

Parallel VLSI Neural System Design

DAVID ZHANG

Springer



TP183
Z63

Parallel VLSI Neural System Design

DAVID ZHANG



E200000525



Springer

Dr. David Zhang
Department of Computer Science
Hong Kong Polytechnic University
Hung Hom Kowloon
Hong Kong

Library of Congress Cataloging-in-Publication Data

Parallel VLSI Neural System Design / David Zhang

p. cm.

Includes bibliographical references and index

ISBN 9813083301

1. Neural Networks (Computer science). 2. Parallel processing (Electronic computers). 3. Integrated circuits--Very large scale integration. I. David Zhang, 1949- . II. Title.

QA76.87.Z473 1998

006.3/2 --dc21

98-41790
CIP

ISBN 981-3083-30-1

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on micro-films or in any other way, and storage in databanks or in any system now known or to be invented. Permission for use must always be obtained from the publisher in writing.

© Springer-Verlag Singapore Pte. Ltd. 1999
Printed in Singapore

The publisher makes no representation, express or implied, with regard to the accuracy of the information contained in this book and cannot accept any legal responsibility or liability for any errors or omissions that may be made.

Typesetting: Camera-ready by author
SPIN 10628591 5 4 3 2 1 0

Preface

Neural information processing has been emerging as a new field. This type of processing is an alternative form of computation that attempts to mimic the functionality of the biological human brain in solving demanding pattern recognition problems. However, researchers and engineers have long been fascinated by how efficient and how fast *artificial neural networks* (ANNs) are capable of performing such complex tasks as recognition. Such networks are capable of recognizing input data from *any of the five senses* with the necessary accuracy and speed to allow living creatures to survive. Machines which perform such complex tasks as recognition, with similar accuracy and speed, were difficult to be implemented until the technological advances of integrated circuits and VLSI systems. Since then, VLSI neural systems have witnessed an exponential growth and a new engineering discipline was born.

There were a number of excellent text and reference books on the subject, each dealing with one or two topics. This book attempts to present a parallel VLSI neural system design methodology for pattern recognition applications. The methodology emphasizes a coordination between model definition, architectural description, and hardware implementation. Depending on the different pattern recognition applications, the methodology provides appropriate ANN models suited to parallel/pipeline processing, mapping the models onto the corresponding VLSI architectures and finally hardware implementation.

A parallel ANN model is a basis of VLSI design because it directly reflects highly parallel, regular and modular VLSI architectures. Such three kinds of parallel ANN models, including an unsupervised learning model for fuzzy clustering, a supervised training model for pattern classification, and a neural-like network model for finite ring computing, are developed in Chapters 2, 3, and 4, respectively.

VLSI technology offers a highly advanced implementation medium both at the fabrication and the CAD level if efficient architectures can be provided. In Chapter 5, the mapping policies from ANNs to systolic arrays are given. Three typical VLSI architectures based on the effective matching policies are introduced in Chapters 6-8. They are a parallel architecture built by systolic arrays, a pipeline architecture based on window operation, and a simplified architecture using a priori knowledge.

The implementation of ANN architectures into silicon as special-purpose hardware is an important step in parallel VLSI neural system design. In Chapter 9, a computational block design for digital ANN based on dynamic pipelines is proposed. Chapter 10 presents a digital array compressor design based on C²PL (complex complementary pass-transistor logic). In Chapter 11, a hybrid programmable ANN design using BiCMOS circuit building blocks is described. As an example of VLSI implementation, a finite ring ANN given in Chapter 4 is also implemented in Chapter 12.

The effectiveness of parallel VLSI neural system design methodology is illustrated by applying the designs to various pattern recognition applications, and analyzing the performances of the given systems.

This book is not a primer in ANN, in that a certain amount of prior knowledge, such as parallel processing and VLSI design, is assumed. It is my hope that this book will contribute to our understanding of this new and exciting discipline: Parallel VLSI Neural System Engineering.

David D. Zhang
The Hong Kong Polytechnic University

Acknowledgements

My sincere thanks goes to Professors M.I. Elmasry and M. Kamel at University of Waterloo, Ontario, Canada, for their support and advice throughout this research. I would like to thank Dr. Harold H. Szu for his useful comments on artificial neural networks. I would also like to express my gratitude to Professor G.A. Jullien and Professor W.C. Miller at University of Windsor, Ontario, Canada, for their previous supports and their encouragement to me to study VLSI design. Special thanks are due to my research assistant, Lei Huang, for his help in the preparation of this book. The support from the University Grants Committee in Hong Kong is appreciated.

