



Catedra de Structura sistemelor de calcul

Recunoasterea cifrelor scrise de mana folosind o retea artificiala neuronală

Dunca Lucian - Marian

Grupa 302219

Profesor Laborator: Vlad Ratiu

Data: 20.11.2020



Content

Content.....	2
1. Rezumat.....	3
2. Introducere	4
3. Fundamentare teoretica	6
4. Proiectare și implementare	8
4.1 Implementarea software a RAN	8
4.2 Implementarea hardware in limbajul VHDL	8
4.2.1. Implementarea modulului de inmultire	8
4.2.2. Implementarea modului functiei de activare reLU si normalizare	9
4.2.3. Implementarea primului strat	9
4.2.4. Implementarea stratului ascuns	9
4.2.5. Implementarea stratului de iesire	9
4.2.6 Implementarea modului principal	9
Rezultate experimentale	10
Testarea modului de inmultire matriceala	10
Testarea modului de normalizare.....	10
Testarea modulului primului strat	11
Testarea modului stratului ascuns.....	11
Testarea rețelei neurale	12
Concluzii.....	13
Bibliography.....	14



1. Rezumat

O rețea neuronală artificială este o colecție de unități numite neuroni, aceștia au ca model neuronii creierului uman. Printr-o rețea neuronală putem rezolva diferite probleme din lumea reală prin aplicarea în practică a modului în care creierul uman funcționează.

S-a folosit această tehnologie pentru a modela un program care va avea capacitatea să recunoască și să clasifice niște imagini cu cifre scrise de mână. Obiectivele principale ale proiectului sunt conturarea unei arhitecturi capabile de a rula pe o placă FPGA și de a prezice cifra din imagine.

Pentru a rezolva această problemă am utilizat mai întâi un program scris în limbajul Python care va antrena un model de recunoaștere a cifrelor dintr-o imagine, setul de date este cel oferit de MNIST. Pentru scrierea propriu-zisă a aplicației care va rula pe placă FPGA am folosit limbajul de descriere hardware VHDL. În acesta am implementat o rețea neuronală cu propagare înainte, deoarece antrenarea modelului se face software și nu este necesară implementarea propagării înapoi.

Rezultatele dezvoltării RAN au fost unele așteptate, precizia rețelei nu este una foarte mare datorită utilizării reprezentării în virgulă mobilă dar nu este nici foarte mică. Așadar această implementare este una fiabilă care duce la creșterea vitezei calculelor necesare în rețeaua neuronală.



2. Introducere

Proiectul urmarește implementarea unei rețele artificiale neurale (RAN) cu propagare înainte în limbajul “Very high speed integrated circuit hardware description language” (VHDL) pentru o placă “Field programmable gate array” (FPGA.) RAN este o tehnică puternică și flexibilă de învățare automată care simulează creierul uman în sensul în care procesează semnale de intrare și le transformă în semnale de ieșire [1]. Limbajul VHDL este un limbaj de descriere hardware folosit pentru a descrie sisteme digitale și mixte cum ar fi FPGA. Plăcuțele FPGA sunt niste circuite integrate foarte rapide care sunt create pentru a fi configurate de către client. Datorită acestei proprietăți aceste circuite integrate sunt des folosite în domeniul RAN.

RAN fac parte din domeniul învățării profunde (deep learning), care face parte la rândul ei dintr-un domeniu mai larg, cel al învățării automate (machine learning). Învățarea automată este o ramură a inteligenței artificiale, al cărui obiectiv este de a dezvolta tehnici care dau posibilitatea calculatoarelor de a învăța. Mai precis, se urmarește să creeze programe capabile de generalizare pe baza unor exemple. RAN sunt niste structuri de procesare a informației prin simularea artificială a fiziologiei și structurii creierului uman. Sistemul neural uman funcționează pe baza legăturilor între neuroni. Semnalul de intrare al unui neuron este compus din nenumărate semnale ale altor neuroni care sunt conectați la el, atunci când suma tuturor semnalelor depășește o anumită limită, neuronal va genera un semnal electric care se va propaga prin sinapse la alți neuroni.[2]

RAN sunt formate din unități elementare numite neuroni combinați după o anumită arhitectură. De exemplu, aceștia pot să fie aranjați în multiple straturi. Arhitecturile multistrat au următoarele componente: [2]

1. Stratul cu datele de intrare, format din n neuroni (un neuron pentru fiecare semnal de intrare)
2. Stratul ascuns, format din unul sau mai multe straturi formate din m neuroni.
3. Stratul de ieșire, format din p neuroni (un neuron pentru fiecare semnal de ieșire)

