de AMD-ATI fiind mediocră. De asemenea, BROOK+ nu oferă acces de nivel scăzut la arhitectura GPU, optimizările fiind lăsate exclusiv în seama compilatorului. Din aceste motive și din cauza inexistenței unui debugger, am renunțat la a mai optimiza varianta BROOK+.

GPU-ul de la NVIDIA este mult mai ușor de programat pentru a obține performanță. Toate deficiențele întâlnite la BROOK+ nu există la limbajul CUDA folosit pentru GPU-urile NVIDIA. Am implementat două variante de optimizare, cea mai bună atingând o creștere a vitezei de 21.95 ori. În cazul comparării timpului de rulare cu cel obținut de varianta inițială, open source, obținem o creștere de performanță și mai mare, de 58.35 ori.

Optimizările efectuate cu ajutorul GPU-urilor au arătat că gradul de accelerare poate fi și mai mare dacă rețeaua neuronală are dimensiuni mai mari. Accelerarea de 21.95 ori a fost limitată de dimensiunea mică a rețelei simulate.

Pentru a putea implementa eficient o rețea neuronală cu ajutorul GPU-urilor, este obligatorie cunoașterea și accesul prin software a arhitecturii interne a procesorului grafic. În caz contrar, performanțele obținute sunt comparabile (sau chiar mai slabe) cu cele obținute utilizând CPU-uri.

**Concluzia cea mai importantă a acestui capitol este că plăcile grafice reprezintă cea mai performantă soluție de optimizare a algoritmilor rețelelor neuronale ce poate fi aplicată. Algoritmii implementați pe GPU sunt de departe mult mai rapizi decât cei implementați pe CPU, o accelerare de peste 20 de ori permițând antrenarea unor rețele mult mai complexe.**

# Capitolul 7

## Concluzii

În acest capitol vor fi prezentate contribuțiile autorului, urmate de sumarizarea concluziilor tezei. La final vor fi expuse principalele direcții de continuare a cercetării.

### 7.1 Contribuții

- **Crearea unui clasificator pentru perechi de cifre parțial suprapuse** și demonstrarea eficienței acestuia [20]. Acest clasificator pune bazele construirii unui sistem de recunoaștere a șirurilor numerice cu număr nedefinit de cifre, unite sau nu.
- **Crearea unui nou algoritm de recunoaștere** [17] a șirurilor numerice scrise de mână. Metoda propusă evită folosirea segmentării prin utilizarea a doi clasificatori, unul la nivel de cifră și unul la nivel de două cifre parțial suprapuse. Ratele de recunoaștere obținute sunt foarte apropiate de cele mai bune existente în domeniu, cu toate că metoda nu folosește nici segmentare, nici antrenare cu exemple negative. Pentru evitarea antrenării cu exemple negative, au fost concepute trei metode simple și eficiente de selecție a clasificatorului potrivit conținutului imaginii. Analiza erorilor arată că în cazul utilizării antrenării cu exemple negative, sistemul de recunoaștere va ajunge la rate de recunoaștere comparabile, probabil chiar mai bune decât cele mai eficiente metode actuale, păstrând, însă, simplitatea dată de neutilizarea segmentării.
- **Accelerarea pasului de propagare înainte** a rețelelor neuronale de peste 20 de ori folosind arhitecturi GPU, cu perspectiva măririi acestui factor pentru rețele și mai mari. Punerea bazelor creării unui algoritm de antrenare pe GPU pentru rețele neuronale practic imposibil de antrenat pe CPU.
- Crearea unui metode ce generează automat perechi de cifre parțial suprapuse.
- Crearea seturilor de imagini cu perechi de cifre parțial suprapuse TDNS.
- Crearea seturilor de imagini cu cifre LNIIST.
- Crearea seturilor de imagini cu șiruri numerice 3DNS, 6DNS și 10DNS.

- Dezvoltarea unui algoritm de extragere automată a câmpurilor din formulare. Algoritmul poate fi adaptat pentru a extrage câmpuri rectangulare din orice tip de formular.
- Conceperea unui algoritm de clusterizare (concatenare) a imaginilor cu părți de cifre în imagini de cifre sau de perechi de cifre.

## 7.2 Sumarul tezei

Teza propune o metodă eficientă pentru recunoașterea șirurilor numerice formate dintr-un număr nedefinit de cifre izolate sau parțial suprapuse.

Întrucât orice metodă trebuie validată, prezentarea începe cu **construirea unor seturi de imagini pentru antrenarea și testarea clasificatorilor**.

