

EXPLORING COMPLEXITY REDUCTION FOR LEARNING IN DEEP
NEURAL NETWORKS

by

Sourya Dey

PhD Dissertation Proposal
UNIVERSITY OF SOUTHERN CALIFORNIA
(ELECTRICAL ENGINEERING)

April 2019

Copyright 2019

Sourya Dey

*Let yourself be silently drawn by the
strange pull of what you really love.
-Rumi*

Table of Contents

List of Figures	vi
List of Tables	ix
List of Acronyms	x
Acknowledgements	xi
Abstract	xii
Related Publications with links	xiii
1 Introduction	1
1.1 Neural networks	1
1.1.1 Complexity of neural networks	1
1.1.2 Model search	2
1.2 Dissertation Contributions	3
1.2.1 Achieved Contributions	3
1.2.2 Expected Contributions	5
1.3 Proposal Organization	7
2 Background	8
2.1 Mathematical Notation	8
2.2 Notation and Basic Operations for Neural Networks	12
2.2.1 Feedforward (FF)	14
2.2.2 Backpropagation (BP)	18
2.2.3 Update (UP)	19
2.3 Training and Inference	20
2.4 Other types of networks	25
2.4.1 Convolutional neural networks	26
2.4.2 Recurrent neural networks	27

3	Pre-Defined Sparsity	28
3.1	Related Work	28
3.2	Structured Pre-defined sparsity	30
3.2.1	Motivation and Preliminary Examples	34
3.2.2	Structured Constraints	35
3.2.3	Modifications to Neural Network Operations	38
3.3	Performance Results, Trends and Guidelines	40
3.3.1	Datasets and Experimental Configuration	41
3.3.2	Dataset Redundancy	46
3.3.3	Individual junction densities	49
3.3.4	‘Large and sparse’ vs ‘small and dense’ networks	54
3.4	Summary	58
4	Hardware Architecture	59
4.1	Junction pipelining and Operational parallelism	63
4.2	Memory organization	66
4.3	Clash-freedom	68
4.4	Batch size	69
4.5	Architectural Constraints	70
4.6	Special Case: Processing a FC junction	72
4.7	FPGA Implementation	74
4.7.1	Network Configuration and Training Setup	75
4.7.2	Bit Width Considerations	76
4.7.3	Implementation Details	82
4.8	Summary	87
5	Connection Patterns	88
5.1	Biadjacency Matrices	88
5.2	Clash-free memory access patterns	91
5.2.1	Types of memory access patterns	93
5.3	Comparison between classes of Pre-defined Sparsity	97
5.4	Comparison to other methods of sparsity	101
5.5	Metrics for Connection Patterns	102
5.5.1	Window biadjacency matrices	103
5.5.2	Scatter	107
6	Dataset Engineering	113
6.1	Generating Algorithm	116
6.1.1	Variations and Difficulty Scaling	119
6.2	Neural Network Results and Analysis	121
6.2.1	Results	121
6.2.2	Results for Pre-Defined Sparse Networks	125

6.3	Metrics for Dataset Difficulty	128
6.3.1	Goodness of the Metrics	131
6.3.2	Limitations of the Metrics	133
6.4	Summary	134
7	Model Search	135
7.1	Architecture Search	136
7.1.1	Evolutionary Algorithms	137
7.1.2	Reinforcement Learning (RL) methods	138
7.1.3	Other methods of Architecture Search	140
7.2	Hyperparameter optimization	140
7.3	Our Proposed Research	142
7.3.1	Proposed research in Architecture Search	142
7.3.2	Proposed research in Hyperparameter Optimization	146
8	Conclusion	147
8.1	Summary	147
8.2	Future Work	151
8.3	Final Word	151
	References	153

List of Figures

2.1	Complete training flow for a single input sample for a simple MLP.	15
2.2	Comparison of ReLU and sigmoid activations and their derivatives.	17
2.3	Flowchart showing the correct way to use data and set parameters and hyperparameters.	22
2.4	Example of a simple CNN.	26
2.5	Example of a RNN.	27
3.1	Illustrating basic concepts of structured pre-defined sparsity.	33
3.2	Histograms of weight values in different junctions for fully connected networks trained on MNIST, and preliminary examples of the performance effects of pre-defined sparsity.	36
3.3	Comparison of classification accuracy for original and changed redundancy datasets.	48
3.4	Comparing histograms of weight values for original and reduced redundancy versions of MNIST, after training fully-connected networks.	49
3.5	Comparison of classification accuracy as a function of ρ_{net} for different ρ_L , where $L = 2$.	51
3.6	Comparison of classification accuracy as a function of ρ_{net} for ρ_2 vs ρ_3 in three-junction MNIST networks, keeping ρ_1 fixed.	51
3.7	Comparison of classification accuracy as a function of ρ_{net} for datasets with varying redundancy.	52
3.8	Comparing ‘large and sparse’ to ‘small and dense’ networks for MNIST networks with different numbers of junctions.	55
3.9	Comparing ‘large and sparse’ to ‘small and dense’ networks for a Reuters network.	56
3.10	Comparing ‘large and sparse’ to ‘small and dense’ networks for TIMIT and CIFAR networks.	57
4.1	Overview of our hardware architecture.	60
4.2	Architecture for parallel operations for an intermediate junction i ($i \neq 1, L$) showing the three operations along with associated inputs and outputs.	64

4.3	An example of hardware processing inside a sparse junction.	67
4.4	An example of hardware processing inside a fully connected junction.	73
4.5	Maximum absolute values (left y-axis) for all weights, biases, and deltas in the network, and percentage classification accuracy (right y-axis), as the network is trained.	77
4.6	Histograms of absolute value of \mathbf{s}_1 with respect to dynamic range for (a) sparse vs. (b) fully connected cases, as obtained from software simulations. Values right of the pink line are clipped.	80
4.7	Comparison of activation functions for \mathbf{a}_1 , as obtained from software simulations.	81
4.8	Performance for different ρ_2 , keeping ρ_1 fixed at 6.25%, as obtained from software simulations.	85
4.9	Breaking up each operation into 3 clock cycles.	86
4.10	Our design working on the Xilinx XC7A100T-1CSG324C FPGA.	87
5.1	Different connection patterns arising in a structured pre-defined sparse network.	89
5.2	Repeating Fig. 4.3 for convenience.	93
5.3	Various types of clash-free left memory access patterns and memory dithering.	95
5.4	An example of biadjacency matrices and equivalent junctions.	103
5.5	Different types of windowing in left neurons.	104
5.6	Window biadjacency matrices and scatter.	106
5.7	Constructing window biadjacency matrices.	107
5.8	Performance vs.scatter.	110
6.1	Generating the Morse codeword $\bullet - \bullet - \bullet$ corresponding to the $+$ symbol.	119
6.2	Classification performance for different variations of the Morse datasets.	122
6.3	Effects of noise leading to spaces getting confused with dots and dashes.	123
6.4	Effects of increasing the size of <i>Morse 3.1</i> by a factor of x on test accuracy after 30 epochs (blue), and (Training Accuracy – Test Accuracy) after 30 epochs (orange).	124
6.5	Effects of imposing pre-defined sparsity on classification performance for different Morse datasets.	126
6.6	Validation performance results for varying ρ_1 and ρ_2 individually so as to keep ρ_{net} fixed at (a) 25%, (b) 50%.	127
6.7	Plotting each metric for dataset difficulty vs. percentage accuracy obtained for different Morse datasets.	132
7.1	Types of model search.	136

8.1	Trends in AI papers on Scopus.	148
8.2	Trends in AI papers on arXiv.	149
8.3	Parameter explosion in deep neural networks over the years.	149

List of Tables

4.1	Hardware Architecture Total Storage Cost Comparison for FC vs. sparse cases.	65
4.2	Implemented Network Configuration	76
4.3	Effect of Bit Width on Performance	79
5.1	Comparison of Clash-free Left Memory Access Types and associated Hardware Cost for a Single Junction i with $(N_{i-1}, N_i, d_i^{\text{out}}, d_i^{\text{in}}, z_i) = (12, 12, 2, 2, 4)$	97
5.2	Comparison of Pre-Defined Sparse Classes	100
6.1	Correlation coefficients between metrics and accuracy	132
7.1	Building on and extending our research to reduce the complexity of architecture search.	144

List of Acronyms

AI	Artificial Intelligence
ASIC	Application Specific Integrated Circuit
ASR	Automatic Speech Recognition
BRAM	Block Random Access Memory
CI	Confidence Interval
CNN	Convolutional Neural Network
DSP	Digital Signal Processing
FC	Fully Connected
FF	Feedforward
FPGA	Field Programmable Gate Array
gcd	Greatest Common Divisor
GPU	Graphics Processing Unit
LDPC	Low Density Parity Check
LED	Light Emitting Diode
LUT	Look Up Table
MFCC	Mel-frequency Cepstral Coefficient
MLP	Multilayer Perceptron
NN	Neural network
PCA	Principal Component Analysis
ReLU	Rectified Linear Unit
RL	Reinforcement Learning
RNN	Recurrent Neural Network
RTL	Register Transfer Level
SGD	Stochastic Gradient Descent
TPC	Test Prediction Comparison
UART	Universal asynchronous receiver-transmitter

Acknowledgements

Pursuing a Ph.D. is a beautiful thing since it allows one to be free to explore uncharted territory instead of being burdened by constraints of banality. In this regard, I would like to thank my co-advisors – Professors Peter Beerel and Keith Chugg – without whose help and guidance I would not have been able to explore the research directions that have led to the material in this dissertation proposal.

I am grateful to past and present members of my research team who have helped and collaborated with me in my research. Special thanks to my primary collaborators – Kuan-Wen Huang, Yinan Shao, Diandian Chen, and Souvik Kundu. I also acknowledge the contribution of secondary collaborators who have since moved on – Zongyang Li, Saad Zafar, Mahdi Jelodari Mamaghani, Zheng Wu, Jiajin Xi, and Venkata Nishanth Narisetty; and current research team members whose constant feedback has been invaluable – Professor Leana Golubchik, Marco Paolieri, Arnab Sanyal, and Kiran Nagendra. I am also indebted to Professor Panayiotis Georgiou and Mohammed Nasir for help in specific efforts.

I would like to thank the National Science Foundation, Software and Hardware Foundations, for funding my research under Grant 1763747. I am grateful to Diane Demetras and Annie Yu for their help in administrative matters related to the progress of my Ph.D. and presentation of my research to the outside world.

Heartfelt thanks to my family members for their constant love and support, despite being located halfway around the world. Nothing would have been possible without them.

Finally, thanks to *you* the reader for picking up this dissertation proposal detailing my research efforts since spring 2016. I hope you have as much enjoyment reading it as I had writing it.

Abstract

Neural networks have proven to be extremely powerful tools for cutting-edge machine learning and artificial intelligence applications such as computer vision and speech recognition. Modern neural networks are usually deep, *i.e.* comprising many trainable parameters across several interconnected layers, which gives rise to significant computational complexity, particularly during their training phase. Such complexity is often prohibitively large and concentrates the deployment of deep neural networks only to the few entities with the means to procure the necessary financial and computational resources. Moreover, the study of neural networks is largely empirical and largely guided by the underlying current of thought “we did this because it works”, instead of “this is *why* it works”.

In this dissertation proposal, we propose *pre-defined sparsity* – a technique to reduce the computational complexity of neural networks during both training (*i.e.* learning) and inference phases. We analyze the resulting sparse connection patterns in an attempt to understand network performance. We also introduce a hardware architecture compatible with pre-defined sparsity to reduce the storage complexity of neural networks and potentially accelerate training. The architecture is flexible in the sense that the complexity of the network to be implemented is decoupled from the available hardware resources. Additionally, we engineer a family of synthetic datasets of algorithmically customizable difficulty for benchmarking neural networks and machine learning methodologies. Finally, we propose expected contributions in *model search*, which involves understanding how and why different network structures and methods work, and consequently selecting appropriate network configurations of low complexity for different applications.

Related Publications with links

Note that the citation counts in the URLs linked below may be inaccurate. Please see the author's [Google Scholar page](#) for a full list of citations.

S. Dey, K. W. Huang, P. A. Beerel and K. M. Chugg, “Pre-defined Sparse Neural Networks with Hardware Acceleration,” in *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* (Early Access), 2019. doi: 10.1109/JET-CAS.2019.2910864.

IEEE: <https://ieeexplore.ieee.org/document/8689061>

S. Dey, Y. Shao, K. M. Chugg and P. A. Beerel, “Accelerating training of deep neural networks via sparse edge processing,” in *26th International Conference on Artificial Neural Networks (ICANN)* Part 1, pp. 273–280, Springer, 2017.

Springer: https://link.springer.com/chapter/10.1007/978-3-319-68600-4_32

S. Dey, P. A. Beerel and K. M. Chugg, “Interleaver design for deep neural networks,” in *51st Annual Asilomar Conference on Signals, Systems, and Computers (ACSSC)*, pp. 1979–1983, Oct 2017.

IEEE: <https://ieeexplore.ieee.org/document/8335713>

S. Dey, K. W. Huang, P. A. Beerel and K. M. Chugg, “Characterizing sparse connectivity patterns in neural networks,” in *2018 Information Theory and Applications Workshop (ITA)*, pp. 1–8, Feb 2018.

IEEE: <https://ieeexplore.ieee.org/document/8502950>

[Awarded Best Paper] **S. Dey**, K. M. Chugg and P. A. Beerel, “Morse Code Datasets for Machine Learning,” in *9th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, pp. 1-7, Jul 2018.

IEEE: <https://ieeexplore.ieee.org/document/8494011>

S. Dey, D. Chen, Z. Li, S. Kundu, K. W. Huang, K. M. Chugg and P. A. Beerel, “A Highly Parallel FPGA Implementation of Sparse Neural Network Training,” in

*2018 International Conference on Reconfigurable Computing and FPGAs (ReCon-
Fig)*, pp. 1–4, Dec 2018. Expanded preprint version available on [arXiv](#).
IEEE: <https://ieeexplore.ieee.org/document/8641739>

Chapter 1

Introduction

1.1 Neural networks

Neural network (NN)s have proven to be ubiquitous in modern machine learning and Artificial Intelligence (AI) applications such as object classification [1], self-driving cars [2], Automatic Speech Recognition (ASR) [3], and so on. Such NNs are usually deep, implying that there are hidden layers between the input (such as a picture of a cat), and the output (such as a probability distribution indicating what animal the picture is). The depth of a NN can vary from just one or two hidden layers, as is common for a Multilayer Perceptron (MLP) [4], to several hundred for a Convolutional Neural Network (CNN) [5]. Throughout this work, we will use the term neural network (NN) to denote a deep neural network, *i.e.* having at least one hidden layer.

1.1.1 Complexity of neural networks

NNs used for classification typically rely on large amounts of labeled data to train – a process during which the numerical values of its internal parameters are tuned

to facilitate better inference – measured as classification performance on unseen data. As more data have become available, the size and complexity of NNs have grown sharply, with modern NNs containing millions [1] or even billions of trainable parameters [6]. These massive NNs come with the cost of large computational and storage demands. The current state of the art is to train large NNs on Graphics Processing Unit (GPU)s in the cloud – a process that can take days to weeks even on powerful GPUs [1, 6, 7] or similar programmable processors with multiply-accumulate accelerators [8]. Therefore, several prominent researchers [9] have identified complexity reduction as a key to NN acceleration, which is a general term referring to speeding up the deployment and operation of a NN. Acceleration is generally performed post-training to reduce complexity of inference only, *e.g.* methods for quantization, compression, and grouping parameters [10–13]. However, reduction in training complexity of NNs remains largely unexplored.

1.1.2 Model search

There are a number of choices which play a crucial role during the deployment of NNs to solve specific problem(s). Some of these are the number of layers, the structure of the layers, how fast the network should learn, and so on. These form the domain of model search. Unfortunately, the design of NNs is a largely empirical process and there is no clear understanding regarding, for example, how many layers and neurons are appropriate for classifying a picture of a handwritten digit

into one out of ten classes. There have been efforts in model search, which can be broadly classified into two categories – a) architecture search to determine the optimal structure of a NN [14–18], and b) hyperparameter search and scheduling [19–22]. The former effort remains computation-intensive, for example, [16] requires 800 GPUs running for 21–28 days to arrive at a good architecture. There have been several algorithms relating to the latter effort; and both are active areas of research.

1.2 Dissertation Contributions

The contributions of this dissertation proposal are divided broadly into two categories – a) achieved, which are focused around complexity reduction, and b) expected, which are focused around model search by leveraging the complexity reduction techniques of the achieved contributions.

1.2.1 Achieved Contributions

This is a list of the contributions we have already achieved and published in various conferences and journals:

1. We proposed the novel technique of **pre-defined sparsity** – a method to reduce the complexity of NNs during both training and inference phases. Our results show that NN complexity, both in terms of the number of parameters

(storage) and the number of operations to be performed (computation), can be reduced by factors greater than 5X without significant performance loss [23–25].

2. We analyzed **trends and design guidelines** for selecting a pre-defined sparse NN [23, 24]. These studies help to accelerate the search for good pre-defined sparse NNs given any problem.
3. We proposed a **hardware architecture** to leverage the benefits of pre-defined sparsity [23, 25]. The architecture is flexible in the sense that the complexity of the NN to be implemented is decoupled from the available hardware resources on the given device (such as Field Programmable Gate Array (FPGA)(s)). Thus, NNs of varying sizes can be supported on various hardware platforms. Moreover, the architecture supports both training and inference phases of a NN. To the best of our knowledge, we are the first to propose such a flexible hardware architecture with the potential to accelerate both training and inference.
4. We developed a working **FPGA implementation** of the hardware architecture [26]. This serves as a proof-of-concept of our ideas.
5. We proposed and analyzed the properties of a class of pre-defined sparse NN **connection patterns** which are suited to our hardware architecture in the sense that they allow maximum throughput [23, 27]. In particular, we

showed that such hardware-friendly sparse patterns result in no performance degradation compared to other sparse patterns.

6. We developed a metric called *scatter* to characterize the ‘goodness’ of a connection pattern in order to **predict NN performance prior to training**. This is helpful in filtering out bad connection patterns without incurring the computational expenses of training.
7. We have also created a **family of synthetic datasets** on classification of Morse code symbols [28]. Owing to their algorithmic generation capability, it is simple to tune the parameters of these datasets. This results in a large number of classification problems of varying difficulty, which can be used to benchmark different NNs and other machine learning algorithms. This work resulted in a ‘Conference Best Paper’ award. The datasets are open-source and available at [29].

1.2.2 Expected Contributions

This is a list of the contributions we expect to make:

1. We will improve the existing hardware implementation in [26]. In particular, the author of this dissertation proposal plans to work on pipelining to **improve operational speed of the hardware implementation**. Other

members of our research team are working on better memory management, and using cloud FPGA resources to support more powerful NNs.

2. We will study techniques for **architecture search** of NNs. This implies deciding on the number of layers, the number of neurons, and different methods of connecting layers in a network. The plan is to leverage our achieved contributions, *i.e.* the underlying theme of the search will be low-complexity (sparse) NNs which lead to good performance.
3. We will analyze the properties of good NN architectures to gain a **deeper understanding of NNs**, such as the crucial parameters and appropriate sparsity levels. This will help us to develop indicators and monitors to predict NN performance. These efforts will expand our achieved contributions of trends in pre-defined sparse connection patterns, and predicting NN performance without having to fully train. We can collectively describe this and the previous expected contribution as **automating the search for good low-complexity NNs**.
4. *Time and resources permitting*, we also plan to work on **hyperparameter configuration** strategies which may lead to faster convergence and better training of NNs. This work will be specifically geared towards hyperparameters determining sparsity and complexity levels in a NN. For example, a

sparse NN needs less regularization as compared to a dense NN (see Chapter 2 for regularization and other related definitions).

1.3 Proposal Organization

This dissertation proposal is organized as follows. Chapter 2 contains background necessary for understanding NNs. Chapter 3 introduces pre-defined sparsity, along with trends and guidelines for selecting connection patterns. Chapter 4 discusses our hardware architecture and its FPGA implementation. Chapter 5 analyzes sparse connection patterns, and metrics to predict NN performance. Chapter 6 discusses our efforts in creating a family of synthetic datasets, along with metrics to characterize their difficulty. Chapter 7 introduces model search of NNs, which pertains to our expected contributions, and finally, Chapter 8 concludes this work.

Chapter 2

Background

In this chapter, we provide the background necessary for understanding the concepts and terms used in the majority of this work. We will begin with mathematical notation, then move on to notation and operations specific to neural networks. We will not provide background on model search in this chapter since it is not necessary to understand our achieved contributions, instead we postpone it to Chapter 7.

2.1 Mathematical Notation

We will use the numerator layout convention for matrix calculus. A complete set of rules for this convention is given in [30]. We summarize some of the key rules below:

- Vectors are written as lower case bold letters, such as \mathbf{x} , and can be either row (dimensions $1 \times n$) or column (dimensions $n \times 1$). Column vectors are the default choice, unless otherwise mentioned. Individual elements are indexed by bracketed superscripts, *e.g.* $x^{(i)}$, where $i \in \{1, \dots, n\}$.

- Matrices are written as upper case bold letters, such as \mathbf{X} . A matrix with dimensions $m \times n$ corresponds to m rows and n columns. Individual elements are indexed by bracketed double superscripts for row and column, respectively, *e.g.* $X^{(i,j)}$, where $i \in \{1, \dots, m\}$, $j \in \{1, \dots, n\}$.
- The **derivative** of f with respect to x is $\frac{\partial f}{\partial x}$, where both x and f can be scalars, vectors, or matrices. The **gradient** of f with respect to x is $\nabla_x f = \left(\frac{\partial f}{\partial x}\right)^T$, *i.e. gradient is transpose of derivative*. The derivative is important because it obeys the chain rule of calculus and helps to derive several results. The gradient is important because it is the form used in updating NN parameters, as will be described in Section 2.2.3.

Some relevant forms of derivatives and gradients are described next.

Scalar-by-scalar

Both f and x are scalars. We do not define the gradient in this case. The derivative is the scalar $\frac{\partial f}{\partial x}$, also written as f' .

Scalar-by-vector

f is a scalar, \mathbf{x} is a vector of dimensions $n \times 1$. Then the derivative is a $1 \times n$ row vector:

$$\frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x^{(1)}} \quad \frac{\partial f}{\partial x^{(2)}} \quad \dots \quad \frac{\partial f}{\partial x^{(n)}} \right] \quad (2.1)$$

and the gradient $\nabla_{\mathbf{x}} f$ is its transposed column vector.

Vector-by-vector of equal size

Both \mathbf{f} and \mathbf{x} are vectors of dimensions $n \times 1$. Then the derivative is the Jacobian matrix of dimensions $n \times n$:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f^{(1)}}{\partial x^{(1)}} & \cdots & \frac{\partial f^{(1)}}{\partial x^{(n)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f^{(n)}}{\partial x^{(1)}} & \cdots & \frac{\partial f^{(n)}}{\partial x^{(n)}} \end{bmatrix} \quad (2.2)$$

Vectorized scalar function

This is a scalar-to-scalar function applied element-wise to a vector, *e.g.* :

$$f \left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \right) = \begin{bmatrix} f(x^{(1)}) \\ f(x^{(2)}) \\ \vdots \\ f(x^{(n)}) \end{bmatrix} \quad (2.3)$$

In this case, both the derivative and gradient are the same $n \times n$ diagonal matrix, given as:

$$\nabla_x \mathbf{f} = \begin{bmatrix} f'(x^{(1)}) & & & 0 \\ & f'(x^{(2)}) & & \\ & & \ddots & \\ 0 & & & f'(x^{(n)}) \end{bmatrix} \quad (2.4)$$

An equivalent form of this is to take the diagonal and express it as a $n \times 1$ vector:

$$\mathbf{f}' = \begin{bmatrix} f'(x^{(1)}) \\ f'(x^{(2)}) \\ \vdots \\ f'(x^{(n)}) \end{bmatrix} \quad (2.5)$$

To realize the effect of this equivalent form, let's say we want to multiply the gradient $\nabla_x \mathbf{f}$ from (2.4) with some n -dimensional (column) vector \mathbf{a} . Achieving

the same result with \mathbf{f}' from (2.5) would require the *Hadamard product* \circ , defined as element-wise multiplication of 2 vectors:

$$(\nabla_{\mathbf{x}} \mathbf{f}) \mathbf{a} = \mathbf{f}'(\mathbf{x}) \circ \mathbf{a} = \begin{bmatrix} f'(x^{(1)}) a^{(1)} \\ f'(x^{(2)}) a^{(2)} \\ \vdots \\ f'(x^{(n)}) a^{(n)} \end{bmatrix} \quad (2.6)$$

2.2 Notation and Basic Operations for Neural Networks

A NN, sometimes referred to as an Artificial Neural Network to distinguish it from biological neural networks found in living organisms, is an interconnected set of **nodes, or neurons**, with the capability to learn mathematical representations of the data fed to it. This ability to learn is referred to as **training**. Once a NN is trained, it may be used for **inference**, where it is operated on new data different from what it was trained on, and its performance based on some metric(s) is noted.

We will use a **Multilayer Perceptron (MLP)**, which is a commonly used NN for classification problems, and assume **supervised learning**, *i.e.* where labels are provided for different training samples and the NN is asked to learn from them. Our work also extends to other networks and forms of learning, however,

to avoid burdening the reader, we will initially introduce notation for MLPs using supervised learning and then extend as required.

A set of nodes in a MLP is referred to as a **layer**. An $(L + 1)$ -layer MLP has N_i nodes in the i^{th} layer, described collectively by the **neuronal configuration** $\mathbf{N}_{\text{net}} = (N_0, N_1, \dots, N_L)$,¹ where layer 0 is the input layer. We use the convention that layer i is to the ‘right’ of layer $i - 1$, or ‘next’ to layer $i - 1$.

There are L **junctions** between layers, with junction i ($i \in \{1, 2, \dots, L\}$) connecting the N_{i-1} nodes of its left layer $i - 1$ with the N_i nodes of its right layer i . For a **Fully Connected (FC)** MLP, all the nodes in a layer connect to all the nodes in its next and previous layers, if present. These connections are defined using **edges**, which have associated **weight** values. Weights in junction i are collected as a matrix \mathbf{W}_i . Note that the rows (1st index) refer to the next layer, and the columns (2nd index) refer to the previous layer. This is to facilitate matrix-vector multiplications, as will be seen in (2.7a). Additionally, individual nodes in all layers except for the input have **bias** values. For layer i , the biases are collected as a vector \mathbf{b}_i .

Note that we denote layer or junction number via subscript. This is why we denote individual elements of matrices and vectors as bracketed superscripts, as

¹We make a distinction between an ordered tuple of numbers providing information, such as \mathbf{N}_{net} , and a vector, such as \mathbf{x} . The former is not meant to be used in algebraic calculations, hence its individual elements are written with round brackets around them. The latter is used in algebra and hence its individual elements are written with the more conventional square brackets around them, as in the equations presented so far in this chapter.

shown in Section 2.1. For example $b_2^{(3)}$ refers to the bias of the 3rd neuron in layer 2, and $W_1^{(4,5)}$ refers to the weight of the edge connecting the 5th node in layer 0 to the 4th node in layer 1. This notation suits us since throughout this work we will need to reference layer numbers far more than individual weight or bias values, and hence benefit from using subscripts for layer number instead of the somewhat inconvenient bracketed superscripts.

Weights and biases are all **trainable parameters**, which means that their values change during training. This is usually done until the values converge, *i.e.* the change in values obtained by training on more data is negligible. In practical scenarios, training is often stopped due to time or computational constraints instead of waiting for convergence. The process of training is affected by some values which the network does not learn, instead, they are adjusted by an entity external to the NN such as a human. These values are known as **hyperparameters**, and describe quantities such as how fast weight values should change, how many inputs should be trained on in parallel, and so on. Next we describe the three operations of a NN. These are summarized in Fig. 2.1.

2.2.1 Feedforward (FF)

This process starts by accepting an external input datum. Each input is represented by a vector of **features** \mathbf{a}_0 , and a ground truth labeling \mathbf{y} denoting

$N_0 \times 1$ $N_L \times 1$

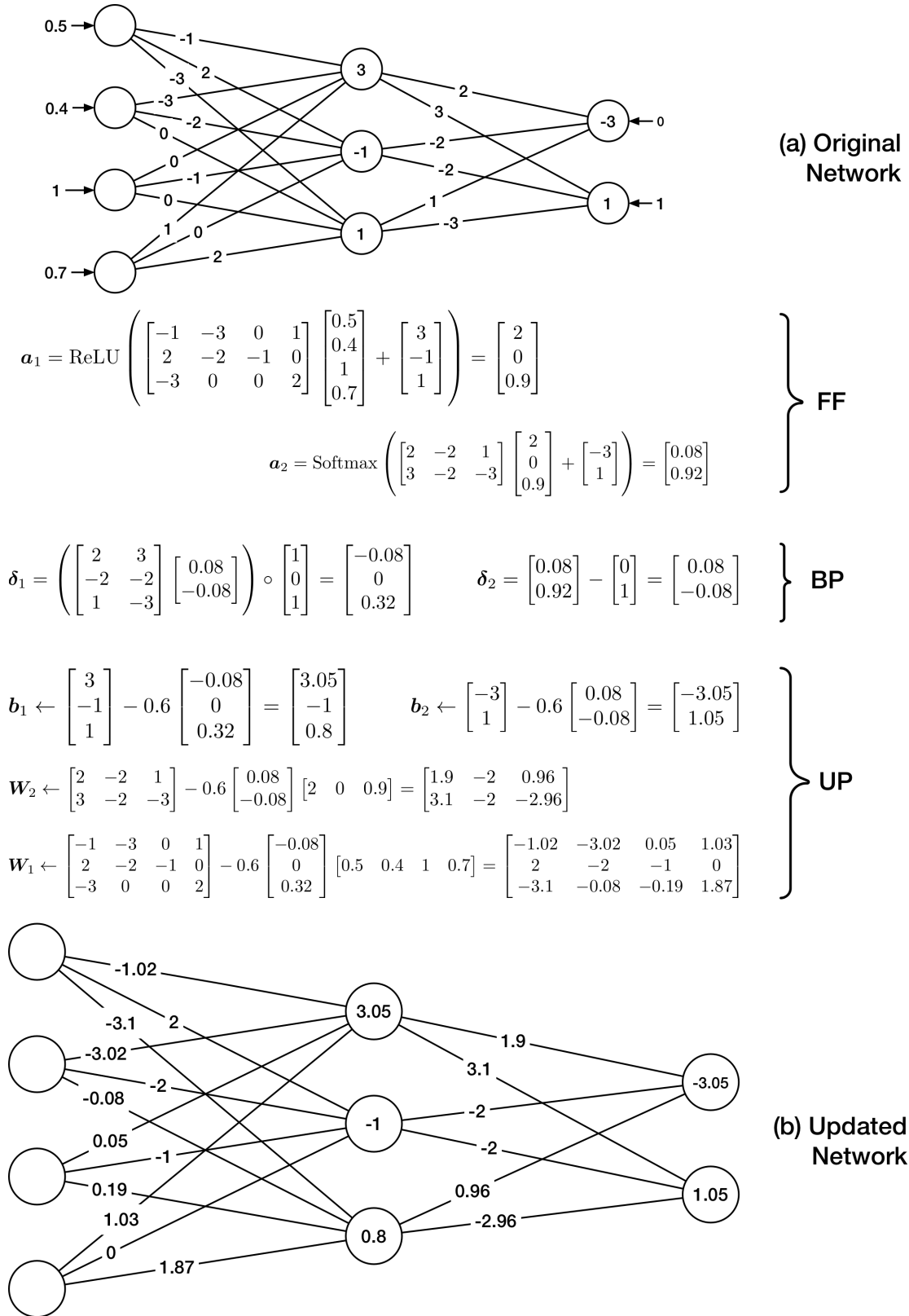


Figure 2.1: Complete training flow for a single input sample for a MLP with $\mathbf{N}_{\text{net}} = (4, 3, 2)$. (a) Original NN, presented with an input $\mathbf{a}_0 = [0.5, 0.4, 1, 0.7]^T$ and corresponding one-hot ground-truth labeling $\mathbf{y} = [0, 1]^T$. FF proceeds with ReLU and Softmax activations for junctions 1 and 2 respectively, BP uses cross-entropy cost, and UP uses $\eta = 0.6$. No regularization is applied and batch size $M = 1$. (b) After updating the weights and biases of the original NN.

the ideal output layer values. Then, $\forall i \in \{1, 2, \dots, L\}$, the Feedforward (FF) operation proceeds as:

$$\mathbf{s}_i = \mathbf{W}_i \mathbf{a}_{i-1} + \mathbf{b}_i \quad (2.7a)$$

$$\mathbf{a}_i = \mathbf{h}(\mathbf{s}_i) \quad (2.7b)$$

$$\mathbf{H}'_i = \frac{\partial \mathbf{a}_i}{\partial \mathbf{s}_i} \quad (2.7c)$$

where \mathbf{s} is the linear output, \mathbf{h} is the non-linear **activation function**, and \mathbf{a} is the **activation** output. Note that (2.7c) is not strictly a part of FF, but it is computed since the derivative values are required for BP.

For most layers, \mathbf{h} is a vectorized scalar function applied to \mathbf{s}_i . In such cases, \mathbf{H}'_i is a diagonal matrix and can be reduced to the vector \mathbf{h}'_i , as shown in Section 2.1.

One such common activation function which we will use frequently is **Rectified Linear Unit (ReLU)** [31], shown below in its scalar-to-scalar form:

$$\text{ReLU}(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (2.8a)$$

$$\text{ReLU}'(x) = \begin{cases} 1, & x > 0 \\ [0, 1], & x = 0 \\ 0, & x < 0 \end{cases} \quad (2.8b)$$

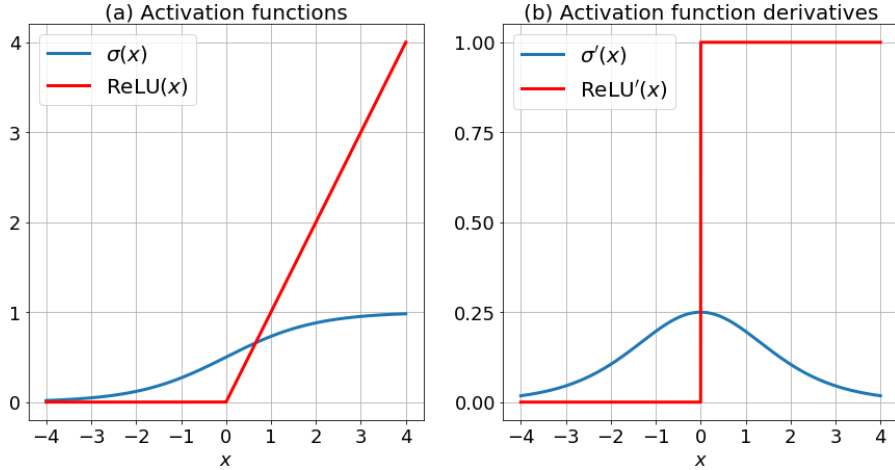


Figure 2.2: Comparison of ReLU and sigmoid (a) activations, and (b) activation derivatives.

An alternative to ReLU is the sigmoid function given as:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (2.9a)$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (2.9b)$$

Fig. 2.2 illustrates ReLU and sigmoid.

For the last layer ($i = L$) however, \mathbf{h} is usually the **softmax** function which outputs a probability distribution. This is a vector-to-vector function where each element of the output vector depends on each element of the input vector, hence \mathbf{H}'_L is not (necessarily) a diagonal matrix.

$$\mathbf{a}_L = \text{Softmax}(\mathbf{s}_L) = \left[\frac{e^{s_L^{(1)}}}{\sum_{i=1}^{N_L} e^{s_L^{(i)}}} \quad \frac{e^{s_L^{(2)}}}{\sum_{i=1}^{N_L} e^{s_L^{(i)}}} \quad \cdots \quad \frac{e^{s_L^{(N_L)}}}{\sum_{i=1}^{N_L} e^{s_L^{(i)}}} \right]^T \quad (2.10)$$

2.2.2 Backpropagation (BP)

The final layer activations \mathbf{a}_L are compared with the ground-truth labels \mathbf{y} to compute a scalar **cost** (or loss) function C . The cost function we use is **cross-entropy**:

$$C = - \sum_{i=1}^{N_L} y^{(i)} \ln a_L^{(i)} \quad (2.11)$$

The **delta** (or error) values for every layer are computed next. These are the gradients of the cost with respect to \mathbf{s} :

$$\boldsymbol{\delta}_L = \left(\frac{\partial C}{\partial \mathbf{s}_L} \right)^T = \mathbf{H}'_L{}^T \nabla_{\mathbf{a}_L} C \quad (2.12a)$$

$$\boldsymbol{\delta}_i = \left(\frac{\partial C}{\partial \mathbf{s}_i} \right)^T = \mathbf{H}'_i{}^T \mathbf{W}_{i+1}^T \boldsymbol{\delta}_{i+1} \quad \forall i \in \{1, 2, \dots, L-1\} \quad (2.12b)$$

The above equations can be derived using the chain rule for derivatives.