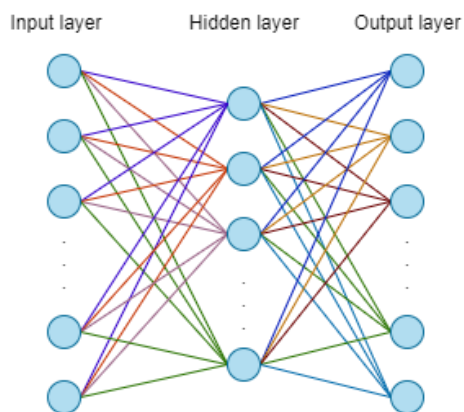


Fig 1. Straturi RAN



Proiectului are ca scop implementarea unei RAN de recunoastere a cifrelor scrise de mana, in limbajul de descriere hardware VHDL. Mai apoi rularea pe o placa de dezvoltare FPGA pentru generarea rezultatelor. Pentru realizarea acestuia se urmaresc urmatoarele obiective:

1. Dezvoltarea modelului RAN in python
2. Importarea setului de imagini pe care se va antrena modelul
3. Antrenarea modelului si generarea greutatilor fiecarui strat al RAN in python
4. Implementarea RAN in limbajul VHDL
5. Testarea RAN cu ajutorul greutatilor obtinute din modelul retelei din python
6. Rularea pe placa de dezvoltare BASYS3

Implementarea aleas pentru RAN in limbajul VHDL presupune multiple compromisuri. Antrenarea modelului unei RAN este un lucru de o dificultate ridicata pentru placutele FPGA deoarece acestea nu au capacitatea de a face calcule in virgula mobila. Asadar, se abordeaza implementarea unei RAN strict cu propagare inainte care nu presupune antrenarea modelului ci doar utilizarea unui model gata antrenat. Rolul acesteia va fi doar de a ne da rezultate in urma aplicarii unor imagini la intrarea retelei.

Pentru obtinerea modelului care va fi folosit pe placuta, se va folosi o implementare software in limbajul python care va crea si antrena modelul. Modelul va fi antrenat pe imagini care se vor importa din baza de date MNIST. Imaginile importate sunt sub forma unei matrici cu marimea de 28 pe 28 de pixeli. Pentru a reduce numarul de neuroni si greutati generate, modelul va fi antrenat cu imagini de 20 pe 20 de pixeli. Imaginile sunt reduse cu ajutorul unei interpolari cubice, astfel imaginea isi mentine forma. Fiecare pixel din imagine va avea o valoare cuprinsa intre 0 si 255 (scala gri). Se va seta fiecare pixel cu valoarea mai mare decat 0 pe 1 pentru a reduce complexitatea calculelor.

Deoarece greutatile generate vor fi reprezentate in virgula mobila si aceasta reprezentare nu este suportata de limbajul VHDL, ne vom folosi de reprezentare in punct fix, prin scalarea tuturor greutatilor cu 10000. Aceasta reprezentare insa va reduce din precizia rezultatelor.

In sectiunile urmatoare se vor prezenta diferite aspecte ale proiectului: fundamentarea teoretica se prezinta baza teoretica pe care se bazeaza proiectul precum si tehnologiile care pot sa fie utilizate, proiectare si implementare reprezinta partea principala a proiectului care contine descrierea fiecarei etape parcurse pentru realizarea obiectivelor proiectului, rezultate experimentale unde se demonstreaza ca proiectul implementat a fost implementat cu success si rezultatele obtinute sunt valide.



3. Fundamentare teoretică

Pentru a înțelege cum funcționează o RAN cu propagare înainte vom porni de la unitatea de bază a acesteia, perceptronul (Fig 2). Perceptronul este precursorul neuronului din RAN cu propagare înapoi, acesta este format doar dintr-un strat de intrare, unul de ieșire, o funcție de însumare și o funcție de activare.

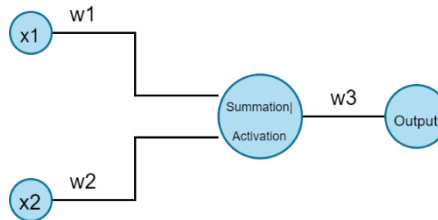


Fig 2. Perceptron

Imaginea de mai sus reprezintă un singur perceptron, dar dacă avem multe intrări și multe date, un singur perceptron nu este de ajuns. Se va mari numărul de perceptroni formându-se modelul de bază a unei RAN cu propagare înainte (Fig 3). Fiecare linie între perceptroni este reprezentată de o greutate, valoarea perceptronilor dintr-un strat poate să fie scrisă ca o înmulțire matriceală între o matrice a greutăților și o matrice coloană a intrărilor (Fig 4).

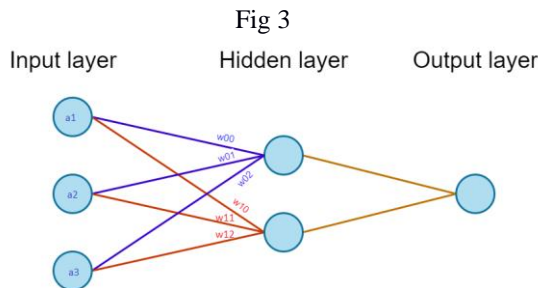


Fig 3

$$\begin{pmatrix} w^0 & w^1 & w^2 \\ w^{10} & w^{11} & w^{12} \end{pmatrix} \cdot \begin{pmatrix} a^1 \\ a^2 \\ a^3 \end{pmatrix}$$

Fig 4

Procesul de antrenare al unui model al unei RAN implică direct aceste greutăți. Prin antrenarea unei RAN se înțelege modificarea greutăților cu scopul ca predicția rețelei să fie aceeași cu predicția așteptată în etapa de antrenare, astfel reducându-se valoarea erorii din rețea. Aceasta “antrenare” se realizează prin mici modificări ale greutăților care sunt reprezentate în general de numere în virgula mobilă.

