神经元网络系统设计方法

# NEURAL NETWORKS SYSTEM DESIGN METHODOLOGY

David D. Zhang

*Department of Computer Science*
*City University of Hong Kong*
*Kowloon, Hong Kong*

# NEURAL NETWORKS SYSTEM DESIGN METHODOLOGY

## 神经元网络系统设计方法

**David D. Zhang**

*Department of Computer Science*
*City University of Hong Kong*
*Kowloon, Hong Kong*

（京）新登字 158 号

## 内 容 简 介

　　本书从讨论神经元网络的模型入手，系统阐述了神经元网络系统设计原理和方法及应用，并着重介绍了如何用 VLSI（大规模集成电路）实现神经元网络。本专著包含了作者在该领域的研究成果，提供了参考资料目录。既是人工智能、计算机和微电子学方面的技术参考书，又可作高年级大学生和研究生的专业教材。

# NEURAL NETWORKS SYSTEM DESIGN METHODOLOGY

**Neural Networks System Design Methodology** presents an integrated design approach to this new engineering discipline. Three topics of the system design methodology, including several model definitions, architectural descriptions, and hardware implementations, are investigated and their coordination are discussed in detail after an overview of the field of ANNs. Engineering applications of ANNs to fuzzy clustering, speech recognition, classification and pattern recognition are covered. The book can be used as a text book or reference for graduate or senior undergraduate courses on the subject. It can also be used by researchers in the field.

**David Zhang** graduated in computer science from Peking University in 1974 and received his M.Sc and Ph.D degrees in computer science and engineering from Harbin Institute of Technology (HIT) in 1983 and 1985, respectively. From 1986 to 1988 he was a postdoctoral fellow at Tsinghua University and became an associate professor at Academia Sinica, Beijing, China. In 1988, he joined the University of Windsor, Ontario, Canada, as a visiting professor in electrical engineering. He received his second Ph.D in electrical and computer engineering at University of Waterloo, Ontario, Canada, in 1994. Currently, he is an associate professor in City University of Hong Kong. He has authored and co-authored near 100 papers including two books, as well as received several recognisable project awards. Dr. Zhang is a senior member of the IEEE.

*To the memory of my beloved supervisor*
*——Professor Tong Chang*

# Preface

Researchers and engineers have long been fascinated by how efficient and how fast biological neural networks 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 VLSI circuits and systems in the late 1980's. Since then, the field of Artificial Neural Networks (ANNs) have witnessed an exponential growth and a new engineering discipline was born. Today, many engineering curriculums have included a course or more on the subject at the graduate or senior undergraduate levels.

This book attempts to present a system design methodology of ANNs 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 VLSI implementation. The book discusses these three phases:

1. Parallel ANN Model: Three types of models, 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. Compared with the conventional approaches, the new models can greatly reduce the complexity of the VLSI implementation.

2. VLSI Architecture: Three typical architectures, including a parallel architecture built by systolic arrays, a pipeline architecture based on window operation and a simplified architecture using a priori knowledge, are designed. They are all easily implemented in VLSI medium.

3. Hardware Implementation: Two design approaches are investigated. One is a digital array compressor design based on a complex complementary pass-transistor logic ($C^2PL$) and the other is a hybrid programmable ANN design using BiCMOS circuit building blocks. As an example, a VLSI implementation for finite ring neural network is developed. Our simulation results show their advantages in power, time and area.

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

It is my hope that this book will contribute to our understanding of this new and exciting discipline; ANNs System Engineering.

David D. Zhang
City University of Hong Kong

# Acknowledgements

# CONTENTS

## LIST OF TABLES

# LIST OF FIGURES

# INTRODUCTION

## 1.1   ANN FOR PATTERN RECOGNITION

Artificial Neural Networks (ANN) 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]. 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 are 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.
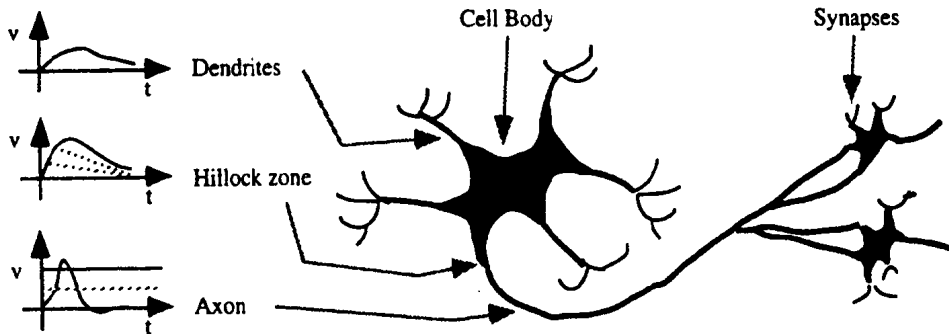
Fig.1.1 A prototype biological neuron

Today 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 McCelland [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 MODEL

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 dentrites are the branch-like structures which provide the sensory input to a cell body (See Fig.1.1). ANNs mimic the functioning of the neural networks of a brain. ANN consists of a large number of simple node. 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.

The connections between the nodes can be bidirectional or unidirectional. The activation function and the activation values to be used in the network are often restricted in the range [0, 1]. In the case of discrete values, the activation values can take only two values – 0 or 1. The number of activation functions can be used to define the propagation law in the network. The commonly used activation functions are shown in Fig.1.4. For a sigmoid function the output at node j is

$$y_j = f(\varphi_j) = \frac{1}{1 + \exp\left(-(\varphi_j + \theta_j)\right)} \qquad (1.1)$$

where $\theta_j$ is a bias and the net input $\varphi_j$ is represented as

$$\varphi_j = \sum_{i=1}^{n} (x_i \, w_{ij}) \qquad (1.2)$$