Classification is often done using **one-hot labels**, *i.e.* all the elements in \mathbf{y} are 0s except for the correct class, which has value 1. For example, when classifying an image of a digit into one out of ten possible classes 0 – 9, if the correct class for a particular input sample is 3, its one-hot ground truth labeling would be $\mathbf{y} = [0, 0, 0, 1, 0, 0, 0, 0, 0, 0]^T$. For the commonly occurring case of one-hot labels, cross-entropy cost, softmax activation for the output layer, and a scalar-to-scalar

activation applied element-wise to all other layers such as ReLU, (2.12) simplifies to:

$$\underset{N_L \times 1}{\boldsymbol{\delta}_L} = \mathbf{a}_L - \mathbf{y} \quad (2.13a)$$

$$\underset{N_i \times 1}{\boldsymbol{\delta}_i} = (\mathbf{W}_{i+1}^T \boldsymbol{\delta}_{i+1}) \circ \mathbf{h}'_i \quad \forall i \in \{1, 2, \dots, L-1\} \quad (2.13b)$$

Note that the sigmoid derivative has a maximum value of 0.25, as shown in Fig. 2.2(b). As a result, the delta values can become very small when backpropagating through deep networks. This phenomenon is known as the vanishing gradient problem [32], and is the reason why sigmoid activations have generally fallen out of favor with MLPs.

2.2.3 Update (UP)

The goal of training a NN is to minimize the cost. An optimization algorithm commonly used for achieving this goal is **Stochastic Gradient Descent (SGD)**, wherein the value of a particular trainable parameter of the NN is updated using the gradient of the cost with respect to it. The gradients for a single input sample are given as $\nabla_{\mathbf{w}_i} C = \boldsymbol{\delta}_i \mathbf{a}_{i-1}^T$ and $\nabla_{\mathbf{b}_i} C = \boldsymbol{\delta}_i$. Knowing the true gradient requires knowledge of the underlying data distribution, which is usually not obtainable. In SGD, the true gradients are approximated by averaging the single sample gradients

over a **batch** (sometimes called minibatch) of M input samples. This gives the SGD update rule $\forall i \in \{1, 2, \dots, L\}$ as:

$$\mathbf{W}_i \leftarrow \mathbf{W}_i - \frac{\eta}{M} \sum_{m=1}^M (\boldsymbol{\delta}_i \mathbf{a}_{i-1}^T)^{[m]} \quad (2.14a)$$

$$\mathbf{b}_i \leftarrow \mathbf{b}_i - \frac{\eta}{M} \sum_{m=1}^M (\boldsymbol{\delta}_i)^{[m]} \quad (2.14b)$$

where η is the **learning rate** hyperparameter determining how fast the network should learn, and square-bracketed superscript m denotes input sample number m (not to be confused with round-bracketed superscripts for individual elements in a vector or matrix).

While trainable parameters refers to weights and biases, we will use the term **network parameters** to collectively refer to $\{\mathbf{a}_i, \mathbf{h}'_i, \boldsymbol{\delta}_i, \mathbf{W}_i, \mathbf{b}_i\} \forall i \in \{1, 2, \dots, L\}$.

² These network parameters account for the total storage cost.

2.3 Training and Inference

A NN is initially trained using all the operations described, *i.e.* FF, BP and UP. Inputs are traversed in batches until all inputs in the training dataset are exhausted, this constitutes an **epoch** of training. A NN can be trained for dozens

²The range of \mathbf{h}'_i is from 1 to $L - 1$ since final layer softmax derivatives are directly combined with cost derivatives to yield $\boldsymbol{\delta}_L$.

or even hundreds of epochs depending on the performance desired. Two metrics for determining performance are a) the cost, which should be minimized, and more commonly, b) the classification accuracy, computed as the fraction of inputs correctly classified, which should be maximized.

Validation is commonly done after every epoch to determine how the NN will behave on unseen data. This validation data is a separate subset of the training data. The hyperparameters are adjusted according to validation performance. For example, if training performance keeps improving but validation performance deteriorates, the NN is **overfitting**. This implies that the NN is learning its training data ‘too well’ and failing to generalize on unseen data. This typically happens when the number of trainable parameters is far more than the number of input data samples. In such cases, training should be stopped, or regularization may prove to be helpful.

Once training is complete, the NN performs inference on test data. Inference performance, such as **classification accuracy on test data**, often serves as the defining measure for how good a NN is. Note that inference only involves computation of (2.7a) and (2.7b), and hence is of much lower complexity than training. The complete cycle of using data for training and inference is summarized in Fig. 2.3. Next we discuss some hyperparameters and design choices we used when performing experiments.

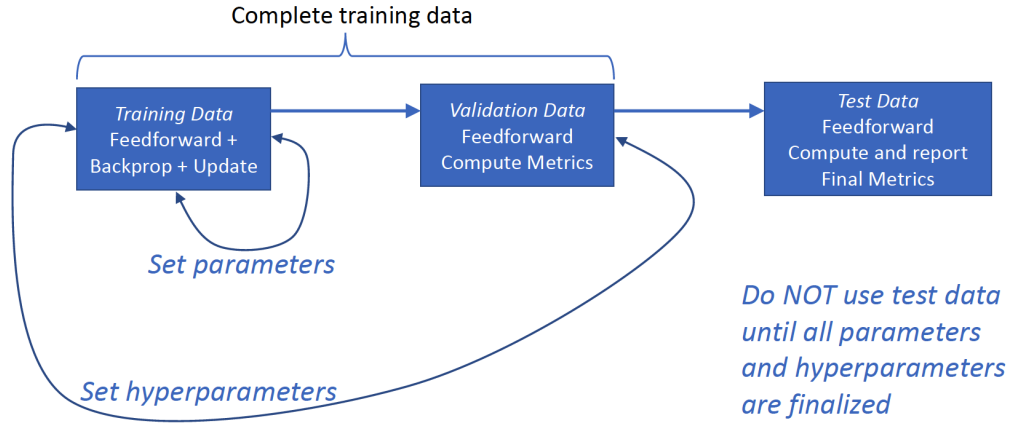


Figure 2.3: Flowchart showing the correct way to use data and set parameters and hyperparameters.

Regularization

The idea behind regularization is to impose an additional constraint on the cost objective to make the NN generalize better to unseen data (*i.e.* validation and test), and combat overfitting. A typical way to do this is by imposing an **L2 penalty** on the magnitudes of the weights. The regularized cost function becomes:

$$C = C_0 + \lambda \|\mathbf{W}\|_2^2 \quad (2.15)$$

where C_0 is the original cost function such as given in (2.11), λ is a hyperparameter denoting the importance of making the weights small relative to minimizing the original cost function, and \mathbf{W} is the vector containing all the weights across all junctions of the network (*i.e.* concatenation of the flattened matrices $W_i \quad \forall i \in \{1, 2, \dots, L\}$). Typical values for λ range between 10^{-3} to 10^{-6} .

Regularization helps to improve performance particularly when NN is over-parametrized, *i.e.* it has a large number of trainable parameters. As we shall see later in this work, imposing pre-defined sparsity is a form of regularization since it reduces the number of trainable parameters.

Initializing trainable parameters

Initializing all weights to a constant such as 0 usually slows down learning. A popular technique to initialize weights is according to a normal distribution, such as **Glorot Normal** [33]

$$W_i \sim \mathcal{N}\left(0, \frac{2}{d_i^{\text{in}} + d_i^{\text{out}}}\right) \quad (2.16)$$

or **He Normal** [34]:

$$W_i \sim \mathcal{N}\left(0, \frac{2}{d_i^{\text{in}}}\right) \quad (2.17)$$

where $\mathcal{N}(\mu, \sigma^2)$ denotes the Normal distribution with mean μ and variance σ^2 .

Biases can be initialized with 0s, or a small positive constant to help ReLU units remain active (*i.e.* not operating on inputs less than 0).

Optimizer

SGD, as given in (2.14), is a baseline optimizer for reducing cost. Other adaptive optimizers extend SGD by incorporating techniques such as momentum and bias correction (see [35] for an overview of these techniques). A popular optimizer is **Adam** [36], which has default values ($\eta = 0.001, \rho_1 = 0.9, \rho_2 = 0.999, \epsilon = \text{machine epsilon}$), and works as shown below (for any arbitrary trainable parameter p).

$$v_1 \leftarrow \rho_1 v_1 + (1 - \rho_1) \nabla_p C \quad (2.18a)$$

$$v_2 \leftarrow \rho_2 v_2 + (1 - \rho_2) (\nabla_p C)^2 \quad (2.18b)$$

$$\tilde{v}_1 = \frac{v_1}{1 - \rho_1^t} \quad (2.18c)$$

$$\tilde{v}_2 = \frac{v_2}{1 - \rho_2^t} \quad (2.18d)$$

$$p \leftarrow p - \frac{\eta}{\sqrt{\tilde{v}_2} + \epsilon} \tilde{v}_1 \quad (2.18e)$$

where t is time, *i.e.* ρ_1^t is ρ_1 raised to the t th power. Time is incremented after each update (*i.e.* after the value of p changes). The time step at which the above equations occur is $t + 1$, and the initial values for v_1 and v_2 are both 0. Note that the final equation (2.18e) is similar to the regular SGD update equation with $M = 1$.

A possible modification to Adam is to add a decay parameter d such that for time step $t+1$, instead of using η , the optimizer uses $\frac{\eta}{1+dt}$. This leads to smoother convergence.

Batch Normalization

Batch normalization [37] has been found to be an effective in training deep NNs. Typically intermediate outputs are normalized before applying the activation function. Any variable x is normalized by subtracting its mean and dividing by its standard deviation over a batch of samples, *i.e.* $x \leftarrow \frac{x - \mu_x}{\sigma_x}$. Following this, the normalized variable is subjected to an affine transformation, *i.e.* $x \leftarrow \gamma x + \beta$, where γ and β are trainable parameters.

2.4 Other types of networks

We introduced several core concepts of NNs using MLPs for classification. This is because the primary focus of our achieved contributions are MLPs. However, we have also experimented on other kinds of networks, and will continue to do so as part of our expected contributions. A quick introduction to some other commonly used forms of NNs are given below. The reader is encouraged to refer to [38] for further exposition of related concepts.

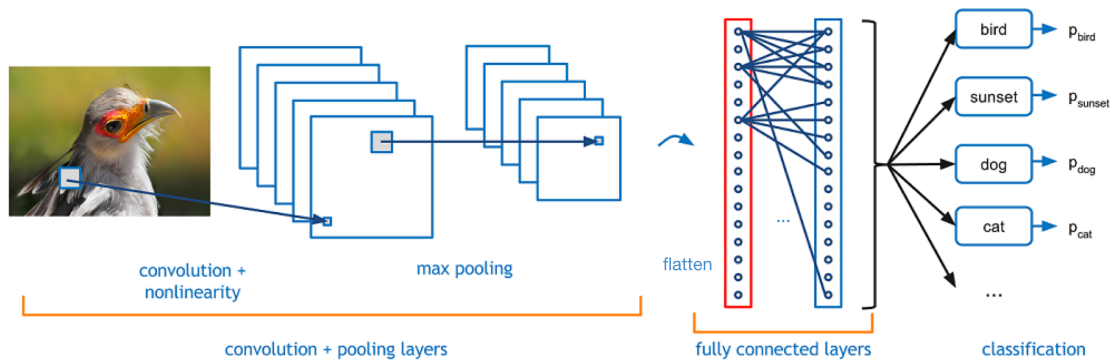


Figure 2.4: Example of a simple CNN for object classification, showing convolution, pooling and FC layers. Figure modified from [39].

2.4.1 Convolutional neural networks

These are mostly used for object detection and classification. Like MLPs, they are also a form of ‘feedforward’ NNs (not to be confused with the FF operation) in the sense that there are no state variables and cycles within the network. CNNs primarily include **convolution layers**, which consist of a number of **filters** of some **window size**, which perform a linear correlation operation on a group of nodes in a left layer, then apply a non-linearity like ReLU to get the value of a single node in the right layer. This is repeated for the next group of left nodes until the whole left layer is covered, and then repeated for a certain number of filters (also known as **channels**). CNNs also include **pooling layers**, which take a group of left nodes and **downsample** them to create a single right node value. Downsampling is done via an operation such as taking the maximum, or averaging. Convolution and pooling layers are typically followed by a MLP, *i.e.* FC layers, when doing classification. A simple CNN is shown in Fig. 2.4.

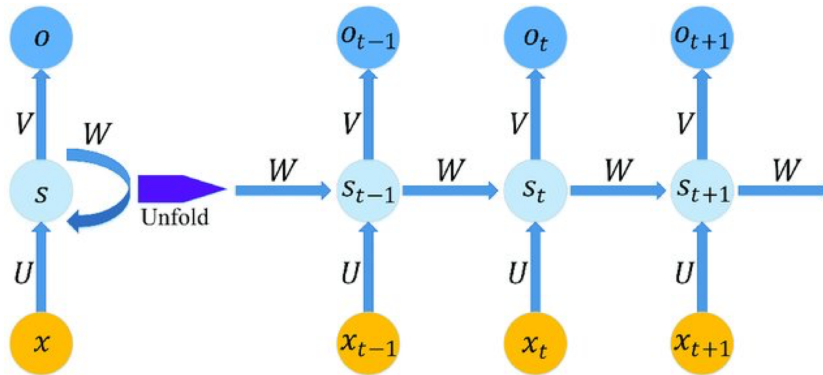


Figure 2.5: Example of a RNN with input x , states s and output o . Time is denoted by t , while U , V and W denote linear operations. Non-linear operations are not shown here for simplicity. Figure courtesy [40, Fig. 4].

2.4.2 Recurrent neural networks

Contrary to MLPs and CNNs, Recurrent Neural Network (RNN)s have internal **states** and may have cycles in their network structure. They are primarily used for learning dependencies over time, such as speech or video signals. Such an input signal is broken up into segments in time and may be operated on in a way similar to MLPs, *i.e.* a linear combination followed by a non-linearity, to generate intermediate states and outputs at different instances of time. A RNN is shown in Fig. 2.5.

Chapter 3

Pre-Defined Sparsity

This chapter discusses the major idea driving our research – pre-defined sparsity – which deletes parameters and simplifies a NN prior to training. We will begin with related work for complexity reduction of NNs. This also includes hardware implementations, which serve as related work for our own hardware efforts described in Chapter 4.

3.1 Related Work

As mentioned in the Introduction, modern NNs suffer from parameter explosion in the sense that the number of trainable parameters ranges from millions to billions [1, 5, 6, 34]. These demand large amounts of memory to store and arithmetic resources to operate on, particularly during training. Moreover, training a network with too many parameters makes it likely to overfit [41], and memorize undesirable noise patterns [42]. However, most forms of complexity reduction aim to simplify and accelerate the inference phase only, not the more complexity-intensive training phase. Some of these are described next.

[10] used vector quantization to compress the weights of deep CNNs. [11] grouped weights into hash buckets, effectively implementing a form of coarser granularity in weight values. This approach is taken to its extreme by using only values $+1$ and -1 for weights [43]. The efforts in [7, 12, 13] use a number of techniques for reducing the storage footprint of the weights, such as deleting weights with low values (pruning), retraining only necessary weights, and Huffman coding for efficient storage. All these methods need to first train the fully connected (FC) NN with all weights present, before deciding compression strategies for inference. Similarly, other pruning and trimming methods post-process a trained FC NN to remove weights [44–46].

One well-known method for reducing the size of NNs is dropout [47], where some parameters of the NN are randomly deleted during each batch of training. However, dropout is an ensemble method which actually increases overall training complexity since many different configurations need to be trained. Moreover, the final NN used for inference is FC, *i.e.* uncompressed. Other mathematical methods employed during training include using low-precision arithmetic [48–50], special matrices to structure the weights [51, 52], and regularizers to discard unimportant weights [53, 54] for inference, however, the latter two method classes lead to more efficient inference at the cost of increased training complexity.

Some custom hardware implementations for NNs have been developed such as Application Specific Integrated Circuit (ASIC)-based [13, 45, 55–58], and FPGA-based [43, 44, 52, 59–61]. We note that while all these methods have achieved excellent results in making inference faster and more efficient **on-device**, training is typically done for the complete uncompressed network off-device, such as on a power-hungry GPU or cloud server. Note that there have been some efforts in on-device training [62–66]. However, these have not been targeted towards complexity reduction and hence some have been limited to implementing NNs with a fairly small number of total nodes – 21 in [64], 83 in [63], and a more impressive 221 in [62] (some of the pipelining ideas used in our hardware architecture described in Chapter 4 were inspired by [62]).

In summary, based on related work and the present state of NNs, we identify two open problems – a) Design methods and algorithms to reduce the complexity of *training* NNs, and b) a flexible hardware architecture to support both low-complexity training and inference on-device. This chapter proposes a method aimed at tackling the first problem.

3.2 Structured Pre-defined sparsity

Pre-defined sparsity refers to a class of methods where a NN is made sparse by removing some of its edges (connections) *prior to training*. This implies that both

training and inference use a NN of lower complexity as compared to a regular non-sparse NN, such as one having FC layers. To the best of our knowledge, we are the first to propose pre-defined sparsity as a technique for complexity reduction in [25]. Note that several other authors have recently proposed pre-defined sparsity [67–69] independent of our work.

To understand pre-defined sparsity, we need to recall and add to the basic definitions from Section 2.2. The **out-degree** of a node j in the left layer $i - 1$ of junction i , $d_i^{\text{out}(j)}$, is the total number of edges connecting it to layer i to its right. Likewise, the **in-degree** of a node j in the right layer i of junction i , $d_i^{\text{in}(j)}$, is the total number of edges connecting it to layer $i - 1$ to its left. When these numbers are constant for all nodes in a layer, the superscript can be omitted and we get d_i^{out} and d_i^{in} for junction i . Note that for a conventional FC junction, $d_i^{\text{out}} = N_i$ and $d_i^{\text{in}} = N_{i-1}$. We refer to the total number of edges in junction i as $|\mathbf{W}_i|$, thus $|\mathbf{W}_i| = N_{i-1}N_i$ for FC.

► **Definition 1: *Pre-defined sparsity*:** A junction i is pre-defined sparse if it does not have all N_iN_{i-1} edges present. A NN is pre-defined sparse if it has at least one junction which is pre-defined sparse.

► **Definition 2: *Structured pre-defined sparsity*:** Structured pre-defined sparsity is pre-defined sparsity with fixed in- and out-degrees for each junction.

Thus, every node in layer $i - 1$ has a fixed $d_i^{\text{out}} \leq N_i$, and every node in layer i has a fixed $d_i^{\text{in}} \leq N_{i-1}$. This leads to $|\mathbf{W}_i| = N_{i-1}d_i^{\text{out}} = N_id_i^{\text{in}}$, thus

$|\mathbf{W}_i| \leq N_{i-1}N_i$. The **density** of junction i is measured relative to FC and denoted as $\rho_i = |\mathbf{W}_i|/(N_{i-1}N_i)$. (Occasionally we may refer to sparsity as the opposite of density, so a junction which is 20% dense is 80% sparse). As will be shown in subsequent chapters, imposing the structured constraint leads to performance improvement and ease of hardware implementation as compared to distributing connections randomly. For the remainder of this work, pre-defined sparsity will always refer to the structured form unless otherwise mentioned.

The **overall density** of a pre-defined sparse NN is defined as:

$$\rho_{\text{net}} = \frac{\sum_{i=1}^L |\mathbf{W}_i|}{\sum_{i=1}^L N_{i-1}N_i} = \frac{\sum_{i=1}^L N_{i-1}d_i^{\text{out}}}{\sum_{i=1}^L N_{i-1}N_i} = \frac{\sum_{i=1}^L N_i d_i^{\text{in}}}{\sum_{i=1}^L N_{i-1}N_i} \quad (3.1)$$

i.e. the total number of edges in the NN as a fraction of the total number of edges in the corresponding FC NN. Thus, specifying $\mathbf{N}_{\text{net}} = (N_0, N_1, \dots, N_L)$ and either of the **out-degree configuration** $\mathbf{d}_{\text{net}}^{\text{out}} = (d_1^{\text{out}}, d_2^{\text{out}}, \dots, d_L^{\text{out}})$ or the **in-degree configuration** $\mathbf{d}_{\text{net}}^{\text{in}} = (d_1^{\text{in}}, d_2^{\text{in}}, \dots, d_L^{\text{in}})$ completely determines the density of each junction and the overall density. Fig. 3.1 illustrates these basic concepts for a simple MLP.

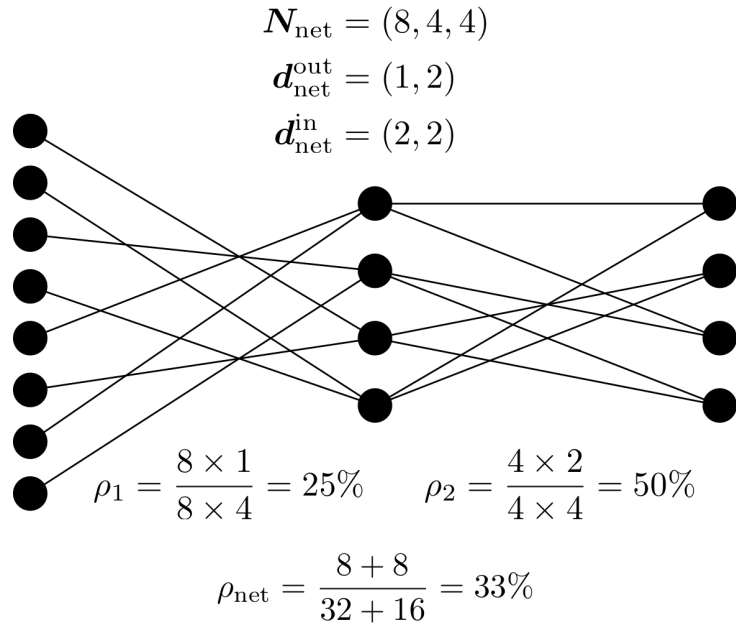


Figure 3.1: Illustrating basic concepts of structured pre-defined sparsity.

The ‘pre’ in pre-defined sparsity

We would like to emphasize that all densities of a pre-defined sparse NN are set prior to training, and then held fixed throughout training and inference.

This means that the **connection pattern**, *i.e.* which node(s) connects to which node(s) and which edges are present and absent, does not change for the NN once it is set prior to training at the beginning.

Note that the concept of pre-defined sparsity as defined by us applies only to MLPs. The connection pattern in CNNs is different – every node in a convolution layer receives inputs from various filters applied to a localized region of the inputs in its previous layer. For RNNs, a node may receive raw inputs and those from adjacent network states, which themselves may be MLPs or CNNs. A complete

investigation of recurrent and convolution layers is beyond the scope of this work. The concepts of pre-defined sparsity presented in this work are applied only to MLP layers, however, as will be shown by our experiments, they extend to MLP layers which may be part of a different network such as a CNN.

3.2.1 Motivation and Preliminary Examples

Pre-defined sparsity can be motivated by inspecting the histogram for trained weight values in a FC NN. There have been previous efforts to study such statistics [7, 70], however, not for individual junctions. Fig. 3.2 shows weight histograms for each junction in both a 2-junction and 4-junction FC NN trained on the MNIST dataset on handwritten digit classification (datasets will be described in more detail in Section 3.3.1). Note that many of the weights are zero or near-zero after training, especially in the earlier junctions. This motivates the idea that some weights in these junctions could be set to zero (*i.e.* the edges excluded).

Even with this intuition, it is unclear that one can pre-define a set of weights to be zero and let the NN learn around this constraint. Fig. 3.2(c) and (h) show that, in fact, this is the case – *i.e.* this shows classification accuracy performance on the test set as a function of the overall density ρ_{net} for structured pre-defined sparsity. The circled point on the extreme right of each subfigure is the FC case. Note that even at $\rho_{\text{net}} = 20\%$, the degradation in classification performance is

within 1%. Since the computational and storage complexity is directly proportional to the number of edges in the NN, operating at an overall density of 20% results in a 5X reduction in complexity both during training and inference. These preliminary examples serve to show the **effectiveness of pre-defined sparsity in reducing complexity of NNs while incurring minimal performance loss**. Detailed numerical experiments in Section 3.3 will further build on these preliminary examples.

3.2.2 Structured Constraints

In terms of the weight matrix, if junction i is pre-defined sparse, then \mathbf{W}_i has zeroes indicating the absence of edges and non-zero elements indicating the weight values of present edges. Note that in the course of random initialization or training updates, weight values of present edges can also become zero. However, the absent edges in a NN will always remain absent, *i.e.* :

$$\mathbb{I}\left(\mathbf{W}_i^{(j,:)} \neq 0\right) \leq d_i^{\text{in}} \quad \forall j \in \{1, \dots, N_i\} \quad (3.2a)$$

$$\mathbb{I}\left(\mathbf{W}_i^{(:,k)} \neq 0\right) \leq d_i^{\text{out}} \quad \forall k \in \{1, \dots, N_{i-1}\} \quad (3.2b)$$

where \mathbb{I} is the indicator function, which has value 1 if its argument is true, otherwise 0.

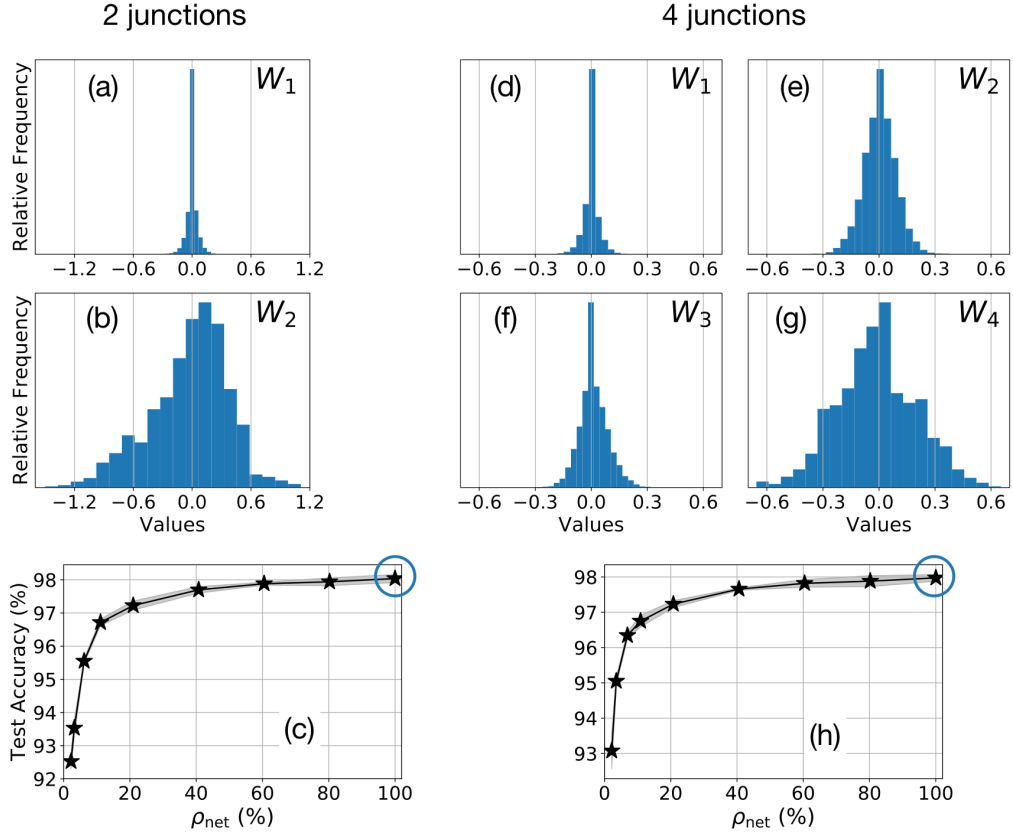


Figure 3.2: Histograms of weight values in different junctions for FC NNs trained on MNIST for 50 epochs, with (a-b) $\mathbf{N}_{\text{net}} = (800, 100, 10)$, and (d-g) $\mathbf{N}_{\text{net}} = (800, 100, 100, 100, 10)$. Test set classification accuracy shown in (c,h) for different NNs with same \mathbf{N}_{net} and varying ρ_{net} . The FC cases ($\rho_{\text{net}} = 100\%$) are circled to provide a baseline. The remaining ρ_{net} values are set by reducing ρ_1 since junction 1 has more weights close to zero in the FC cases.

Imposing the structured constraint restricts the number of possible density values. As an example, the NN shown in Fig. 2.1 with $\mathbf{N}_{\text{net}} = (4, 3, 2)$ cannot have $d_1^{\text{out}} = 2$ (*i.e.* $\rho_1 = 2/3$) since that would lead to $d_1^{\text{in}} = 8/3$, and a node cannot have a fractional number of connections. This leads to:

Theorem 1

For a given junction i , the total number of different ρ_i values possible is equal to the Greatest Common Divisor (gcd) of N_i and N_{i-1} .

Proof: Consider a NN junction i . Its density ρ_i cannot be arbitrary, since $\rho_i = d_i^{\text{out}}/N_i = d_i^{\text{in}}/N_{i-1}$, where d_i^{out} and d_i^{in} are natural numbers satisfying the equation $N_{i-1}d_i^{\text{out}} = N_i d_i^{\text{in}}$. Therefore, the number of possible ρ_i values is the same as the number of $(d_i^{\text{out}}, d_i^{\text{in}})$ values satisfying the structured pre-defined sparsity constraints:

$$d_i^{\text{out}} = \frac{N_i d_i^{\text{in}}}{N_{i-1}}, \quad d_i^{\text{in}} \leq N_{i-1}, \quad d_i^{\text{out}}, d_i^{\text{in}} \in \mathbb{N} \quad (3.3)$$

where \mathbb{N} denotes the set of natural numbers.

The smallest value of d_i^{in} which satisfies $d_i^{\text{out}} \in \mathbb{N}$ is $N_{i-1}/\text{gcd}(N_{i-1}, N_i)$, and other values are its integer multiples. Since d_i^{in} is upper bounded by N_{i-1} , the total number of possible $(d_i^{\text{out}}, d_i^{\text{in}})$ is $\text{gcd}(N_{i-1}, N_i)$. Thus, the set of possible ρ_i is

$$\left\{ \rho_i \in (0, 1] \mid \rho_i = \frac{k}{\text{gcd}(N_{i-1}, N_i)}, k \in \mathbb{N} \right\}. \quad (3.4)$$

which is a set of cardinality equal to the gcd of N_i and N_{i-1} . ■

As a concrete example, consider a NN with $\mathbf{N}_{\text{net}} = (117, 390, 13)$ (this will arise in our experiments on the TIMIT dataset described in Section 3.3). The

number of possible junction densities are $\gcd(117, 390) = 39$ and $\gcd(390, 13) = 13$.

Therefore, the possible junction densities are:

$$\rho_1 \in \left\{ \frac{1}{39}, \frac{2}{39}, \dots, \frac{39}{39} \right\}, \quad \rho_2 \in \left\{ \frac{1}{13}, \frac{2}{13}, \dots, \frac{13}{13} \right\}. \quad (3.5)$$

The above proof and example for Theorem 1 was achieved in collaboration with Kuan-Wen Huang.

3.2.3 Modifications to Neural Network Operations

The basic operations of a NN – FF, BP and UP – get modified for the case of structured pre-defined sparsity. The primary change is that only the present weights are used for computation, which leads to reduction in complexity. Now we will present the modified equations for the commonly occurring case of one-hot labels, cross-entropy cost, softmax activation for the output layer, and a scalar-to-scalar activation applied element-wise to all other layers. The matrix-vector products are split into summations to highlight usage of only the present weights.

The FF equations (2.7) become:

$$s_i^{(j)} = \sum_{f=1}^{d_i^n} W_i^{(j,k_f)} a_{i-1}^{(k_f)} + b_i^{(j)} \quad (3.6a)$$

$$a_i^{(j)} = h \left(s_i^{(j)} \right) \quad (3.6b)$$

$$h'_i{}^{(j)} = \frac{\partial a_i^{(j)}}{\partial s_i^{(j)}} \quad (3.6c)$$

where (3.6a) and (3.6b) are $\forall i \in \{1, 2, \dots, L\}$, and (3.6c) is $\forall i \in \{1, 2, \dots, L-1\}$. This is because softmax activation derivatives are directly combined with cross-entropy cost derivatives in the final layer to give (3.7a), hence are not separately required to be computed.

The BP equations (2.13) become:

$$\delta_L^{(j)} = a_L^{(j)} - y^{(j)} \quad (3.7a)$$

$$\delta_i^{(j)} = h_i^{(j)} \left(\sum_{f=1}^{d_i^{\text{out}}} W_{i+1}^{(k_f, j)} \delta_{i+1}^{(k_f)} \right) \quad (3.7b)$$

where (3.7b) is $\forall i \in \{1, 2, \dots, L-1\}$.

The UP equations (2.14), assuming batch size $M = 1$ to reduce clutter, become:

$$b_i^{(j)} \leftarrow b_i^{(j)} - \eta \delta_i^{(j)} \quad (3.8a)$$

$$W_i^{(j, k)} \leftarrow W_i^{(j, k)} - \eta a_{i-1}^{(k)} \delta_i^{(j)} \quad (3.8b)$$

where both (3.8a) and (3.8b) are $\forall i \in \{1, 2, \dots, L\}$, and (3.8b) is only for those (j, k) node pairs which have an edge connecting them.

For all the equations (3.6)–(3.8) above, single superscript j refers to nodes in a layer and hence has range $\{1, 2, \dots, N_i\}$. Double superscripts denote *only those edges which are present*. For example, in (3.6a), the summation for a particular row j of \mathbf{W}_i is carried out over d_i^{in} different values from $W_i^{(j, k_1)}$ to $W_i^{(j, k_{d_i^{\text{in}}})}$. Considering Fig. 3.1 as a concrete example, when doing the computation for the 1st right

neuron in junction 1 (*i.e.* $i = 1, j = 1$), we get $k_1 = 5$ and $k_2 = 7$. This is because nodes 5 and 7 from layer 0 connect to node 1 of layer 1.

This concludes the theory relating to pre-defined sparsity. In the next section, we present performance results of pre-defined sparsity.

3.3 Performance Results, Trends and Guidelines

Fig. 3.2(c) and (h) showed some preliminary examples regarding the potential of pre-defined sparsity to reduce complexity with minimal degradation in performance. This section analyzes results and trends observed when experimenting with several different classification datasets via software simulations. We intend the following to provide guidelines on designing pre-defined sparse NNs.

Guidelines for designing pre-defined sparse NNs

1. The performance of pre-defined sparsity is better on datasets that have more inherent redundancy (Section 3.3.2).
2. Junction density should increase to the right, *i.e.* junctions closer to the output should generally have more density than junctions closer to the input (Section 3.3.3).
3. Larger and more sparse NNs are better than smaller and denser NNs, given the same number of layers and trainable parameters. Specifically, ‘larger’ refers to more hidden neurons (Section 3.3.4).

The reader is reminded that when we refer to pre-defined sparsity, we are actually referring to structured pre-defined sparsity, *i.e.* fixed in- and out-degrees. The remainder of this section first describes the datasets we experimented on, and then examines these trends in detail.

3.3.1 Datasets and Experimental Configuration

NNs for classification need significant amounts of quality labeled data to train and perform inference. Here we discuss some widely used datasets made publicly available for this purpose, and the training configuration we designed for each of them.

MNIST handwritten digits [71]

The Modified National Institute of Standards and Technology (MNIST) database has 70,000 images of handwritten digits from 0 – 9, split into 60,000 for training and 10,000 for test. We further split the training set into 50,000 for actual training and 10,000 for validation. Each image is 28×28 pixels of varying shades of gray values, from 0 signifying completely black to 255 signifying completely white. We rasterized each input image into a single layer of 784 features. This is the *permutation-invariant* format, wherein spatial information from the inputs are not used, implying that the ordering of the 784 pixels is not important. On certain occasions we added 16 input features which are always trivially 0 so as to get 800

features for each input. This leads to easier selection of different sparse network configurations. In Fig. 3.2(c) for example, we used 800 input neurons and 100 hidden neurons since $\text{gcd}(800, 100) > \text{gcd}(784, 100)$, so more values of ρ_{net} could be simulated. We verified that adding extra always-0 input features did not alter performance. Also note that no data augmentation was applied.

Reuters RCV1 corpus of newswire articles [72]

Reuters Corpus Volume I (RCV1) is an archive of newswire stories, each assigned categories based on its content. The classification categories are grouped in a tree structure, *e.g.* a top-level category like politics can be subdivided into second level categories for international, national and state, international can be further subdivided into third level categories for Asia, Africa and so on. An article can have multiple categories, however, we used preprocessing techniques inspired by [4] to isolate only those articles which belonged to a single second level category. We finally obtained 328,669 articles in 50 categories, split into 50,000 for validation, 100,000 for test, and the remaining 178,669 for training.

The original data has a list of token strings for each story, for example, a story on finance would frequently contain the token ‘financ’ (which can refer to ‘finance’, ‘financial’, ‘financier’ and so on). We chose the most common 2000 tokens across all articles and computed counts for each of these in each article. Each count x

was transformed into $\log(1 + x)$ to form the final 2000-dimensional feature vector for each input.