Pentru implementarea acestor greutăți în limbajul VHDL se va folosi reprezentarea în punct fix pentru a exprima numeral în virgula mobilă obținute prin antrenarea modelului RAN. Prin folosirea reprezentării în punct fix se va reduce precizia RAN și îi va crește complexitatea calculelor necesare datorită înmulțirii unor numere cu valori mari, însă este necesară abordarea acesteia deoarece reprezentarea în virgula mobilă nu este acceptată în limbajul de descriere VHDL fără a fi creat de către noi.

O altă componentă vitală într-o RAN este funcția de activare care are ca rol convertirea sumei ponderate a semnalelor de intrare a unui neuron într-un semnal de ieșire. Există nenumărate funcții de activare care se folosesc în funcție de ce avem nevoie. Cea mai simplă funcție de activare din punct de vedere al nevoii computaționale este funcția RELU : $f(x) = \max(0, x)$ (Fig 5). Alte funcții folosite sunt funcția sigmoidă, folosită frecvent deoarece simplifică calculele ulterioare prin maparea semnalelor de ieșire în intervalul (0,1). La fel și funcția tanh care mapează valorile în intervalul (-1,1). Funcțiile prezentate precedent sunt folosite mai mult în straturile ascunse datorită proprietăților lor. O altă



functie de activare importanta, folosita pentru obtinerea rezultatelor in ultimul strat este functia SoftMax. Aceasta functie este o normalizare exponentiala care mapeaza rezultatele intr-o distributie de probabilitati, avand suma componentelor 1.

Fig 5

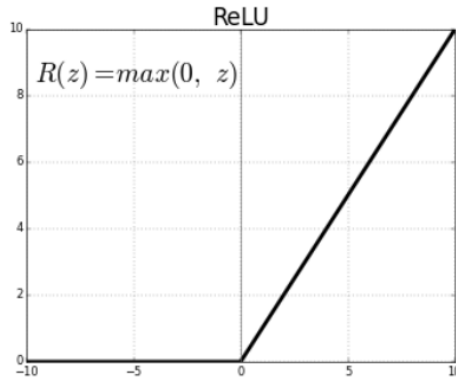
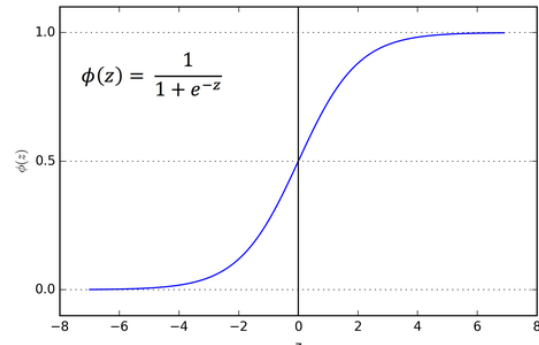


Fig 6



RAN sunt capabile de a invata si din acest motive le trebuie antrenate mai intai. Exista cateva metode de invatare:

- **Invatarea supravegheta:** implica un professor si RAN insasi, aceasta va ghici raspunsul dorit iar profesorul ii va spune ulterior raspunsul corect. Mai apoi retea va compara raspunsul sau cu cel oferit de catre profesor si va face ajustarile necesare in functie de eroarea existenta.
- **Invatarea nesupravegheta:** este necesara atunci cand nu exista un set de date cu raspuns cunoscut. In acest caz, se vor imparti datele in grupuri in functie de un criteriu bazat pe datele prezente.
- **Invatare consolidate:** acest tip de invatare se bazeaza pe observatie, RAN face o decizie prin observarea mediului. Daca decizia este falsa, retea isi va ajusta greutatea pentru a fi capabila sa faca o decizie diferita data viitoare.

Crearea unei RAN pentru recunoasterea cifrelor de mana nu este un subiect nou abordat. Exista nenumarate lucrari despre aceasta tema si multe abordari diferite [1][4]. O abordare posibila este folosind retele neurale convolutionale, abordare prezentata in lucrarile [5] si [6] care au una dintre cele mai bune performante in procesarea imaginilor. Aceasta abordare presupune introducerea unor straturi suplimentare numite straturi de convolutie si pooling care au rolul de a extrage caracteristici dintr-o imagine astfel crescand precizia si viteza.

Creare RAN poate avea o abordare software folosind librariile keras/tensorFlow in limbajul de programare python, se pot gasi numeroase articole pe aceasta tema. Pe de alta parte, contruirea unei RAN intr-un limbaj de descriere hardware nu este un subiect asa de des intalnit comparat cu varianta software. O analiza a metodei de implementare a unei RAN pentru o placa FPGA se poate gasi in articolul [7].