Contents

<i>Preface</i>	v
----------------------	---

CHAPTER 1	<i>VLSI Neural System Design Methodology</i>.....	1
1.1	INTRODUCTION	1
1.2	NEURAL NETWORK MODELS	3
1.2.1	Biological Neural Networks	3
1.2.2	Artificial Neural Networks	6
1.3	ANN ARCHITECTURES	10
1.4	HARDWARE IMPLEMENTATIONS	11
1.4.1	Analog and Mixed Implementations	12
1.4.2	Digital Implementations	13
1.5	VLSI SYSTEM DESIGN METHODOLOGY	14

<u>PART I</u>	<u><i>PARALLEL ANN MODELS</i></u>.....	21
----------------------	---	-----------

CHAPTER 2	<i>An Unsupervised Learning Model</i>.....	23
2.1	INTRODUCTION	23
2.2	FUZZY CLUSTERING NEURAL NETWORKS ..	28
2.2.1	Network Architecture	28
2.2.2	Learning Algorithm	28
2.3	PARALLEL FCNN MODEL	32
2.4	EXPERIMENTAL RESULTS	32
2.5	SUMMARY	40

CHAPTER 3	<i>A Supervised Training Model</i>	41
3.1	INTRODUCTION	41
3.2	LINEAR SEPARABILITY ANALYSIS	44
3.2.1	Definitions	44
3.2.2	Layered Perceptrons	46

3.3 LAYER ADAPTATION APPROACH	48
3.3.1 Linear Separability Principle	49
3.3.2 Adaptation Approach	50
3.4 EXPERIMENT: PATTERN RECOGNITION	53
3.5 COMPARISONS	55
3.6 SUMMARY	57
 CHAPTER 4 <i>A Neural-Like Network Model</i>	59
4.1 INTRODUCTION	59
4.1.1 Notation	60
4.1.2 Residue Number System (RNS)	61
4.1.3 Neural Network Architecture	61
4.2 FRNN COMPUTING MODEL	63
4.3 FRNN ARCHITECTURE	66
4.4 CASE STUDIES	67
4.4.1 A Multiplier	68
4.4.2 RNS to Binary Converter	69
4.5 SUMMARY	70
 <u>PART II</u> <i>VLSI ARCHITECTURES</i>	71
 CHAPTER 5 <i>Mapping ANN onto Systolic Arrays</i>	73
5.1 INTRODUCTION	73
5.1.1 Systolic Arrays	74
5.1.2 Mapping Algorithm to Systolic Architecture	77
5.2 MAPPING POLICIES	79
5.3 DESIGN APPROACH	82
5.3.1 Typical SA Structures	82
5.3.2 Pipeline Matching	84
5.3.3 Iteration Processing	84
5.4 CASE STUDY	86
5.4.1 Hamming Net	86
5.4.2 Simulations and Experiments	90
5.5 SUMMARY	91

CHAPTER 6	<i>A Parallel Architecture Implemented by Systolic Arrays</i>	93
6.1	INTRODUCTION	93
6.2	FCNN ARCHITECTURE	94
6.3	PERFORMANCE ANALYSIS	96
6.4	MAPPING FCNN ONTO SA	100
6.5	SUMMARY	106
CHAPTER 7	<i>A Pipelined Architecture Based on Window Operation</i>	107
7.1	INTRODUCTION	107
7.2	PIPELINED ARCHITECTURE	109
7.2.1	Window Model	109
7.2.2	Pipelined Architecture	111
7.2.3	Building Unit Design	112
7.3	WINDOW IMPLEMENTATION	116
7.3.1	Parallel Data Flow Window	116
7.3.2	Serial Data Flow Window	117
7.3.3	Window Computation Element	118
7.4	CASE STUDIES	120
7.5	PERFORMANCE ANALYSIS	126
7.6	SUMMARY	129
CHAPTER 8	<i>A Simplified Architecture Using A Priori Knowledge</i>	131
8.1	INTRODUCTION	131
8.2	TYPICAL STRUCTURE MODELS	133
8.2.1	Output ROM Model	133
8.2.2	Input ROM Model	135
8.2.3	Learning ROM Model	135
8.3	ROM LAYER IN VLSI	136
8.4	EXAMPLES	136
8.5	SUMMARY	147