TIMIT speech corpus [73]

TIMIT (presumably abbreviated from Texas Instruments, Massachusetts Institute of Technology) is a speech dataset comprising approximately 5.4 hours of 16 kHz audio commonly used in ASR. A modern ASR system has three major components: (i) preprocessing and feature extraction, (ii) acoustic model, and (iii) dictionary and language model. A complete study of an ASR system is beyond the scope of this work. Instead we focus on the acoustic model which is typically implemented using a NN. The input to the acoustic model is feature vectors and the output is a probability distribution on phonemes (*i.e.* speech sounds), *i.e.* the classification targets for our experiments are different phonemes. We used a phoneme set of size 39 as defined in [74].

We extracted 25ms speech frames with 10ms shift, as in [4], and computed a feature vector of 39 Mel-frequency Cepstral Coefficient (MFCC)s for each frame. MFCCs are obtained from a sequence of operations – windowing, Fourier transform, mapping powers on to mel scale [75], taking logarithm, discrete cosine transform, measuring amplitude. We used the complete training set of 818,837 training samples (462 speakers), 89,319 validation samples (50 speakers), and 212,093 test samples (118 speakers).

CIFAR-100 images [76]

The Canadian Institute For Advanced Research (CIFAR) datasets come in two forms – 10 classes, and 100 classes. Each has 50,000 training samples (which we split into 40,000 for actual training and 10,000 for validation), and 10,000 test samples. The samples are $32 \times 32 \times 3$ (for 3 color channels – red, green and blue) images of objects such as different vehicles and animals. We experimented on the CIFAR-100 dataset with a CNN. The CNN has 3 blocks, each block has 2 convolution layers with window size 3×3 followed by a max pooling layer of pool size 2×2 . The number of filters for the six convolution layers are (from left to right) 60, 60, 125, 125, 250 and 250. This results in a total of approximately one million trainable parameters in the convolutional portion of the network. Batch normalization is applied before activations. The output from the 3rd block, after flattening into a vector, has 4000 features. This is the input to a MLP. Typically dropout is applied in the MLP portion, however we omitted it there since pre-defined sparsity is an alternate form of parameter reduction. Instead we found that a dropout probability of 0.5, *i.e.* half the nodes and weights dropped for each batch of training, applied to the convolution blocks improved performance. No data augmentation was applied.

Experimental Configuration

For each dataset, we performed classification using one-hot labels and measured accuracy on the test set as a performance metric.¹ We also calculated the top-5 test set classification accuracy for CIFAR-100, *i.e.* the percentage of samples for which the ground-truth label was among the top 5 softmax outputs of the network.

We found the optimal training configuration and hyperparameters for each FC setup by doing a grid search using validation performance as a metric. This resulted in choosing ReLU activations for all layers except for the final softmax layer. He initialization worked best for the weights; while for biases, we found that an initial value of 0.1 worked best in all cases except for Reuters, for which zeroes worked better. The Adam optimizer was used with all parameters set to default, except that we set the decay parameter to 10^{-5} for best results. We used a batch size of 1024 for TIMIT and Reuters since the number of training samples is large, and 256 for MNIST and CIFAR. All experiments were run for 50 epochs of training and all

¹The NN in a complete ASR system would be a ‘soft’ classifier and feed the phoneme distribution outputs to a decoder to perform ‘hard’ final classification decisions. Therefore for TIMIT, we computed another performance metric called Test Prediction Comparison (TPC), measured as KL divergence between predicted test output probability distributions of sparse vs the respective FC case. In other words,

$$\text{TPC} = \frac{1}{212093} \sum_{i=1}^{212093} \sum_{j=1}^{39} A_{ij} \log \left(\frac{A_{ij}}{B_{ij}} \right) \quad (3.9)$$

where $A_{212093 \times 39}$ and $B_{212093 \times 39}$ are the complete output test matrices for the 212,093 test samples with 39 labels each for the FC and sparse cases, respectively. Lesser values of TPC are better as they indicate minimal performance degradation due to sparsification. We found that performance results obtained using the TPC metric were qualitatively very similar to those obtained from the test accuracy metric, hence we omit the TPC results.

hyperparameters listed so far were kept the same when sparsifying the network to maintain consistency.

The only exception is regularization. We found from validation that sparser networks need lesser regularization than denser networks. This confirms our belief that **pre-defined sparse NNs are less prone to overfitting due to having fewer trainable parameters**. Accordingly we applied an L2 penalty to the weights, but reduced the coefficient λ as ρ_{net} decreased.

Each experiment was run at least five times to average out randomness and the 90% Confidence Interval (CI)s for each metric are shown as shaded regions (this also holds for the results in Fig. 3.2(c,h)).

3.3.2 Dataset Redundancy

Many machine learning datasets have considerable redundancy in their input features. For example, one may not need information from the ~ 800 input features of MNIST to infer the correct image class. We hypothesize that pre-defined sparsity takes advantage of this redundancy, and will be less effective when the redundancy is reduced. To test this, we changed the feature vector for each dataset as follows:

- MNIST: Principal Component Analysis (PCA) is a technique to compute variances in each feature across all inputs in the dataset, and assign more ‘importance’ to features with large variance. This is because such features are more discriminatory, *i.e.* observing their values will provide the network

with more information. We used PCA to reduce the feature count of MNIST to the most important (least redundant) 200, *i.e.* a factor of $\sim 4X$.

- Reuters: Size of the feature vector is the number of most frequent tokens considered. We reduced this from 2000 to 400, *i.e.* by $5X$.
- TIMIT: We both reduced and increased the number of MFCCs by $3X$ to 13 and 117, respectively. Note that the latter increases redundancy.
- CIFAR-100: In the CNN used for CIFAR, a source of redundancy is the depth of the convolution+pooling portion which extracts features and discriminates between classes before the MLP performs final classification. In other words, the convolution blocks ease the burden of the MLP. So a way to reduce redundancy and increase the classification burden of the MLP is to lessen the effectiveness of the convolution layers by reducing their number. Accordingly, we used a single convolution layer with 250 filters of window size 5×5 followed by a 8×8 max pooling layer. This results in the same number of features, 4000, at the input of the MLP as the original network, but has reduced redundancy for the MLP.

Classification performance results are shown in Fig. 3.3 as a function of ρ_{net} . For MNIST and CIFAR-100, the performance degrades more sharply with reducing ρ_{net} for the networks using the reduced redundancy datasets. To explore this further, we recreated the histograms from Fig. 3.2 for the reduced redundancy

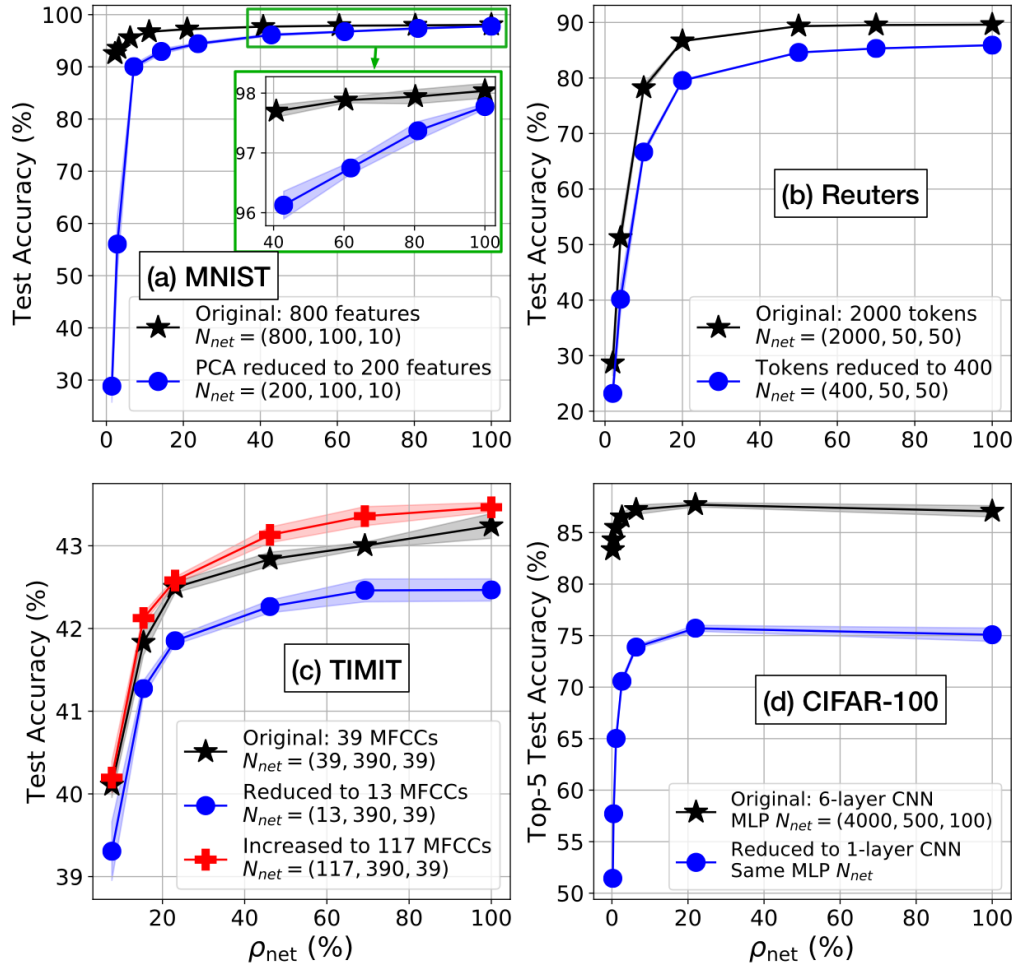


Figure 3.3: Comparison of classification accuracy as a function of ρ_{net} for different versions of datasets – original, reduced in redundancy by reducing feature space (MNIST, Reuters, TIMIT) or performing less processing prior to the MLP (CIFAR-100), and increasing redundancy by enlarging feature space (TIMIT). Higher density points for MNIST are magnified.

datasets, *i.e.* a FC NN with $N_{\text{net}} = (200, 100, 10)$ training on MNIST after PCA.

We observed a wider spread of weight values, implying less opportunity for sparsification (*i.e.* fewer weights were close to zero), as shown in Fig. 3.4. Similar trends are less discernible for Reuters and TIMIT, however, reducing redundancy led to worse performance overall.

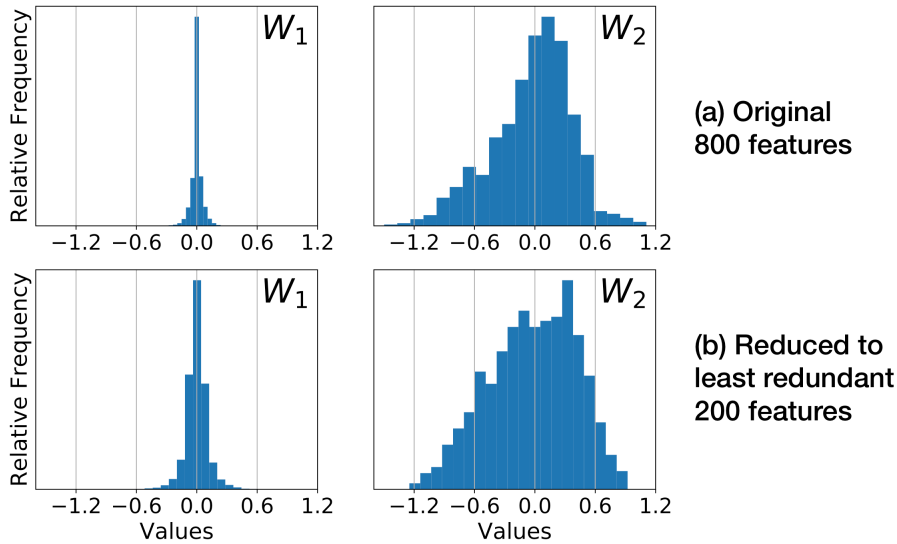


Figure 3.4: Comparing histograms of weight values for (a) original MNIST with $N_{\text{net}} = (800, 100, 10)$, and (b) MNIST reduced to least redundant 200 features such that $N_{\text{net}} = (200, 100, 10)$, after training FC NNs for both.

The results in Fig. 3.3 further demonstrate the effectiveness of pre-defined sparsity in greatly reducing network complexity with negligible performance degradation. For example, even the reduced redundancy problems perform well when operating with half the number of connections. For CIFAR in particular, *FC performs worse than an overall MLP density of around 20%*. Thus, in addition to reducing complexity, structured pre-defined sparsity may be viewed as an alternative to dropout in the MLP for the purpose of improving classification.

3.3.3 Individual junction densities

The weight value histograms in Figs. 3.2 and 3.4 indicate that latter junctions, particularly junction L closest to the output, have a wide spread of weight values.

This suggests that a good strategy for reducing ρ_{net} would be to use lower densities in earlier junctions, *i.e.* $\rho_1 < \rho_L$.

This is demonstrated in Fig. 3.5 for the cases of MNIST, CIFAR-100 and Reuters, each with $L = 2$ junctions in their MLPs. Each curve in each subfigure is for a fixed ρ_2 , *i.e.* reducing ρ_{net} across a curve is done solely by reducing ρ_1 . For a fixed ρ_{net} , the performance improves as ρ_2 increases. For example, the circled points in Reuters both have $\rho_{\text{net}} = 4\%$, but the starred point with $\rho_2 = 100\%$ has approximately 40% better test accuracy than the pentagonal point with $\rho_2 = 2\%$. The trend clearly holds for MNIST and Reuters, and is also discernible for CIFAR-100.

This trend extends to NNs with more than two junctions. Fig. 3.6 compares classification performance by trading off ρ_2 and ρ_3 (keeping ρ_1 fixed) for different three-junction NNs training on MNIST. For each individual plotted curve (same color), ρ_3 is kept constant and ρ_{net} is varied by varying ρ_2 . Note that all the results are qualitatively similar to Fig. 3.5, *i.e.* for the same ρ_{net} , better results are obtained for higher values of ρ_3 . Note that ρ_1 is fixed for each subfigure at fairly low values, as given in the caption of Fig. 3.6. We also experimented with keeping ρ_1 fixed at higher values, such as 50%. Since the input layer has the most neurons, this resulted in junction 1 dominating the total number of weights. As a result, changing ρ_2 and ρ_3 had little effect on accuracy. Hence we decided to include the more interesting low ρ_1 cases in Fig. 3.6.

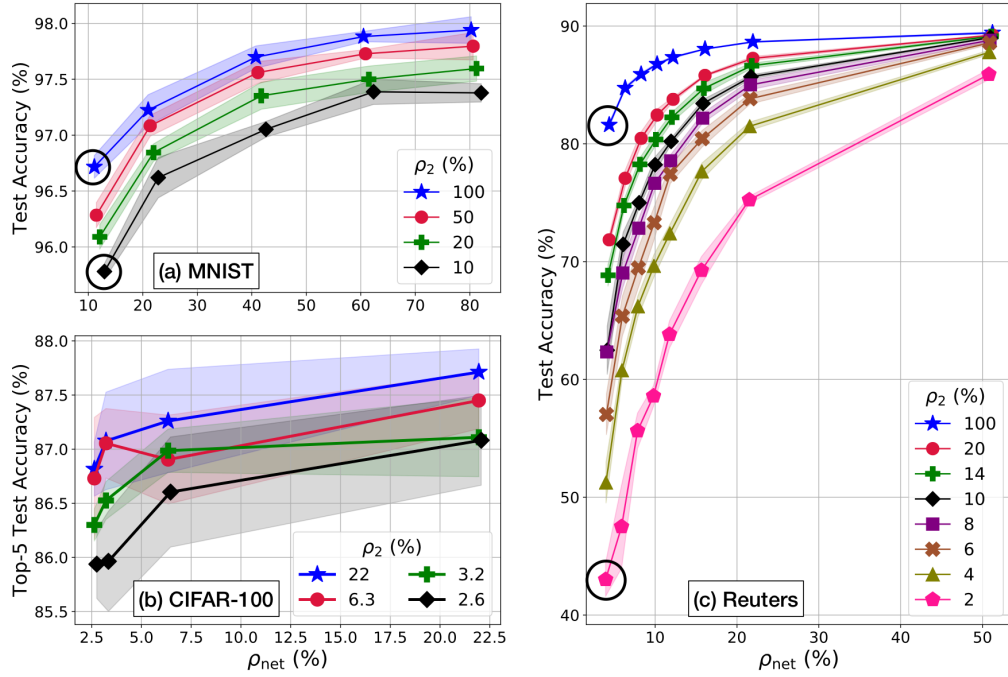


Figure 3.5: Comparison of classification accuracy as a function of ρ_{net} for different ρ_L , where $L = 2$. Black-circled points show the effects of ρ_2 when ρ_{net} is the same. N_{net} values are (800, 100, 10) for MNIST, (2000, 50, 50) for Reuters, and (4000, 500, 100) for the MLP portion in CIFAR-100.

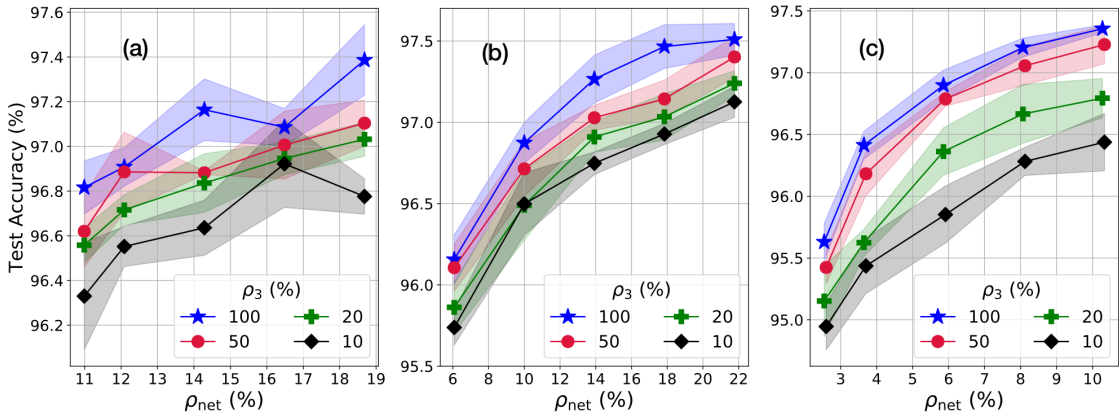


Figure 3.6: Comparison of classification accuracy as a function of ρ_{net} for ρ_2 vs ρ_3 in three-junction MNIST networks, keeping ρ_1 fixed. (a) $N_{\text{net}} = (800, 100, 100, 10)$, $\rho_1 = 10\%$, (b) $N_{\text{net}} = (800, 100, 200, 10)$, $\rho_1 = 4\%$, (c) $N_{\text{net}} = (800, 200, 100, 10)$, $\rho_1 = 1\%$.

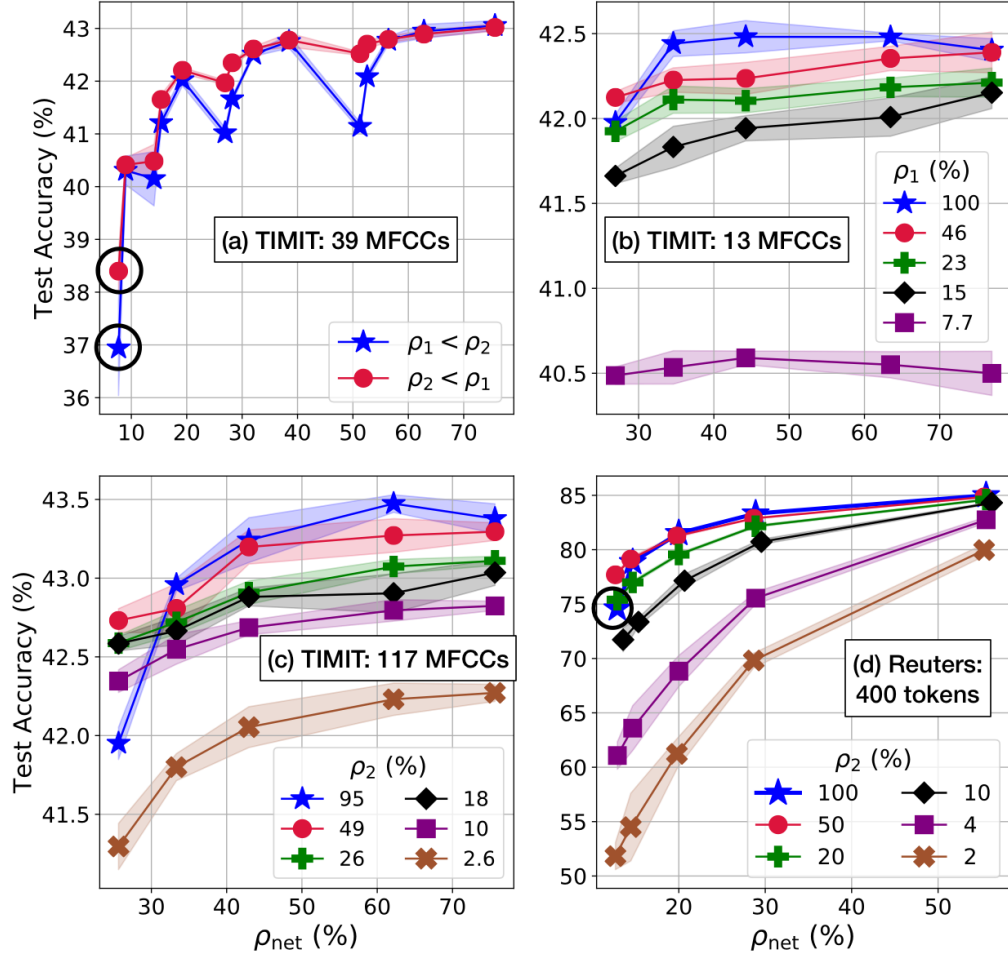


Figure 3.7: Comparison of classification accuracy as a function of ρ_{net} for: (a) TIMIT with 39 MFCCs for the two cases where one junction is always sparser than the other and vice-versa. Black-circled points show how reducing ρ_1 degrades performance to a greater extent. (b) TIMIT with 13 MFCCs for different ρ_1 . (c,d) TIMIT with 117 MFCCs, and Reuters reduced to 400 tokens, for different ρ_2 . N_{net} values are (a) (39, 390, 39), (b) (13, 390, 39), (c) (117, 390, 39), (d) (400, 50, 50).

We further observed that this trend of $\rho_{i+1} > \rho_i$ improving performance is related to the redundancy inherent in the dataset, and may not hold for datasets with very low levels of redundancy. To explore this, results analogous to those in Fig. 3.5 are presented in Fig. 3.7 for TIMIT, but with varying sized MFCC feature vectors – *i.e.* datasets corresponding to larger feature vectors will contain

more redundancy. The results in Fig. 3.7(c) are for 117-dimensional MFCCs and are consistent with the trend in Fig. 3.5. However, for a MFCC dimension of 13, this trend actually reverses – *i.e.* junction 1 should have higher density for better performance. This is shown in Fig. 3.7(b), where each curve is for a fixed ρ_1 . This reversed trend is also observed for the case of 39 dimensional feature vectors, considered in Fig. 3.7(a), where $\mathbf{N}_{\text{net}} = (39, 390, 39)$. Due to this symmetric neuronal configuration, for each value of ρ_{net} on the x-axis in Fig. 3.7(a), the two curves have complementary values of ρ_1 and ρ_2 ($\rho_1 \neq \rho_2$) – *e.g.* the two curves at $\rho_{\text{net}} = 7.69\%$ have (ρ_1, ρ_2) values of $(2.56\%, 12.82\%)$ and $(12.82\%, 2.56\%)$. We observe that the curve for $\rho_1 < \rho_2$ is generally worse than the curve for $\rho_2 < \rho_1$, which indicates that junction 1 should have higher density for improving performance in this case as well.

Fig. 3.7(d) depicts the results for Reuters with the feature vector size reduced to 400 tokens. While junction 2 is still more important (as in Fig. 3.5(c) for the original Reuters dataset), notice the circled star-point at the very left of the $\rho_2 = 100\%$ curve. This point has very low ρ_1 . Unlike Fig. 3.5(c), it crosses below the other curves, indicating that it is more important to have higher density in the first junction with this less redundant set of features.

In summary, if an individual junction density falls below a certain value, referred to as the *critical junction density*, it will adversely affect performance regardless of the density of other junctions. This explains why some of the curves

cross in Fig. 3.7. The critical junction density is much smaller for earlier junctions than for later junctions in most datasets with sufficient redundancy. However, the critical density for earlier junctions increases for datasets with low redundancy.

3.3.4 ‘Large and sparse’ vs ‘small and dense’ networks

As mentioned previously, NNs with lower values of ρ_{net} need lower values of the L2 penalty coefficient λ for regularization. This raises the question – why not design a small FC NN, *i.e.* one with a small number of nodes in its layers? This would also have a low number of parameters to begin with and would thus have low complexity.

We experimented with this aspect and observed that when keeping the total number of trainable parameters the same, sparser NNs with larger hidden layers (*i.e.* more neurons) generally performed better than denser networks with smaller hidden layers. This is true as long as the larger NN is not so sparse that individual junction densities fall below the critical density, as explained in Sec. 3.3.3. While the critical density is problem-dependent, it is usually low enough to obtain significant complexity savings above it. Thus, ‘large and sparse’ is better than ‘small and dense’ for many practical cases, including NNs with more than one hidden layer (*i.e.* $L > 2$).

Fig. 3.8 shows this for networks having one (two) and three (four) hidden layers (junctions) trained on MNIST. For the four-junction network, all hidden layers

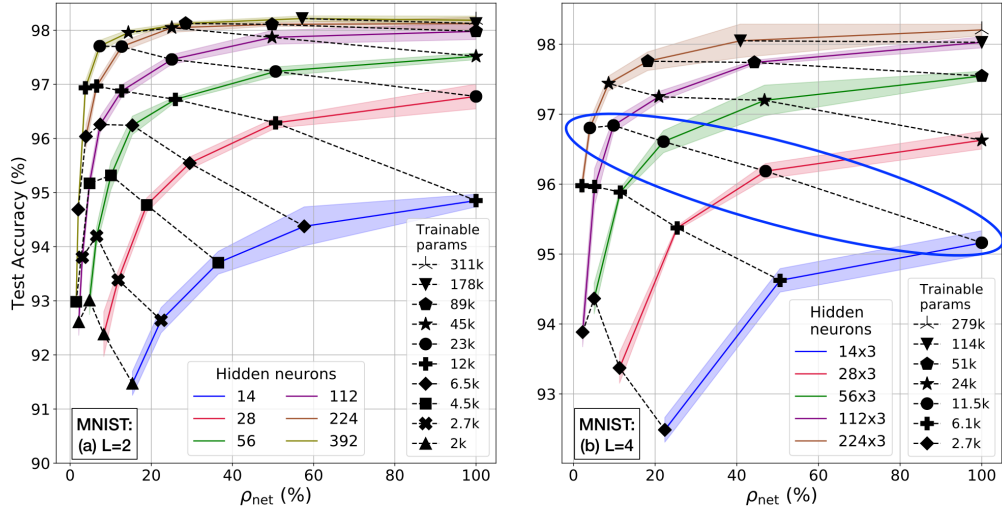


Figure 3.8: Comparing ‘large and sparse’ to ‘small and dense’ networks for MNIST with 784 features, with (a) $\mathbf{N}_{\text{net}} = (784, x, 10)$ (on the left), and (b) $\mathbf{N}_{\text{net}} = (784, x, x, x, 10)$ (on the right). Solid curves (with the shaded CIs around them) are for constant x , black dashed curves with same marker are for same number of trainable parameters. The final junction is always FC. Intermediate junctions for the $L = 4$ case have d^{out} values similar to junction 1.

have the same number of neurons. Each solid curve shows classification performance vs ρ_{net} for a particular \mathbf{N}_{net} , while the black dashed curves with identical markers are configurations that have approximately the same number of trainable parameters. As an example, the points with circular markers (with a big blue ellipse around them) in Fig. 3.8(b) all have the same number of trainable parameters and indicate that the larger, more sparse NNs perform better. Specifically, the network with $\mathbf{N}_{\text{net}} = (784, 112, 112, 112, 10)$ and $\mathbf{d}_{\text{net}}^{\text{out}} = (10, 10, 10, 10)$ corresponding to $\rho_{\text{net}} = 9.82\%$ performs significantly better than the FC network with $\mathbf{N}_{\text{net}} = (784, 14, 14, 14, 10)$, and other smaller and denser networks, despite each having 11,500 trainable parameters. Increasing the network size further to

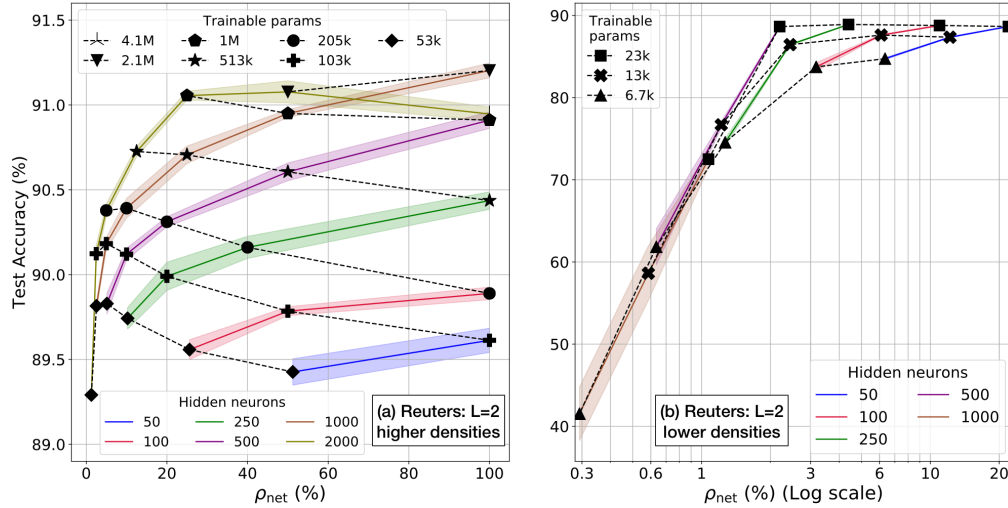


Figure 3.9: Comparing ‘large and sparse’ to ‘small and dense’ networks for Reuters with 2000 tokens, with $N_{\text{net}} = (2000, x, 50)$. The x-axis is split into higher values on the left (a), and lower values on the right in log scale (b). Solid curves (with the shaded CIs around them) are for constant x , black dashed curves with same marker are for same number of trainable parameters. Junction 1 is sparsified first until its number of total weights is approximately equal to that of junction 2, then both are sparsified equally.

$N_{\text{net}} = (784, 224, 224, 224, 10)$, and reducing ρ_{net} to 4% to fix the number of trainable parameters at 11,500 leads to performance degradation. This is because this ρ_{net} was achieved by setting $\rho_2 = \rho_3 = 2.68\%$, which appears to be below the critical density.

Fig. 3.9 summarizes the analogous experiment on Reuters with similar conclusions. Both subfigures are for the same results with the x-axis split into higher and lower density range (on log scale), to show more detail. Observe that the trend of ‘large and sparse’ being better than ‘small and dense’ holds for subfigure (a), but reverses for (b) since densities are very low (the black dashed curves have positive slope instead of negative). This is due to the critical density effect.

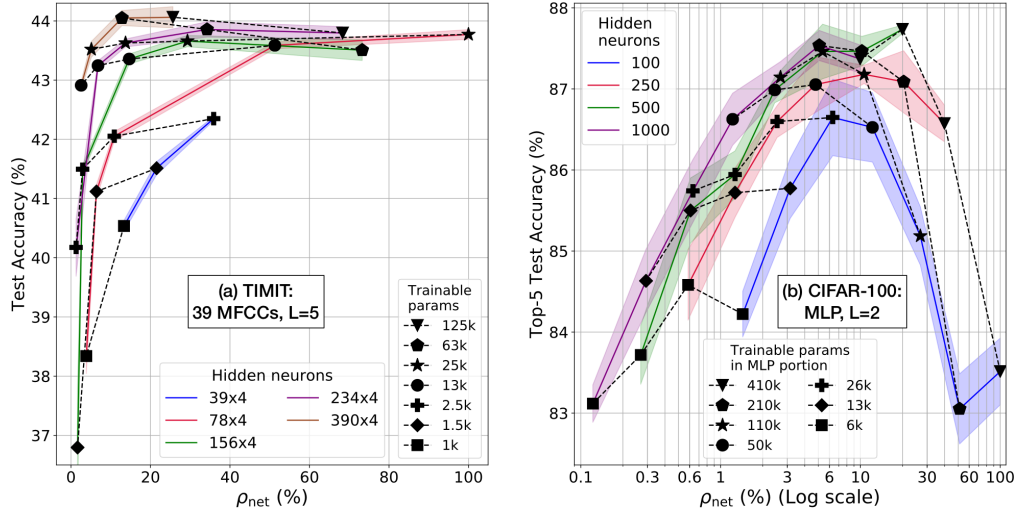


Figure 3.10: Comparing ‘large and sparse’ to ‘small and dense’ networks for (a) TIMIT with 39 MFCCs and $N_{net} = (39, x, x, x, x, 39)$ (on the left), and (b) CIFAR-100 with the deep 9-layer CNN and MLP $N_{net} = (4000, x, 100)$ with log scale for the x-axis (on the right). Solid curves (with the shaded CIs around them) are for constant x , black dashed curves with same marker are for same number of trainable parameters (in the MLP portion only for CIFAR). Since TIMIT has symmetric junctions, we tried to keep input and output junction densities as close as possible and adjusted intermediate junction densities to get the desired ρ_{net} . CIFAR-100 is sparsified in a way similar to Reuters in Fig. 3.9.

Fig. 3.10(a) shows the result for the same experiment on TIMIT with four hidden layers. The trend is less clearly discernible, but it exists. Notice how the black dashed curves have negative slopes at appreciable levels of ρ_{net} , indicating ‘large and sparse’ being better than ‘small and dense’, but high positive slopes at low ρ_{net} , indicating the rapid degradation in performance as density is reduced beyond the critical density. This is exacerbated by the fact that TIMIT with 39 MFCCs is a dataset with low redundancy, so the effects of very low ρ_{net} are better observed.

Fig. 3.10(b) for the MLP portion of CIFAR-100 shows similar results as TIMIT, but on a log x-scale for more clarity. As noted in Sec. 3.3.2, the best performance for a given N_{net} occurs at an overall density less than 100%. It appears that for any N_{net} for CIFAR-100, peak performance occurs at around 10–20% overall MLP density. We are intrigued by this trend, and plan to investigate further.

Some results have been omitted due to being qualitatively similar to those already shown. These are experiments on TIMIT with $L = 2$ (similar to 3.10(a)), Reuters with $L = 3$ (similar to 3.9), and CIFAR100 with the reduced redundancy having a single convolutional layer (similar to 3.10(b)).

3.4 Summary

This chapter introduced the concept of pre-defined sparsity as a method of complexity reduction and presented performance results indicating its effectiveness and delineating trends to design such sparse NNs. The next chapter will present a concrete application of complexity reduction – a hardware architecture we designed and implemented to leverage pre-defined sparsity.

Chapter 4

Hardware Architecture

This chapter describes our proposed flexible hardware architecture which can leverage structured pre-defined sparsity to both train and test NNs of any complexity on-device. We begin with an overview of the key features, which are encapsulated in Fig. 4.1.

► **Definition 3: Degree of parallelism:** The degree of parallelism z_i is the number of edges processed in parallel in junction i .

Thus, we call our architecture **edge-based**. z_i determines the number of clock **cycles** taken to process a complete junction. Fig. 4.1(a) shows an example junction i with 6 edges processed using $z_i = 3$. The $z_i = 3$ blue edges are processed in cycle 0 (we begin numbering from 0), and the remaining $z_i = 3$ pink edges in cycle 1. A given hardware device can support some largest values of $\{z_i\}$, so NNs with more edges will simply take more cycles to process.