Pentru antrenarea clasificatorului la nivel de cifră am creat seturile din LNIST (66214 imagini pentru setul de antrenare și 45398 de imagini pentru setul de test) prin extragerea de imagini din NIST SD 19. Acestea sunt construite identic cu cele folosite de autorii din [44, 45] pentru a permite o comparație cât mai riguroasă a rezultatelor.

Pentru antrenarea clasificatorului la nivel de două cifre parțial suprapuse (2RNC) am creat seturile de antrenare și de test TDNS. Din NIST SD 19 au fost extrase 20000 de imagini de cifre (10000 pentru antrenare și 10000 pentru test) din care s-au creat câte 200000 de imagini cu perechi de cifre parțial suprapuse atât pentru setul de antrenare, cât și pentru setul de test. Pentru formarea imaginilor cu perechi de cifre parțial suprapuse am dezvoltat un algoritm care unește cifrele în funcție de morfologia acestora.

Dacă imaginile setului LNIST au fost extrase pentru a fi folosite de către clasificatorul la nivel de cifră, pentru cazul general, cel al șirurilor numerice formate din mai multe cifre, au fost construite seturile de test 3DNS, 6DNS și 10DNS. Și în construirea acestor seturi am urmărit obținerea unor seturi cât mai apropiate de cele folosite de alți cercetători pentru evaluarea performanțelor metodelor de recunoaștere. Cele trei seturi sunt identice cu cele folosite în domeniu, în limita informațiilor furnizate de autorii care le-au utilizat anterior.

Pentru extragerea din NIST SD 19 a imaginilor necesare construirii seturilor 3DNS, 6DNS și 10DNS am elaborat un algoritm care folosește o variantă a transformatei Hough pentru detectarea liniilor din formularele NIST. După detecția acestor linii, sunt detectate dreptunghiurile ce încadrează câmpurile dorite. Apoi are loc extragerea propriu-zisă a imaginilor și se aplică un algoritm de înlăturare automată a posibilelor părți din dreptunghiul delimitator care au fost extrase în același timp cu imaginile șirurilor numerice. La final a fost aplicat un proces de inspectare și curățare manuală a imaginilor. Algoritmul de extragere automată a imaginilor din câmpurile formularelor a diminuat considerabil timpul necesar construirii pe cale manuală a celor trei seturi de test. De asemenea, algoritmul este destul de general, putând fi adaptat facil pentru alte tipuri de formulare.

După ce seturile pentru antrenare și testare au fost create, se trece la prezentarea **sistemului de recunoaștere**. Se începe cu antrenarea celor doi clasificatori. Ratele de eroare obținute sunt de 0.86% pentru o cifră, respectiv de 5.35% pentru două cifre, ceea ce arată că RNC reprezintă o bună alegere atât pentru clasificatorul dedicat cifrelor, cât și pentru clasificatorul dedicat perechilor de cifre parțial suprapuse.

Folosind cei doi clasificatori și trei metode distincte de selecție a lor în funcție de conținutul imaginii, am implementat sistemul de recunoaștere a șirurilor numerice.

Ratele de recunoaștere obținute sunt: 93.77% pentru șiruri de trei cifre (setul 3DNS), 92.52% pentru șiruri de șase cifre (setul 6DNS) și 92.82% pentru șiruri de zece cifre (setul 10DNS).

Deși RNC nu au fost antrenate cu exemple negative, rezultatele obținute pentru seturile de test 3DNS și 6DNS (pentru 10DNS nu există date de comparație în literatură) sunt mult mai bune decât cele obținute anterior în domeniu, cu rețele neuronale artificiale. Diferențele ajung până la câteva zeci de procente în cazul rețelelor antrenate fără exemple negative.

Utilizarea combinației de clasificatori s-a dovedit a fi eficientă, rata de recunoaștere crescând față de sistemul ce utilizează numai clasificatorul de cifre cu 4.47% pentru setul 3DNS, respectiv cu 0.61% pentru setul 6DNS.

Metoda propusă oferă rate de recunoaștere cu numai câteva procente (3.03% pentru 3DNS și 4.42% pentru 6DNS) mai mici decât cele mai bune rezultate obținute anterior [44]. Având în vedere că rezultatele sunt obținute fără a aplica segmentarea și fără antrenare cu exemple negative, ratele de recunoaștere pot fi considerate foarte bune.