PART III HARDWARE IMPLEMENTATIONS 149

CHAPTER 9 *Computational Blocks Design for Digital ANN . 151*

9.1 INTRODUCTION.....	151
9.2 PIPELINED SWITCHING TREES	152
9.3 GRAPH BASED REDUCTION.....	153
9.3.1 Graph Reduction Rules	154
9.3.2 Minimization Considerations.....	157
9.3.3 Example Results.....	159
9.4 CIRCUIT CONSIDERATIONS	161
9.4.1 Worst Case Test	161
9.4.2 Reduction of Charge Sharing.....	162
9.4.3 Reduction of Tree Height.....	165
9.4.4 Transistor Sizing	169
9.5 PRELIMINARY FABRICATION RESULTS	170
9.6 SUMMARY	171

CHAPTER 10 *Digital ANN Compressor Design 173*

10.1 INTRODUCTION.....	173
10.2 C^2 PL MODEL	175
10.3 3-2 COMPRESSOR DESIGN	176
10.3.1 Basic Structure	177
10.3.2 Comparison with CPL.....	181
10.4 DNN APPLICATIONS	183
10.5 SUMMARY	189

CHAPTER 11 *Hybrid Programmable ANN Design 191*

11.1 INTRODUCTION.....	191
11.2 ANALYSIS AND DESIGN FOR PRNN	193
11.3 IMPROVED PRNN CIRCUIT	197
11.3.1 Synapse Building Block.....	197
11.3.2 Neuron Building Block	199
11.3.3 Connection Network	201

11.4 EXPERIMENTAL RESULTS201

11.5 SUMMARY203

APPENDIX A203

CHAPTER 12 *VLSI Implementation for Finite Ring ANN.....* 207

12.1 INTRODUCTION207

12.2 FRRR Architecture208

 12.2.1 Modulo Reduction208

 12.2.2 MSB Carry Iteration209

 12.2.3 Feedforward Processing212

12.3 VLSI IMPLEMENTATION213

 12.3.1 Carry Look-Ahead Adder214

 12.3.2 ROM Implementation216

 12.3.3 ROM Logic Cell217

12.4 COMPARISON218

12.5 SUMMARY223

APPENDIX B225

CHAPTER 13 *Conclusions and Prospects.....* 227

Bibliography235

Index253

1

VLSI Neural System Design Methodology

In this chapter, we begin by introducing the background for artificial neural networks (ANNs). Then, the existing ANN models, architectures and hardware implementation techniques are briefly reviewed. This leads to the definition of the research objective and provides insight into a parallel VLSI neural system design methodology.

1.1 INTRODUCTION

ANNs are massively parallel interconnected networks of simple (usually adaptive) nodes which are intended to interact with objects of the real world in the same way as biological nervous systems do [1].

The interest in these networks is due to the general opinion that they are able to perform some complicated and creative tasks, such as pattern recognition, similar to the way they are performed by human brains [2,10,35]. The implementations of these tasks by traditional computing methods have only reached relatively low performances in some limited aspects or environments. Nevertheless, as neural systems show some properties, like association, generalization, parallel searching, and adaptation to changes in the environment, which are analogous to human brain properties, they promise improved results.

The usage of ANNs for pattern recognition may be traced back to the perceptron models originated by Rosenblatt in 1950 [2]. The perceptron models used the concept of reward and punishment. In late 1960s, the progress in ANN models slowed down due to the limited capabilities of the early single layer perceptron models. In the mid-1970s and early 1980s, with the availability of enhanced computing power the progress in the development of ANN models accelerated. Researchers were able to model and test their theories about the functioning of the brain.

Today a number of well-developed theories and models of ANNs are available [3,34-36,40-47,55-64,198-201]. These networks consist of a large number of simple processing elements called nodes that represent the neurons. These nodes are interconnected by the synaptic connections. These models are capable of learning and making decisions; and are suitable for a variety of pattern recognition tasks [36,39].

Pattern recognition techniques can be grouped into two classes: supervised and unsupervised techniques. In supervised methods, certain number of samples is available for each category and these samples are used to train the classifier. In the case of unsupervised classification no training samples are available, and the network learns by detecting the similarity between the input patterns.

Now many ANN models and algorithms for pattern recognition applications are available. They include Back-Propagation (BP) learning, Competitive learning, Kohonen learning, Adaptive Resonance Theory, Neocognitron models, Hopfield Networks, and Boltzman machines [5-13]. Applications of ANNs include character recognition, human face identification, speech recognition, multispectral image analysis, and expert systems [14-18]. The BP learning is essentially of the supervised type and the network learns with the help of training sets. The BP networks have been successfully used for many pattern recognition problems [5,15-17].