For a given processing operation in junction i , there are z_i computational units to perform it and z_i memories to store each quantity associated with it. This is shown in Fig. 4.1(b) for the FF operation, *i.e.* there are $z_i = 3$ computational

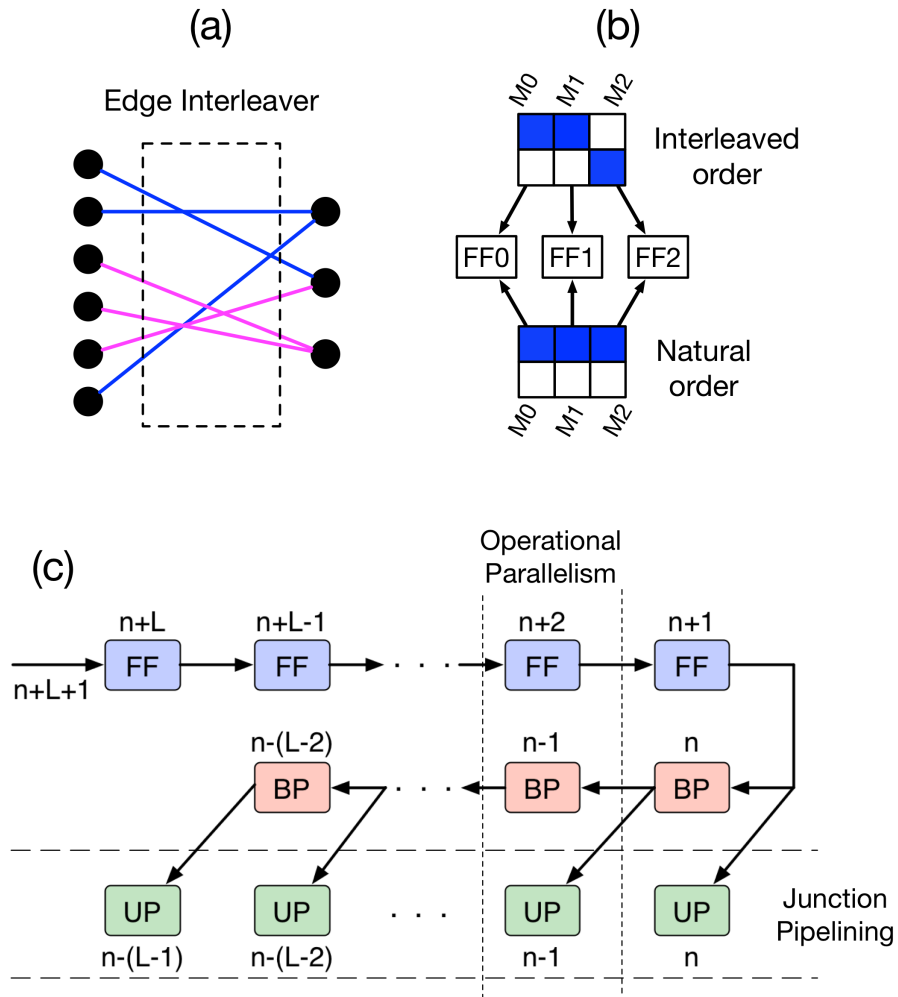


Figure 4.1: (a) Processing $z_i = 3$ edges in each cycle (blue in cycle 0, pink in cycle 1) for some junction i with a total of 6 edges. (b) Accessing $z_i = 3$ memories – M_0 , M_1 and M_2 shown as columns – from two separate banks, one in natural order (same address from each memory), the other in interleaved order. Clash-freedom is achieved by accessing only one element from each memory. The accessed values are fed to $z_i = 3$ processors to perform FF simultaneously. (c) Operational parallelism in each junction (vertical dotted lines denote processing for one junction), and junction pipelining of each operation across junctions (horizontal dashed lines) for different inputs.

units FF0–FF2. A challenge with this architecture is that in order to achieve high-throughput without memory duplication, the parallel memories must be accessed in two manners: **natural order** and **interleaved order**. Natural order accesses

memory elements sequentially. Interleaved order leads to the possibility of memory contention. This leads us to define:

► **Definition 4: *Clash*:** A clash arises if the same memory needs to be accessed multiple times in the same cycle, which may lead to stalls and/or wait states.¹

► **Definition 5: *Clash-free*:** The clash-free condition occurs when the NN connection pattern and hardware architecture is designed such that no clashes occur during memory accesses.

Thus, the z_i computational units must access the memories such that no clashes and memory contention occur. Natural and interleaved order accesses are shown for $z_i = 3$ memories on the bottom and top of Fig. 4.1(b), respectively. Notice that there is only one shaded cell per column in cycle 0. The unshaded cells will be accessed in cycle 1. The concept of clash-freedom is tied to the connection patterns in the junctions, as will be expounded in Chapter 5.

Apart from clash-freedom to maximize throughput, some other characteristics regarding memory requirements guided us in developing the proposed architecture. Firstly, since weight memories are the largest (*i.e.* the number of weights is more than any other network parameter), their number should be minimized. Secondly, having a few deep memories is more efficient in terms of power and area than

¹Some multi-ported memories can be accessed more than once in a cycle, however, clashes can still occur when trying to access more elements than the number of ports. To be more specific, a) for single-ported memories, attempting two reads or two writes or a read and a write in the same cycle is a clash, and b) for simple dual-ported memories with one port exclusively for reading and the other exclusively for writing, a read and a write can be performed in the same cycle, but attempting to perform two reads or two writes in the same cycle will result in a clash.

having many shallow memories [77]. These help to alleviate the key concern of large storage requirement when implementing NNs in hardware. Memory organization is described in Section 4.2.

In addition to edge processing in a given junction, our architecture is **pipelined across junctions**. Thus, the z_i values are selected to set the number of cycles required to process a layer to a constant – *i.e.* junctions with more weights have larger z_i so that the computation time of all junctions is the same. Furthermore, the **three operations associated with training are performed in parallel**, *i.e.* FF, BP, and UP. These concepts are shown in Fig. 4.1(c) and elaborated in Section 4.1. An $(L + 1)$ -layer NN has L junction pipeline stages so that the **throughput**, *i.e.* the frequency of processing input samples, is determined by the time taken to perform a single operation in a single junction.

Summary of architectural features

1. Edge-based, *i.e.* not tied to a specific number of nodes in a layer.
2. Flexible, *i.e.* the amount of logic is determined by the degree of parallelism which trades size for speed. This means that the architecture is also compatible with conventional FC junctions, as will be shown in Section 4.6.

3. Fully pipelined for the parallel operations associated with NN training.

The architecture can also operate in inference only mode by eliminating the logic and memory associated with BP and UP, and the \mathbf{h}' computations in (3.6c).

4.1 Junction pipelining and Operational parallelism

Our edge-based architecture is motivated by the fact that all three operations – FF, BP, UP – use the same weight values for computation. Since z_i edges are processed in parallel in a single cycle, the time taken to complete an operation in junction i is $C_i = |\mathbf{W}_i|/z_i$ cycles. The **degree of parallelism configuration** $\mathbf{z}_{\text{net}} = (z_1, \dots, z_L)$ is chosen to achieve $C_i = C \quad \forall i \in \{1, \dots, L\}$. This allows efficient junction pipelining since each operation takes exactly C cycles to be completed for each input in each junction, which we refer to as a **junction cycle**. This determines throughput.

The following is an analysis of Fig. 4.1(c) in more detail for an example NN with $L = 2$. While a new training input numbered $n + 3$ is getting loaded as \mathbf{a}_0 , junction 1 is processing the FF stage for the previous input $n + 2$ and computing \mathbf{a}_1 . Simultaneously, junction 2 is processing FF and computing cost δ_L via cost

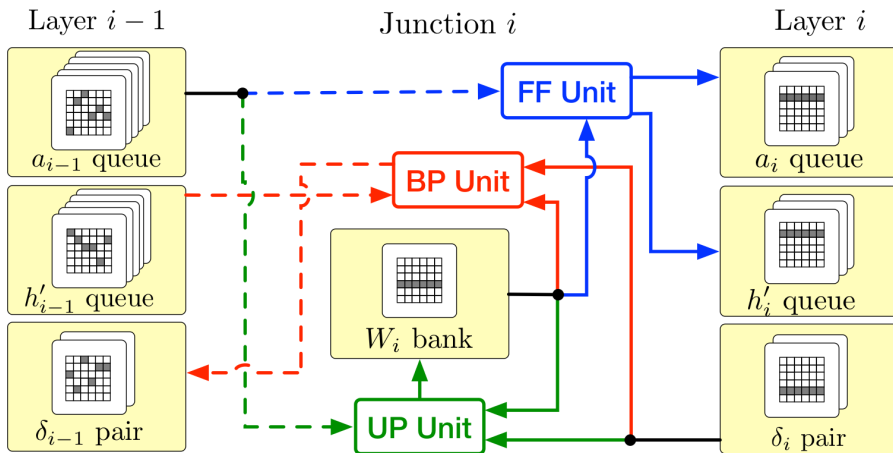


Figure 4.2: Architecture for parallel operations for an intermediate junction i ($i \neq 1, L$) showing the three operations along with associated inputs and outputs. Natural and interleaved order accesses are shown using solid and dashed lines, respectively. The \mathbf{a} and \mathbf{h}' memory banks occur as queues, the $\boldsymbol{\delta}$ memory banks as pairs, while there is a single weight memory bank.

derivatives for input $n + 1$. It is also doing BP on input n to compute $\boldsymbol{\delta}_1$, as well as updating (UP) its parameters from the finished $\boldsymbol{\delta}_L$ computation of input n . Simultaneously, junction 1 is performing UP using $\boldsymbol{\delta}_1$ from the finished BP results of input $n - 1$.² This results in *operational parallelism* in each junction, as shown in Fig. 4.2. The combined speedup is approximately a factor of $3L$ as compared to doing one operation at a time for a single input.

Notice from Fig. 4.2 that there is only one weight memory bank which is accessed for all three operations. However, UP in junction 1 needs access to \mathbf{a}_0 for input $n - 1$, as per the weight update equation (3.8b). This means that there need to be $2L + 1 = 5$ left activation memory banks for storing \mathbf{a}_0 for inputs $n - 1$ to $n + 3$,

²Note that BP does not occur in the first junction because there are no $\boldsymbol{\delta}_0$ values to be computed

Table 4.1: Hardware Architecture Total Storage Cost Comparison for $\mathbf{N}_{\text{net}} = (800, 100, 10)$ FC vs sparse with $\mathbf{d}_{\text{net}}^{\text{out}} = (20, 10)$, $\rho_{\text{net}} = 21\%$

Parameter	Expression	Count (FC)	Count (sparse)
\mathbf{a}	$\sum_{i=0}^{L-1} (2(L-i) + 1) N_i$	4300	4300
\mathbf{h}'	$\sum_{i=1}^{L-1} (2(L-i) + 1) N_i$	300	300
δ	$2\sum_{i=1}^L N_i$	220	220
\mathbf{b}	$\sum_{i=1}^L N_i$	110	110
\mathbf{W}	$\sum_{i=1}^L N_i d_i^{\text{in}}$	81000	17000
TOTAL	Σ (All above)	85930	21930

i.e. a **queue**-like structure. Similarly, UP in junction 2 will need $2(L-1) + 1 = 3$ queued banks for each of its left activation \mathbf{a}_1 and its derivative \mathbf{h}'_1 memories – for inputs from n (for which values will be read) to $n+2$ (for which values are being computed and written). There also need to be 2 banks for all δ memories – 1 for reading and the other for writing. Thus junction pipelining requires multiple memory banks, but only for layer parameters \mathbf{a} , \mathbf{h}' and δ , *not* for weights. The number of layer parameters is insignificant compared to the number of weights for practical networks. This is why pre-defined sparsity leads to significant storage savings, as quantified in Table 4.1 for the circled FC point vs the $\rho_{\text{net}} = 21\%$ point from Fig. 3.2(c). Specifically, memory requirements are reduced by 3.9X in this case. Furthermore, the computational complexity, which is proportional to the number of weights for a MLP, is reduced by 4.8X. For this example, these complexity reductions come at a cost of degrading the classification accuracy from 98.0% to 97.2%.

4.2 Memory organization

For the purposes of memory organization, edges are numbered sequentially from top to bottom on the right side of the junction. Other network parameters such as \mathbf{a} , \mathbf{h}' and $\boldsymbol{\delta}$ are numbered according to the neuron numbers in their respective layer. Consider Fig. 4.3 as an example, where junction i is flanked by $N_{i-1} = 12$ left neurons with $d_i^{\text{out}} = 2$ and $N_i = 8$ right neurons, leading to $|\mathbf{W}_i| = 24$ and $d_i^{\text{in}} = 3$. The three weights connecting to right neuron 0 are numbered 0, 1, 2; the next three connecting to right neuron 1 are numbered 3, 4, 5, and so on. A particular right neuron connects to some subset of left neurons of cardinality d_i^{in} .

Each type of network parameter is stored in a **bank of memories**. The example in Fig. 4.3 uses $z_i = 4$, *i.e.* 4 weights are accessed per cycle. We designed the weight memory bank to have the minimum number of memories to prevent clashes, *i.e.* z_i , and their depth equals C_i . Weight memories are read in natural order – 1 row per cycle (shown in same color).

Right neurons are processed sequentially due to the weight numbering. The number of right neuron parameters of a particular type needing to be accessed in a cycle is upper bounded by $\lceil z_i/d_i^{\text{in}} \rceil$, which leads to $z_{i+1} \geq \lceil z_i/d_i^{\text{in}} \rceil$ in order to prevent clashes in the right memory bank. (This does not limit most practical designs, as will be shown in Section 4.5). For FF in Fig. 4.3 for example, cycles 0 and 1 finish computation of $a_i^{(0)}$ and $a_i^{(1)}$ respectively, while cycle 2 finishes computing both $a_i^{(2)}$ and $a_i^{(3)}$. For BP or UP, everything remains same except for

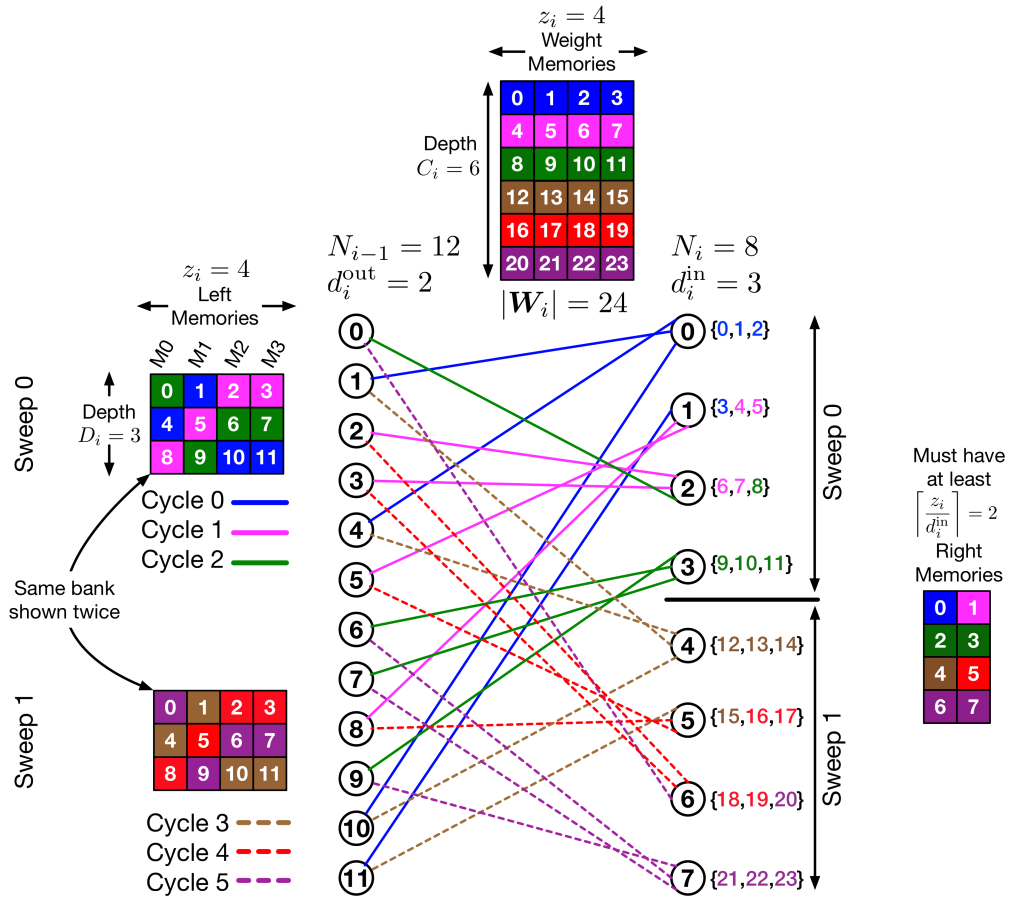


Figure 4.3: An example of processing inside junction i with $z_i = 4$ memories in the weight and left banks, and $z_{i+1} = 2$ memories in the right bank. The banks are represented as numerical grids, each column is a memory, and the number in each cell is the number of the edge / left neuron / right neuron whose parameter value is stored in it. Edge are sequentially numbered on the right (shown in curly braces). $z_i = 4$ weights are read in each of the six cycles, with the first three colored blue, pink and green, respectively. These represent sweep 0, while the next 3 (using dashed lines) colored brown, red and purple, respectively, represent sweep 1. Clash-freedom leads to at most one cell from each memory in each bank being accessed each cycle. Weight and right memories are accessed in natural order, while left memories are accessed in interleaved order.

the right memory accesses. Now $\delta_i^{(0)}$ and $\delta_i^{(1)}$ are used in cycle 0, $\delta_i^{(1)}$ and $\delta_i^{(2)}$ in cycle 1, and $\delta_i^{(2)}$ and $\delta_i^{(3)}$ in cycle 2. Thus the maximum number of right neuron parameters ever accessed in a cycle is $\lceil z_i/d_i^{\text{in}} \rceil = 2$.

Since edges are interleaved on the left, in general, the z_i edge processing logic units will need access to z_i parameters of a particular type from layer $i - 1$. So all the left memory banks have z_i memories, each of depth $D_i = N_{i-1}/z_i$, which are accessed in interleaved order. For example, after D_i cycles, N_{i-1} edges have been processed – *i.e.* $(D_i \times z_i) = N_{i-1}$. We require that each of these edges be connected to a different left neuron to eliminate the possibility of duplicate edges. This completes a **sweep**, *i.e.* one complete access of the left memory bank. Since each left neuron connects to d_i^{out} edges, d_i^{out} sweeps are required to process all the edges, *i.e.* each left activation is read d_i^{out} times in the whole junction cycle. The reader can verify that D_i cycles multiplied by d_i^{out} sweeps results in C_i total cycles, *i.e.* one junction cycle.

4.3 Clash-freedom

We define a **clash** as attempting to perform a particular operation more than once on the same memory at the same time, which would stall processing. The idea of **clash-freedom** is to pre-define a pattern of connections and z values such that no operation in any junction of the NN results in a clash. Section 4.2 described how z values should be designed to prevent clashes in the weight and right memory banks.

This subsection analyzes the left memory banks, which are accessed in interleaved order. Their memory access pattern should be designed so as to prevent

clashes. Additionally, the following properties are desired for practical clash-free patterns. Firstly, it should be easy to find a pattern that gives good performance. Secondly, the logic and storage required to generate the left memory addresses should be low complexity.

We generate clash-free patterns by initially specifying the left memory addresses to be accessed in cycle 0 using a **seed vector** $\phi_i \in \{0, 1, \dots, D_i - 1\}^{z_i}$. Subsequent addresses are cyclically generated. Considering Fig. 4.3 as an example, $\phi_i = (1, 0, 2, 2)$. Thus in cycle 0, we access addresses $(1, 0, 2, 2)$ from memories $(M0, M1, M2, M3)$, *i.e.* left neurons $(4, 1, 10, 11)$. In cycle 1, the accessed addresses are $(\phi_i + 1) \% D_i = (2, 1, 0, 0)$, and so on. Since $D_i = 3$, cycles 3–5 access the same left neurons as cycles 0–2.

We found that this technique results in a large number of possible connection patterns, and randomly sampling from this set results in performance comparable with non-clash-free NNs. These ideas are further discussed in Chapter 5.

4.4 Batch size

As shown in (2.14), for a batch size of M , the UP operation in (3.8) is performed only once for M inputs by using the average over the M gradients. Our architecture performs an UP for every input and therefore may be viewed as having batch size one. However, the processing in our architecture differs from a typical software

implementation with $M = 1$ due to the pipelined and parallel operations. Specifically, in our architecture, FF and BP for the same input use different weights, as implied by Fig. 4.1(c). However, this deviation from the conventional notion of batch size did not result in performance degradation in our initial hardware implementation (presented in Section 4.7). We would also like to point out that there is considerable ambiguity in the literature regarding ideal batch sizes (cf. [21, 78]).

The concept of batch size for our architecture can be modified in several different ways. Firstly, by performing UP once every ($M > 1$) samples. As an example, a practical case for a deep MLP could be $M = 64$ and $L = 5$. For all such cases where $M \gg L$, FF and BP will use the same weights for most inputs in a batch. Secondly, a more conventional minibatch update can also be obtained by completely removing the UP logic from the junction pipeline. After performing FF and BP on M samples, the pipeline will be flushed and the averaged gradients over M samples used to update all parameters in all junctions simultaneously. This would eliminate the UP arithmetic units from the pipeline, at the cost of increased storage for accumulating intermediate values from (3.8). However, we have not implemented these ideas in actual hardware yet.

4.5 Architectural Constraints

The depth of left memories in our hardware architecture is $D_i = N_{i-1}/z_i$. Thus N_{i-1} should preferably be an integral multiple of z_i . This is not a burdening

constraint since the choice of z_i is independent of network parameters and depends on the capacity of the device. In the unusual case that this constraint cannot be met, the extra cells in memories can be filled with dummy values such as 0.

There are also 2 conditions placed on the z_i values to eliminate stalls in processing: for all layers $i \in \{1, \dots, L\}$, (i) $|\mathbf{W}_i|/z_i = C$, and (ii) $z_{i+1} \geq \lceil z_i/d_i^{\text{in}} \rceil$. Let us examine how much of a burden these constraints impose on choosing a network configuration. Without loss of generality, (i) can be written as:

$$\begin{aligned}
\frac{|\mathbf{W}_{i+1}|}{z_{i+1}} &= \frac{|\mathbf{W}_i|}{z_i} \\
\Rightarrow z_{i+1} &= \frac{|\mathbf{W}_{i+1}| z_i}{|\mathbf{W}_i|} \\
\Rightarrow z_{i+1} &= \frac{z_i d_{i+1}^{\text{out}}}{d_i^{\text{in}}} \tag{4.1}
\end{aligned}$$

Then (ii) becomes:

$$\begin{aligned}
\frac{z_i d_{i+1}^{\text{out}}}{d_i^{\text{in}}} &\geq \left\lceil \frac{z_i}{d_i^{\text{in}}} \right\rceil \\
\Rightarrow d_{i+1}^{\text{out}} &\geq \frac{d_i^{\text{in}}}{z_i} \left\lceil \frac{z_i}{d_i^{\text{in}}} \right\rceil \tag{4.2}
\end{aligned}$$

which needs to be satisfied $\forall i \in \{1, \dots, L-1\}$.

In practice, it is desirable to design z_i/d_i^{in} to be an integer so that an integral number of right neurons finish processing every cycle. This simplifies hardware implementation by eliminating the need for additional storage, for example, of the

intermediate activation values during FF. In this case, (4.2) reduces to $d_{i+1}^{\text{out}} \geq 1$, which is always true.

For non-integral z_i/d_i^{in} , there are two cases. If $z_i > d_i^{\text{in}}$, (4.2) reduces to $d_{i+1}^{\text{out}} \geq 2$. On the other hand, if $z_i < d_i^{\text{in}}$, there is no bound on the right hand side of (4.2). In general, note that (4.2) becomes a burdening constraint only if d_i^{in} is large, and d_{i+1}^{out} and z_i are both desired to be small. This corresponds to earlier junctions being denser than later, which is typically not desirable according to the observations in Sec. 3.3.3, or to very limited hardware resources. We thus conclude that (4.2) is not a limiting constraint in most practical cases.

4.6 Special Case: Processing a FC junction

Fig. 4.4 shows the FC version of the junction from Fig. 4.3, which has 96 edges to be accessed and operated on. This can be done keeping the same junction cycle $C_i = 6$ by increasing z_i to 16, *i.e.* using more hardware. On the other hand, if hardware resources are limited, one can use the same $z_i = 4$ and pay the price of a longer junction cycle $C_i = 24$, as shown in Fig. 4.4. *This demonstrates the flexibility of our architecture.*

Note that FC junctions are clash-free in all practical cases due to the following reasons. Firstly, the left memory accesses are in natural order just like the weights, which ensures that no more than one element is accessed from each memory per cycle. Secondly, $\lceil z_i/d_i^{\text{in}} \rceil = 1$ for all practical cases since $z_i \leq N_{i-1}$, as discussed

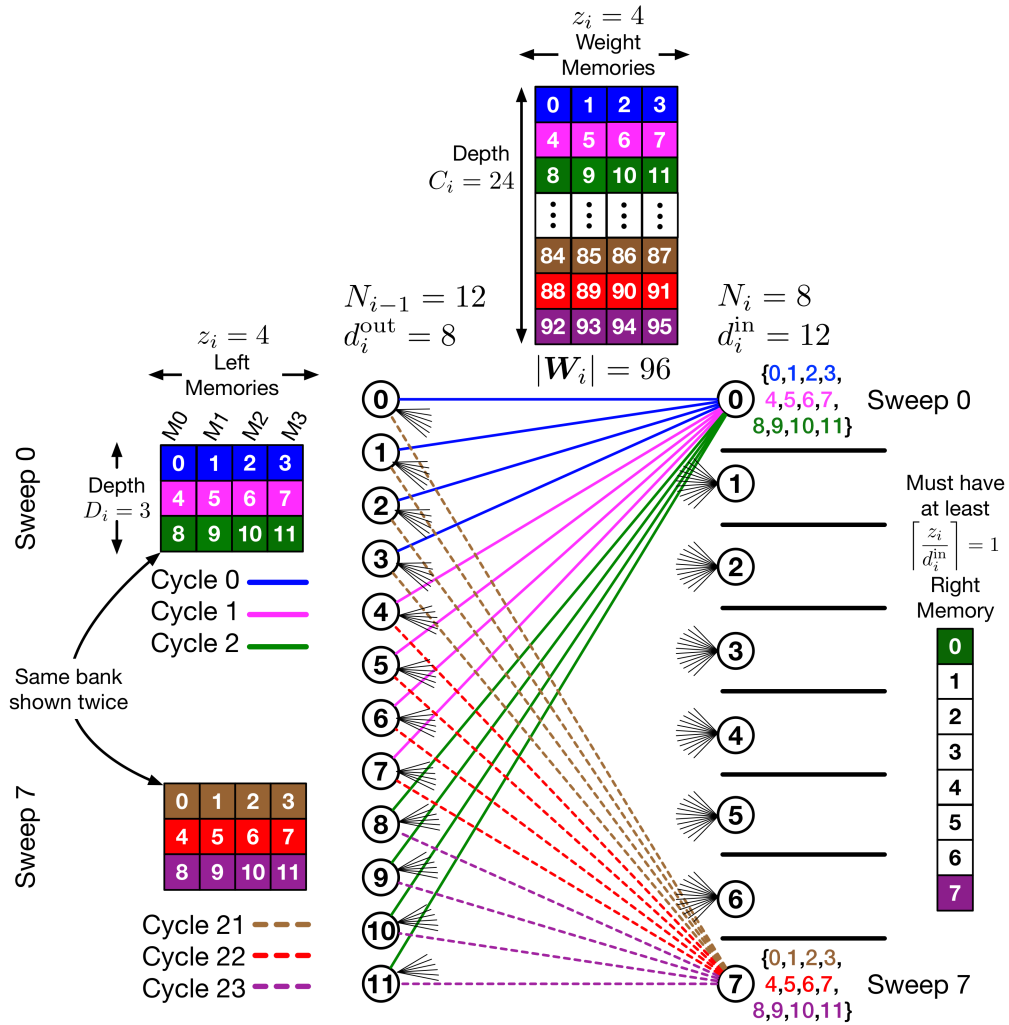


Figure 4.4: Processing the FC version of the junction from Fig. 4.3. For clarity, only the first 12 and last 12 edges (dashed) are shown, corresponding respectively to right neurons 0 and 7, sweeps 0 and 7, cycles 0–2 and 21–23.

in Section 4.5, and $d_i^{\text{in}} = N_{i-1}$ always for FC junctions. This means that at most one right neuron is processed in a cycle, so clashes will never occur when accessing the right memory bank. As an example, in Fig. 4.4, one right neuron finishes processing every 3rd cycle.

Note that compared to Fig. 4.3, the weight memories in Fig. 4.4 are deeper since C_i has increased from 6 to 24. However, the left layer memories remain the same size since $N_{i-1} = 12$ and $z_i = 4$ are unchanged, but the left memory bank is accessed more times since the number of sweeps has increased from 2 to 8. Also note that even if cycle 0 (blue) accesses some other clash-free subset of left neurons, such as $\{4, 5, 6, 7\}$ instead of $\{0, 1, 2, 3\}$, the connection pattern would remain unchanged. This implies that *different memory access patterns do not necessarily lead to different connection patterns*. This point will be discussed further in Chapter 5.

4.7 FPGA Implementation

This section describes our implementation of the architecture described thus far. Certain efforts in this implementation were achieved by Diandian Chen, – such as Universal asynchronous receiver-transmitter (UART) interfacing and observing outputs via Light Emitting Diode (LED)s – and hence have not been presented in detail in this dissertation proposal.

We developed our implementation on a Digilent Nexys-4 DDR board containing a FPGA from the **Xilinx Artix-7** family with part number XC7A100T-1CSG324C [79]. This is a fairly modest FPGA in terms of capacity and performance, cf. the Xilinx Virtex UltraScale+ series [80] or the cloud-based FPGAs from Amazon Web Services [81] for more powerful alternatives. However, our purpose was to obtain

an initial proof-of-concept in which we could explore efficient design styles and optimize the Verilog Register Transfer Level (RTL) code to make it more robust and scalable. The Artix-7 device serves this purpose. We are currently in the process of extending our initial implementation to more powerful platforms.

Our initial proof-of-concept implementation used the MNIST dataset. We used powers of 2 for most network configuration values to simplify the hardware realization. Accordingly we padded each input with 0s to extend the number of features from 784 to 1024 (this does not alter performance, as mentioned in Section 3.3.1), and each ground-truth output with 0s to extend the number of labels from 10 to 32. Note that the output encoding remains one-hot.

4.7.1 Network Configuration and Training Setup

We implemented a structured pre-defined sparse NN with $\mathbf{N}_{\text{net}} = (1024, 64, 32)$ and $\mathbf{d}_{\text{net}}^{\text{out}} = (4, 16)$, leading to $\rho_{\text{net}} = 7.576\%$. The complete configuration is given in Table 4.2. The values shown were chosen on the basis of hardware constraints and experimental results, which will be described in Sections 4.7.2 and 4.7.3.

We selected 12,544 inputs to comprise 1 epoch of training. Learning rate (η) is initially 2^{-3} , decay is implemented by halving after the first 2 epochs, then after every 4 epochs until its value became 2^{-7} . Designing η to always be a power of 2 leads to the multiplications with η in eq. (3.8) getting reduced to bit shifts. Pre-defined sparsity leads to a total number of trainable parameters = $(|\mathbf{W}_1| = 4096) +$

Table 4.2: Implemented Network Configuration

Junction Number (i)	1	2
Left Neurons (N_{i-1})	1024	64
Right Neurons (N_i)	64	32
Out-degree (d_i^{out})	4	16
Number of Weights ($ \mathbf{W}_i = N_{i-1} \times d_i^{\text{out}}$)	4096	1024
In-degree ($d_i^{\text{in}} = W_i/N_i$)	64	32
z_i	128	32
Junction cycle ($ \mathbf{W}_i /z_i$) ^a	32	32
Density ($\rho_i = W_i/(N_{i-1}N_i)$)	6.25%	50%
Overall Density	$\rho_{\text{net}} = 7.576\%$	

^aIn terms of number of clock cycles. Not considering the additional clock cycles needed for memory accesses.

$(|\mathbf{W}_2| = 1024) + (|\mathbf{b}_1| = N_1 = 64) + (|\mathbf{b}_2| = N_2 = 32) = 5216$, which is much less than 12,544, so we theorized that overfitting was not an issue. We verified this using software simulations, and hence did not apply any regularization.

4.7.2 Bit Width Considerations

Parameter Initialization

We initialized weights using the Glorot Normal technique (2.16), which translates to a three standard deviation range of ± 0.51 for junction 1 and ± 0.61 for junction 2 in our NN configuration described in Table 4.2.

The biases in our architecture are stored along with the weights as an augmentation to the weight memory banks. So we initialized biases in the same manner

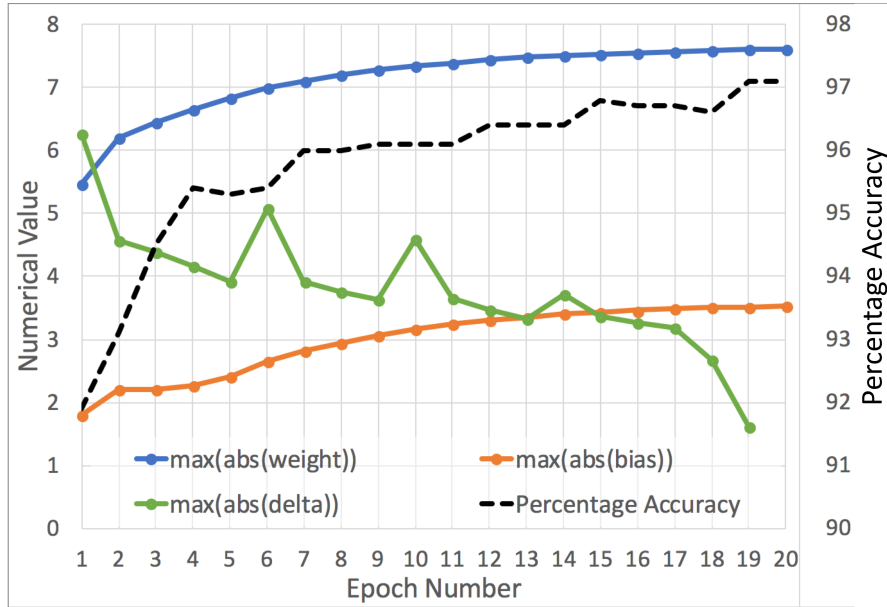


Figure 4.5: Maximum absolute values (left y-axis) for all weights, biases, and deltas in the NN, and percentage classification accuracy (right y-axis), as the NN is trained.

as weights. Software simulations showed that this led to no degradation in performance from the conventional method of initializing biases with constant values such as 0 or 0.1. This makes sense since the maximum absolute value from initialization is much closer to 0 than their final values when training converges, as shown in Fig. 4.5.

To simplify the RTL, we used the same set of $|\mathbf{W}_i|/z_i$ unique values to initialize all weights and biases in junction i . Again, software simulations showed that this led to no degradation in performance as compared to initializing all of them randomly. This is not surprising since an appropriately high value of initial learning rate will drive each weight and bias towards its own optimum value, regardless of similar values at the start.

Fixed Point Configuration

We recreated the aforementioned initial conditions in software and trained our configuration to study the range of values for network variables until convergence. The results for are shown in Fig. 4.5 for all weights, biases and delta values. We used the **sigmoid** activation function for our hardware implementation, which has output range = $[0, 1]$.

To keep the hardware optimal, we decided on the **same fixed point bit configuration** for all computed values of network parameters — $\{\mathbf{a}_i, \mathbf{h}'_i, \boldsymbol{\delta}_i, \mathbf{W}_i, \mathbf{b}_i\}$, $i \in \{1, 2\}$. This uses less hardware resources compared to implementations such as a) [50], which uses fixed point adders, but more resource-intensive floating point multipliers and floating-to-fixed-point converters, or b) [49], which incurs additional Digital Signal Processing (DSP) resources due to its use of probabilistic fixed point rounding techniques.

Our configuration is characterized by the bit triplet (b_w, b_n, b_f) , which are respectively the total number of bits, integer bits, and fractional bits, with the constraint $b_w = b_n + b_f + 1$, where the 1 is for the sign bit. This gives a numerical range of $[-2^{b_n}, 2^{b_n} - 2^{-b_f}]$ and precision of 2^{-b_f} . Fig. 4.5 shows that the maximum absolute values of various network parameters during training stays within 8. Accordingly we set $b_n = 3$. We then experimented with different values for the bit triplet and obtained different utilizations for the FPGA Look Up Table (LUT)

Table 4.3: Effect of Bit Width on Performance

b_w	b_n	b_f	FPGA LUT Utilization %	Accuracy after 1 epoch	Accuracy after 15 epochs
8	2	5	37.89	78	81
10	2	7	72.82	90.1	94.9
10	3	6	63.79	88	93.8
12	3	8	83.38	90.3	96.5
16	4	11	112	91.9	96.5

(*i.e.* logic) resources. The results are shown in Table 4.3. For this hardware implementation, we measured accuracy as classification performance on the last 1000 training samples³. Noting the diminishing returns and impractical utilization of hardware resources for high bit widths, we chose the bit triplet (12, 3, 8) as being the optimal case, *i.e.* this is a **12-bit implementation**.

Dynamic Range Reduction due to Sparsity

We found that sparsity leads to reduction in the dynamic range of network parameters, since the summations in (3.6) and (3.7) are over smaller ranges. This motivated us to use a special form of adder and multiplier which preserves the bit triplet between inputs and outputs by clipping large absolute values of output to either the positive or negative maximum allowed by the range. For example, 10 would become 7.996 and -10 would become -8 .

³This differs from the conventional method of measuring performance as classification accuracy on the test set. We chose the last 1000 training samples to measure performance on in order to simplify the measurement techniques for this initial proof-of-concept implementation.