Concluzia experimentelor efectuate este că unul din modurile în care rata de recunoaștere poate fi îmbunătățită este antrenarea clasificatorilor cu imagini mai mari. Acest lucru implică rețele neuronale mai mari, antrenarea cărora va dura mult mai mult timp. Folosind programele existente deja, timpul de antrenare s-ar întinde pe mai multe săptămâni. Din acest motiv devine necesară căutarea unor **soluții mai rapide**.

Pentru că scrierea tuturor programelor pentru fiecare tip de optimizare ar fi durat multe luni de zile, am ales un mod mai simplu de a testa diferitele moduri de optimizare. Mai exact, am ales o parte reprezentativă a algoritmului rețelelor neuronale și am implementat-o în diverse moduri. Pentru experimentele de optimizare am ales stratul L3 al 1RNC (rețeaua neuronală la nivel de cifră). Acest strat este cel mai mare din ambii clasificatori folosiți.

Am efectuat două tipuri de optimizări folosind arhitectura x86 și apoi folosind GPU-urile. Prima variantă optimizată folosind numai procesorul a implicat rescrierea în C++ a stratului L3 și aranjarea datelor de intrare astfel încât să fie cât mai ușor accesibile. S-a obținut o îmbunătățire semnificativă față de algoritmul inițial, de 2.63 ori. Această versiune de optimizare a fost folosită drept referință pentru celelalte tipuri de optimizări. Deși ar fi fost avantajos din prisma mărimii cifrelor, nu am folosit drept referință varianta inițială deoarece este un algoritm open source care este scris într-un mod didactic, nu având ca principiu un timp cât mai mic de rulare. Al doilea motiv pentru care am ales algoritmul C++ optimizat pentru a fi referință este că doream să compar arhitecturi diferite și nu este normal a compara o variantă neoptimizată scrisă pe o arhitectură cu una optimizată scrisă pe cealaltă arhitectură.

Pe arhitectura x86 am mai scris și o variantă în limbaj de asamblare ce folosește instrucțiuni SIMD SSE3. Aceasta s-a dovedit a fi mai rapidă ca varianta C++ optimizată, dar nu am folosit-o ca referință deoarece nu intenționez să scriu în limbaj de asamblare variantele pentru plăcile grafice și s-ar fi ajuns la comparații nefondate.

Celălalt tip de arhitectură folosit a fost reprezentat de plăci grafice, așa numitele GPU-uri. Am efectuat experimente cu GPU-uri atât de la AMD-ATI, cât și de la NVIDIA. Varianta scrisă în BROOK+ pe placa AMD-ATI este foarte lentă deoarece nu este optimizată, documentația oferită de AMD-ATI fiind mediocră. De asemenea, BROOK+ nu oferă acces de nivel scăzut la arhitectura GPU, optimizările fiind lăsate exclusiv în seama compilatorului. Din aceste motive



și din cauza inexistenței unui debugger, am renunțat la a mai optimiza varianta BROOK+.

GPU-ul de la NVIDIA este mult mai ușor de programat pentru a obține performanță. Toate deficiențele întâlnite la BROOK+ nu există la limbajul CUDA folosit pentru GPU-urile NVIDIA. Am implementat două variante de optimizare, cea mai bună atingând o creștere a vitezei de 21.95 ori. În cazul comparării timpului de rulare cu cel obținut de varianta inițială, open source, obținem o creștere de performanță și mai mare, de 58.35 ori.

Optimizările efectuate cu ajutorul GPU-urilor au arătat că gradul de accelerare poate fi și mai mare dacă rețeaua neuronală are dimensiuni mai mari. Accelerarea de 21.95 ori a fost limitată de dimensiunea mică a rețelei simulate.

Pentru a putea implementa eficient o rețea neuronală cu ajutorul GPU-urilor, este obligatorie cunoașterea și accesul prin software a arhitecturii interne a procesorului grafic. În caz contrar, performanțele obținute sunt comparabile (sau chiar mai slabe) cu cele obținute utilizând CPU-uri.

Observația cea mai importantă este că plăcile grafice reprezintă cea mai performantă soluție de optimizare a algoritmilor rețelelor neuronale. Algoritmii implementați pe GPU sunt de departe mult mai rapizi decât cei implementați pe CPU, o accelerare de peste 20 de ori permițând antrenarea unor rețele mult mai complexe.