4. Proiectare și implementare

Pentru implementarea RAN în limbajul de descriere hardware VHDL s-a dezvoltat mai întâi modelul acestuia printr-o metodă software în limbajul python cu ajutorul framework-ului “keras”. A fost ales să se dezvolte o RAN cu propagare înapoi care are 3 straturi, o altă alternativă prezentată în ‘Fundamentarea teoretică’ ar fi fost o RAN convolutivă. S-a ales această metodă deoarece complexitatea imaginilor nu este ridicată, fiind posibilă procesarea acestora folosind o RAN simplă, fără tehnicile prezente în rețelele convolutive.

4.1 Implementarea software a RAN

Pentru crearea RAN s-a folosit framework-ul ‘keras’ care ne ajută să construim straturi cu un număr de neuroni variabili și în același timp să alegem funcția de activare care va fi aplicată pe fiecare strat în parte. Mai întâi s-au importat imaginile folosite pentru a antrena și testa rețeaua. Imaginile folosite sunt cele din baza de date MNIST, acestea au fost încărcate sub formă de matrici cu mărimea de 28 pe 28 de pixeli. Setul de date a fost împărțit în 2 părți, 60.000 de poze pentru antrenarea sistemului și 10.000 de poze pentru testarea acestuia.

Inițial s-a ales ca stratul de intrare să aibă 36 de neuroni, stratul ascuns un număr de 5 neuroni iar implicit stratul de ieșire un număr de 10 neuroni. Datorită reprezentării greutăților în reprezentarea în punct fix, rezultatele acestei rețele au fost scăzute. Asadar, s-a ales creșterea numărului de neuroni chiar dacă complexitatea calculului va fi mai ridicată. Prin urmare, RAN are în stratul de intrare un număr de 400 de neuroni (indicate de o imagine de 20 pe 20 de pixeli), în stratul ascuns un număr de 10 de neuroni, iar în stratul de ieșire tot 10 neuroni.

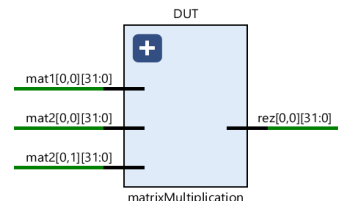
Imaginile importate din baza de date MNIST au o mărime de 28 pe 28 de pixeli. Pentru a obține imagini cu mărimea dorită vom folosi o interpolare cubică care ne va reduce imaginea la mărimea dorită, în același timp nedegradând imaginile. În continuare vom prelucra imaginile pentru a putea fi direct trimise către modelul RAN. În primul rând se vor transforma matricile de mărime 20 pe 20 într-o matrice cu o singură linie cu mărimea de 1 pe 400. Mai apoi se vor reprezenta variabile de ieșire cu ajutorul “one hot encoding” pentru a fi posibilă procesarea acestora de către RAN. Crearea straturilor se va realiza utilizând metode din framework-ul ‘keras’. Stratul ascuns va avea 10 neuroni și funcția de activare va fi ReLU. Stratul de ieșire va avea tot 10 neuroni (reprezentând cele 10 cifre posibile) și funcția de activare va fi o normalizare simplă. La final se va activa antrenarea modelului rețelei prin apelarea unei funcții din keras.

4.2 Implementarea hardware în limbajul VHDL

Pentru modelarea RAN în limbajul de descriere hardware VHDL este nevoie de a proiecta mai multe module cu scopuri diferite. Pentru a stoca informațiile într-un mod eficient și de a ne ajuta în momentul în care trebuie făcute calcule, se va declara o structură nouă sub formă de matrice. Această structură are o mărime generică care se poate schimba și va avea rolul de a memora toate datele RAN care sunt reprezentate de către semnalele de intrare, imagini, greutăți și semnale de ieșire intermediare. Se va folosi un pachet pentru a declara această structură de date împreună cu toate constantele aferente necesare reprezentate de către greutățile și imaginile RAN.

4.2.1. Implementarea modulului de înmulțire

Pentru a simplifica înmulțirea greutăților cu cea a semnalelor de intrare, ne vom ajuta de către structura folosită în stocarea datelor și se va implementa un modul de înmulțire matriceală generic. Acest modul va putea primi ca intrări 2 matrici de mărime $m \times n$ respectiv $n \times p$ și va crea o matrice ca semnal de ieșire de mărimea $m \times p$





4.2.2. Implementarea modului funcției de activare reLU și normalizare

Se vor implementa 2 module pentru definirea celor 2 funcții de activare. Definirea funcției reLU este trivială, necesitând alegerea maximului între 0 și valoarea de input. Pe de altă parte, pentru stratul de ieșire se va folosi o funcție de normalizare simplă deoarece în urma deciziilor de implementare alese nu putem aplica funcția softmax asupra valorilor obținute în ultimul strat.