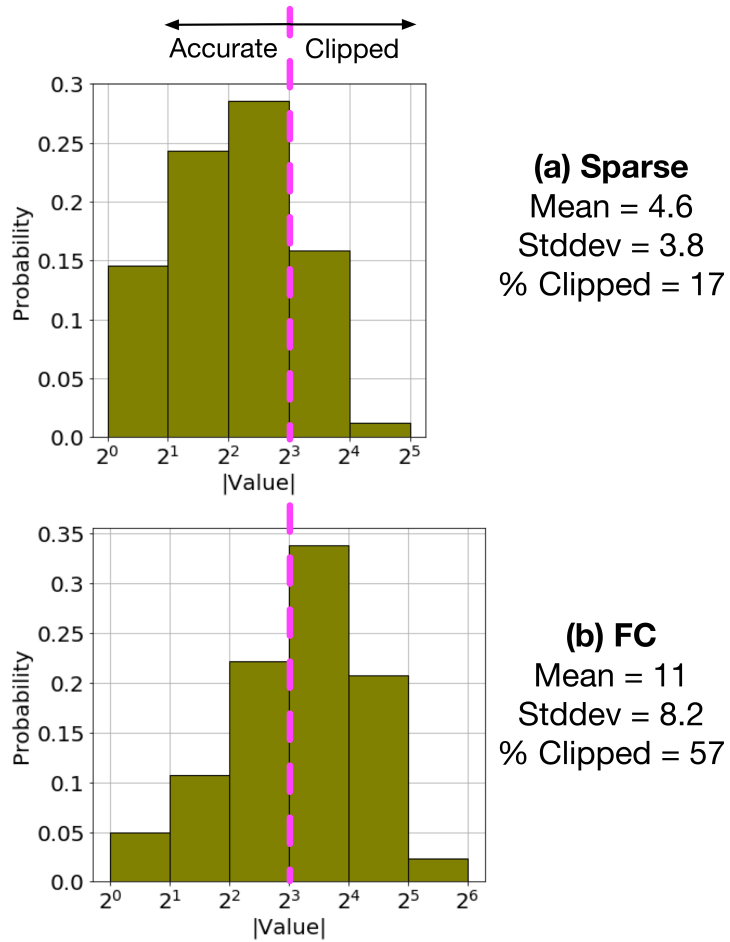


Figure 4.6: Histograms of absolute value of s_1 with respect to dynamic range for (a) sparse vs. (b) fully connected cases, as obtained from software simulations. Values right of the pink line are clipped.

Fig. 4.6 analyzes the worst clipping errors by comparing the absolute values of the argument of the sigmoid function in the hidden layer, *i.e.* $s_1 = \sum \mathbf{W}_1 \mathbf{a}_0 + \mathbf{b}_1$ from (3.6a), for our sparse case vs. the corresponding FC case (which would have $d_1^{\text{out}} = 64$, $d_2^{\text{out}} = 32$). Notice that the sparse case only has 17% of its values clipped due to being outside the dynamic range afforded by $b_n = 3$, while the FC case has 57%. The sparse case also has a smaller variance. This implies that the

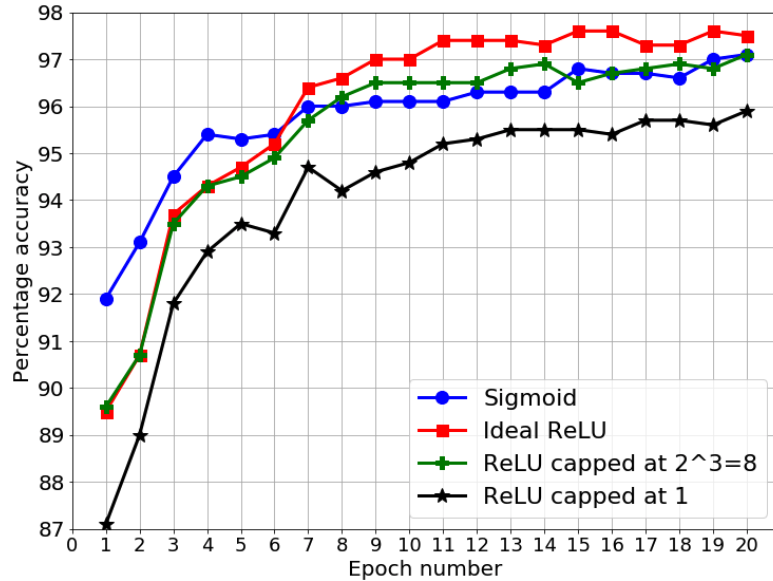


Figure 4.7: Comparison of activation functions for α_1 , as obtained from software simulations.

hardware errors introduced due to finite bit-width effects are less pronounced for our pre-defined sparse configuration as compared to FC.

Why we chose sigmoid activations?

As has been amply demonstrated in literature (*e.g.* [1, 9, 82]), the ideal ReLU activation function ((2.8)) is more widely used than sigmoid ((2.9)) due to the former's better performance. This is primarily due to elimination of the vanishing gradient problem. Note also that ReLU has a tendency towards generating sparse outputs, *i.e.* several nodes output zero values for activation.

However, ideal ReLU is not practical for hardware due to its unbounded range. We experimented with a modified form of ReLU where the outputs were clipped to a) 8, which is the maximum supported by $b_n = 3$, and b) 1, to preserve bit width

consistency in the multipliers and adders and ensure compatibility with sigmoid activations. Fig. 4.7 shows software simulations comparing sigmoid with these cases. Note that ReLU clipped at 8 converges similar to sigmoid, but sigmoid has better initial performance. Also, there is no need to promote extra sparsity by using ReLU because our configuration is already sparse. Finally, in our hardware implementation, sigmoid does not suffer from vanishing gradient problems due to just having two junctions. We therefore concluded that sigmoid activation is the best choice for *all* layers. For the output layer, we simply picked the correct class as corresponding to the output node with the maximum sigmoid output value.

4.7.3 Implementation Details

Sigmoid Activation

The sigmoid function uses exponentials, which are computationally infeasible to obtain in hardware. So we pre-computed the values of $\sigma(\cdot)$ and $\sigma'(\cdot)$ and stored them in LUTs. Interpolation was not used, instead we computed sigmoid for all 4096 possible 12-bit arguments up to the full 8 fractional bits of accuracy. On the other hand, its derivative values were computed to 6 fractional bits of accuracy since they have a range of $[0, 2^{-2}]$. The number of sigmoid LUTs required is

$$\sum_{i=1}^L z_i/d_i^{\text{in}} = 3.$$

Arithmetic Units

As indicated in Table 4.2, we designed z_i/d_i^{in} to be an integer for both junctions, which ensures that in every cycle, an integral number of right neurons finish their FF processing. This implies that the FF summations in (3.6a) can occur in a single cycle, and eliminates the need for storing partial sums. The total number of multipliers required for FF is $\sum_{i=1}^L z_i = 160$. The summations also use a tree adder of depth $= \log_2(d_i^{\text{in}})$ for every neuron processed in a cycle.

The BP summation in (3.7b) will need several cycles to complete for a single left neuron since weight numbering is permuted on the left. This necessitates storing $\sum_{i=2}^L z_i$ partial sums, however, tree adders are no longer required. Note that (3.7b) has 2 multiplications, and moreover, BP does not occur in the first junction, so the total number of multipliers required for BP is $2 \sum_{i=2}^L z_i = 64$.

The UP operation in each junction i requires z_i adders for the weights and z_i/d_i^{in} adders for the biases, since that many right neurons are processed every cycle. Only the weight update requires multipliers, so their total number is $\sum_{i=1}^L z_i = 160$.

Our FPGA device has 240 DSP blocks. Accordingly, we implemented the 224 FF and BP multipliers using 1 DSP for each, while the other 160 UP multipliers and all adders were implemented using logic.

Memories and Data

All memories were implemented using Block Random Access Memory (BRAM). The memories for \mathbf{a} and \mathbf{h}' never need to be read from and written into in the same cycle, so they are single-port. The δ memories are true dual-port, i.e. both ports support reads and writes. This is required due to their read-modify-write nature, since they accumulate partial sums. The ‘weight+bias’ memories are simple dual-port, with one port used exclusively for reading the k th cell in cycle k , and the other for simultaneously writing the $(k - 1)$ th cell. These were initialized using Glorot normal values while all other memories were initialized with zeroes.

Note that regardless of access order, the logic to compute memory addresses simply consists of z_i incrementers. This is because natural order accesses occur sequentially, while interleaved order accesses also occur sequentially according to the seed vector. A more complete description of clash freedom is deferred to Chapter 5.

The ground truth one-hot encoding for all 12,544 inputs were stored in a single-port BRAM, and initialized with word size = 10 to represent the 10 MNIST outputs. After reading, the word was padded with 0s to make it 32-bit long. On the other hand, the input data was too big to store on-chip. Since the native MNIST images are $28 \times 28 = 784$ pixels, the total input data size is $12544 \times 784 \times 8 = 78.68$ Mb, while the total device BRAM capacity is only 4.86 Mb. So the input data was fed from a host computer using the UART interface.

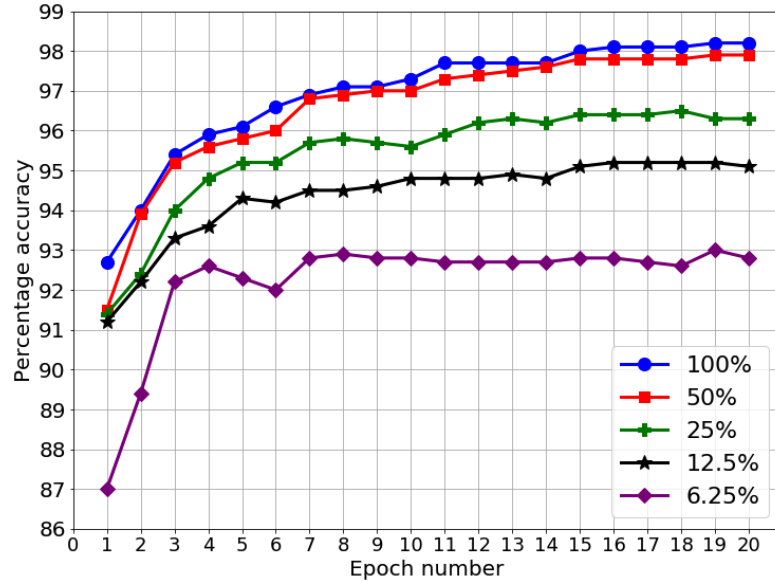


Figure 4.8: Performance for different ρ_2 , keeping ρ_1 fixed at 6.25%, as obtained from software simulations.

Network Configuration

Here we explain the choice of NN configuration in Table 4.2. We initially picked $N_2 = 16$, which is the minimum power of 2 above 10. Since later junctions need to be denser than earlier ones to improve performance, we experimented with junction 2 density and show its effects on performance in Fig. 4.8. We concluded that 50% density is optimum for junction 2. Attempting to meet the architectural constraints in Section 4.5 led to $z_1 = 256$, which was beyond the capacity of our FPGA. So we increased N_2 to 32 and set z_2 to the minimum value of 32, leading to $z_1 = 128$. We also experimented with $d_1^{\text{out}} = 8$, but the resulting accuracy was within 1 percentage point of our final choice of $d_1^{\text{out}} = 4$.

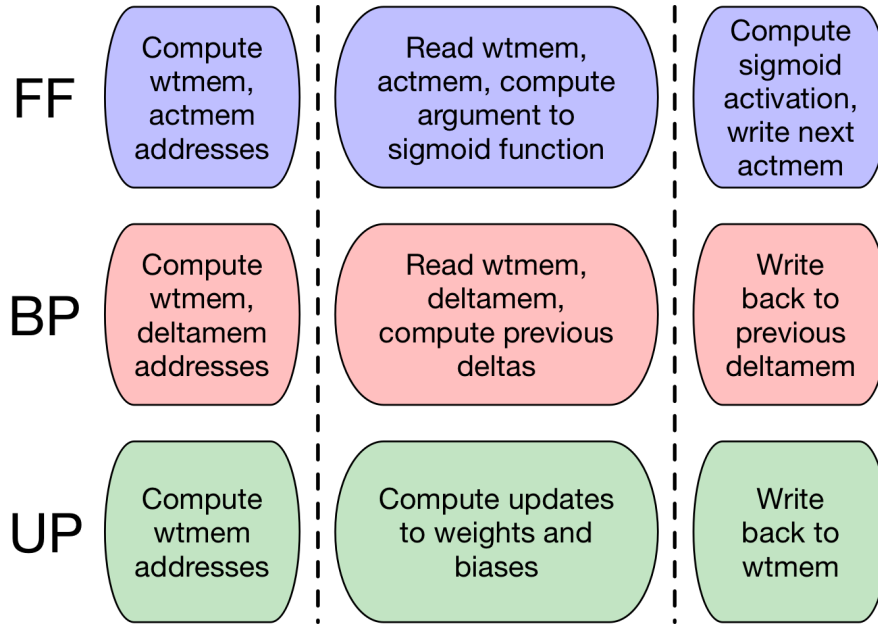


Figure 4.9: Breaking up each operation into 3 clock cycles.

Timing and Results

A junction cycle in our design is actually $(|\mathbf{W}_i|/z_i + 2) = 34$ clock cycles since each set of z_i weights in a junction need a total of 3 clock cycles for each operation, as explained in Fig. 4.9. To summarize, the first and third are used to compute memory addresses, while the second is the most time-consuming since it performs most of the arithmetic computations. This determines our clock frequency, which is 15MHz.

We stored the results of several training inputs and fed them out to 10 LEDs on the board, each representing an output from 0-9. The FPGA implementation

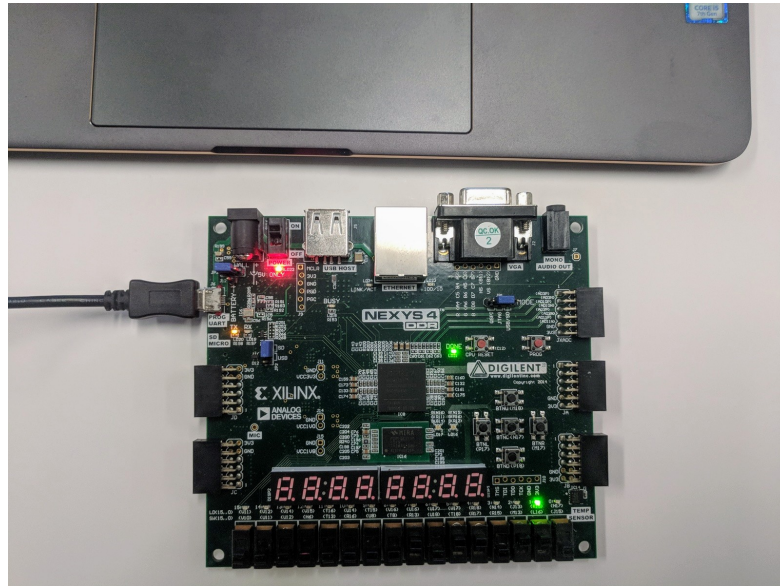


Figure 4.10: Our design working on the Xilinx XC7A100T-1CSG324C FPGA.

performed according to RTL simulations and close to the ideal floating point software simulations, giving 96.5% accuracy in 14 epochs of training. Fig. 4.10 shows our FPGA implementation in action.

4.8 Summary

This chapter described our proposed hardware architecture and detailed an initial FPGA implementation we developed. While this is a good initial proof-of-concept, there are significant engineering tasks required to demonstrate state-of-the-art training speeds. This is an area of ongoing work.

Chapter 5

Connection Patterns

A fully connected neural network consists of fully connected junctions, *i.e.* every node in a layer connects to every node in its previous and next layer. This gives rise to a single possible **connection pattern**. On the other hand, a structured pre-defined sparse NN with a fixed N_{net} and $\mathbf{d}_{\text{net}}^{\text{out}}$ has several different ways of connecting the nodes in each junction. Fig. 5.1 illustrates this for a single junction with 4 left nodes and 2 right nodes. Subfigure (a) shows the FC case, while subfigures (b)–(g) show the six possible ways to arrange the connections for the 50% density case.

In practice, we desire connection patterns which give good performance and are hardware friendly, *i.e.* lead to clash-free memory accesses and are simple to implement. The remainder of this chapter will elaborate on these ideas.

5.1 Biadjacency Matrices

First we will introduce a useful concept. The biadjacency matrix for a bipartite graph is a rectangular matrix with 1 and 0 denoting the presence and absence of an edge, respectively [83]. Accordingly:

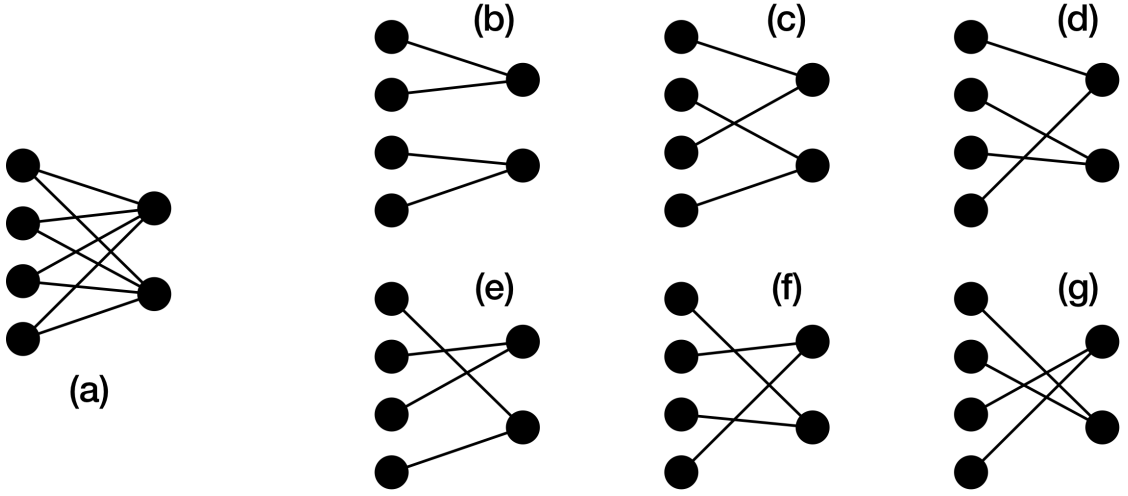


Figure 5.1: Different connection patterns for a single-junction network with $N_{\text{net}} = (4, 2)$. (a) The FC case, with just one possible connection pattern. (b)–(g) Six different ways of arranging connections when $d_{\text{net}}^{\text{out}} = (1)$, *i.e.* $d_{\text{net}}^{\text{in}} = (2)$. Notice that in each case, each left node has exactly one outgoing connection and each right node has exactly two incoming connections.

► **Definition 6: Biadjacency matrix:** The biadjacency matrix \mathcal{B}_i for $N_i \times N_{i-1}$ junction i is defined such that $\mathcal{B}_i^{(j,k)} = 1$ if an edge exists between node k in layer $i - 1$ and node j in layer i , otherwise $\mathcal{B}_i^{(j,k)} = 0$.

Then, if junction i is structured pre-defined sparse, \mathcal{B}_i has exactly d_i^{in} 1s in each row and exactly d_i^{out} 1s in each column, *i.e.* :

$$\mathbf{1}^T \mathcal{B}_i = d_i^{\text{out}} \mathbf{1}^T \quad (5.1a)$$

$$\mathcal{B}_i \mathbf{1} = d_i^{\text{in}} \mathbf{1} \quad (5.1b)$$

where $\mathbf{1}$ is the vector of all 1s with size set by the context. Moreover, the density can be computed as the fraction of 1s in the biadjacency matrix:

$$\rho_i = \frac{\sum_{j=1}^{N_i} \sum_{k=1}^{N_{i-1}} \mathcal{B}_i^{(j,k)}}{N_i N_{i-1}} \quad (5.2)$$

Now we attempt to tackle the question of the number of possible connection patterns in a junction i , which we denote as G_i^{CP} . We saw that the single junction in Fig. 5.1 with the structured constraints $d_i^{\text{out}} = 1$ and $d_i^{\text{in}} = 2$ had $G_i^{\text{CP}} = 6$. These are the six possible (2×4) $\{0, 1\}$ -only matrices with each row summing to 2 and each column summing to 1, as given below:

$$\left\{ \begin{array}{l} \left[\begin{array}{cccc} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{array} \right], \left[\begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array} \right], \left[\begin{array}{cccc} 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{array} \right], \\ \left[\begin{array}{cccc} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{array} \right], \left[\begin{array}{cccc} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{array} \right], \left[\begin{array}{cccc} 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{array} \right] \end{array} \right\}$$

Generalizing this, we can state:

Theorem 2

Given a junction i with N_{i-1} , N_i and ρ_i (from which d_i^{out} and d_i^{in} can be computed), the number of possible connection patterns it can have, G_i^{CP} , is the same as the number of possible matrices of dimensions $N_i \times N_{i-1}$

consisting only of $\{0, 1\}$ entries such that sum of each row equals d_i^{in} and sum of each column equals d_i^{out} .

Moreover, the number of possible connection patterns in a complete L -junction NN is

$$G_{\text{net}}^{\text{cp}} = \prod_{i=1}^L G_i^{\text{cp}} \quad (5.3)$$

Finding the exact value of G_i^{cp} for any junction may not be feasible, cf. [84]. However, computing the exact value is usually of little interest. As will be discussed subsequently, finding a class of connection patterns with certain desirable properties is of more importance.

5.2 Clash-free memory access patterns

Section 4.3 introduced the idea of a clash as attempting to perform a particular operation more than once on the same memory at the same time, which would stall processing. For convenience in the upcoming discussion, Fig. 4.3 is repeated as Fig. 5.2.

Our edge-based architecture described in Chapter 4 is inspired by architectures proposed for iterative decoding of modern sparse-graph-based error correction codes (*i.e.* Turbo and Low Density Parity Check (LDPC) codes) (cf., [85, 86]).

For modern codes, the clash-free property of the memories is ensured by defining clash-free interleavers (*i.e.* permutations) [87], or clash-free parity check matrices [86]. Similarly in our architecture, the connection pattern in a junction is defined by an interleaver, *i.e.* a mapping from sequentially numbered edges on the right to permuted numbering on the left. This leads to interleaved order access of the left memory banks, and the possibility of clashes. We defined a simple algorithm to prevent this from happening – start with a seed vector ϕ_i specifying the left memory addresses to be accessed in cycle 0, then perform $(\phi_i + 1) \% D_i$ for every subsequent cycle. This defines a certain left **memory access pattern** for junction i , and can be visualized using the shaded left memories in Fig. 5.2. Similar to connection patterns, we define G_i^{map} and $G_{\text{net}}^{\text{map}}$ for the number of possible left memory access patterns in junction i and a complete NN, respectively (*i.e.* $G_{\text{net}}^{\text{map}} = \prod_{i=1}^L G_i^{\text{map}}$).

When $z_i \geq d_i^{\text{in}}$, which is expected to be true for practical cases of implementing sparse NNs on powerful hardware devices, we state a result without formal proof: G_i^{cp} is lower bounded by G_i^{map} . Thus, for a given NN with \mathbf{N}_{net} and $\mathbf{d}_{\text{net}}^{\text{out}}$ and a given hardware device with \mathbf{z}_{net} , $G_{\text{net}}^{\text{map}}$ **is a measure of the number of possible pre-defined sparse connection patterns which are clash-free and hardware friendly.**

For the case of $z_i < d_i^{\text{in}}$ however, it could happen that $G_i^{\text{cp}} < G_i^{\text{map}}$. In particular, a FC junction has only one possible connection pattern, as shown in Fig.

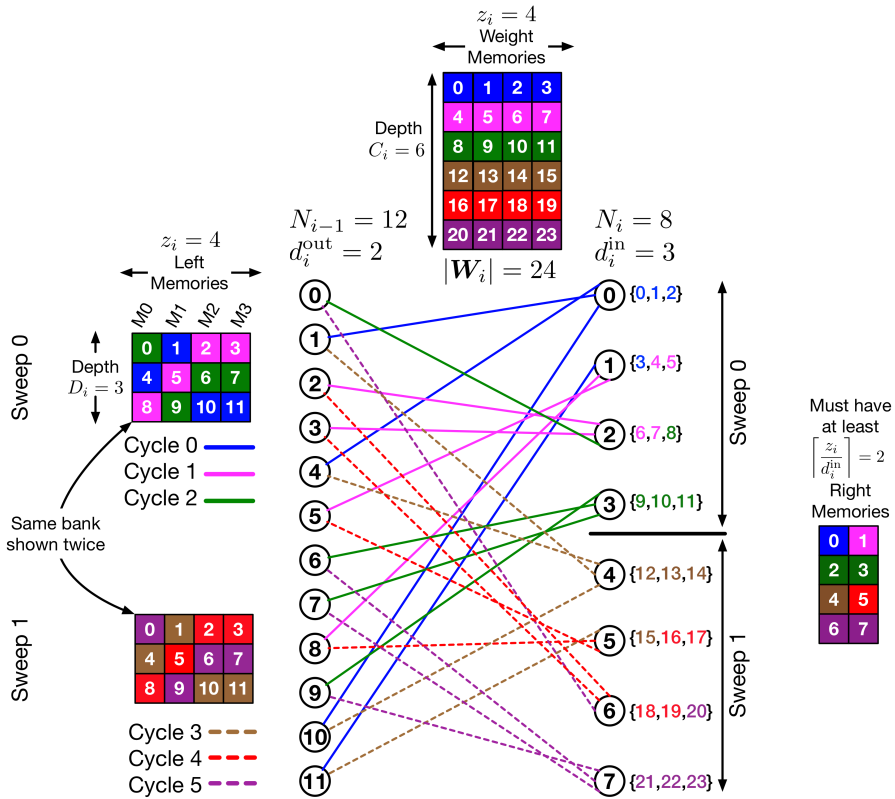


Figure 5.2: Repeating Fig. 4.3 for convenience.

5.1(a), but there can be many possible left memory access patterns. As an example, consider the left memory access pattern for the FC junction in Fig. 4.4. Cycle 0 accesses left nodes $\{0, 1, 2, 3\}$, while cycle 1 accesses left nodes $\{4, 5, 6, 7\}$. These can be switched to get a different clash-free left memory access pattern, but the connection pattern would not change.

5.2.1 Types of memory access patterns

We refer to the seed vector technique discussed previously as **Type 1** memory accesses. This is recapitulated in Fig. 5.3(a). This approach only requires storing

the z_i -dimensional seed vector ϕ_i , and uses z_i incrementers to generate subsequent addresses. We can compute G_i^{map} as the number of possible ways of designing ϕ_i , *i.e.* $G_i^{\text{map}} = D_i^{z_i}$. This is the simplest and most hardware friendly method for memory accesses we could come up with. However, if desired, there can be other types of memory accessing which increase G_i^{map} (and therefore $G_{\text{net}}^{\text{map}}$ and the number of possible connection patterns) at the cost of more hardware. These are discussed next.

In **Type 2**, a new seed vector ϕ_i is defined for every sweep. Considering the example in Fig. 5.3(b), $\phi_i = (1, 0, 2, 2)$ for sweep 0, but $(2, 0, 0, 0)$ for sweep 1. Since there are d_i^{out} sweeps, there will be d_i^{out} different ϕ_i vectors for each junction. This results in $G_i^{\text{map}} = D_i^{z_i d_i^{\text{out}}}$. This approach requires storing d_i^{out} different ϕ_i s, and uses z_i incrementers to generate subsequent addresses. Our hardware implementation (described in Section 4.7) used type 2 memory accesses. The seed vectors for all sweeps were pre-generated and hard-coded into FPGA logic.

In **Type 3**, the constraint of cyclically accessing the left memories is also eliminated. Instead, any cycle can access any cell from each of the memories. This means that storing ϕ_i is not enough, the entire sequence of memory accesses needs to be stored as a matrix $\Phi_i \in \{0, 1, \dots, D_i - 1\}^{D_i \times z_i}$. In Fig. 5.3(c) for example, $\Phi_i = ((1, 0, 2, 2), (0, 2, 1, 0), (2, 1, 0, 1))$ for sweep 0. Every sweep would also have a different Φ_i , resulting in $G_i^{\text{map}} = (D_i!)^{z_i d_i^{\text{out}}}$. Thus, this approach requires storing

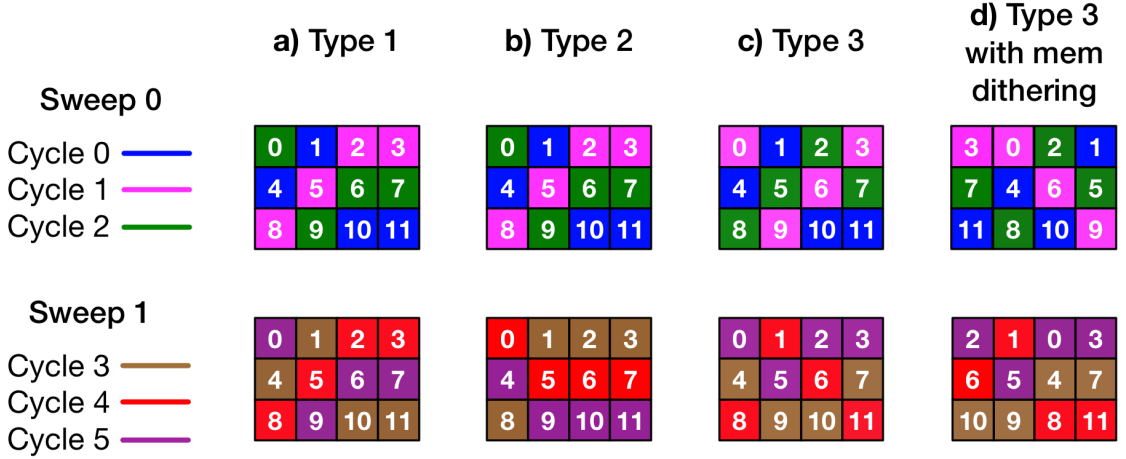


Figure 5.3: (a-c) Various types of clash-free left memory access patterns, and (d) memory dithering for Type 3, using the same left neuronal structure from Fig. 5.2 as an example. The grids represent different access patterns for the same memory bank. The number in each cell represents the left neuron number whose parameter is stored in that cell. Cells sharing the same color are read in the same cycle.

d_i^{out} different Φ_i s, but removes the need for having z_i incrementers to compute subsequent addresses.

The most general kind of clash-freedom resulting in the biggest G_i^{map} arises from ignoring the effect of sweeps. We call this **Type 4**. This means that each address in each left memory of junction i can be accessed in any cycle, as long as it is accessed a total of d_i^{out} times during the junction cycle (consisting of $C = D_i d_i^{\text{out}}$ cycles). However, this may lead to connectivity patterns where a pair of neurons in adjacent layers have more than one connection between them, which is invalid due to redundancy. Hence we skip the details, but state the resulting count of possible NNs: $G_i^{\text{map}} = \left(\frac{C!}{(d_i^{\text{out}}!)^{D_i}} \right)^{z_i}$.

A technique that can be applied to all the types of clash-freedom is **memory dithering**, which is a permutation of the z_i memories (*i.e.* the columns) in a

bank. This permutation can change every sweep, as shown in Fig. 5.3(d). Memory dithering incurs no additional incrementers, but incurs an additional storage cost because of the z_i permutation, which goes up to possibly d_i^{out} different z_i permutations for Types 2 and 3. However, G_i^{map} is increased by a factor K_i , as described next.

- If d_i^{in}/z_i is an integer, an integral number of cycles are required to process each right neuron. Since a cycle accesses all memories, dithering has no effect and $K_i = 1$.
- If z_i/d_i^{in} is an integer greater than 1, the effects of dithering on connectivity patterns are only observed when switching from one right neuron to the next within a cycle. This results in

$$K_i = \left(\frac{z_i!}{d_i^{\text{in}}! \frac{z_i}{d_i^{\text{in}}}} \right)^{d_i^{\text{out}}} \quad (5.4)$$

for Types 2 and 3, and the d_i^{out} exponent is omitted for Type 1 since the access pattern does not change across sweeps.

- When either of z_i or d_i^{in} does not perfectly divide the other, an exact value of K_i is hard to arrive at since some proper or improper fraction of right neurons are processed every cycle. In such cases, K_i is upper-bounded by $(z!)^{d_i^{\text{out}}}$.

Table 5.1: Comparison of Clash-free Left Memory Access Types and associated Hardware Cost for a Single Junction i with $(N_{i-1}, N_i, d_i^{\text{out}}, d_i^{\text{in}}, z_i) = (12, 12, 2, 2, 4)$

Type	Memory Dithering	K_i	G_i^{map}	Cost to Compute Memory Addresses	
				Storage	Incrementers
1	No	N/A	$D_i^{z_i} = 81$	$z_i = 4$	$z_i = 4$
	Yes	6	$K_i D_i^{z_i} = 486$	$2z_i = 8$	$z_i = 4$
2	No	N/A	$D_i^{z_i d_i^{\text{out}}} = 6561$	$z_i d_i^{\text{out}} = 8$	$z_i = 4$
	Yes	36	$K_i D_i^{z_i d_i^{\text{out}}} \approx 236\text{k}$	$2z_i d_i^{\text{out}} = 16$	$z_i = 4$
3	No	N/A	$(D_i!)^{z_i d_i^{\text{out}}} \approx 1.68\text{M}$	$N_{i-1} d_i^{\text{out}} = 24$	None
	Yes	36	$(D_i!)^{z_i d_i^{\text{out}}} \approx 60\text{M}$	$(N_{i-1} + z_i) d_i^{\text{out}} = 32$	None

Table 5.1 compares the count of possible left memory access patterns and associated hardware cost for computing memory addresses for Types 1–3, with and without memory dither. The junction used is the same as in Fig. 5.2, except N_i is raised to 12 such that d_i^{in} becomes 2 and allows us to better show the effects of memory dithering. K_i is calculated as given in (5.4).

5.3 Comparison between classes of Pre-defined Sparsity

At this point, the reader might be wondering about the performance of such clash-free patterns. To study this, we group pre-defined sparsity into three classes in order of least to most restrictive, or equivalently, from most to least number of possible connection patterns, or equivalently, from ‘difficult to realize’ to ‘simple to realize’ on our hardware architecture.

- Random, wherein connections are randomly distributed without regard for fixed in- and out-degrees. Given N_i and N_{i-1} , a member of this class is completely defined by ρ_i .
- Structured, which imposes the constraints of fixed in- and out-degrees. Given N_i and N_{i-1} , a member of this class is completely defined by d_i^{out} .
- Clash-free, which further imposes the constraint of clash-free memory access patterns. Based on hardware friendliness, we can further expand clash-free into subclasses – from type 3 with memory dithering to type 1 with no memory dithering (*i.e.* the opposite order of Table 5.1). Given N_i and N_{i-1} , a member of this class is completely defined by d_i^{out} and z_i .

Table 5.2 compares performance on different datasets for these three classes. For the clash-free class, *we chose the most restrictive Type 1 with no memory dithering*, and experimented with different z_{net} settings to simulate different hardware environments:

- Reuters: One junction cycle is 50 cycles for all the different densities. This is because we scale z_{net} accordingly, corresponding to a more powerful hardware device being used for each NN as ρ_{net} increases.
- CIFAR-100 and MNIST: These simulate cases where hardware choice is limited, such as a high-end, a mid-range and a low-end device being available. Thus three different z_{net} values are used for CIFAR depending on ρ_{net} .

- TIMIT: We keep z_{net} constant for different densities. Junction cycle length varies from 90 cycles for $\rho_{\text{net}} = 7.69\%$ to 810 for $\rho_{\text{net}} = 69.23\%$. This shows that when limited to a single low-end hardware device, denser NNs can be processed in longer time by simply changing z_{net} .

Table 5.2 confirms that **the most restrictive and hardware friendly clash-free pre-defined sparse connection patterns do not lead to any statistically significant performance degradation**. Thus, for most practical purposes, one can opt for maximum hardware simplicity and design for Type 1 clash-freedom with no memory dithering. For certain pathological cases of particularly small junctions, one may opt for the other Types and incorporate memory dithering to incorporate more flexibility in getting a particularly desired connection pattern. Finally, we would like to mention that while clash-free connection patterns are focused on ease of hardware implementation, they can also be good for software implementation – *e.g.* the clash-free class could serve as the basis for pattern selection in software packages.

We also observed that random pre-defined sparsity performs poorly for very low density networks, as shown by the blue values. This is possibly because there is non-negligible probability of neurons getting completely disconnected, leading to irrecoverable loss of information.