## 7.3 Continuări posibile

Pentru mărirea ratei de recunoaștere trebuie îmbunătățiți clasificatorii:

- 1RNC:
  - adăugarea în setul de antrenare de exemplare ale cifrei 1 cu linie orizontală în partea de jos. În prezent setul de antrenare nu conține niciun astfel de exemplar și toate cifrele 1 scrise în acest mod sunt recunoscute ca fiind 2. Verificarea existenței a suficiente exemplare din fiecare stil de scriere a cifrelor [42] și completarea setului de antrenare, dacă este cazul
  - generarea sau colectarea de exemple negative
  - reantrenarea clasificatorului
- 2RNC:
  - mărirea dimensiunii imaginilor setului de antrenare de la  $21 \times 13$  pixeli la  $41 \times 29$  pixeli
  - rafinarea algoritmului de formare automată a perechilor de cifre parțial suprapuse prin implementarea cazurilor particulare  $\bar{1}\bar{X}$  și  $X\bar{1}$
  - studierea efectului aplicării deformărilor înaintea procesului de unire a cifrelor
  - generarea sau colectarea de exemple negative
  - reantrenarea clasificatorului

Reantrenarea clasificatorilor (sau cel puțin a unuia dintre ei) cu exemple negative este foarte importantă deoarece ar îmbunătăți drastic acuratețea alegerii clasificatorului corect. Din analiza manuală a erorilor rezultă că dacă totdeauna s-ar alege clasificatorul potrivit

imaginii ( $RNC_1$  pentru cifre, respectiv  $RNC_2$  pentru perechi de cifre), rata de recunoaștere ar deveni mai bună decât cea mai bună existentă în momentul de față în domeniu.

Pentru a putea antrena în timp rezonabil (câteva ore, maxim o zi, și nu săptămâni) clasificatorul 2RNC cu imagini cu dimensiunea de  $41 \times 29$  pixeli este necesară finalizarea implementării pe GPU a algoritmului de antrenare a rețelei neuronale.

Întrucât există, deși într-un număr destul de redus, și imagini care conțin mai mult de două cifre unite, trebuie dezvoltată o metodă de rezolvare a acestui caz.

Algoritmul de extragere a câmpurilor din formulare poate fi îmbunătățit pentru a diminua cazurile în care trebuie rezolvate manual. Acest aspect este important numai dacă sistemul de clasificare va trebui în viitor să extragă imagini de test din alte formulare. Bineînțeles, etapele de extragere a câmpurilor și cea de recunoaștere a informațiilor ar putea fi concatenate și introduse într-un sistem unitar care să proceseze automat formulare cu informații numerice.

Studierea implementării clasificatorilor cu SVM (Support Vector Machine) sau alte metode de clasificare. Sunt importante atât rata de recunoaștere, cât și gradul de paralelism al algoritmului clasificatorului în ideea implementării pe GPU.

O continuare interesantă și foarte utilă ar fi extinderea metodelor de clasificare și la literele scrise de mână. Întrucât literele sunt de aproximativ trei ori mai multe decât cifrele, clasificatorul la nivel de literă ar avea de trei ori mai multe ieșiri, devenind mai complex. Cel mai mult va crește gradul de complexitate al clasificatorului la nivel de două litere parțial suprapuse (sau scrise legat) deoarece numărul claselor va ajunge la 900. Din fericire, gradul acesta de complexitate este potrivit cu capacitatea de procesare a plăcilor grafice, și problema recunoașterii literelor și cuvintelor scrise de mână ar reuși să utilizeze mult mai eficient resursele care acum sunt doar parțial folosite din cauza dimensiunilor prea mici ale problemei recunoașterii cifrelor și a perechilor de cifre parțial suprapuse.

Multe din aspectele prezentate se alfă deja în faza de dezvoltare.

# Anexa 1 - Formatul fișierelor IDX

Acest format este utilizat pentru a stoca tablouri unidimensionale și multidimensionale. Tipul elementelor poate fi ales din cele mai uzuale tipuri numerice. Elementele dintr-un fișier trebuie să aibă toate același tip.

Fișierul este structurat în modul următor:

- cod,
- numărul de elemente din dimensiunea 0,
- numărul de elemente din dimensiunea 1,
- numărul de elemente din dimensiunea 2,
- ...
- numărul de elemente din dimensiunea N,
- date.

Codul este un întreg pe 4 octeți. Primii doi octeți sunt nuli. Al treilea octet codifică tipul de date memorate:

- 0x08: unsigned byte,
- 0x09: signed byte,
- 0x0B: short (2 bytes),
- 0x0C: int (4 bytes),
- 0x0D: float (4 bytes),
- 0x0E: double (8 bytes).