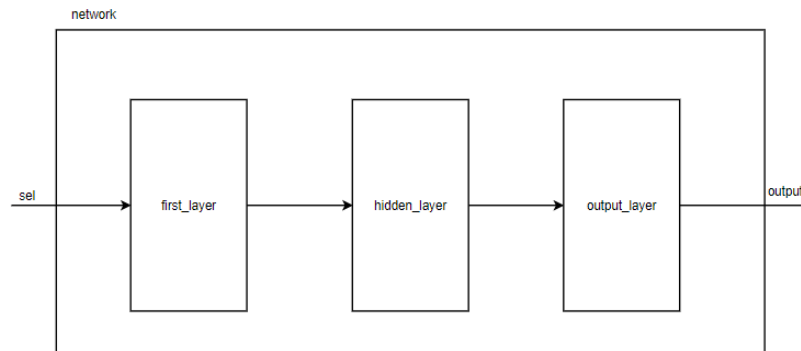
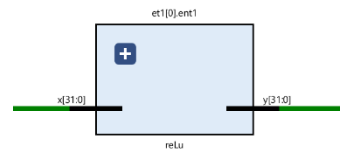


Fig 7 - Schema bloc a sistemului

4.2.3. Implementarea primului strat

Un prim modul necesar este primul strat al rețelei care are rolul de a procesa imaginile sub forma de matrice și de a le transforma într-o matrice cu un singur rând. Această nouă matrice obținută va fi folosită ca și semnal de intrare de către următorul strat al RAN (stratul ascuns).

4.2.4. Implementarea stratului ascuns

Stratul ascuns va primi ca input imaginile procesate de către stratul de intrare. Acesta va folosi modulul de multiplicare pentru a genera ieșirile prin înmulțirea greutăților cu cea a input-ului. Înainte de a trimite spre ieșire valorile obținute din multiplicare, se va aplica funcția de activare reLU asupra fiecărei valori.

4.2.5. Implementarea stratului de ieșire

Stratul de ieșire va primi ca input valorile trimise de către stratul ascuns. Acesta va folosi modulul de multiplicare pentru a genera ieșirile prin înmulțirea greutăților cu cea a input-ului. Înainte de a trimite spre ieșire valorile obținute din multiplicare, se va aplica funcția de activare pentru a normaliza rezultatele și de a le aduce în intervalul (0,100), aceste valori reprezentând șansa ca o cifră specifică să fie în imagine.

4.2.6 Implementarea modului principal

Rolul modului principal este de a instanția restul entităților RAN și de a le lega între ele cu ajutorul unor semnale intermediare. Se folosește un semnal generic pentru declararea mărимii imaginii (se consideră imaginea ca fiind o matrice pătratică cu latura de lungime dată de semnalul generic), iar ca semnal de intrare avem un stimul care selectează imaginea care va trece prin RAN. Semnalul de ieșire va fi o matrice cu 10 coloane în care vom avea o distribuție de valori care reprezintă probabilitatea ca imaginea să reprezinte cifra de pe care poziție se afla valoarea respectivă.

Imaginile sunt niste constante declarate într-un pachet care conține toate constantele necesare sistemului. O dată cu selecția imaginii, primul strat primește această imagine pe care o aplatizează și o trimite stratului ascuns printr-un semnal intermediar. Stratul ascuns primește această matrice pentru procesare, totodată se declară semnalele generice necesare acestuia (numărul de neuroni și mărimea matricii aplatizate). Rezultatul este trimis către ultimul strat. Stratul de ieșire va face operațiile necesare și va aplica funcția de normalizare, ulterior transmitând matricea finală către ieșire. Această matrice este conectată la semnalul de ieșire al modului principal.

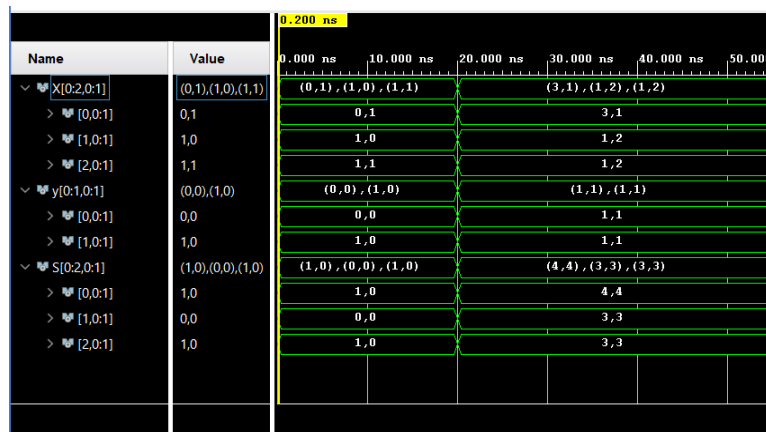


Rezultate experimentale

Pentru implementare și testarea s-a folosit limbajul VHDL, mediul de dezvoltare utilizat pentru implementare este 'Vivado Design Suite 2019.2'. Pentru a testa implementarea proiectului s-au folosit numeroase banci de test pentru fiecare modul în parte. Un banc de test este un cod scris în VHDL care ne permite să avem un set de stimuli repetitive și portabili pentru nenumărate simulări. Acesta poate să fie un simplu fișier cu un semnal de ceas și date de intrare sau un fișier mai complex cu verificări de erori și testare condițională. Pentru testare s-a ales simpla introducere a unor stimuli și verificarea corectitudinii datelor de ieșire.