Table 5.2: Comparison of Pre-Defined Sparse Classes

$d_{\text{net}}^{\text{out}}$	$\rho_{\text{net}}(\%)$	z_{net}	Test Accuracy Performance (%)		
			Clash-free	Structured	Random
MNIST: $N_{\text{net}} = (800, 100, 100, 100, 10)$, FC test accuracy = $(98 \pm 0.1)\%$					
(80, 80, 80, 10)	80.2	(200, 25, 25, 4)	97.9 ± 0.2	97.9 ± 0.2	97.8 ± 0.2
(60, 60, 60, 10)	60.4	(200, 25, 25, 4)	97.6 ± 0.1	97.8 ± 0.1	97.6 ± 0.2
(40, 40, 40, 10)	40.6	(200, 25, 25, 5)	97.5 ± 0.1	97.7	97.6 ± 0.1
(20, 20, 20, 10)	20.8	(200, 25, 25, 10)	97.2 ± 0.2	97.2 ± 0.1	97.1 ± 0.1
(10, 10, 10, 10)	10.9	(200, 25, 25, 25)	96.7 ± 0.1	96.8 ± 0.2	96.7 ± 0.2
(5, 10, 10, 10)	6.9	(100, 25, 25, 25)	96.3 ± 0.1	96.3 ± 0.1	96.2 ± 0.1
(2, 5, 5, 10)	3.6	(80, 25, 25, 50)	95 ± 0.2	95.1 ± 0.1	95 ± 0.3
(1, 2, 2, 10)	2.2	(80, 20, 20, 100)	93.3 ± 0.3	93.1 ± 0.5	92 ± 0.3
Reuters: $N_{\text{net}} = (2000, 50, 50)$, FC test accuracy = $(89.6 \pm 0.1)\%$					
(25, 25)	50	(1000, 25)	89.4 ± 0.1	89.3	89.4
(10, 10)	20	(400, 10)	87 ± 0.1	86.7 ± 0.1	86.5 ± 0.1
(5, 5)	10	(200, 5)	78.5 ± 0.5	78.2 ± 0.7	77.5 ± 0.6
(2, 2)	4	(80, 2)	53.3 ± 1.8	51.2 ± 1.7	46.8 ± 2.9
(1, 1)	2	(40, 1)	28.4 ± 2.4	28.7 ± 2.3	28 ± 1.9
TIMIT: $N_{\text{net}} = (39, 390, 39)$, FC test accuracy = $(43.2 \pm 0.2)\%$					
(270, 27)	69.2	(13, 13)	43 ± 0.1	43	43 ± 0.1
(180, 18)	46.2		42.7 ± 0.1	42.8 ± 0.1	42.9 ± 0.1
(90, 9)	23.1		42.1 ± 0.1	42.5 ± 0.1	42.4 ± 0.1
(60, 6)	15.4		41.5 ± 0.1	41.8 ± 0.2	41.9 ± 0.1
(30, 3)	7.7		40.5 ± 0.2	40.1 ± 0.2	39.4 ± 0.8
CIFAR-100 ^a : $N_{\text{net}} = (4000, 500, 100)$, FC top-5 test accuracy = $(87.1 \pm 0.6)\%$					
(100, 100)	22	(2000, 250)	87.5 ± 0.2	87.7 ± 0.2	87.4 ± 0.3
(29, 29)	6.4		86.8 ± 0.3	87.2 ± 0.5	87.1 ± 0.2
(12, 12)	2.6	(400, 50)	86.3 ± 0.2	86.5 ± 0.4	86.6 ± 0.4
(5, 5)	1.1		85.3 ± 0.5	85.5 ± 0.5	85.7 ± 0.3
(2, 2)	0.4	(80, 10)	84.1 ± 0.5	84.3 ± 0.3	83.8 ± 0.3
(1, 1)	0.2		83 ± 0.5	83.3 ± 0.4	81.7 ± 0.7

^aFor CIFAR-100, given values of N_{net} , $d_{\text{net}}^{\text{out}}$, z_{net} and ρ_{net} are just for the MLP portion, which follows convolution and pooling layers to form the complete NN, as described in Section 3.3.1. Reported values are top-5 test accuracies obtained from training on the complete NN.

5.4 Comparison to other methods of sparsity

As we have seen, pre-defined sparsity is a method of complexity reduction ‘from the get-go’, *i.e.* prior to training. It can be thought of as a preemptive method since it foresees that most NNs are over-parametrized, hence, most of the weights can be absent from the start without appreciable performance loss. In this sense, **pre-defined sparsity should not be compared on an apples-to-apples basis with methods which train the complete FC NN and then reduce complexity based on training results.** Nevertheless, for the purposes of thoroughness, our previous work [23] compared our method of pre-defined sparsity with the most restrictive Type 1 clash freedom against two other methods of sparsifying NNs.

This comparison and the consequent results were achieved by Kuan-Wen Huang, hence are not presented in this dissertation proposal. We state the primary conclusion – even though clash-free patterns are highly structured and pre-defined, there is no significant performance degradation when compared to advanced methods for producing sparse models which exploit specific properties of the dataset or learn sparse patterns during training. In fact, the performance of our method is comparable to FC NNs at most values of ρ_{net} . The reader is encouraged to refer to Section V of [23] for details.

5.5 Metrics for Connection Patterns

This section discusses possible methods and metrics to judge the ‘goodness’ of a connection pattern before training a NN and computing inference accuracies. These methods can be helpful in filtering out bad connection patterns without incurring the time and computational expense of training. We will use the concept of biadjacency matrices introduced earlier in Section 5.1.

Biadjacency matrices for individual junctions can be multiplied to yield the effective connection pattern for an **equivalent junction** spanning any two junctions x and y , *i.e.* $\mathbf{B}_{x:y}^{N_y \times N_{x-1}} = \prod_{i=y}^x \mathbf{B}_i$, where element $\mathbf{B}_{x:y}^{(j,k)}$ denotes the number of paths from the k th neuron in layer $x - 1$ to the j th neuron in layer y . Thus, these elements are integers ≥ 0 . For the special case where $x = 1$ and $y = L$, we obtain the input-output biadjacency matrix $\mathbf{B}_{1:L}$, *i.e.* the equivalent junction is the whole NN. Each element of $\mathbf{B}_{1:L}$ is the number of paths from a certain input neuron to a certain output neuron.

As an example, consider the NN with $\mathbf{N}_{\text{net}} = (6, 4, 2)$ shown in Fig. 5.4(a). The input-output biadjacency matrix is $\mathbf{B}_{1:2} = \mathbf{B}_2 \mathbf{B}_1$, the equivalent out-degree is $d_{1:2}^{\text{out}} = d_1^{\text{out}} d_2^{\text{out}} = 2$, and equivalent in-degree is $d_{1:2}^{\text{in}} = d_1^{\text{in}} d_2^{\text{in}} = 6$. The equivalent input-output junction 1 : 2 is shown in Fig. 5.4(b). Note how the equivalent biadjacency matrix has double connections between some pairs of neurons. This indicates that there are multiple paths connecting certain inputs to certain outputs.

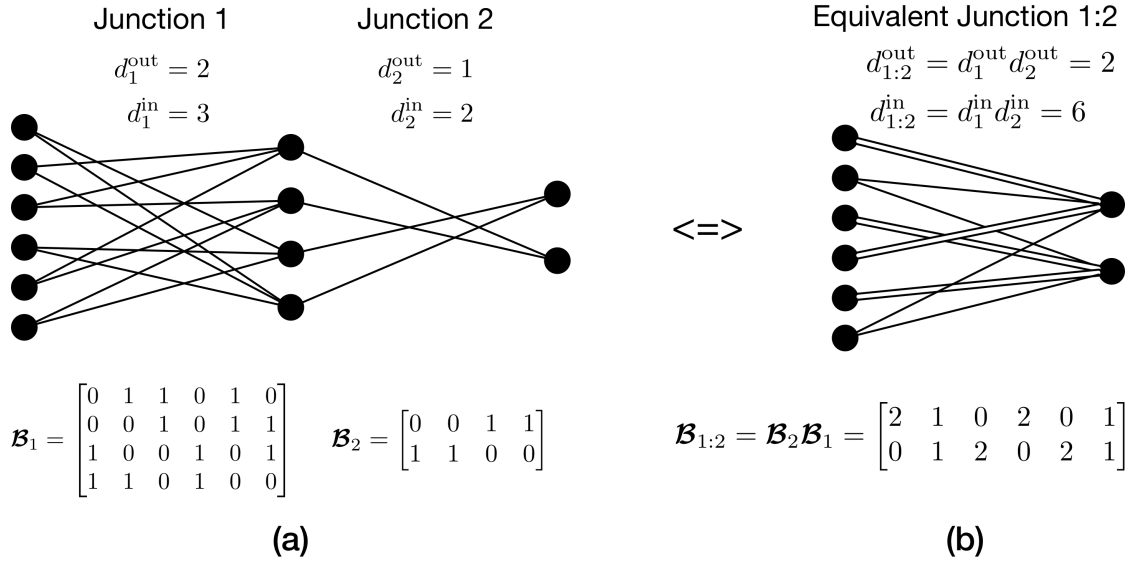


Figure 5.4: (a) Biadjacency matrices for a NN with $N_{\text{net}} = (6, 4, 2)$, $\mathbf{d}_{\text{net}}^{\text{out}} = (2, 1)$, $\mathbf{d}_{\text{net}}^{\text{in}} = (3, 2)$. (b) Equivalent input-output junction and its biadjacency matrix. Double lines indicate double connections between the same pair of neurons.

Note that an equivalent junction is only an abstract concept that aids visualizing how neurons connect across junctions. The reader should not mistake it for an actual junction; such an understanding would be incorrect because it would ignore the effects of the non-linearity in the hidden layer. However, constructing equivalent junctions is helpful in formalizing some metrics for connection patterns, as will be discussed in the remainder of this section.

5.5.1 Window biadjacency matrices

We now attempt to characterize the quality of a pre-defined sparse connection pattern, *i.e.* we try to find the best possible way to connect neurons to optimize

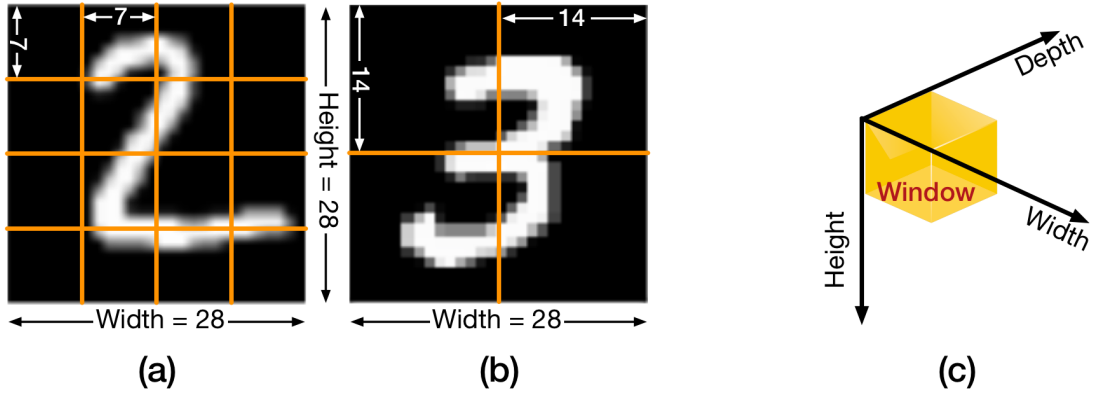


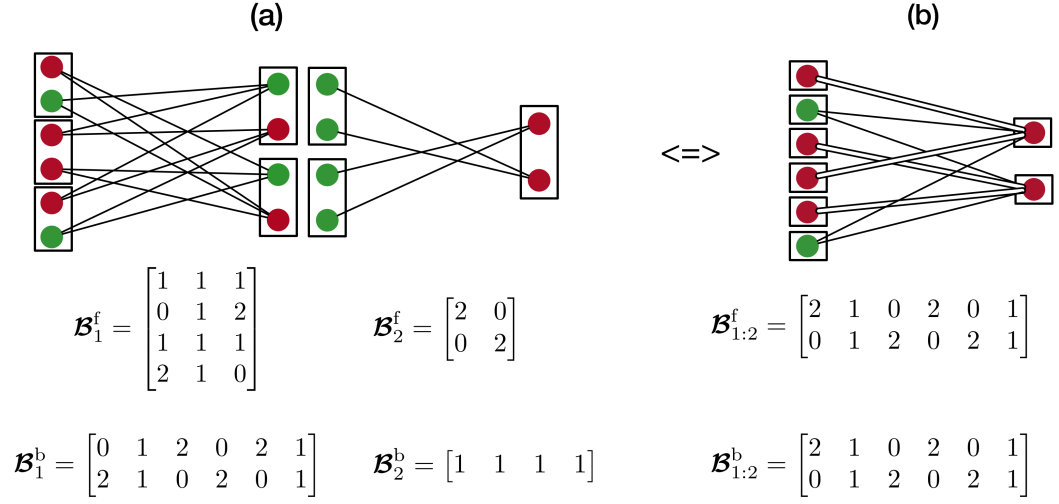
Figure 5.5: (a–b) Images can be divided into 2D windows corresponding to width and height at the output of the input layer of a NN. For example, a 28×28 MNIST image can be divided into (a) 16 windows of size 7×7 each, or (b) 4 windows of size 14×14 each. (c) Example of a 3D window where there is a 3rd dimension for depth, such as color images like CIFAR, or the output of a convolutional layer.

performance. As indicated in Chapter 3, one reason behind the success of pre-defined sparsity is that there exists redundancy / correlated information in typical NN problems. Intuitively, we assume that left neurons of a junction can be grouped into **windows** depending on the **dimensionality** of the left layer output. For example, the input layer for a NN training on MNIST would have 2D windows, each of which might correspond to a fraction of the image, as shown in Fig. 5.5(a–b). When outputs from a NN have an additional depth dimension, such as color channels in CIFAR, each window is a cuboid capturing fractions of width, height and depth, as shown in Fig. 5.5(c).

The idea behind constructing windows is that *neurons in a layer should get information from all portions of the entity its adjacent layer is attempting to represent*. Thus, we will try to maximize the number of left windows to which each

right neuron connects. To realize the importance of this, consider the MNIST output neuron representing digit 2. Let's say the sparse connection pattern is such that when the connections to output 3 are traced back to the input layer, they all come from the top half of the image. This would be undesirable since the top half of an image of a 2 can be mistaken for a 3, as in Fig. 5.5(a-b). A good pre-defined sparse connection pattern will try to avoid such scenarios by spreading the connections to any right neuron across as many input windows as possible. The problem can also be mirrored so that every left neuron connects to as many different right windows as possible. This ensures that local information from left neurons is spread to different parts of the right layer. The grouping of right windows will depend on the dimensionality of the input to the right layer.

The window size is chosen to be the minimum possible such that the ideal number of connections from or to it remains integral. In order to achieve the minimum window size, we let the number of left windows be d_i^{in} and the number of right windows be d_i^{out} . So the number of neurons in each left and right window is N_{i-1}/d_i^{in} and N_i/d_i^{out} , respectively. The example from Fig. 5.4 is reproduced in Fig. 5.6. For example, since $d_1^{\text{in}} = 3$, the inputs must be grouped into 3 windows of size 2 so that ideally one connection from each reaches every hidden neuron. The 1st and 3rd hidden neurons achieve this and are colored green, while the other two do not and are colored red.



$$\mathcal{S}^{\text{net}} = (S_1^f = 0.83, S_1^b = 0.67, \quad S_2^f = \mathbf{0.5}, S_2^b = 1, \quad S_{1:2}^f = 0.67, S_{1:2}^b = 0.67)$$

Figure 5.6: Window biadjacency matrices and scatter (to be discussed in Section 5.5.2) for the NN in Fig. 5.4. Green neurons indicate ideal connectivity. The hidden layer is split into two parts to show separate constructions of \mathcal{B}_1^f and \mathcal{B}_2^b . (a) shows the complete NN, while (b) shows the equivalent input-output junction 1 : 2. The final scatter value S is bolded within the scatter vector \mathcal{S}^{net} .

Then we construct the forward window biadjacency matrix \mathcal{B}_i^f and backward window biadjacency matrix \mathcal{B}_i^b by summing up entries of \mathcal{B}_i as shown in Fig. 5.7. These window biadjacency matrices describe connectivity between windows in a layer and neurons on the opposite side of the junction. Forward indicates windows on the left to neurons on the right, and vice-versa for backward. Ideally, *every window biadjacency matrix for a single junction should be the all 1s matrix, which signifies exactly 1 connection from every window to every neuron on the opposite side.*

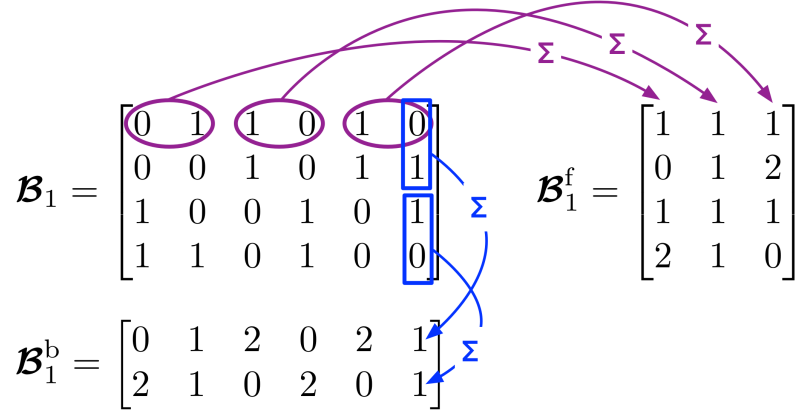


Figure 5.7: Constructing window biadjacency matrices for junction 1 of the NN in Fig. 5.4(a). The forward window biadjacency matrix \mathcal{B}_1^f is constructed by summing (Σ) across $d_1^{\text{in}} = 3$ windows of each row of \mathcal{B}_1 , each window is of size $N_0/d_1^{\text{in}} = 2$. The backward window biadjacency matrix \mathcal{B}_1^b is constructed by summing (Σ) down $d_1^{\text{out}} = 2$ windows of each column of \mathcal{B}_1 , each window is of size $N_1/d_1^{\text{out}} = 2$.

Note that these matrices can also be constructed for equivalent junctions, *i.e.* $\mathcal{B}_{x:y}^f$ and $\mathcal{B}_{x:y}^b$, by multiplying matrices for individual junctions. It could happen that the effective $d_{x:y}^{\text{in}}$ for layer y becomes more than N_{x-1} , or the effective $d_{x:y}^{\text{out}}$ for layer $x-1$ becomes more than N_y . In such cases, the number of neurons in each window is rounded up to 1. This scenario signifies that there are multiple paths connecting neurons in the equivalent junction, *i.e.* $\mathcal{B}_{x:y}^f$ and/or $\mathcal{B}_{x:y}^b$ have entries greater than 1. In such cases, a good connection pattern should try to distribute the connections as evenly as possible. as quantified in Section 5.5.2.

5.5.2 Scatter

We propose **scatter** S as a proxy for the goodness of a pre-defined sparse MLP connection pattern. In other words, it can be useful in predicting NN performance

before training. To compute scatter, we count the number of entries greater than or equal to 1 in the window biadjacency matrix. If a particular window gets more than its fair share of connections to a neuron on the opposite side, then it is depriving some other window from getting its fair share. This should not be encouraged, so we treat entries greater than 1 the same as 1. Scatter is the average of the count, *i.e.* for junction i :

$$S_i^f = \frac{1}{N_i d_i^{\text{in}}} \sum_{j=1}^{N_i} \sum_{k=1}^{d_i^{\text{in}}} \mathbb{I} \left(\mathcal{B}_i^{f(j,k)} \geq 1 \right) \quad (5.5a)$$

$$S_i^b = \frac{1}{d_i^{\text{out}} N_{i-1}} \sum_{j=1}^{d_i^{\text{out}}} \sum_{k=1}^{N_{i-1}} \mathbb{I} \left(\mathcal{B}_i^{b(j,k)} \geq 1 \right) \quad (5.5b)$$

The subscript and superscript notation for scatter is the same as the biadjacency matrices, *i.e.* f and b denote forward (left windows to right neurons) and backward (right neurons to left windows), respectively. As an example, the NN in Fig. 5.6(a) has $S_1^f = 10/12 = 0.83$ since 10 out of the 12 elements in \mathcal{B}_1^f are ≥ 1 . The other scatter values can be computed similarly to form the scatter vector $\mathbf{S}^{\text{net}} = (S_1^f, S_1^b, S_2^f, S_2^b, S_{1:2}^f, S_{1:2}^b)$. Notice that \mathbf{S}^{net} will be all 1s for FC NNs, which is the ideal case. Incorporating sparsity leads to reduced \mathbf{S}^{net} values. Finally, we define:

► **Definition 7: Scatter:** Scatter is a measure of how well-distributed the connections from neurons in any layer of a NN are to its adjacent layers. It is measured as the minimum element in the scatter vector, *i.e.* :

$$S = \min \mathbf{S}^{\text{net}} \quad (5.6)$$

In Fig. 5.6 for example, $\mathbf{S}^{\text{net}} = (0.83, 0.67, 0.5, 1, 0.67, 0.67)$, so $S = 0.5$ (shown in bold). Our experiments indicate that any low value in \mathbf{S}^{net} leads to bad performance, so we picked the critical minimum value. These experiments are described next.

Experimental Results on Scatter

We ran experiments to evaluate scatter using structured pre-defined sparse (non-clash-free) NNs training on three datasets: a) Morse code symbol recognition, b) MNIST, and c) CIFAR-10 images. (a) will be described in Chapter 6, while (c) has the same image format as CIFAR-100, but 10 classes instead of 100 (see Section 3.3.1 for more details on MNIST and CIFAR). The reader is encouraged to refer to our previous work [24] for complete details on the network and training configuration. We skip those details here and focus on the results instead. These are shown in Fig. 5.8, which plots classification performance on validation data vs. scatter S . Subfigures (a–b) show that performance gets better with increasing scatter. Subfigure (c) does not show the trend clearly, we hypothesize that this is

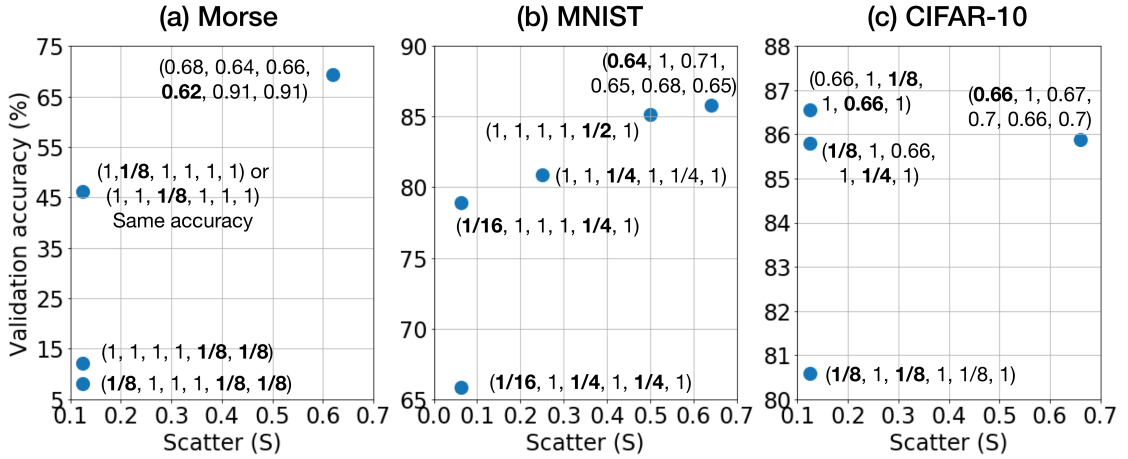


Figure 5.8: Performance vs. scatter for structured pre-defined sparse (non-clash-free) NNs training on (a) Morse (b) MNIST, and (c) CIFAR-10, part of which is convolutional. All minimum values that need to be considered to differentiate between connection patterns are bolded.

due to using convolution layers in the NN for CIFAR-10, leading to reduced impact of the MLP portion.

\mathbf{S}^{net} is shown alongside each point. When S is equal for different connection patterns, the next minimum value in \mathbf{S}^{net} needs to be considered to differentiate the networks, and so on. Considering the Morse results, the leftmost three points all have $S = 1/8$, but the number of occurrences of $1/8$ in \mathbf{S}^{net} is 3 for the lowest point (8% accuracy), 2 for the second lowest (12% accuracy) and 1 for the highest point (46% accuracy). For the MNIST results, both the leftmost points have a single minimum value of $1/16$ in \mathbf{S}^{net} , but the lower has two occurrences of $1/4$ while the upper has one.

We draw several insights from these results. Firstly, although we defined S as a single value for convenience, there may arise cases when other (non-minimum)

elements in \mathbf{S}^{net} are important. Secondly, our results indicate that *scatter is a sufficient metric for performance, not necessary*. In other words, a network with a high S value should perform well, but a network with a slightly lower S value than another cannot be conclusively dismissed as being worse. But if a network has multiple low values in \mathbf{S}^{net} , it should be rejected. Finally, carefully choosing which neurons to group in a window will increase the predictive power of scatter. A priori knowledge of the dataset will lead to better window choices.

The rightmost point in each subfigure of Fig. 5.8 corresponds to structured pre-defined sparsity as we have described so far, *i.e.* the only constraints are fixed in- and out-degree. To be more specific, the connections can arrange themselves in any way they want as long as these constraints are obeyed. On the other hand, all other points were generated by specifically planning individual connections, *i.e.* designing the complete biadjacency matrices manually. These perform poorly because when we designed a particular junction to have high values in \mathbf{S}^{net} , it invariably led to low values for another junction, leading to a low value for the final S . This explains why the unplanned patterns perform the best.

In summary, we would like to state that our experiments on scatter are not as exhaustive as those performed to discover trends and guidelines in pre-defined sparsity, as detailed in Section 3.3. Hence, the results regarding scatter require

additional study before being definitive, as we plan to do in our expected contributions. Also note that the trends discussed in Section 3.3 also help in indicating how well a pre-defined sparse NN can be expected to perform.

Chapter 6

Dataset Engineering

Thus far, this document has discussed methods and architectures for complexity reduction of neural networks primarily aimed towards classification problems. An underlying assumption was that there is a large amount of high quality labeled data available to train and benchmark the performance of different NNs. The effects of dataset size on network performance has been explored in [88], in particular, more data is beneficial in reducing overfitting and improving robustness and generalization capabilities of NNs [38, 89]. However, it is often a challenge to obtain adequate amounts of quality labeled data, and several entities [90–92] have identified this as a bottleneck to pushing the frontiers of NN performance.

A possible solution is to obtain data by synthetic instead of natural means. **Synthetic data** is generated using computer algorithms instead of being collected from real-world scenarios. The advantages are that a) computer algorithms can be tuned to mimic real-world settings to desired levels of accuracy, and b) a theoretically unlimited amount of data can be generated by running the algorithm long enough. Synthetic data has been successfully used in problems such as 3D

imaging [93], point tracking [94], breaking Captchas on popular websites [95], and augmenting real world datasets [96].

We conducted some investigations on synthetic data with the underlying motive of creating **several low-redundancy datasets**, *i.e.* which would challenge a pre-defined sparse NN more than datasets with considerable redundancy such as discussed in Section 3.3.2. We came up with a family of synthetic datasets on **classifying Morse codewords**. Morse code is a system of communication where each letter, number or symbol in a language is represented using a sequence of dots and dashes, separated by spaces. It is widely used to communicate in situations where voice is not possible, such as helping people with disabilities talk [97–99], or where message transmission needs to be achieved using only 2 states [100], or in rehabilitation and education [101]. Morse code is a useful skill to learn and there exist cellphone apps designed to train people in its usage [102, 103].

Our classification problem groups Morse codewords into 64 character classes comprising letters, numbers and symbols. This is different from previous works with just two classes for dots and dashes [97, 98, 100], and other approaches using fuzzy logic [100], time series to decode English words in Morse code [104, 105], or radial basis functions [106]. We also investigated **metrics to characterize the difficulty of a dataset**. In summary, we refer to our efforts discussed in this chapter as *dataset engineering*. We published these efforts in [28], which won a

‘Conference Best Paper’ award. The codebase is **available on Github** [29]. The key contributions are as follows:

Contributions in dataset engineering

1. To the best of our knowledge, we are the first to develop an algorithm for generating machine learning datasets of varying difficulty – measured as test set classification accuracy of a NN training on it. Harder datasets lead to lower accuracy, and vice-versa. Encountering harder datasets leads to aggressive exploration of hyperparameters and learning algorithms, which ultimately increases the robustness of NNs training on them. Some particularly hard datasets are useful for testing the limits of our sparse NNs.
2. We introduce metrics to quantify the difficulty of a dataset before having a NN train on them. While some of these arise from information theory, we also come up with a new metric which achieves a high correlation coefficient with the eventual accuracy obtained after training and inference.
3. This work is one of few to introduce a spatially 1-dimensional dataset. This is in contrast to the wide array of image and character recognition datasets which are usually 2D such as MNIST, where each image has width and height, or 3D such as CIFAR, where each image has width,

height and depth (color channels). The dimensionality of the inputs is important when benchmarking pre-defined sparse connection patterns, since some of the metrics discussed in Chapter 5 depend on dimension.

6.1 Generating Algorithm

We picked **64 class labels** for our dataset – the 26 English letters, the 10 Arabic numerals, and 28 other symbols such as (, +, :, etc. Each of these is represented by a sequence of dots and dashes in Morse code, for example, + is represented as • — • — •. So as to mimic a real-world scenario in our algorithm, we imagined a human or a Morse code machine writing out this sequence within a frame of fixed size. Wherever the pen or electronic instrument touches is darkened and has a high intensity, indicating the presence of dots and dashes, while the other parts are left blank, *i.e.* spaces.

Step 1 – Frame Partitioning

For our algorithm, each Morse codeword lies in a frame which is a vector of 64 values, *i.e.* **64 features** per input. Within the frame, the length of a sequence having consecutive similar values is used to differentiate between a dot and a dash. In the baseline dataset, a dot can be 1-3 values wide and a dash 4-9. This is in accordance with international Morse code regulations [107] where the size or

duration of a dash is around three times that of a dot. The space between a dot and a dash can have a length of 1-3 values. The exact length of a dot, dash or space is chosen from these ranges according to a uniform probability distribution. This is to mimic the human writer who is not expected to make each symbol have a consistent length, but can be expected to make dots and spaces around the same size, and dashes longer than them. The baseline dataset has no leading spaces before the 1st dot or dash, *i.e.* the codeword starts from the left edge of the frame. There are trailing spaces to fill up the right side of the frame after all the dots and dashes are complete.

Step 2 – Assigning Values for Intensity Levels

All values in the frame are initially real numbers in the range $[0, 16]$ and indicate the intensity of that point in the frame. For dots and dashes, the values are drawn from a normal distribution with mean $\mu = 12$ and standard deviation $\sigma = 4/3$. The idea is to have the ‘six-sigma’ range from $(12 - 3 \times 4/3) = 8$ to $(12 + 3 \times 4/3) = 16$. This ensures that any value making up a dot or a dash will lie in the upper half of possible values, *i.e.* in the range $[8, 16]$. The value of a space is exactly 0. Once again, these conditions mimic the human or machine writer who is not expected to have consistent intensity for every dot and dash, but can be expected to not let the writing instrument touch portions of the frame which are spaces.

Step 3 – Noising

Noise in input samples is often deliberately injected as a means of avoiding overfitting in NNs [38], and has been shown to be superior to other methods of avoiding overfitting [108]. This was, however, the secondary reason behind our experimenting with noise. The primary reason was to deliberately make the data hard to classify and test the limits of different NNs processing it. Noise can be thought of as a human accidentally varying the intensity of writing the Morse codeword, or a Morse communication channel having noise. The baseline dataset has no noise, while others have additive noise from a mean-zero normal distribution applied to them. Fig. 6.1 shows the 3 steps up to this point. Finally, all the values are normalized to lie within the range $[0, 1]$ with precision of 3 decimal places.

Step 4 – Mass Generation

Steps 1-3 describe the generation of one input sample corresponding to some particular class label. **This can be repeated as many times as required for each of the 64 class labels.** This demonstrates a key advantage of synthetic over real-world data – the ability to generate an arbitrary amount of data having an arbitrary prior probability distribution over its classes. The baseline dataset has 7000 examples for each class, for a total of 448,000 examples.

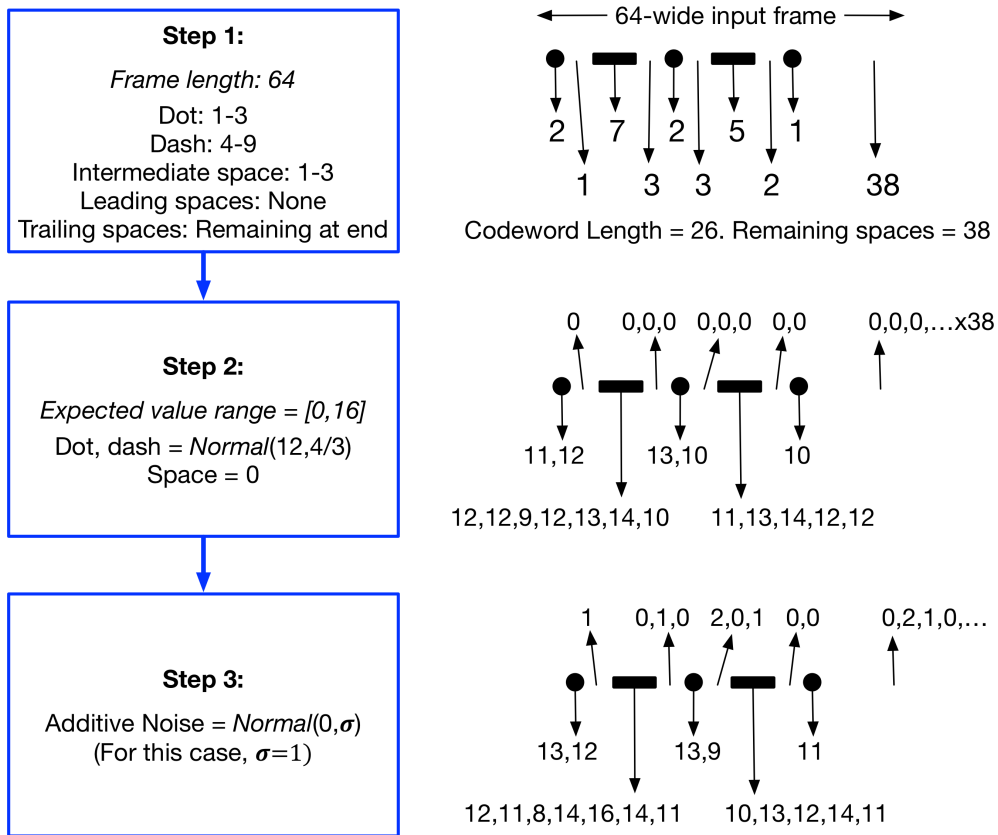


Figure 6.1: Generating the Morse codeword $\bullet - \bullet - \bullet$ corresponding to the $+$ symbol. The first 3 steps, prior to normalizing, are shown. Only integer values are shown for convenience, however, the values can and generally will be fractional. $Normal(\mu, \sigma)$ denotes a normal distribution with mean = μ , standard deviation = σ . For this figure, $\sigma = 1$.

6.1.1 Variations and Difficulty Scaling

The baseline dataset is as described so far, except that $\sigma = 0$, *i.e.* it has no additive noise. We experimented with the following variations in datasets:

1. Baseline with **additive noise** = $Normal(0, \sigma)$, $\sigma \in \{0, 1, 2, 3, 4\}$. These are called *Morse 1.σ*, *i.e.* 1.0 to 1.4, where 1.0 is the baseline.

2. Instead of having the codeword start from the left edge of the frame, we introduced a random number of **leading spaces**. For example, in Fig. 6.1, the codeword occupies a length of 26 values. The remaining 38 space values can be randomly divided between leading and trailing spaces. This increases the difficulty of the dataset since no particular set of neurons are expected to be learning dots and dashes as the actual codeword could be anywhere in the frame. Just like variation 1, we added noise and call these datasets *Morse 2. σ* , $\sigma \in \{0, 1, 2, 3, 4\}$.
3. There is no overlap between the lengths of dots and dashes in the datasets described so far. The difficulty can be increased by making **dash length = 3-9 values**, which is exactly according to the convention of having dash length thrice of dot length. This means that dashes can masquerade as dots and spaces, and vice-versa. This is done on top of introducing leading spaces. These datasets are called *Morse 3. σ* , σ being as before.
4. The Morse datasets only have 64 input features, which is quite small compared to some others such as MNIST (784 features) or CIFAR (3072 features). To decrease dataset difficulty, we introduced **dilation** by a factor of 4. This is done by scaling all lengths in variation 3 by a factor of 4, *i.e.* frame length (number of features) becomes 256, dot sizes and space sizes are 4-12, and dash size is 12-36. These datasets are called *Morse 4. σ* , σ being as before.

5. Increasing the number of training examples, *i.e.* the **size of the dataset**, makes it easier to classify since a NN has more labeled training examples to learn from. Accordingly we chose *Morse 3.1* and scaled the number of examples to obtain *Morse Size* x , $x \in \{1/8, 1/4, 1/2, 2, 4, 8\}$. For example, *Morse Size* $1/2$ has 3500 examples for each class, for a total of 224,000 examples.

6.2 Neural Network Results and Analysis

6.2.1 Results

We initially designed a FC MLP NN to benchmark the different variations of the Morse datasets. N_0 always matches the frame length, *i.e.* 256 for the *Morse 4.σ* datasets and 64 for all others, while $N_L = 64$ to match the number of classes. We used a single hidden layer with 1024 neurons, *i.e.* $\mathbf{N}_{\text{net}} = (256, 1024, 64)$ for *Morse 4.σ* and $\mathbf{N}_{\text{net}} = (64, 1024, 64)$ for all other datasets. The performance, *i.e.* accuracy, generally increases on adding more hidden neurons, however, we stuck with 1024 since values above that yielded diminishing returns. The hidden layer has ReLU activations, while the output is a softmax probability distribution. We used the Adam optimizer with default parameters, He normal initialization for the weights, and trained for 30 epochs using a minibatch size $M = 128$. We used $5/7$ th of the total examples for training, and $1/7$ th each for validation and testing.