Another important class of neural networks is self-organizing neural networks. The networks with competitive learning algorithms are self-organizing networks. Early models of competitive learning were developed by Malsburg in his study of visual cortex [19]. Rumelhart and Zipser have suggested an algorithm for competitive learning [9]. The main disadvantage of competitive learning is that the network forgets its earlier learning with new learning and the network may get set into an unstable state with the spurious input patterns. To overcome this drawback, Grossberg developed an adaptive resonance architecture [10-11] and Fukushima proposed the Neocognitron models [12]. Kohonen developed a learning paradigm for self-organizing networks known as Kohonen learning [6-7]. These algorithms can be used for a variety of tasks in pattern recognition.

ANNs consist of parallel distributed processing (PDP) models. The PDP models are well described in the work of Rumelhart and McClelland [4]. The functional synthesis of these models consists of establishing a relationship between the several inputs and one or more outputs. In ANN, the nodes are connected to each other by the synaptic connections or the links. There is an

associated synaptic strength or a weight with each connection. During the learning, the weights which represent the knowledge stored in the network are updated. The ANNs consist of two or several layers of nodes and each layer contains several nodes. The observed feature vector is presented to the input nodes. The input values may represent the probability that the discrete feature is present. Each possible decision or outcome can be represented by a node in the output layer.

1.2 NEURAL NETWORK MODELS

ANNs are information processing systems. Due to the variety of fields of interest and applications, a broad range of ANN models has been emerged. In all these fields, the term “neural networks” is characterized by a combination of adaptive learning algorithms and parallel distributed implementations. Although ANNs are biologically motivated, their resemblance to the brain models is not straightforward. Since the basic concepts of biological neural networks provide a common ground for understanding the ANN models, we will first introduce basic biological neural networks in order to appreciate the differences and similarities with ANN.

1.2.1 Biological Neural Networks

As you know, the von Neumann computer performs poorly on certain tasks, such as pattern recognition that humans handle routinely. It is interesting to compare the human brain with a serial modern von Neumann computer from the information processing aspects as shown in Table 1.1. Although the neuron’s switching time (a few milliseconds) is about a million times slower than modern computer elements, they have a thousand-fold greater connectivity than today’s supercomputers. Furthermore, the brain is very powerfully efficient – it consumes less than 100 watts – by contrast a supercomputer may dissipate 10^5 watts.

The basic element of neural networks of a brain is a neuron. The neurons consist of four basic parts: cell body, synapses, axons, and dendrites. The cell body essentially sums the membrane potential provided by the synapses. The synapses provide an output. Axons are the connections between the neurons that carry charge, and the dendrites are the branch-like structures, which provide the sensory input to a cell body (See Fig.1.1).

In fact, the synapse represents the junction between an axon and a dendrite. How two or more neurons interact remains largely mysterious, and complexities of different neurons vary greatly. Generally speaking, a neuron sends its output to other neurons via its axon. An axon carries information through a series of action potentials, or waves of current, which depends on the neuron’s voltage potential. More precisely, the membrane generates the action potential and propagates down the axon and its branches, where axonal insulators restore and amplify the signal as it propagates, until it arrives at a synaptic junction.

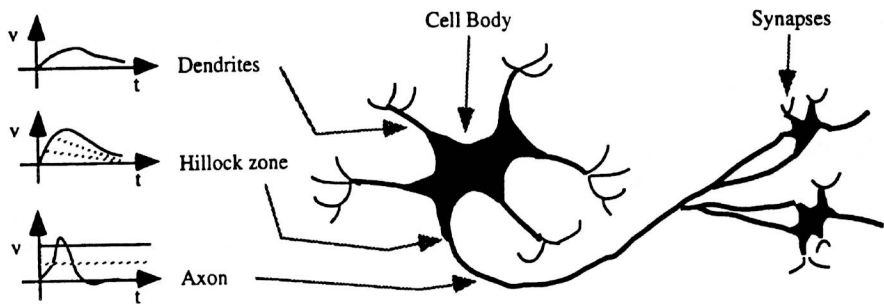


Fig.1.1 A prototype biological neuron

Table 1.1 Comparison between human brain and modern computer

	<i>Human Brain</i>	<i>Modern Computer</i>
<i>Element</i>	10^{11} neurons 10^{14} synapses	10^9 transistors
<i>Fanout</i>	10^3	3
<i>Implementation</i>	Analog	Digital
<i>Processing</i>	passively parallel	largely serial
<i>Switching Time</i>	10^{-3} s	10^{-9} s
<i>Power</i>	100 watts	10^5 watts