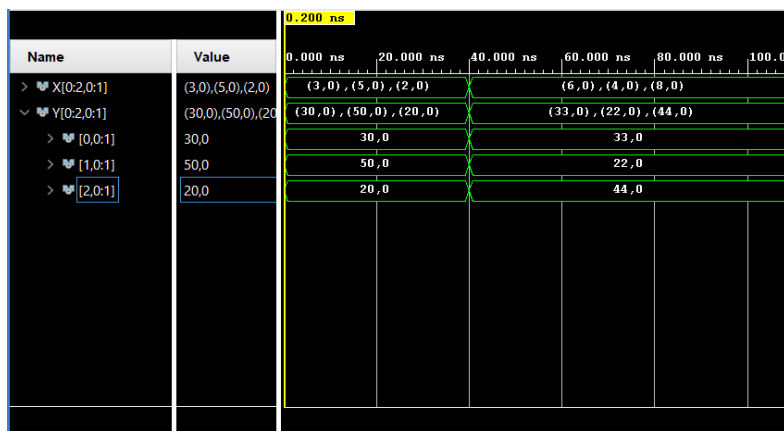
Testarea modului de înmulțire matriceală

Pentru a testa modulul de înmulțire matriceală s-a creat un banc de test separat în care s-au instantiat toate semnalele necesare instantierii modulului de înmulțire. S-au declarat 2 semnale care vor fi folosite ca stimuli de intrare al modulului, X și Y. Se poate observa că în urma înmulțirii celor 2 matrici X și Y, rezultatul este corect din punct de vedere matematic.



Testarea modului de normalizare

Pentru a testa modulul de înmulțire matriceală s-a creat un banc de test separat în care s-au instantiat toate semnalele necesare instantierii modulului de înmulțire. S-au declarat un semnal care va fi folosit ca stimul de intrare al modulului. Ca ieșire se va putea observa un alt semnal care va avea un procent care reprezintă rata de acoperire a valorii de pe poziția respectivă în suma totală.





Testarea modului primului strat

Modulul primului strat s-a testat într-un mod similar cu modulele precedente. S-a creat un banc de test și s-au declarat toate semnalele necesare instanțierii entității primului strat. Stimulul de intrare al modului este o matrice cu mărimea de 7 pe 7. Ieșirea obținută este o matrice cu mărimea de 49 pe 1 care este obținută din aplatizarea matricii de la intrare.

Name	Value	499,998 ps	499,999 ps
✓ X[0:6,0:6]	(1,2,3,4,5,6,7),(1,2,3,4,5,6,7)	(1,2,3,4,5,6,7), (1,2,3,4,5,6,7)	(1,2,3,4,5,6,7), (1,2,3,4,5,6,7)
> X[0,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[1,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[2,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[3,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[4,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[5,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
> X[6,0:6]	1,2,3,4,5,6,7	1,2,3,4,5,6,7	1,2,3,4,5,6,7
✓ y[0:48,0:1]	(1,0),(2,0),(3,0),(4,0),(5,0),(6,0),(7,0)	(1,0), (2,0), (3,0), (4,0), (5,0), (6,0), (7,0)	(1,0), (2,0), (3,0), (4,0), (5,0), (6,0), (7,0)
> y[0,0:1]	1,0	1,0	1,0
> y[1,0:1]	2,0	2,0	2,0
> y[2,0:1]	3,0	3,0	3,0
> y[3,0:1]	4,0	4,0	4,0
> y[4,0:1]	5,0	5,0	5,0
> y[5,0:1]	6,0	6,0	6,0
> y[6,0:1]	7,0	7,0	7,0
> y[7,0:1]	1,0	1,0	1,0
> y[8,0:1]	2,0	2,0	2,0
> y[9,0:1]	3,0	3,0	3,0
> y[10,0:1]	4,0	4,0	4,0
> y[11,0:1]	5,0	5,0	5,0