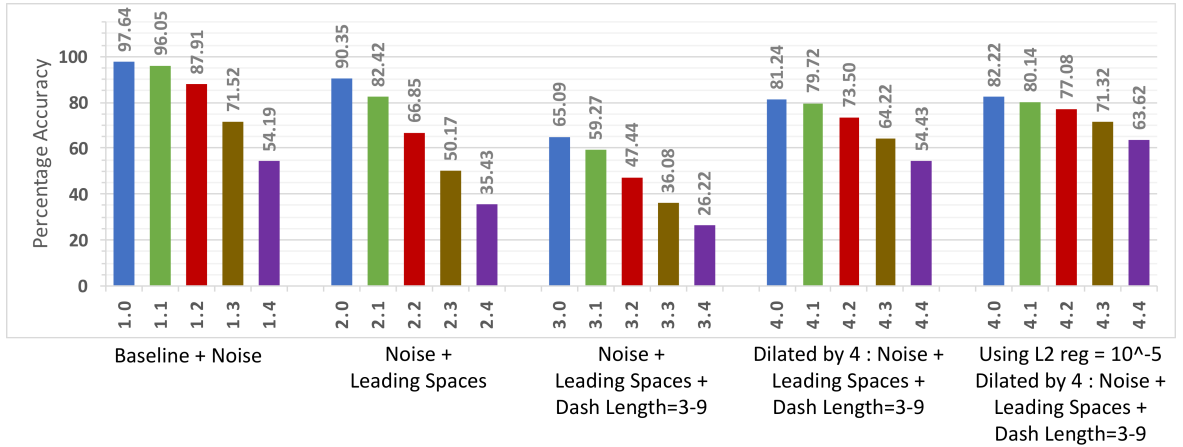


Figure 6.2: Percentage classification accuracies obtained by the FC NN (described in Section 6.2.1) on the test set of different Morse datasets. The rightmost set of bars corresponds to *Morse* $4.\sigma$ with L2 regularization with $\lambda = 10^{-5}$.

All reported accuracies are those obtained on the test samples unless otherwise mentioned.

Validation results suggested that no regularization was required when training on *Morse* $1.\sigma$, $2.\sigma$ and $3.\sigma$. However, the NN for *Morse* $4.\sigma$ is more prone to overfitting due to having more input neurons, leading to more weight parameters. Accordingly we applied L2 regularization with $\lambda = 10^{-5}$, which was the best value as determined from validation.

Test accuracy results after training the NN on the different Morse datasets are shown in Fig. 6.2. As expected, increasing the standard deviation of noise results in drop in performance. This effect is not felt strongly when $\sigma = 1$ since the 3σ range can take spaces to a value of 3 (on a scale of $[0, 16]$, *i.e.* before normalizing to $[0, 1]$), while dots and dashes can drop to $8 - 3 = 5$, so the probability of a space

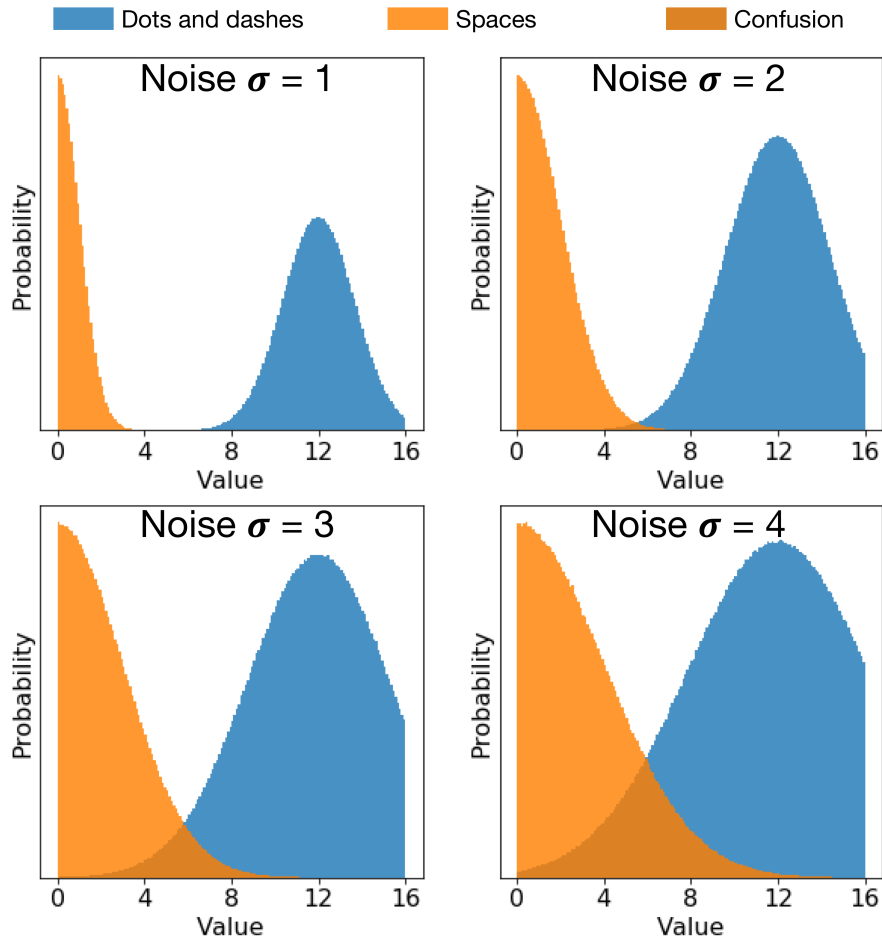


Figure 6.3: Effects of noise leading to spaces (orange) getting confused (brown) with dots and dashes (blue). Higher values of noise σ lead to increased probability of the brown region, making it harder for the NN to discern between ‘dots and dashes’ and spaces. The x-axis in each plot shows values in the range $[0, 16]$, *i.e.* before normalizing to $[0, 1]$.

being confused with a dot or dash is basically 0. Confusion can occur for $\sigma \geq 2$, and gets worse for higher values, as shown in Fig. 6.3.

Since the codeword lengths do not often stretch beyond 32, the first half of neurons usually encounter high input intensity values corresponding to dots and dashes during training. This means that the latter half of neurons mostly encounter

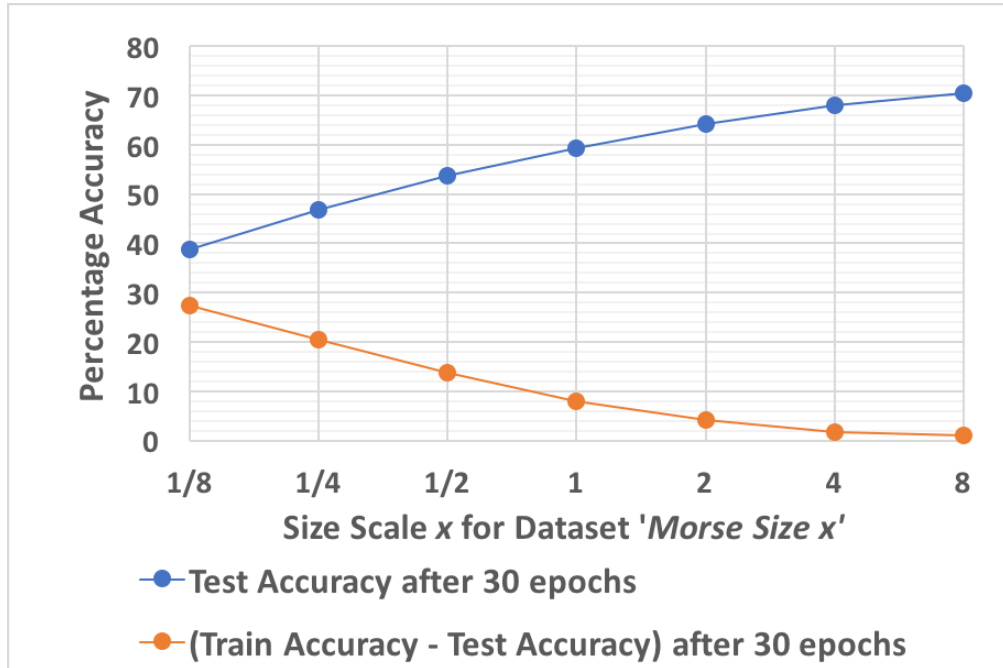


Figure 6.4: Effects of increasing the size of *Morse 3.1* by a factor of x on test accuracy after 30 epochs (blue), and (Training Accuracy – Test Accuracy) after 30 epochs (orange).

lower input values corresponding to spaces. This aspect changes when introducing leading spaces, which become inputs to some neurons in the first half. The result is an increase in the variance of the input to each neuron. As a result, accuracy drops. The degradation is worse when dashes can have a length of 3-9. Since the lengths are drawn from a uniform distribution, $1/7$ th of dashes can now be confused with $1/3$ rd of dots and $1/3$ rd of intermediate spaces. As an example, for the + codeword which has 2 dashes, 3 dots and 4 intermediate spaces, there is a $(2/9 \times 1/7 + 3/9 \times 1/3 + 4/9 \times 1/3) = 29\%$ chance of this confusion occurring. Dilating by 4, however, reduces this chance to $(2/9 \times 1/25 + 3/9 \times 1/9 + 4/9 \times 1/9) = 9.5\%$. Accuracy is better as a result.

Increasing dataset size has a beneficial effect on performance. Giving the NN more examples to train from is akin to training on a smaller dataset for more epochs, with the important added advantage that overfitting is reduced. This is shown in Fig. 6.4, which shows improving test accuracy as the dataset is made larger. At the same time, the difference between final training accuracy and test accuracy reduces, which implies that the network is generalizing better and not overfitting. Note that *Morse Size 8* has 3 million labeled training examples – a beneficial consequence of being able to cheaply generate large quantities of synthetic data.

6.2.2 Results for Pre-Defined Sparse Networks

Fig. 6.5 shows classification performance for 4 different Morse datasets when using our method of structured pre-defined sparse NNs. Both junctions are sparsified equally. Note how the baseline dataset only has mild performance degradation even when ρ_{net} is reduced to $1/4$, while performance drops off much more rapidly when dataset variations are introduced. These variations lead to increased information content per neuron, *i.e.* reduced redundancy. This enforces the trend observed in Section 3.3.2 of reduced redundancy datasets being less robust to sparsification.

Also note that as density is reduced, *Morse 4.2* has the best performance out of the non-baseline models tested in Fig. 6.5. This is because it has more weights to begin with, due to the increased number of input neurons. For example, the

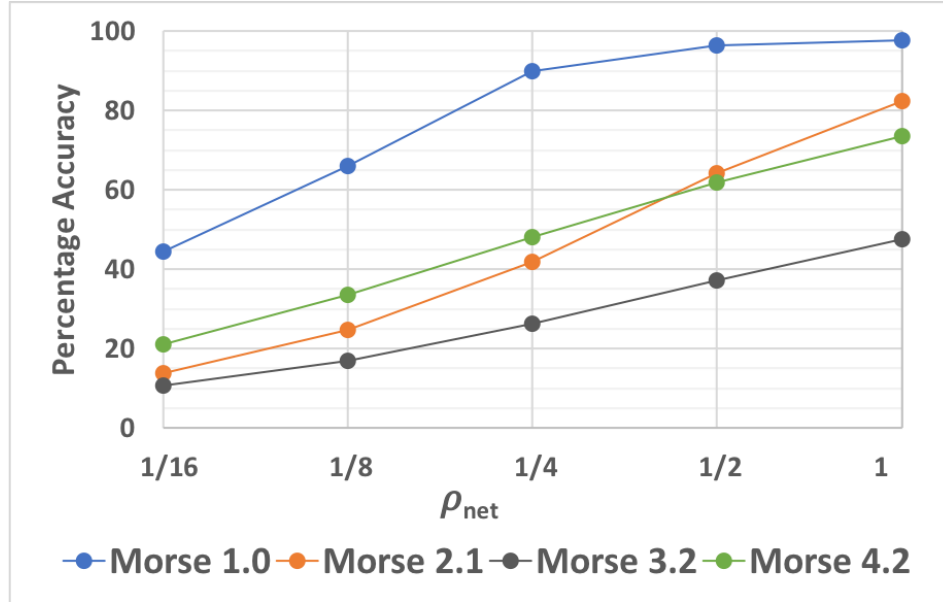


Figure 6.5: Effects of imposing pre-defined sparsity on classification performance for different Morse datasets.

performance at $\rho_{\text{net}} = 1/4$ for *Morse 4.2* is better than FC for *Morse 2.1*. This enforces the trend of ‘large and sparse’ NNs performing better than ‘small and dense’ NNs, as discussed in Section 3.3.4. However, there is a subtle difference, while Section 3.3.4 experimented on varying numbers of hidden neurons, here it is the number of input neurons which varies.

We also experimented with the trend of latter junctions requiring a higher density than former, as discussed in Section 3.3.3, on the dataset *Morse 1.0*. Similar to the case of $\mathbf{N}_{\text{net}} = (39, 390, 39)$ for TIMIT, the NN used for Morse code has symmetric junctions since $\mathbf{N}_{\text{net}} = (64, 1024, 64)$. We set up the experiment somewhat differently than TIMIT. We swept over several possible (ρ_1, ρ_2) value pairs such that they all led to the same ρ_{net} value. We did this for two different ρ_{net}

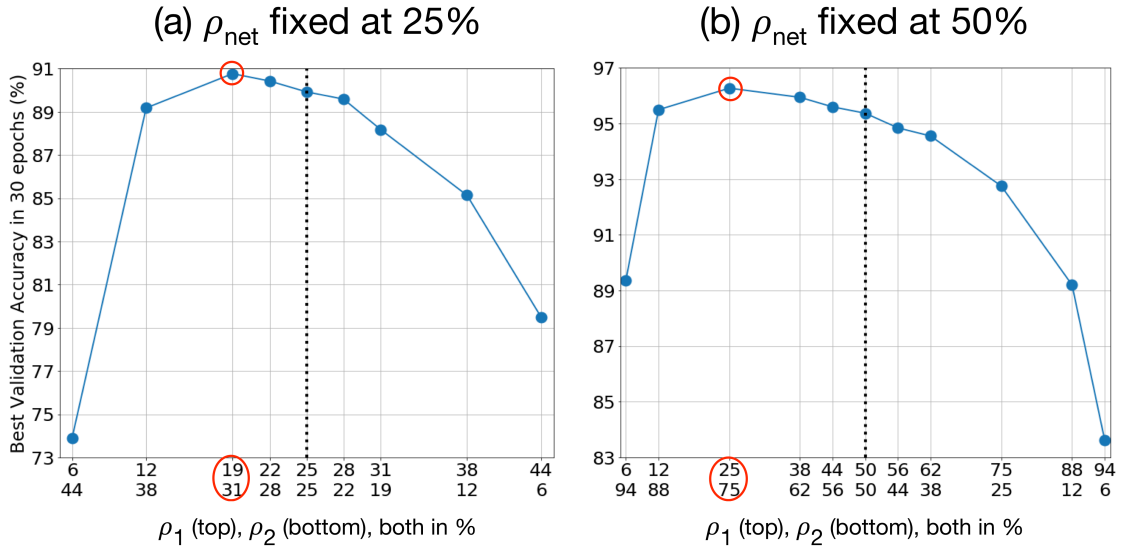


Figure 6.6: Validation performance results for varying ρ_1 (x-axis top) and ρ_2 (x-axis bottom) individually so as to keep ρ_{net} fixed at (a) 25%, (b) 50%. The black dashed line is when $\rho_1 = \rho_2$, while the red circles indicate peak performance. The NN has $N_{\text{net}} = (64, 1024, 64)$ and dataset used is *Morse 1.0*.

values – 25% in Fig. 6.6(a) and 50% in Fig. 6.6(b), and plot the peak validation accuracy obtained in 30 epochs. The black dashed line is for the case $\rho_1 = \rho_2$, and should be the point where the performance peaks provided both junctions have equal importance. The fact that performance peaks *to the left* of the black dashed line for both figures (red circles) indicates that a higher value for ρ_2 is beneficial, as indicated by the tend in Section 3.3.3. For example when $\rho_{\text{net}} = 50\%$, performance peaks at $(\rho_1, \rho_2) = (25\%, 75\%)$. Also note that points to the extreme left and right go lower than the critical density for junctions 1 and 2, respectively.

6.3 Metrics for Dataset Difficulty

This section discusses possible metrics for quantifying how difficult a dataset is to classify. Each sample in a dataset is a point in an N_0 -dimensional space. For the Morse datasets (not considering dilation), $N_0 = 64$. There are N_L classes of points, which is also 64 in this case. The classification problem is essentially finding the class of any new point. Any machine learning classifier will attempt to learn and construct **decision boundaries** between the classes by partitioning the whole space into N_L regions. The samples of a particular class m are clustered in the m th region. Suppose a particular input sample actually belongs to class m . The classifier commits an error if it ranks some class j , $j \neq m$, higher than m when deciding where that input sample belongs. The probability of this happening is $P_{PW}(j|m)$, where subscript PW stands for pairwise and indicates that the quantity is specific to classes j and m . The overall probability of error $P(E)$ would also depend on the prior probability $P(m)$ of the m th class occurring. Considering all classes in the dataset, $P(E)$ is given as:

$$\sum_{m=1}^{N_L} P(m) \left[\max_{\substack{j \in \{1, 2, \dots, N_L\} \\ j \neq m}} P_{PW}(j|m) \right] \leq P(E) \leq \sum_{m=1}^{N_L} P(m) \sum_{\substack{j=1 \\ j \neq m}}^{N_L} P_{PW}(j|m) \quad (6.1)$$

This is a standard result in decision theory, cf. [109].

The pairwise probabilities can be approximately computed by assuming that the locations of samples of a particular class m are from a Gaussian distribution

with mean located at the centroid c_m , which is the average of all samples for the class. To simplify the math, we take the average variance across all N_0 dimensions within a class – this gives the variance σ_m^2 for class m . The distance between 2 classes m and j is the L2-norm of the displacement vector between their centroids, *i.e.* $d(m, j) = \|c_m - c_j\|_2$. A particular class will be more prone to errors if it is close to other classes. This can be quantified by looking at $\frac{d_{\min}(m)}{\sigma_m}$, where the numerator is given as:

$$d_{\min}(m) = \min_{\substack{j \in \{1, 2, \dots, N_L\} \\ j \neq m}} d(m, j) \quad (6.2)$$

With the Gaussian assumption, (6.1) simplifies to $V_{\text{lower}} \leq P(E) \leq V_{\text{upper}}$, where:

$$V_{\text{lower}} = \sum_{m=1}^{N_L} P(m) Q \left(\sqrt{\frac{d_{\min}(m)^2}{4\sigma_m^2}} \right) \quad (6.3a)$$

$$V_{\text{upper}} = \sum_{m=1}^{N_L} P(m) \sum_{\substack{j=1 \\ j \neq m}}^{N_L} Q \left(\sqrt{\frac{d(m, j)^2}{4\sigma_m^2}} \right) \quad (6.3b)$$

where $Q(\cdot)$ is the tail function for a standard Gaussian distribution. See [109] for a further exposition of these concepts.

The **lower and upper bounds for error** V_{lower} and V_{upper} can thus be used as metrics for dataset difficulty, since higher values for them imply higher probabilities of error, *i.e.* lower accuracy. A simpler metric can be obtained by just considering

$\frac{\sigma_m}{d_{\min}(m)}$. Higher values for this indicate that a) class m is close to some other class and the NN will have a hard time differentiating between them, and b) Variance of class m is high, so it is harder to form a decision boundary to separate inputs having labels m from those with other labels. Since $\frac{\sigma_m}{d_{\min}(m)}$ is different for every class, we experimented with ways to reduce it to a single measure such as taking the minimum, the average and the median. The average worked best, which gives our 3rd metric, which we call the **distance metric** V_{dist} :

$$V_{\text{dist}} = \frac{\sum_{m=1}^{N_L} \frac{\sigma_m}{d_{\min}(m)}}{N_L} \quad (6.4)$$

Therefore, high values of V_{dist} lead to low accuracy.

The 4th and final metric is the **threshold metric** V_{thresh} , which we came up with. To obtain this, we first compute the class centroids just as before. Then we compute the L1-norm between every pair of centroids and average over N_0 , *i.e.* :

$$d_1(m, j) = \frac{\|c_m - c_j\|_1}{N_0} \quad (6.5)$$

Since all N_0 features in each input sample are normalized to $[0, 1]$, all the elements in all the centroid vectors also lie in the range $[0, 1]$. So the d_1 number for every pair of classes is always between 0 and 1, in fact, it is proportional to the absolute

distance between the 2 classes. Then, we simply count how many of the d_1 numbers are less than a threshold, which we empirically set to 0.05. This gives V_{thresh} , *i.e.* :

$$V_{\text{thresh}} = \sum_{m=1}^{N_L} \sum_{\substack{j=1 \\ j \neq m}}^{N_L} \mathbb{I}(d_1(m, j) < 0.05) \quad (6.6)$$

The higher the value of V_{thresh} , the lower the accuracy. Note that the total number of d_1 values will be $\binom{N_L}{2}$, so the count for V_{thresh} will typically be higher for datasets that have more classes. This is a desired property since more number of classes usually makes a dataset harder to classify. Note that the maximum value of V_{thresh} for the Morse datasets is $\binom{64}{2} = 2016$.

6.3.1 Goodness of the Metrics

We computed V_{lower} , V_{upper} , V_{dist} and V_{thresh} values for all the Morse datasets and plotted these with the test set classification accuracy results obtained from the FC NN in Section 6.2.1. The results are shown in Fig. 6.7, while the Pearson correlation coefficient r of each metric with the accuracy is given in Table 6.1. Note that the metrics are an indicator of dataset difficulty, so they are negatively correlated with accuracy. It is apparent that the V_{upper} and V_{thresh} metrics are the best since their r values have the highest magnitude.

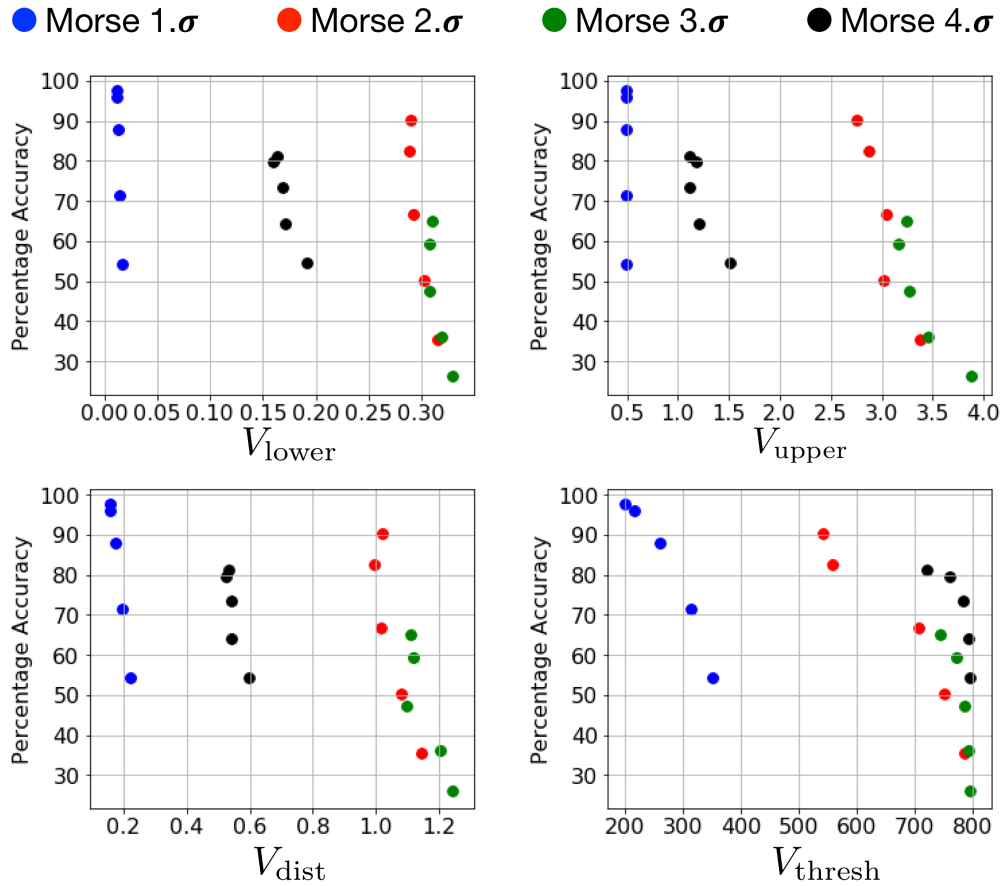


Figure 6.7: Plotting each metric for dataset difficulty vs. percentage accuracy obtained for datasets *Morse 1.σ* (blue), *2.σ* (red), *3.σ* (green) and *4.σ* (black). The accuracy results are using the FC NN, as reported in Section 6.2.1. Color coding is just for clarity, the r values in Table 6.1 take into account all the points regardless of color.

Table 6.1: Correlation coefficients between metrics and accuracy

Metric	r
V_{lower}	-0.59
V_{upper}	-0.64
V_{dist}	-0.63
V_{thresh}	-0.64

6.3.2 Limitations of the Metrics

As mentioned, each class has a single variance value which is the average variance across dimensions. This is a reasonable simplification to make because our experiments indicate that the variance of the variance values for different dimensions is small. However, this simplification possibly leads to the error bounds V_{lower} and V_{upper} not being sufficiently tight. A possible improvement, involving significantly more computation, would be to compute the $N_0 \times N_0$ covariance matrix K_m for each class.

It is worthwhile noting that **all these metrics are a function of the dataset only and are independent of the machine learning algorithm or training setup used.** On the other hand, percentage accuracy depends on the learning algorithm and training conditions. As shown in Fig. 6.4, increasing dataset size leads to accuracy improvement, *i.e.* the dataset becoming easier, since the NN has more training examples to learn from. However, increasing dataset size drives all the metric values towards indicating higher difficulty. This is because the occurrence of more examples in each class increases its standard deviation σ_m and also makes samples of a particular class more scattered, leading to reduced values for d and d_1 . We hypothesize that these shortcomings of the metrics are due to the fact that most variations of the Morse datasets have low redundancy, while the metrics (the error bounds in particular) are designed for high signal-to-noise ratio (SNR) problems, *i.e.* high redundancy.

6.4 Summary

This chapter presented an algorithm to generate datasets of varying difficulty on classifying Morse code symbols. While the results have been shown for NNs, any machine learning algorithm can be tried and the challenge arising from more difficult datasets used to fine tune it. The datasets are synthetic and consequently may not completely represent reality unless statistically verified with real-world tests. However, the different aspects of the generating algorithm help to mimic real-world scenarios which can suffer from noise or other inconsistencies. This chapter highlights one of the biggest advantages of synthetic data – the ability to easily produce large amounts of it and thereby improve the performance of learning algorithms. The given Morse datasets are also useful for testing the limits of various learning algorithms and identifying when they fail or possibly overfit/underfit. Finally, the metrics discussed, while not perfect, can be used to understand the inherent difficulty of the classification problem on any dataset before applying learning algorithms to it.

With the end of this chapter, we come to the end of our achieved contributions. The remainder of this dissertation proposal will discuss proposed research, *i.e.* our expected contributions.

Chapter 7

Model Search

Neural networks form a fascinating branch of study which has the potential to transform the way we lead our everyday lives and cause paradigm shifts in society. Some ‘big’ ideas which employ NNs at their core are AI and smart systems. Despite the explosion of interest in them, NNs offer little to human practitioners in terms of revealing their secrets. For example, why does a particular network with a particular number of layers and hyperparameter configuration perform well for some given task? Is there a general way or methodology to choosing the structure of a NN and its hyperparameters? These are the questions the field of **model search** hopes to answer. Model search forms the core of our proposed research, *i.e.* expected contributions.

We will broadly classify model search into two categories, each of which can be further classified into two types. These are summarized in Fig. 7.1. The remainder of this chapter will first explain these categories along with related work, then discuss our proposed research.

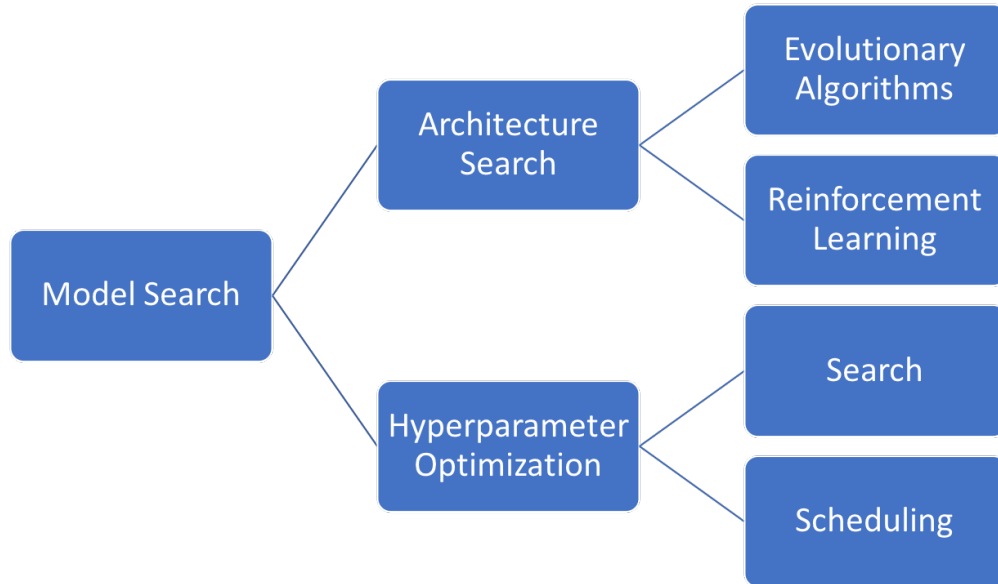


Figure 7.1: Types of model search.

7.1 Architecture Search

► **Definition 8: *Architecture search*:** Architecture search, also referred to in some works as Neural Architecture Search (NAS), is the category of model search concerned with the structure of NNs, including but not limited to the number, type and connection patterns across layers.

Such structural decisions are typically made by training a NN for some time, then using its validation performance as a metric for deciding on architectural modifications with the potential to improve performance. The strategy for making modifications can be implemented in several ways, as described next.

7.1.1 Evolutionary Algorithms

Evolutionary algorithms observe performance of many existing NN architectures, then **mutate** or evolve those with good performance and use these new derived architectures to replace other poorly-performing architectures. This is similar to how a species evolves, *i.e.* desirable traits are kept and bettered, and consequently replace undesirable traits.

Some previous works employing evolution are [14, 15, 110, 111]. The approach followed in most of these is to treat a layer in a NN as a node in a graph, and operations between different layers as edges of the graph. Examples of operations include different kinds of convolutions such as normal, pointwise (using filter size 1×1), depthwise separable [112], and dilated [113], different kinds of pooling layers such as max and average, different configurations of convolution filters such as varying filter size, number and stride (displacement between two consecutive groups of left nodes to which the convolution filter is applied), and so on. In particular, a) [14] uses a blueprint graph for the complete NN, and the actual NN traverses the blueprint and replaces individual graph nodes (*i.e.* NN layers) with new layers which perform better, b) [110] employs genetic traits such as selection, mutation and crossover to generate more competitive structures (*i.e.* performing better on the validation set) which can replace weaker ones, c) [111] treats a CNN *motif* – module consisting of a small number of convolution and pooling layers – as a node in a graph, and assembles well-performing motifs in a hierarchical way to get

the complete CNN (see [111, Fig. 1]), and d) [15] proposes a regularized variant of usual evolution by having the mutated architectures replace the oldest architectures instead of the worst-performing ones, which is similar to elderly humans in a society dying and new babies taking their place.

7.1.2 Reinforcement Learning (RL) methods

Reinforcement Learning (RL) differs from supervised learning in the sense that ground-truth labels are not provided during training. Instead there is a notion of **reward** for doing well, *e.g.* a NN being trained to play chess will get a positive reward for winning a game, making it more likely for it to repeat moves from the won game in future games. The entity being trained is the **agent**, which follows a **policy** which tells it to take certain **actions** depending on the **state** it is in. These actions and states have associated **values** (which are related to the rewards), and it is the goal of the learning algorithm to maximize these values by modifying the actions and/or policy accordingly. The reader is encouraged to refer to [114] for a complete exposition of RL concepts. We provide a brief overview here for the purposes of understanding model search.

Architecture search methods employing RL use validation performance, such as accuracy, as reward. The agent is often a controller NN, typically a RNN, while states are existing architectures for the NN whose architecture search is being performed. These states are also called *child networks*. The actions are modifications

to existing architectures, and the policy determines what modifications to make, *e.g.* adding a max-pooling layer after two convolution layers. Previous works using this technique are [16, 115, 116]. In particular, a) [16] introduced the approach of using a controller RNN, b) [116] reduced complexity by sharing parameters between different child networks, and c) [115] used transfer learning to select best child networks for a small dataset and applied these networks to bigger datasets, which saves time as compared to directly performing architecture search for bigger datasets. These works use variations of the *policy gradient* algorithm [117], which parametrizes the policy for choosing architectural modifications, then improves it by gradient ascent on its parameters using validation accuracy as reward. Thus, the reward plays a role similar to the cost in conventional BP, except that reward should be maximized, unlike cost.

An alternative to policy gradient is *Q-learning*, which picks the action to maximize some function estimating future rewards at every state. This is the technique used in [17, 118], who also use an ϵ -greedy exploration strategy. This means that the modifications to existing architectures can be completely new ones with probability ϵ , or alternatively with probability $1 - \epsilon$ can be an old architecture known to perform well. This is the *exploration vs exploitation* tradeoff commonly encountered in model search, *i.e.* how much should sticking with tried and trusted modifications known to perform well take precedence over trying out a completely new architecture?

7.1.3 Other methods of Architecture Search

This subsection lists some other notable efforts in architecture search of NNs. [119] used a *hypernetwork* to generate NN architectures and set their weights, however, the weights are restricted to a low rank space. [120] improvise on early stopping by observing several metrics during training, such as the first- and second-order differences between validation accuracies across epochs. [18] used surrogate functions in an attempt to approximate the performance of NNs prior to training. This is similar in spirit to our scatter metric, described in Section 5.5.2. Finally, [121] parametrizes architectures and uses gradient descent using the validation cost to better them.

7.2 Hyperparameter optimization

► **Definition 9: *Hyperparameter Optimization:*** Hyperparameter optimization is the problem of finding good values for non-trainable hyperparameters, and dynamically varying their values as required to improve NN performance.

Note that in general, any parameter which is not trainable, *i.e.* is chosen and adjusted by an entity external to the NN, is a hyperparameter. These also include the number, type and connection pattern layers, however, for the purposes of research focus, we have included those under the architecture search section. The hyperparameters we are concerned with are variables used in the optimization

algorithm such as learning rate η and momentum coefficients, those used in regularization such as the coefficient λ (there may be multiple λ values, *e.g.* for L1 and L2), and those used in the training process such as batch size M . Note that this list is not exhaustive.

There are two aspects to dealing with hyperparameters. Firstly, the **search** for an appropriate value. For example in our experiments in Section 3.3, we found from validation performance that L2 coefficient $\lambda = 1.6 \times 10^{-4}$ worked best for the FC MLP with $\mathbf{N}_{\text{net}} = (784, 224, 10)$, but when imposing pre-defined sparsity with $\mathbf{d}_{\text{net}}^{\text{out}} = (12, 10)$, the best value was $\lambda = 0$. There have been several previous efforts in searching for good values of hyperparameters. A celebrated result is the discovery in [20] that **random search** – randomly sampling hyperparameter values within reasonable ranges – worked better than **grid search** – choosing values of different hyperparameters from a finite, fixed list of values. This is because random search samples more values of the important hyperparameters. This raises the topic of NN performance being more sensitive to certain important hyperparameters, as will be discussed in the next paragraph. Other techniques for hyperparameter search include Bayesian optimization [122, 123], tree-structured algorithms [19, 124], and the successive halving algorithm [125].

The second aspect in dealing with hyperparameters is **scheduling**, *i.e.* dynamically changing the values of hyperparameters to optimize NN performance. Scheduling is typically done for important hyperparameters, *i.e.* those which affect

NN performance more than others. One previous work [126] investigated the importance of individual and groups of hyperparameters. In general, it has been found that the learning rate is possibly the most important hyperparameter. Values that are too high tend to make the cost value oscillate about its optimum point without reaching it, while values that are too low slow down convergence. Some optimization algorithms automatically reduce the effective learning rate over time, such as (2.18) for Adam. However, it has been found that periodically increasing the learning rate may also be beneficial, and previous works have proposed triangular [22], cosine [127], or batch-size dependent [21] schemes for doing so.

7.3 Our Proposed Research

Having discussed related work in model search, this section details the research we plan to do in this domain. This section will expand the expected contributions enumerated in Section 1.2.2 (except the one on improving hardware implementation).