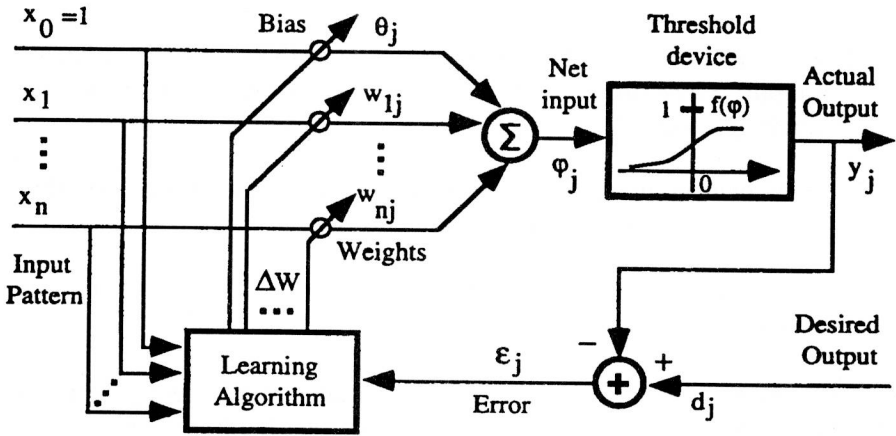


Fig.1.2 A typical artificial neuron model

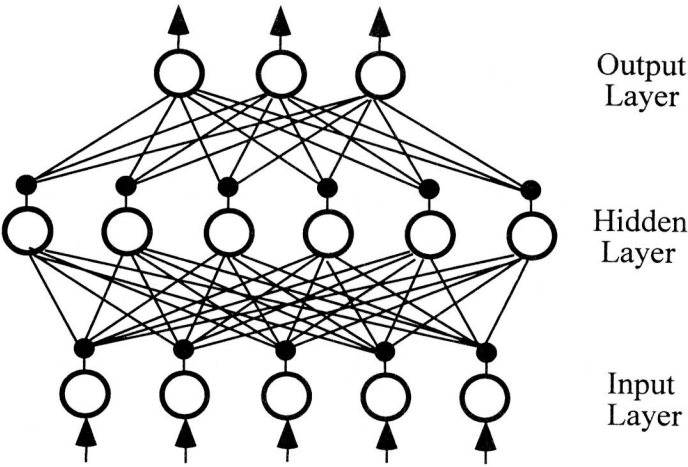


Fig.1.3 A three-layer feedforward network

There are several important features obtained from biological neural networks, which apply to all ANNs:

- 1) Each neuron acts independently of all others – each neuron's output relies only on its constantly available inputs from the abutting connections.
- 2) Each neuron relies only on local information – the information that is provided by the adjoining connections.
- 3) The large number of connections provides a large amount of redundancy and facilitates a distributed representation.

First two features allow neural networks to operate efficiently in parallel. The last one provides neural networks with inherent fault-tolerance.

1.2.2 Artificial Neural Networks

ANNs mimic the functioning of the neural networks of a brain. ANN consists of a large number of simple nodes. Each of the nodes is connected to another node(s) through a synaptic connection or a link [34].

Information processing takes place through the interaction between the nodes. Each node is associated with an activation value $\varphi_j(t)$. The activation value passes through an activation function $f(\varphi_j)$ to provide an actual output $y_j(t)$. These outputs pass through the unidirectional synaptic connections. There is an associated number, w_{ij} , called the weight or the connection strength, that determines the amount of effect node i can have on node j . For each node all the inputs are combined, and the total input, along with the current activation, determines the new activation value (Fig.1.2).

Usually ANNs consist of a number of layers and nodes in each layer. The most general model assumes the complete interconnections between all the nodes, and resolves the cases of the nonconnected nodes (i, j) by setting the weights $w_{ij} = 0$. A simple three-layer feedforward network is shown in Fig.1.3. The networks can be synchronous or asynchronous. The synchronous networks are controlled by clock pulses; whereas in asynchronous networks the nodes respond instantaneously to the incoming inputs.