Al patrulea octet specifică numărul de dimensiuni ale tabloului.

Datele sunt stocate în forma tablourilor C, cu indicele ultimei dimensiuni schimbându-se cel mai des.

Toate numerele întregi sunt reprezentate pe patru octeți în format MSB (non Intel).

Pentru detalii se poate consulta [40].

# Anexa 2 - Formatul fișierelor MNIST

Este un caz particular al tipului de fișier idx descris în anexa precedentă. Toate numerele întregi sunt reprezentate în format MSB first (cel mai semnificativ octet este stocat la adresa cea mai mică). Cele patru fișiere din MNIST sunt structurate în felul următor:

- train-images-idx3-ubyte

offset	tip	valoare	descriere
0000	32 bit integer	0x00000803(2051)	tip date
0004	32 bit integer	60000	număr de imagini
0008	32 bit integer	28	număr de rânduri
0012	32 bit integer	28	număr de coloane
0016	unsigned byte	??	valoare pixel
0017	unsigned byte	??	valoare pixel
.....			
xxxx	unsigned byte	??	valoare pixel

- train-labels-idx1-ubyte

offset	tip	valoare	descriere
0000	32 bit integer	0x00000801(2049)	tip date
0004	32 bit integer	60000	număr de elemente
0008	unsigned byte	??	etichetă
0009	unsigned byte	??	etichetă
.....			
xxxx	unsigned byte	??	etichetă

- t10k-images-idx3-ubyte

offset	tip	valoare	descriere
0000	32 bit integer	0x00000803(2051)	tip date
0004	32 bit integer	10000	număr de imagini
0008	32 bit integer	28	număr de rânduri
0012	32 bit integer	28	număr de coloane
0016	unsigned byte	??	valoare pixel
0017	unsigned byte	??	valoare pixel
.....			
xxxx	unsigned byte	??	valoare pixel

- t10k-labels-idx1-ubyte

offset	tip	valoare	descriere
0000	32 bit integer	0x00000801(2049)	tip date
0004	32 bit integer	10000	număr de elemente
0008	unsigned byte	??	etichetă
0009	unsigned byte	??	etichetă
.....			
xxxx	unsigned byte	??	etichetă

Etichetele au valori între 0 și 9, iar pixelii între 0 (alb) și 255 (negru).

# Bibliografie

- [1] Einstein@home - distributed computing. <http://einstein.phys.uwm.edu/>.
- [2] Folding@home - distributed computing. <http://folding.stanford.edu/>.
- [3] Neuralnet recognition project on the codeproject page, [www.codeproject.com/kb/library/neuralnetrecognition.aspx](http://www.codeproject.com/kb/library/neuralnetrecognition.aspx).
- [4] Nist special database 19 - nist handprinted forms and characters database. [www.nist.gov/srd/nistsd19.htm](http://www.nist.gov/srd/nistsd19.htm).
- [5] MSDN Run-Time Library Reference: `_aligned_malloc`. [http://msdn.microsoft.com/en-us/library/8z34s9c6\(vs.80\).aspx](http://msdn.microsoft.com/en-us/library/8z34s9c6(vs.80).aspx). *Visual C++ Developer Center*.
- [6] AMD. Ati radeon hd 4800 series - gpu specifications. <http://ati.amd.com/products/Radeonhd4800/specs.html>.
- [7] AMD. Brook+. <http://ati.amd.com/technology/streamcomputing/AMD-Brookplus.pdf>.
- [8] Sven Behnke. *Hierarchical Neural Networks for Image Interpretation*, volume LNCS 2766. Springer, 2002.
- [9] Christopher M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 13th edition, 2005.
- [10] Christopher M. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [11] R. G. Casey and E. Lecolinet. A survey of methods and strategies in character segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(7):690–706, July 1996.
- [12] K Chellapilla, S. Puri, and Simard P. High performance convolutional neural networks for document processing. In *International Workshop on Frontiers in Handwriting Recognition*, 2006.
- [13] K Chellapilla, M. Shilman, and Simard P. Combining multiple classifiers for faster optical character recognition. In *Proc. of the 7th International Workshop on Document Analysis Systems*, volume LNCS 3872, pages 358–367, 2006.
- [14] C.E. Cheong, K. Ho-Yon, S. Jang-Won, and J.H. Kim. Handwritten numeral string recognition with stroke grouping. *Proc. of the Fifth International Conference on Document Analysis and Recognition (ICDAR'99)*, (1):745–748, September 1999.