Testarea modului stratului ascuns

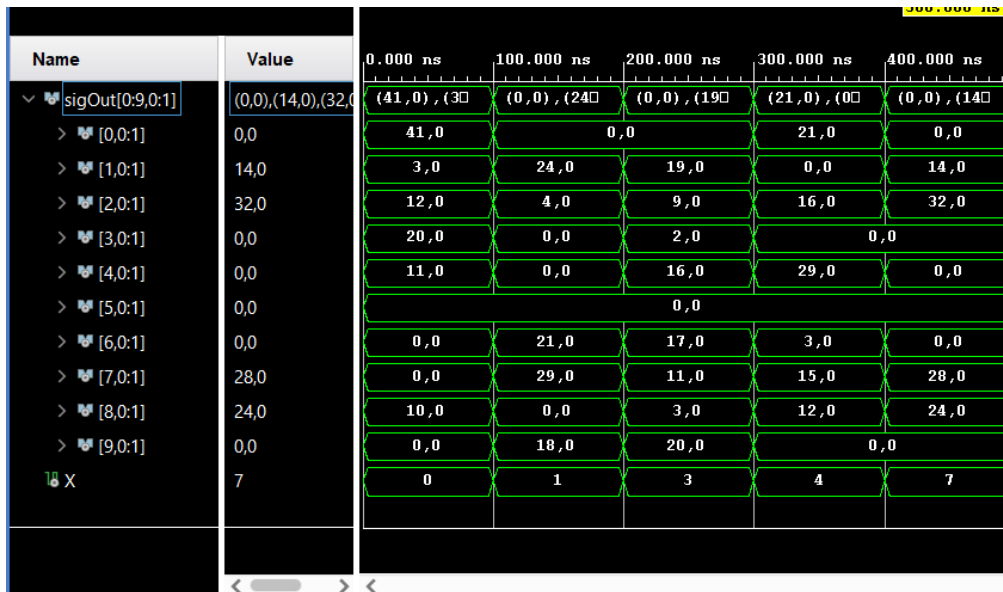
Pentru testarea stratului ascuns se vor instanția entitățile primului strat și al celui ascuns. Ca stimuli de intrare pentru primul strat se va folosi o matrice de 2 pe 2 care va fi trimisă către primul strat. Acesta va procesa matricea și o va aplatiza pentru a putea să fie folosită de către stratul ascuns. Stimulii de intrare pentru stratul ascuns va fi matricea aplatizată împreună cu matricea greutăților. Rolul modului celui de-al 2-lea strat este de a înmulți matricea greutăților cu cea a valorilor din matricea aplatizată, mai apoi să aplice funcția de activare reLU asupra valorilor obținute înainte să le trimită spre ieșire.

Name	Value	499,995 ps	499,996 ps	499,997 ps	499,998 ps	499,999 ps	500,000 ps
✓ X[0:1,0:1]	(3,2),(6,4)	(3,2), (6,4)	(3,2), (6,4)	(3,2), (6,4)	(3,2), (6,4)	(3,2), (6,4)	(3,2), (6,4)
> X[0,0:1]	3,2	3,2	3,2	3,2	3,2	3,2	3,2
> X[1,0:1]	6,4	6,4	6,4	6,4	6,4	6,4	6,4
✓ y[0:3,0:1]	(3,0),(2,0),(6,0),(4,0)	(3,0), (2,0), (6,0), (4,0)	(3,0), (2,0), (6,0), (4,0)	(3,0), (2,0), (6,0), (4,0)	(3,0), (2,0), (6,0), (4,0)	(3,0), (2,0), (6,0), (4,0)	(3,0), (2,0), (6,0), (4,0)
> y[0,0:1]	3,0	3,0	3,0	3,0	3,0	3,0	3,0
> y[1,0:1]	2,0	2,0	2,0	2,0	2,0	2,0	2,0
> y[2,0:1]	6,0	6,0	6,0	6,0	6,0	6,0	6,0
> y[3,0:1]	4,0	4,0	4,0	4,0	4,0	4,0	4,0
✓ weights[0:1,0:3]	(1,1,0,0),(0,0,1,1)	(1,1,0,0), (0,0,1,1)	(1,1,0,0), (0,0,1,1)	(1,1,0,0), (0,0,1,1)	(1,1,0,0), (0,0,1,1)	(1,1,0,0), (0,0,1,1)	(1,1,0,0), (0,0,1,1)
> weights[0,0:3]	1,1,0,0	1,1,0,0	1,1,0,0	1,1,0,0	1,1,0,0	1,1,0,0	1,1,0,0
> weights[1,0:3]	0,0,1,1	0,0,1,1	0,0,1,1	0,0,1,1	0,0,1,1	0,0,1,1	0,0,1,1
✓ output[0:1,0:1]	(5,0),(10,0)	(5,0), (10,0)	(5,0), (10,0)	(5,0), (10,0)	(5,0), (10,0)	(5,0), (10,0)	(5,0), (10,0)
> output[0,0:1]	5,0	5,0	5,0	5,0	5,0	5,0	5,0



Testarea rețelei neuronale

Pentru testarea rețelei neuronale a fost creat un banc de test unde s-a instantiat modulul principal al rețelei 'network'. S-a ales sa se testeze rețeaua cu ajutorul unor imagini de 20 pe 20 de pixeli. A fost declarat un stimul pentru selectarea imaginii transmise RAN, acest stimul va fi schimbat pe parcursul simulării pentru a vedea comportamentul RAN pentru mai multe imagini. Se vor folosi 5 imagini cu 5 cifre diferite si anume cifrele: zero, unu, trei, patru si sapte.