7.3.1 Proposed research in Architecture Search

Several previous works on architecture search were very computationally intensive. For example, the mean total **GPU days** for finding a good architecture for CIFAR-10 are 90 in [17], 96 in [118], 150 in [18], 300 in [111], 1575 in [115] and 19,600 in [16]

(source [116]). This implies that the best architecture presented in [16] – achieving a test error of 3.65% on CIFAR-10 – would have taken a single GPU more than 50 years to discover! In comparison, the technique of parameter sharing between child networks proposed by [116] requires only 0.45 GPU days, or about 11 hours on a single GPU to achieve 2.89% test error on CIFAR-10. This example demonstrates the potential to accelerate the process of architecture search by designing better algorithms, and in our case, which also target lower complexity NNs.

Table 7.1 lists two primary reasons as to why most existing methods for architecture search are time and computation intensive, and how we can build on and extend our existing research to reduce the complexity of architecture search of NNs.

Lower complexity neural networks

Our efforts in pre-defined sparsity have thus far being restricted to MLPs. However, this still led to complexity reduction by a factor of $\sim 5X$ on a typical CNN architecture for CIFAR-100 without performance loss, as was shown in Fig. 3.3(d). We believe that our methods can be extended to more general NN architectures, such as CNNs and RNNs. Thus, our underlying theme in architecture search will be to discover low complexity NNs with minimal performance degradation. This is different from many existing works, where the goal was to discover the best possible NN for a given task. Our goal will be to find NNs that perform well on

Table 7.1: Building on and extending our research to reduce the complexity of architecture search.

Why architecture search is resource intensive?	Leveraging our existing research	Proposed research to simplify architecture search
The NNs being discovered via architecture search are quite complex , involving a large number of interconnected layers (often on the order of a hundred), hence take many GPU hours to train.	Our method of pre-defined sparsity results in NNs with lower complexity which train faster.	We plan to extend complexity reduction methods beyond MLPs and develop low complexity connection patterns for other kinds of NNs.
The decisions regarding how to modify existing architectures can only be taken after training until convergence (so as to get rewards for RL methods for example). Training on large datasets for hundreds of epochs needs many GPU hours.	The trends and design guidelines we discovered, along with methods to predict performance prior to training like scatter, will result in a net reduction in training time since several sparse networks discovered by architecture search can be filtered out quickly.	We plan to improve performance predicting mechanisms by gaining a deeper understanding of NNs , such as appropriate sparsity levels, and crucial parameters which can be monitored to provide information leading to early stopping of a training run.

a given task, while not taking several days of manual tuning of architectures or several weeks worth of GPU hours performing architectural search.

Deeper understanding of neural networks

Our second goal is to attempt to open the ‘black box’ of NNs and try to understand which parameters and hyperparameters are crucial to performance. As was demonstrated in Fig. 3.2, most FC layers have weights which are close to zero. This was an insight which led us to believe that most weights can be safely removed without

affecting performance. The results in Section 3.3 demonstrated the importance of properly choosing sparse architectures, such as designing ‘large and sparse’ networks and increasing the density of later junctions. We also plan to dynamically monitor parameters, such as the rate of change in the values of certain trainable parameters, how many units are ‘dead’ (zero activations) as a result of applying ReLU, and so on. These, along with further investigation into performance predicting measures and surrogate functions like scatter, will help us perform more informed early stopping and better filter out underperforming models. These traits are crucial for more efficient model search *i.e.* using less time and computational resources.

Summarizing proposed research in architectural search

In summary, we plan to automate the process of designing well-performing, low complexity sparse NNs for various applications. While we have not formed a clear consensus regarding evolutionary algorithms vs. RL methods, our approach seems to be better suited for evolutionary algorithms. This is because we plan to make specific modifications at regular training intervals, such as reducing density or changing connection patterns in the NN. This is better suited to evolutionary techniques such as mutation, which is a more natural approach to try as compared to designing rewards for RL methods.

7.3.2 Proposed research in Hyperparameter Optimization

We have already discovered some results regarding hyperparameters, *e.g.* pre-defined sparse NNs require less regularization. There are some other ideas we have at the moment, such as customized learning rates for more crucial parameters in sparse NNs. However, we have decided to focus on architecture search first, and turn to hyperparameter optimization based on availability of time and resources.

Chapter 8

Conclusion

This chapter concludes this dissertation proposal. We provide a summary of the field of neural networks, our achieved contributions, and the work discussed in this document, before ending with concrete plans for future work.

8.1 Summary

Despite the recent spurt of interest in NNs, they are not a completely new concept. The perceptron, which can be considered a very simple NN, was created in 1958 by Rosenblatt [128]. After a few decades of research stagnation, interest in NNs was again piqued in 1986 when Rumelhart et al. proposed backpropagation [129]. However, computers in the 1980s and 90s were not equipped to handle machine learning and NNs on a large scale as is done today. It was much more recently in the 2010s, perhaps with the publication of AlexNet [1], that the research interest and usage of NNs have skyrocketed. As a result, NNs have become the backbone of many AI applications nowadays, and are the major driving force behind the wide variety of ‘smart’ technologies that have become ubiquitous in our lives. This exponential growth in interest is summarized by Figs. 8.1 and 8.2, which show the

Number of AI papers on Scopus by subcategory (1998–2017)
Source: Elsevier

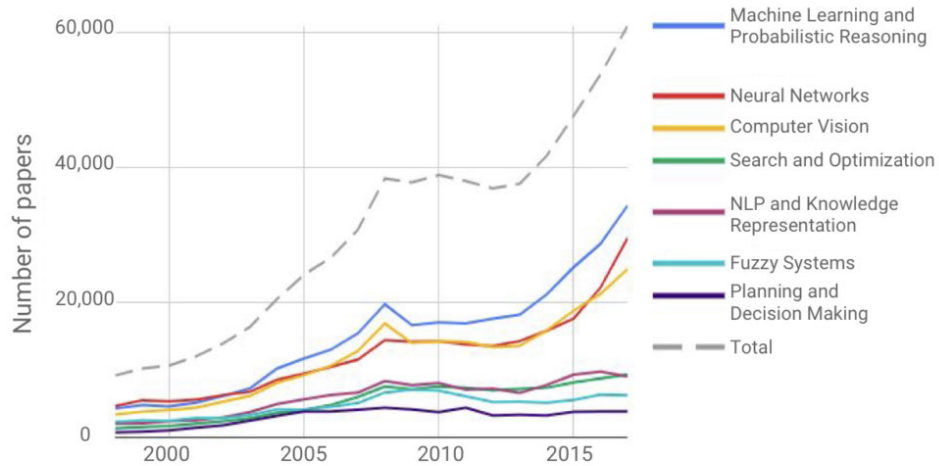


Figure 8.1: Number of AI papers on Scopus from 1998–2017, grouped by subcategory. Figure courtesy [130].

number of papers related to AI published in two widely used research databases – Elsevier’s Scopus, and arXiv, respectively. Notice how NN papers on Scopus have tripled in the last decade, while net AI papers on arXiv have grown by more than 10X from 2010 to 2017.

The growth in interest in NNs has also translated to a growth in their complexity. In particular, there has been a tendency in industry players to drive the field of NNs forward by exponentially increasing the number of parameters, as shown in Fig. 8.3. While this parameter explosion has led to pushing the frontiers of what NNs can achieve, it has often sacrificed **understandability**, *i.e.* being able to understand why a NN is performing in a certain way and how it works. This

Number of AI papers on arXiv by subcategory (2010–2017)

Source: arXiv

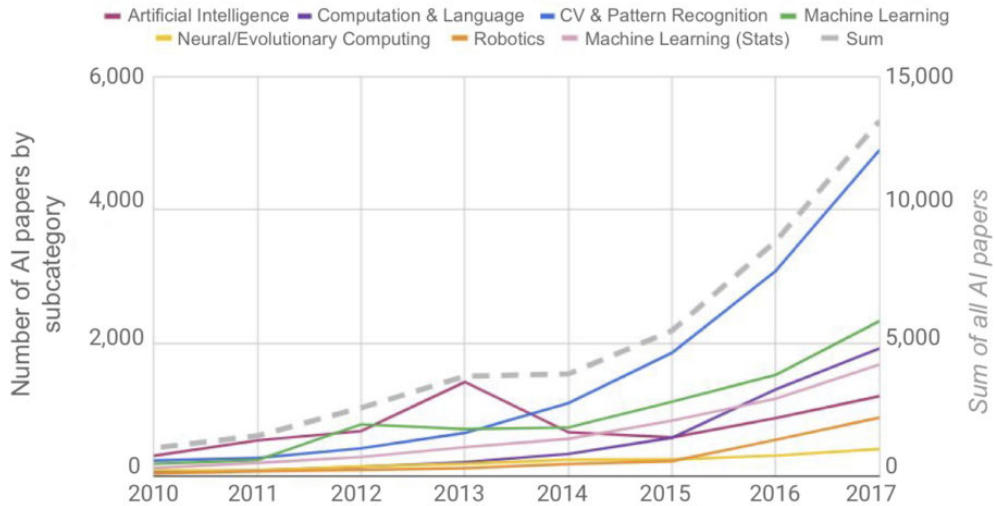


Figure 8.2: Number of AI papers on arXiv from 2010–2017, grouped by subcategory. Figure courtesy [130].

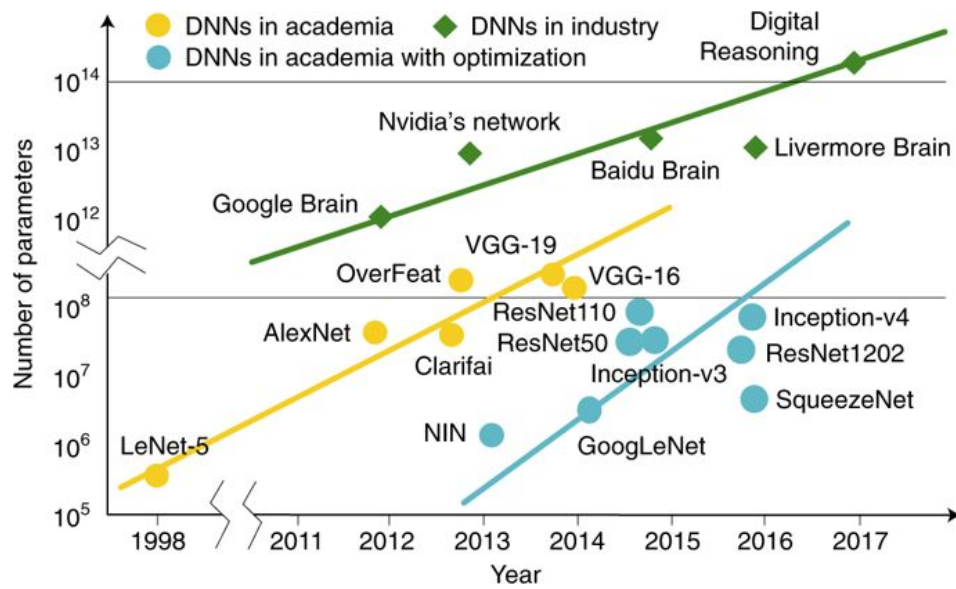


Figure 8.3: Total number of trainable parameters in some popular deep NN architectures over the years. Figure courtesy [131, Fig. 1].

status quo is concerning because it puts the reins of NN research into the sole hands of a few entities with the financial wherewithal to deploy GPUs and other computational resources at a scale well beyond the reach of most other entities. In other words, the NN parameter explosion has resulted in an oligarchy of sorts, where the barrier to entry in terms of required computational resources can push many prospective researchers away.

Our Contributions

The work presented in this dissertation proposal attempts to democratize the development of NNs by lowering their complexity. In particular, our proposed method of pre-defined sparsity simplifies the computational burden of NNs throughout their development cycle, *i.e.* training and inference. The computational complexity reduction also translates to storage complexity reduction when implementing NNs on custom hardware, such as we have developed. Moreover, our goal has not only been the simplification of NNs, but also attempts to understand why they perform in certain ways by analyzing their connection patterns. This will continue on into our expected contributions regarding model search, *i.e.* understanding how and why NNs work the way they do, attempting to predict their performance, and the consequent search for well-performing low-complexity NNs. As an additional effort, we have also developed a family of datasets of varying difficulty for benchmarking machine learning algorithms and applications.

8.2 Future Work

Our first future effort will be pipelining the hardware implementation described in Section 4.7 to improve its operational speed. This is estimated to take a month, *i.e.* May 2019, and will likely be continued by other researchers in our team in the realm of expansion to bigger NNs, possibly using cloud FPGA resources.

After that, we will focus on model search. Our first target in this regard will be to automate the process of discovering good sparse NNs, based on the trends and guidelines discussed in Section 3.3, and performance predicting measures building on the concepts discussed in Section 5.5. As we discover good sparse NNs and connection patterns, we will gain a deeper understanding of NNs and the widespread applicability of sparsifying techniques. This will lead to architecture search of low complexity NNs, as described in Section 7.1. Depending on our results, we expect this to take several months to a year, and will culminate in the final dissertation.

Depending on time and resources beyond that, we will look into hyperparameter optimization and more general techniques for model search of NNs.

8.3 Final Word

We believe our research has the potential to enable NN **exploration** on a wide scale. For more information on our research group and mission, the reader is

encouraged to visit [\[132\]](#) – the website for Hardware Accelerated Learning (HAL) at the University of Southern California (USC).

References

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” in *Proc. Advances in Neural Information Processing Systems 25 (NeurIPS)*, 2012, pp. 1097–1105.
- [2] M. Bojarski, D. D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba, “End to end learning for self-driving cars,” *arXiv preprint arXiv:1604.07316*, 2016.
- [3] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, Nov 2012.
- [4] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors,” *arXiv preprint arXiv:1207.0580*, 2012.
- [5] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” *arXiv preprint arXiv:1602.07261*, 2016.
- [6] A. Coates, B. Huval, T. Wang, D. J. Wu, A. Y. Ng, and B. Catanzaro, “Deep learning with COTS HPC systems,” in *Proc. 30th Int. Conf. Machine Learning (ICML)*, vol. 28, 2013, pp. III–1337–III–1345.
- [7] S. Han, J. Pool, J. Tran, and W. Dally, “Learning both weights and connections for efficient neural networks,” in *Proc. Advances in Neural Information Processing Systems 28 (NeurIPS)*, 2015, pp. 1135–1143.
- [8] N. P. Jouppi, C. Young, N. Patil *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *2017 ACM/IEEE 44th Annu. Int. Symp. Computer Architecture (ISCA)*, June 2017.

- [9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1–9.
- [10] Y. Gong, L. Liu, M. Yang, and L. D. Bourdev, “Compressing deep convolutional networks using vector quantization,” *arXiv preprint arXiv:1412.6115*, 2014.
- [11] W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen, “Compressing neural networks with the hashing trick,” in *Proc. 32nd Int. Conf. Machine Learning (ICML)*, 2015.
- [12] S. Han, H. Mao, and W. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2016.
- [13] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, “EIE: Efficient inference engine on compressed deep neural network,” in *Proc. 43rd Int. Symp. Computer Architecture (ISCA)*, 2016.
- [14] R. Miikkulainen, J. Liang, E. Meyerson, A. Rawal, D. F. O. Francon, B. Raju, H. Shahrzad, A. Navruzyan, N. Duffy, and B. Hodjat, “Evolving deep neural networks,” *arXiv preprint arXiv:1703.00548*, 2017.
- [15] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, “Regularized evolution for image classifier architecture search,” *arXiv preprint arXiv:1802.01548*, 2018.
- [16] B. Zoph and Q. V. Le, “Neural architecture search with reinforcement learning,” *arXiv preprint arXiv:1611.01578*, 2016.
- [17] B. Baker, O. Gupta, N. Naik, and R. Raskar, “Designing neural network architectures using reinforcement learning,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2017.
- [18] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, “Progressive neural architecture search,” *arXiv preprint arXiv 1712.00559*, 2017.
- [19] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kegl, “Algorithms for hyperparameter optimization,” in *Proc. Advances in Neural Information Processing Systems 24 (NeurIPS)*, 2011, pp. 2546–2554.
- [20] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *Journal of Machine Learning Research (JMLR)*, vol. 13, pp. 281–305, 2012.

- [21] P. Goyal, P. Dollar, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, “Accurate, large minibatch SGD: training ImageNet in 1 hour,” *arXiv preprint arXiv:1706.02677*, 2017.
- [22] L. N. Smith, “Cyclical learning rates for training neural networks,” in *Proc. IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2017.
- [23] S. Dey, K.-W. Huang, P. A. Beerel, and K. M. Chugg, “Pre-defined sparse neural networks with hardware acceleration,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 9, no. 2, pp. 332–345, June 2019.
- [24] S. Dey, K.-W. Huang, P. A. Beerel, and K. M. Chugg, “Characterizing sparse connectivity patterns in neural networks,” in *Proc. 2018 Information Theory and Applications Workshop (ITA)*, Feb 2018, pp. 1–9.
- [25] S. Dey, Y. Shao, K. M. Chugg, and P. A. Beerel, “Accelerating training of deep neural networks via sparse edge processing,” in *Proc. 26th Int. Conf. Artificial Neural Networks (ICANN)*. Springer, Sep 2017, pp. 273–280.
- [26] S. Dey, D. Chen, Z. Li, S. Kundu, K.-W. Huang, K. M. Chugg, and P. A. Beerel, “A highly parallel FPGA implementation of sparse neural network training,” in *Proc. Int. Conf. ReConfigurable Computing and FPGAs (ReConFig)*, Dec 2018, pp. 1–4, expanded preprint version available at <https://arxiv.org/abs/1806.01087>.
- [27] S. Dey, P. A. Beerel, and K. M. Chugg, “Interleaver design for deep neural networks,” in *Proc. 51st Asilomar Conf. Signals, Systems, and Computers (ACSSC)*, Oct 2017, pp. 1979–1983.
- [28] S. Dey, K. M. Chugg, and P. A. Beerel, “Morse code datasets for machine learning,” in *Proc. 9th Int. Conf. Computing, Communication and Networking Technologies (ICCCNT)*, Jul 2018, pp. 1–7.
- [29] S. Dey, “Github repository: usc-hal/morse-dataset,” <https://github.com/usc-hal/morse-dataset>.
- [30] —, “Matrix calculus,” University of Southern California, Tech. Rep., 2019, available online at https://www.researchgate.net/publication/332131671_Matrix_Calculus.
- [31] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proc. 27th Int. Conf. Machine Learning (ICML)*, 2010, pp. 807–814.

- [32] S. Hochreiter, “Untersuchungen zu dynamischen neuronalen netzen,” Master’s thesis, TU Munich, 1991.
- [33] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proc. 13th Int. Conf. Artificial Intelligence and Statistics*, 2010, pp. 249–256.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification,” in *Proc. IEEE Int. Conf. Computer Vision (ICCV)*, 2015, pp. 1026–1034.
- [35] S. Ruder, “An overview of gradient descent optimization algorithms,” *arXiv preprint arXiv:1609.04747*, 2016.
- [36] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2014.
- [37] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proc. 32nd Int. Conf. Machine Learning (ICML)*, 2015.
- [38] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [39] A. Deshpande, “A beginner’s guide to understanding convolutional neural networks,” <https://adeshpande3.github.io/A-Beginner%27s-Guide-To-Understanding-Convolutional-Neural-Networks/>, Jul 2016.
- [40] W. Bao, J. Yue, and Y. Rao, “A deep learning framework for financial time series using stacked autoencoders and long-short term memory,” *PLoS ONE*, vol. 12, Jul 2017.
- [41] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. D. Freitas, “Predicting parameters in deep learning,” in *Proc. Advances in Neural Information Processing Systems 26 (NeurIPS)*, 2013, pp. 2148–2156.
- [42] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2017.
- [43] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, “Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1,” *arXiv preprint arXiv:1602.02830*, 2016.

- [44] J. Albericio, P. Judd, T. Hetherington, T. Aamodt, N. E. Jerger, and A. Moshovos, “Cnvlutin: Ineffectual-neuron-free deep neural network computing,” in *Proc. 2016 ACM/IEEE 43rd Annu. Int. Symp. Computer Architecture (ISCA)*, 2016, pp. 1–13.
- [45] B. Reagen, P. Whatmough, R. Adolf *et al.*, “Minerva: Enabling low-power, highly-accurate deep neural network accelerators,” in *Proc. 2016 ACM/IEEE 43rd Annu. Int. Symp. Computer Architecture (ISCA)*, 2016, pp. 267–278.
- [46] A. Aghasi, A. Abdi, N. Nguyen, and J. Romberg, “Net-trim: Convex pruning of deep neural networks with performance guarantee,” in *Proc. Advances in Neural Information Processing Systems 30 (NeurIPS)*, 2017.
- [47] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [48] M. Courbariaux, Y. Bengio, and J.-P. David, “Training deep neural networks with low precision multiplications,” *arXiv preprint arXiv:1412.7024*, 2014.
- [49] S. Gupta, A. Agarwal, K. Gopalakrishnan, and P. Narayanan, “Deep learning with limited numerical precision,” *arXiv preprint arXiv:1502.02551*, 2015.
- [50] K. Kara, D. Alistarh, G. Alonso, O. Mutlu, and C. Zhang, “FPGA-accelerated dense linear machine learning: A precision-convergence trade-off,” in *Proc. IEEE Int. Symp. Field-Programmable Custom Computing Machines*, 2017.
- [51] V. Sindhvani, T. Sainath, and S. Kumar, “Structured transforms for small-footprint deep learning,” in *Proc. Advances in Neural Information Processing Systems 28 (NeurIPS)*, 2015, pp. 3088–3096.
- [52] S. Wang, Z. Li, C. Ding, B. Yuan, Y. Wang, Q. Qiu, and Y. Liang, “C-LSTM: Enabling efficient LSTM using structured compression techniques on FPGAs,” in *Proc. ACM/SIGDA Int. Symp. Field-Programmable Gate Arrays*, 2018.
- [53] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, “Learning structured sparsity in deep neural networks,” in *Proc. Advances in Neural Information Processing Systems 29 (NeurIPS)*, 2016, pp. 2074–2082.
- [54] S. Srinivas, A. Subramanya, and R. V. Babu, “Training sparse neural networks,” in *IEEE Conf. Computer Vision and Pattern Recognition Workshops (CVPRW)*, July 2017, pp. 455–462.

- [55] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, “Dian-
nao: A small-footprint high-throughput accelerator for ubiquitous machine-
learning,” in *Proc. 19th Int. Conf. Architectural Support for Programming
Languages and Operating Systems (ASPLOS)*, 2014, pp. 269–284.
- [56] S. Zhang, Z. Du, L. Zhang *et al.*, “Cambricon-X: An accelerator for sparse
neural networks,” in *Proc. 2016 49th Annu. IEEE/ACM Int. Symp. Microar-
chitecture (MICRO)*, 2016, pp. 1–12.
- [57] Y. Chen, T. Krishna, J. S. Emer, and V. Sze, “Eyeriss: An energy-efficient
reconfigurable accelerator for deep convolutional neural networks,” *IEEE
Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, 2017.
- [58] D. Shin, J. Lee, J. Lee, and H. Yoo, “14.2 DNPU: An 8.1TOPS/W reconfig-
urable CNN-RNN processor for general-purpose deep neural networks,” in
IEEE Int. Solid-State Circuits Conf. (ISSCC), Feb 2017, pp. 240–241.
- [59] Y. Ma, N. Suda, Y. Cao, S. Vrudhula, and J. Seo, “ALAMO: FPGA acceler-
ation of deep learning algorithms with a modularized RTL compiler,” *Inte-
gration, the VLSI Journal*, 2018.
- [60] N. Suda, V. Chandra, G. Dasika, A. Mohanty, Y. Ma, S. Vrudhula, J.-s. Seo,
and Y. Cao, “Throughput-optimized OpenCL-based FPGA accelerator for
large-scale convolutional neural networks,” in *Proc. 2016 ACM/SIGDA Int.
Symp. Field-Programmable Gate Arrays*. ACM, 2016, pp. 16–25.
- [61] S. Han, J. Kang, H. Mao, Y. Hu, X. Li, Y. Li, D. Xie, H. Luo, S. Yao,
Y. Wang, H. Yang, and W. B. J. Dally, “ESE: Efficient speech recognition
engine with sparse LSTM on FPGA,” in *Proc. ACM/SIGDA Int. Symp.
Field-Programmable Gate Arrays*, 2017, pp. 75–84.
- [62] R. G. Gironés, R. C. Palero, J. C. Boluda, and A. S. Cortés, “FPGA imple-
mentation of a pipelined on-line backpropagation,” *Journal of VLSI signal
processing systems for signal, image and video technology*, vol. 40, no. 2, pp.
189–213, Jun 2005.
- [63] N. Izeboudjen, A. Farah, H. Bessalah, A. Bouridene, and N. Chikhi,
“Towards a platform for FPGA implementation of the MLP based back
propagation algorithm,” in *Computational and Ambient Intelligence*, 2007,
pp. 497–505.
- [64] A. Gomperts, A. Ukil, and F. Zurfluh, “Development and implementation
of parameterized FPGA-based general purpose neural networks for online
applications,” *IEEE Transactions on Industrial Informatics*, vol. 7, no. 1,
pp. 78–89, Feb 2011.

- [65] Y. Chen, T. Luo, S. Liu, S. Zhang, L. He, J. Wang, L. Li, T. Chen, Z. Xu, N. Sun, and O. Temam, “Dadiannao: A machine-learning supercomputer,” in *47th Annu. IEEE/ACM Int. Symp. Microarchitecture*, Dec 2014, pp. 609–622.
- [66] C. Wang, L. Gong, Q. Yu, X. Li, Y. Xie, and X. Zhou, “DLAU: A scalable deep learning accelerator unit on FPGA,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 36, no. 3, pp. 513–517, Mar 2017.
- [67] A. Bourely, J. P. Boueri, and K. Choromonski, “Sparse neural network topologies,” *arXiv preprint arXiv:1706.05683*, 2017.
- [68] A. Prabhu, G. Varma, and A. M. Namboodiri, “Deep expander networks: Efficient deep networks from graph theory,” *arXiv preprint arXiv:1711.08757*, 2017.
- [69] D. C. Mocanu, E. Mocanu, P. Stone, P. H. Nguyen, M. Gibescu, and A. Liotta, “Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science,” *Nature Communications*, vol. 9, 2018.
- [70] J. Yosinski and H. Lipson, “Visually debugging restricted boltzmann machine training with a 3D example,” in *Proc. 29th Int. Conf. Machine Learning (ICML)*, 2012.
- [71] Y. LeCun, C. Cortes, and C. J. Burges, “The MNIST database of handwritten digits,” <http://yann.lecun.com/exdb/mnist/>.
- [72] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li, “RCV1: A new benchmark collection for text categorization research,” *Journal of machine learning research*, vol. 5, pp. 361–397, Apr 2004.
- [73] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, N. L. Dahlgren, and V. Zue, “TIMIT acoustic-phonetic continuous speech corpus,” <https://catalog.ldc.upenn.edu/LDC93S1>.
- [74] K.-F. Lee and H.-W. Hon, “Speaker-independent phone recognition using hidden markov models,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 37, no. 11, pp. 1641–1648, Nov 1989.
- [75] S. S. Stevens, J. Volkman, and E. B. Newman, “A scale for the measurement of the psychological magnitude pitch,” *The Journal of the Acoustical Society of America*, vol. 8, no. 185, 1937.

- [76] A. Krizhevsky, “Learning multiple layers of features from tiny images,” Master’s thesis, University of Toronto, 2009.
- [77] N. H. E. Weste and D. M. Harris, *CMOS VLSI Design: A Circuits and Systems Perspective*, 4th ed. Pearson, 2010.
- [78] D. Masters and C. Luschi, “Revisiting Small Batch Training for Deep Neural Networks,” *arXiv preprint arXiv:1804.07612*, 2018.
- [79] Digilent, “Nexys 4 DDR,” <https://reference.digilentinc.com/reference/programmable-logic/nexys-4-ddr/start>.
- [80] Xilinx, “VIRTEX ultra-scale+,” <https://www.xilinx.com/products/silicon-devices/fpga/virtex-ultrascale-plus.html>.
- [81] Amazon Web Services, “Amazon EC2 F1 instances,” <https://aws.amazon.com/ec2/instance-types/f1/>.
- [82] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proc. Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [83] A. S. Asratian, T. M. J. Denley, and R. Haggkvist, *Bipartite Graphs and Their Applications*. Cambridge University Press, 1998.
- [84] A. Barvinok, “On the number of matrices and a random matrix with prescribed row and column sums and 0–1 entries,” *Advances in Mathematics*, vol. 224, no. 1, pp. 316–339, 2010.
- [85] G. Masera, G. Piccinini, M. R. Roch, and M. Zamboni, “VLSI architectures for turbo codes,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 7, no. 3, pp. 369–379, 1999.
- [86] T. Brack, M. Alles, T. Lehnigk-Emden *et al.*, “Low complexity LDPC code decoders for next generation standards,” in *Design, Automation & Test in Europe Conf. & Exhibition (DATE)*, 2007, pp. 1–6.
- [87] S. Crozier and P. Guinand, “High-performance low-memory interleaver banks for turbo-codes,” in *Vehicular Technology Conf., 2001. VTC 2001 Fall. IEEE VTS 54th*, vol. 4, 2001, pp. 2394–2398.
- [88] G. M. Weiss and F. Provost, “Learning when training data are costly: The effect of class distribution on tree induction,” *Journal of Artificial Intelligence Research*, vol. 19, pp. 315–354, 2003.

- [89] I. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2015.
- [90] A. Schumacher, “Problems with imagenet and its solutions,” https://planspace.org/20170911-problems_with_imagenet_and_its_solutions/, September 2017.
- [91] J. Prendki, “The curse of big data labeling and three ways to solve it,” <https://aws.amazon.com/blogs/apn/the-curse-of-big-data-labeling-and-three-ways-to-solve-it/>, November 2018.
- [92] M. Gregory and K. Ocasio, “The [hidden] challenges of ML series: Quadrant 2 data preparation,” <https://www.ntconcepts.com/the-hidden-challenges-of-ml-series-quadrant-2-data-preparation/>, February 2019.
- [93] X. Peng, B. Sun, K. Ali, and K. Saenko, “Learning deep object detectors from 3D models,” in *Proc. IEEE Int. Conf. Computer Vision (ICCV)*. IEEE Computer Society, 2015, pp. 1278–1286.
- [94] D. DeTone, T. Malisiewicz, and A. Rabinovich, “Toward geometric deep SLAM,” *arXiv preprint arXiv:1707.07410*, 2017.
- [95] T. Anh Le, A. G. Baydin, R. Zinkov, and F. Wood, “Using synthetic data to train neural networks is model-based reasoning,” *arXiv preprint arXiv:1703.00868*, 2017.
- [96] N. Patki, R. Wedge, and K. Veeramachaneni, “The synthetic data vault,” in *Proc. IEEE Int. Conf. Data Science and Advanced Analytics (DSAA)*, 2016, pp. 399–410.
- [97] C.-H. Luo and C.-H. Shih, “Adaptive morse-coded single-switch communication system for the disabled,” *Int. Journal of Bio-Medical Computing*, vol. 41, no. 2, pp. 99–106, 1996.
- [98] C.-H. Yang, C.-H. Yang, L.-Y. Chuang, and T.-K. Truong, “The application of the neural network on morse code recognition for users with physical impairments,” *Proc. Institution of Mechanical Engineers, Part H: Journal of Engineering in Medicine*, vol. 215, no. 3, pp. 325–331, 2001.
- [99] C.-H. Yang, L.-Y. Chuang, C.-H. Yang, and C.-H. Luo, “Morse code application for wireless environmental control systems for severely disabled individuals,” *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 11, no. 4, pp. 463–469, Dec 2003.

- [100] C. P. Ravikumar and M. Dathi, “A fuzzy-logic based morse code entry system with a touch-pad interface for physically disabled persons,” in *Proc. IEEE Annu. India Conf. (INDICON)*, Dec 2016.
- [101] T. W. King, *Modern Morse Code in Rehabilitation and Education: New Applications in Assistive Technology*, 1st ed. Allyn and Bacon, 1999.
- [102] R. Sheinker, “Morse code - apps on google play,” <https://play.google.com/store/apps/details?id=com.dev.morsecode&hl=en>, Aug 2017.
- [103] F. Bonnin, “Morse-it on the app store,” <https://itunes.apple.com/us/app/morse-it/id284942940?mt=8>, Mar 2018.
- [104] D. Hill, “Temporally processing neural networks for morse code recognition,” in *Theory and Applications of Neural Networks*. Springer London, 1992, pp. 180–197.
- [105] G. N. Aly and A. M. Sameh, “Evolution of recurrent cascade correlation networks with distributed collaborative species,” in *Proc. 1st IEEE Symp. Combinations of Evolutionary Computation and Neural Networks*, 2000, pp. 240–249.
- [106] R. Li, M. Nguyen, and W. Q. Yan, “Morse codes enter using finger gesture recognition,” in *Proc. Int. Conf. Digital Image Computing: Techniques and Applications (DICTA)*, Nov 2017.
- [107] *International Morse Code*, Radiocommunication Sector of International Telecommunication Union, Oct 2009, available at <http://www.itu.int/rec/R-REC-M.1677-1-200910-I/>.
- [108] R. M. Zur, Y. Jiang, L. L. Pesce, and K. Drukker, “Noise injection for training artificial neural networks: A comparison with weight decay and early stopping,” *Medical Physics*, vol. 36, no. 10, pp. 4810–4818, Oct 2009.
- [109] K. Chugg, A. Anastasopoulos, and X. Chen, *Iterative Detection: Adaptivity, Complexity Reduction, and Applications*. Springer Science & Business Media, 2012, vol. 602.
- [110] L. Xie and A. Yuille, “Genetic CNN,” in *IEEE Int. Conf. Computer Vision (ICCV)*, Oct 2017, pp. 1388–1397.
- [111] H. Liu, K. Simonyan, O. Vinyals, C. Fernando, and K. Kavukcuoglu, “Hierarchical representations for efficient architecture search,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2018.

- [112] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” *arXiv preprint arXiv:1610.02357*, 2016.
- [113] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2016.
- [114] R. S. Sutton and A. G. Barto, *Reinforcement Learning*, 2nd ed. MIT Press, 2018.
- [115] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning transferable architectures for scalable image recognition,” in *IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [116] H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, and J. Dean, “Efficient neural architecture search via parameter sharing,” *arXiv preprint arXiv:1802.03268*, 2018.
- [117] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine Learning*, vol. 8, no. 3, pp. 229–256, May 1992.
- [118] Z. Zhong, J. Yan, W. Wu, J. Shao, and C.-L. Liu, “Practical block-wise neural network architecture generation,” *arXiv preprint arXiv:1708.05552*, 2017.
- [119] A. Brock, T. Lim, J. Ritchie, and N. Weston, “SMASH: One-shot model architecture search through hypernetworks,” *arXiv preprint arXiv:1708.05344*, 2017.
- [120] B. Baker, O. Gupta, R. Raskar, and N. Naik, “Accelerating neural architecture search using performance prediction,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2018.
- [121] H. Liu, K. Simonyan, and Y. Yang, “DARTS: Differentiable architecture search,” *arXiv preprint arXiv:1806.09055*, 2018.
- [122] J. Snoek, H. Larochelle, and R. P. Adams, “Practical bayesian optimization of machine learning algorithms,” in *Proc. Advances in Neural Information Processing Systems 25 (NeurIPS)*, 2012, pp. 2951–2959.
- [123] J. Snoek, O. Rippel, K. Swersky, R. Kiros, N. Satish, N. Sundaram, M. M. A. Patwary, Prabhat, and R. P. Adams, “Scalable bayesian optimization using deep neural networks,” *arXiv preprint arXiv:1502.05700*, 2015.

- [124] J. Bergstra, D. Yamins, and D. D. Cox, “Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures,” in *Proc. 30th Int. Conf. Machine Learning (ICML)*, 2013.
- [125] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, “Hyperband: Bandit - based configuration evaluation for hyperparameter optimization,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2017.
- [126] F. Hutter, H. Hoos, and K. Leyton-Brown, “An efficient approach for assessing hyperparameter importance,” in *Proc. 31st Int. Conf. Machine Learning (ICML)*, 2014, pp. I-754–I-762.
- [127] I. Loshchilov and F. Hutter, “SGDR: Stochastic gradient descent with warm restarts,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2017.
- [128] F. Rosenblatt, “The perceptron: A probabilistic model for information storage and organization in the brain,” *Psychological Review*, pp. 65–386, 1958.
- [129] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [130] D. Enskat, “AI index 2018: Europe – china – united states; AI outpaces CS,” <https://warrenenskat.com/ai-index-2018-europe-china-united-states-ai-outpaces-cs/>, Dec 2018.
- [131] X. Xu, Y. Ding, S. X. Hu, M. Niemier, J. Cong, Y. Hu, and Y. Shi, “Scaling for edge inference of deep neural networks,” *Nature Electronics*, vol. 1, no. 4, pp. 216–222, 2018.
- [132] USC HAL team, “Hardware accelerated learning,” <https://hal.usc.edu/>.