- [15] Dan Ciresan. Analiza și prelucrarea imaginii în vederea procesului de recunoaștere cu rețele neuronale artificiale. Technical report, Referat 1, 2004.
- [16] Dan Ciresan. Image segmentation methods based on natural clustering algorithms, usable for unconstrained handwriting recognition. *Proc. of the 6th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing*, pages 131–140, 2004.
- [17] Dan Ciresan. Avoiding segmentation in multi-digit numeral string recognition by combining single and two-digit classifiers trained without negative examples. *Proc. of 10th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing*, pages 197–203, 2008.
- [18] Dan Ciresan and Cosmin Cernazanu. Linguistic modelling for automatic speech recognition in romanian language. *Proc. of the 7th International Conference on Development and Application Systems*, pages 490–497, 2004.
- [19] Dan Ciresan and Dan Pescaru. Using character moment based invariant features to improve off-line handwriting recognition. *Proc. of the 8th International Conference on Development and Application Systems*, pages 471–476, 2006.
- [20] Dan Ciresan and Dan Pescaru. Off-line recognition of handwritten numeral strings composed from two-digits partially overlapped using convolutional neural networks. *Proc. of IEEE 4th International Conference on Intelligent Computer Communication and Processing*, pages 53–60, 2008.
- [21] L.P. Cordella and M. Vento. Symbol recognition in documents: a collection of techniques? *International Journal on Document Analysis and Recognition*, 3(2):73–88, December 2000.
- [22] CEDAR DATABASE. <http://www.cedar.buffalo.edu/databases/cdrom1/>. *USPS Office of Advanced Technology Database of Handwritten Cities, States, ZIP Codes, Digits, and Alphabetic Characters*.
- [23] Richard O. Duda and Peter E. Hart. Use of the hough transformation to detect lines and curves in pictures. *Communications of the ACM*, 15:11–15, January 1972.
- [24] G. Dzuba, A. Filatov, and A. Volgunin. Handwritten zip code recognition. *Proc. of the 4th International Conference on Document Analysis and Recognition (ICDAR'97)*, 2:766–770, August 1997.
- [25] Streaming SIMD Extensions. [http://en.wikipedia.org/wiki/streaming\\_simd\\_extensions](http://en.wikipedia.org/wiki/streaming_simd_extensions). *Wikipedia*.
- [26] Rob Farber. Cuda, supercomputing for the masses. <http://www.ddj.com/architect/207200659>.
- [27] Jurgen Franke. Expanding the performance of polynomial classifiers by iterative learning. *Advances in Handwriting Recognition*, 34:378–386, 1999.
- [28] Kunihiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4):193–202, April 1980.

- [29] T. M. Ha, M. Zimmermann, and H. Bunke. Off-line handwritten numeral string recognition by combining segmentation-based and segmentation-free methods. *Journal of Pattern Recognition*, 31(3):257–272, 1998.
- [30] T.M. HA, J. Zimmermann, and H. Bunke. Off-line handwritten numeral string recognition by combining segmentation-based and segmentation-free methods. *Pattern Recognition*, 31(3):257–272, 1998.
- [31] K. Kim, Y. Chung, J. Kim, and C.Y. Suen. Recognition of unconstrained handwritten numeral strings using decision value generator. *Proc. Sixth Int'l Conf. Document Analysis and Recognition*, pages 14–17, 2001.
- [32] Kye Kyung Kim, Yun Koo Chung, and C.Y. Suen. Post-processing scheme for improving recognition performance of touching handwritten numeral strings. *Proc. of the 16th International Conference on Pattern Recognition*, 3:327–330, November 2002.
- [33] Kye Kyung Kim, Yun Koo Chung, and C.Y. Suen. Recognition of unconstrained handwritten numeral strings by composite segmentation method. *Proc. of the 15th International Conference on Pattern Recognition*, 2:594–597, September 2002.
- [34] G. Koch, Heutte L., and Paquet T. Automatic extraction of numerical sequences in handwritten incoming mail documents. *Pattern Recognition Letters*, 26(8):1118–1127, June 2005.
- [35] F. Bortolozzi L. S. Oliveira, R. Sabourin and C. Y. Suen. Impacts of verification on a numeral string recognition system. *Pattern Recognition Letters*, 27(7):1023–1031, April 2003.
- [36] F. Lauer, C.Y. Suen, and G. Bloch. A trainable feature extractor for handwritten digit recognition. *Journal of Pattern Recognition*, 40(6):1816–1824, June 2007.
- [37] Y. LeCun and Y. Bengio. Convolutional networks for images, speech, and time-series. In M. A. Arbib, editor, *The Handbook of Brain Theory and Neural Networks*. MIT Press, 1995.
- [38] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.
- [39] Y. LeCun, O. Matan, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, L. D. Jackel, and H. S. Baird. Handwritten zip code recognition with multilayer networks. In IAPR, editor, *Proc. of the International Conference on Pattern Recognition*, volume II, pages 35–40, Atlantic City, 1990. IEEE. invited paper.
- [40] Yann LeCun. The mnist database of handwritten digits <http://yann.lecun.com/exdb/mnist/>.
- [41] Kenneth Levenberg. A method for the solution of certain non-linear problems in least squares. *The Quarterly of Applied Mathematics*, (2):164–168, 1944.
- [42] Z. C. Li and C. Y. Suen. Distinctiveness and similarities of handwritten numerals. *Advances in Handwriting Recognition*, 34:387–396, 1999.