Valoarea lui X reprezinta o imagine care continue cifra indicata de numarul semnalului. Se poate observa din imagine procentele obtinute pentru fiecare imagine in parte. Precizia sistemului nu este una ridicata, avand o precizie de aproximativ 40% datorita reprezentarii greutatilor in punct fix. Totodata distributia acestora este una obtinuta prin normalizare simpla si nu printr-o functie care ne ofera o distributie logaritmica.



Concluzii

Scopul proiectului a fost de a implementa o RAN capabila sa recunoasca intr-o imagine o cifra scrisa de mana si sa ne spuna sub forma unui procent care este sansa ca in imagine sa fie cifra respective. Obiectivele principale fiind dezvoltarea modelului RAN printr-o implementarea software pentru a obtine greutatile unui model antrenat si a le putea utiliza pentru dezvoltarea variantei hardware a RAN. Utilizarea unei retele neurale pe o placuta FPGA duce la o imbunatatire a vitezei de calcul. Din acest motiv RAN sunt de obicei preantrenate si apoi puse sa ruleze pe placutele de dezvoltare deoarece antrenarea sistemului se face o singura data si nu are rost sa se dezvolte o implementare hardware pentru a antrena reseaua direct pe placa de dezvoltare.

Utilizarea reprezentarii in punct fix a greutatilor a dus la o scadere a performantei retelei, acesta fiind un compromise datorita lipsei reprezentarii numerelor in virgula mobile in limbajul VHDL. Asadar reseaua lucreaza doar cu numere intregi reprezentate pe 4 octeti si nu este necesara folosirea numerelor in virgula mobile.

Exista multe posibilitati de dezvoltare ulterioara a retelei: o mare imbunatatire ar fi dezvoltarea unei reprezentari in virgula mobile in limbajul VHDL si utilizarea acesteia pentru greutati. O data cu posibilitatea reprezentarii in virgula mobile se poate utiliza si functia softMax care nu a putut fi implementat datorita reprezentarii in punct fix. O alta posibila dezvoltare ar putea fi schimbarea tipului retelei neuronale in una convolutionala care ar putea sa imbunatateasca viteza si sa scada complexitatea calculelor necesare.



Bibliografie

- [1] Z. Zhang, "A gentle introduction to artificial neural networks," *Annals of Translational Medicine*, p. 7, October 2016 .
- [2] C. Gallo, "Artificial Neural Network Tutorial," *Encyclopedia of Information and Technology, Third Edition*, p. 13, January 2015 .
- [3] J. Brownlee, "How to Develop a CNN for MNIST Handwritten Digit Classification," 24 August 2020. [Online]. Available: <https://machinelearningmastery.com/how-to-develop-a-convolutional-neural-network-from-scratch-for-mnist-handwritten-digit-classification/>.
- [4] S. Ahlawat, "Improved Handwritten Digit Recognition Using Convolutional Neural Networks (CNN)," 25 May 2020 .
- [5] I.-J. K. a. X. Xie, "Handwritten Hangul recognition using deep convolutional," p. 29, 2016 .
- [6] "Image Classification in 10 Minutes with MNIST Dataset," [Online]. Available : <https://towardsdatascience.com/image-classification-in-10-minutes-with-mnist-dataset-54c35b77a38d>.
- [7] F. W. Wibowo, "An Analysis of FPGA Hardware Platform Based Artificial Neural Network," 2019
- [8] M. Nielsen, "Using neural nets to recognize handwritten digits," December 2019. [Online]. Available: <http://neuralnetworksanddeeplearning.com/chap1.html> .
- [9] "7 Types of Neural Network Activation Functions: How to Choose?," [Online]. Available: <https://missinglink.ai/guides/neural-network-concepts/7-types-neural-network-activation-functions-right/>
- [10] "Deep Learning Project – Handwritten Digit Recognition using Python," [Online]. Available: <https://data-flair.training/blogs/python-deep-learning-project-handwritten-digit-recognition/> .
- [11] "Perceptrons and Multi-Layer Perceptrons: The Artificial Neuron at the Core of Deep Learning," [Online]. Available: <https://missinglink.ai/guides/neural-network-concepts/perceptrons-and-multi-layer-perceptrons-the-artificial-neuron-at-the-core-of-deep-learning/> .
- [12] Udemy, "Machine learning from A to Z," [Online]. Available: <https://www.udemy.com/course/machinelearning/>.
- [13] 3Blue1Brown, "But what is a Neural Network? | Deep learning, chapter 1," 2017. [Online]. Available: <https://www.youtube.com/watch?v=aircArvvnKk> .
- [14] "The Complete Guide to Artificial Neural Networks: Concepts and Models," 2018. [Online]. Available: <https://missinglink.ai/guides/neural-network-concepts/complete-guide-artificial-neural-networks/> .
- [15] 3. SERIES, "Gradient descent, how neural networks learn | Deep learning, chapter 2," [Online]. Available: <https://www.youtube.com/watch?v=IHZwWFHWa-w&t=590s> .
- [16] "Neural Networks for Image Recognition: Methods, Best Practices, Applications," [Online]. Available: <https://missinglink.ai/guides/computer-vision/neural-networks-image-recognition-methods-best-practices-applications/> .