- [43] Cheng-Lin Liu and K. Marukawa. Handwritten numeral string recognition: character-level vs string-level classifier training. *Proc. of the 17th International Conference on Pattern Recognition*, 1:405–408, August 2004.
- [44] Cheng-Lin Liu, Hiroshi Sako, and Hiromichi Fujisawa. Effects of classifier structures and training regimes on integrated segmentation and recognition of handwritten numeral string. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 26(11):1395–1407, November 2004.
- [45] L.C. Liu, K. Nakashima, H. Sako, and H. Fujisawa. Integrated segmentation and recognition of handwritten numerals: Comparison of classification algorithms. *Proc. of the International Workshop on Frontiers in Handwritten Recognition (IWFHR'02)*, 8:303–308, 2002.
- [46] George F. Luger. *Artificial Intelligence - Structures and strategies for complex problem solving*. Addison-Wesley, fifth edition, 2005.
- [47] W. Ohyama M. Shi, T. Wakabayashi and F. Kimura. Mirror image learning for handwritten numeral recognition. *Proc. of the Second International Workshop on Machine Learning and Data Mining in Pattern Recognition*, LNCS 2123:239–248, 2001.
- [48] Ranzato Marc'Aurelio, Christopher Poultney, Sumit Chopra, and Yann LeCun. Efficient learning of sparse representations with an energy-based model. In J. Platt et al., editor, *Advances in Neural Information Processing Systems (NIPS 2006)*. MIT Press, 2006.
- [49] O. Matan, J. Bromley, C. Burges, J. Denker, L. Jackel, Y. LeCun, E Pednault, W. Satterfield, C Stenard, and T. Thompson. Reading handwritten digits: A zip code recognition system. *IEEE Computer*, 25(7):59–63, July 1992.
- [50] Ofer Matan, Christopher J. C. Burges, Yann LeCun, and John S. Denker. Multi-digit recognition using a space displacement neural network. In J. M. Moody, S. J. Hanson, and R. P. Lippman, editors, *Neural Information Processing Systems*, volume 4. Morgan Kaufmann Publishers, San Mateo, CA, 1992.
- [51] Michael Negnevitsky. *Artificial Intelligence - A Guide to Intelligent Systems*. Addison-Wesley, second edition, 2005.
- [52] Intel Software Network. Using streaming simd extensions 3 in algorithms with complex arithmetic. [http://cache-www.intel.com/cd/00/00/06/67/66715\\_66715.pdf](http://cache-www.intel.com/cd/00/00/06/67/66715_66715.pdf).
- [53] Daisuke Nishiwaki and Keiji Yamada. A new numeral string recognition method using character touching type verification. *Advances in Handwriting Recognition*, 34:416–425, 1999.
- [54] NVIDIA. Cuda compute unified device architecture - programming guide. [www.nvidia.com](http://www.nvidia.com).
- [55] NVIDIA. Cuda home. [http://www.nvidia.com/object/cuda\\_home.html#](http://www.nvidia.com/object/cuda_home.html#).
- [56] NVIDIA. Cuda overview. <http://www.isi.edu/~ddavis/GPU/Course/Slides/MonoChrome/02 - Introduction to CUDA.pdf>.

- [57] L.S. Oliveira and R. Sabourin. Support vector machine for handwritten numerical string recognition. *9th International Workshop on Frontiers in Handwriting Recognition (IWFHR-9)*, pages 39–44, October 2004.
- [58] L.S. Oliveira, R. Sabourin, F. Bortolozzi, and Suen C.Y. Automatic recognition of handwritten numerical strings: A recognition and verification strategy. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 24(11):1438–1454, November 2002.
- [59] S. Ouchtati, M. Bedda, and A. Lachouri. Segmentation and recognition of handwritten numeric chains. *Journal of Computer Science*, 3(4):242–248, 2007.
- [60] Marc'Aurelio Ranzato, Fu-Jie Huang, Y-Lan Boureau, and Yann LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *Proc. Computer Vision and Pattern Recognition Conference (CVPR'07)*. IEEE Press, 2007.
- [61] Stuart J. Russell and Peter Norvig. *Artificial Intelligence - A Modern Approach*. Prentice Hall, second edition, 2003.
- [62] Stuart C. Shapiro, editor. *Encyclopedia of artificial intelligence*, volume 1. John Wiley & Sons, 1990.
- [63] Z. Shi, N. Srihari, C. Y. Shin, and A. V. Ramanaprasad. A system for segmentation and recognition of totally unconstrained handwritten numeral strings. *Proc. of the 4th International Conference on Document Analysis and Recognition (ICDAR'97)*, 2:445–458, August 1997.
- [64] Patrice Y. Simard, Dave Steinkraus, and John Platt. Best practices for convolutional neural networks applied to visual document analysis. *Proc. of International Conference on Document Analysis and Recognition*, pages 958–962, 2003.
- [65] ASPLOS 2008 Tutorial. Gpu architecture and applications. <http://www.gpgpu.org/asplos2008/ASPLOS08-4-GPU-architecture.pdf>.
- [66] Stanford University. Brookgpu. <http://graphics.stanford.edu/projects/brookgpu/>.
- [67] Xian Wang, Venu Govindaraju, and Sargur Srihari. Holistic recognition of touching digits. *Advances in Handwriting Recognition*, 34:359–367, 1999.
- [68] R. Allen Wilkinson, Michael D. Garris, and Jon Geist. Machine-Assisted Human Classification of Segmented Characters for OCR Testing and Training. volume 1906. SPIE, San Jose, 1993.
- [69] J.F. Wang Y.K.Chen. Segmentation of single- or multiple-touching handwritten numeral string using background and foreground analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(11):1304–1317, November 2000.
- [70] Jie Zhou, Qiang Gan, Adam Krzyzak, and Ghing Y. Suen. Quantum neural network in recognition of handwritten numerals. *Advances in Handwriting Recognition*, 34:368–377, 1999.

# Titluri recent publicate în colecția „TEZE DE DOCTORAT” seria 10: Știința Calculatoarelor

---

1. **Rodica Țirtea** – *Contribuții la îmbunătățirea dependabilității și securității informației*, ISBN 978-973-625-422-2, (2007);
2. **Ionel Muscalagiu** – *Contribuții la implementarea, evaluarea și îmbunătățirea performanțelor tehnicilor de căutare asincrone în cadrul programării bazate pe constrângeri distribuite*, ISBN 978-973-625-592-2, (2007);
3. **Daniel Cioi** – *Contribuții la utilizarea realității virtuale în proiectarea asistată de calculator*, ISBN 978-973-625-613-4, (2008);
4. **Sorin Babii** – *Cercetări privind creșterea performanțelor rețelelor neuronale într-un mediu de calcul distribuit*, ISBN 978-973-625-559-5, (2008);
5. **Norbert Neidenbach** - *Das Service-Management eines IT-Outsourcing-Projektes durch ITIL-Best-Practices, IT-Outsourcing kostenoptimiert planen und steuern*, ISBN 978-973-625-660-8, (2008);
6. **Edwin Hans Wolf** - *Das Geschäftsmodell (Business model) MDS (Managed Desktop Support) im IT-Outsourcing, Leistungserbringung im Rahmen des MDS-Geschäftsmodells*, ISBN 978-973-625-661-5, (2008);
7. **Adrian Zafiu** – *Minimizarea sistemelor decizionale multivalente deterministe și nedeterministe*, ISBN 978-973-625-678-3, (2008);
8. **Daniel Iercan** – *Contributions to the Development of Real-Time Programming Techniques and Technologies*, ISBN 978-973-625-719-3, (2008);
9. **Laurenția Timar** – *Contribuții referitoare la configurarea optimală prin prisma performanță-fiabilitate a unor rețele de dispozitive de achiziția datelor cu aplicabilitate la excavatoarele cu cupe*, ISBN 978-973-625-775-9, (2008).



EDITURA POLITEHNICA