

ABSTRACT

BAKER, LEE B. Design of a 3DIC system to aid in the acceleration of edge systems that employ multiple instances of disparate artificial neural networks. (Under the direction of Paul Franzon.)

This dissertation explores employing a Three-Dimensional Integrated Circuit (3DIC) in the acceleration of Artificial Neural Networks (ANNs) for systems deployed in customer facing applications. Assuming ANNs fulfill their potential, it is this works belief that these systems will employ ANNs for various functions, such as engine monitoring, anomaly detection, navigation etc. and that the various system functions are implemented with a set of disparate ANNs. A further assumption is that these customer facing systems may not have access to cloud servers or the cloud servers do not provide the necessary turn-around time for processing the ANN.

Although ANNs have been known about for many decades, it hasn't been until the last few years that they have demonstrated efficacy in applications such as image recognition and voice recognition. These ANNs have demonstrated significant improvements over what was considered to be state-of-the-art algorithms.

Artificial Neurons (ANes) take their inspiration from neuron behavior observed in the mammalian brain, although implementations are simplifications of what actually exists in the brain. These simplifications range from attempts to emulate the actual spiking behavior of real neurons to ANes that simply encode the spiking behavior in the form of a number or rate.

The ANNs that have demonstrated most efficacy are a family of neural networks that can be described as Deep Neural Networks (DNNs). These DNNs are created by cascading layers of rate-based ANes to form a large, layered ANN employing ten's of thousands or more of ANes.

Considering the storage required for the input and the ANN parameters, the storage requirements result in gigabytes of memory. When these ANNs are required to be solved in fractions of a second, the processing and memory bandwidth becomes prohibitive.

Unfortunately, to achieve a high performance, existing implementations rely on processing a batch of inputs, such as processing a batch of images or voice recordings which all use the same ANN

or by employing a sub-family of DNNs, known as Convolutional Neural Networks (CNNs), which reuse portions of the ANN parameters. These techniques allow these implementations to hold and reuse data in fast local Static Random Access Memory (SRAM). With this works target application, the assumption is there is little opportunity for batch processing or reuse therefore data must be drawn constantly from main memory, which generally is Dynamic Random Access Memory (DRAM).

One area of integrated circuit technology that hasn't been widely used in ANNs is 3DICs. 3DICs have the potential to increase connectivity, and thus bandwidth and keep power dissipation to within acceptable levels.

This work demonstrates how a customized 3DIC DRAM can be combined with application-specific layers to produce a system meeting the required level of performance in systems with multiple instances of disparate ANNs.

© Copyright 2017 by Lee B. Baker

All Rights Reserved

Design of a 3DIC system to aid in the acceleration of edge systems that
employ multiple instances of disparate artificial neural networks

by
Lee B. Baker

A dissertation submitted to the Graduate Faculty of
North Carolina State University
in partial fulfillment of the
requirements for the Degree of
Doctor of Philosophy

Electrical Engineering

Raleigh, North Carolina

2017

APPROVED BY:

Winser Alexander

Gregory Byrd

Richard Warr

Paul Franzon
Chair of Advisory Committee

DEDICATION

To my wife Mandy, my children Adam, Rachel and Paul and my parents Joan and Barry.

BIOGRAPHY

The author was born in the United Kingdom. After performing poorly in high school he took a job in a local electronic engineering firm under a vocational program. After seeing the white coated "engineers" being called down from upstairs to solve the "big" problems, he decided he wanted to wear one of those white coats (his dress sense was wanting). The journey took him to Brighton Polytechnic, now Brighton University and a First Class Honours Degree in Electrical Engineering. After working in the UK for a couple of years, he moved to the United States. The journey included a family with a daughter and two sons. The education continued with a Masters in Engineering from Villanova University and a Masters in Business Administration from North Carolina State University.

With the family now being somewhat independent, he decided to make a career change which would hopefully include teaching.

That career change included enrolling in the Electrical Engineering PhD program at North Carolina State University. This stage of the education journey has resulted in this dissertation.

And remember:

"do not stand still."

"do not let your past dictate your future."

ACKNOWLEDGEMENTS

At a personal level, I would like to thank my wife Mandy and my children Adam, Rachel and Paul for their encouragement.

I would like to thank my advisor, Paul Franzon for his help in making this possible.

I would also like to thank my fellow students, especially Jong Beom, Josh, Sumon and Weifu for their healthy discussions and, being an older student, referring to me as Lee and not Sir or Mr. Baker.

TABLE OF CONTENTS

LIST OF TABLES	vii
LIST OF FIGURES	viii
Chapter 1 Introduction	1
1.1 Overview	1
1.2 Abbreviations	3
Chapter 2 Artificial Neural Networks	5
2.1 Artificial Neural Networks	7
2.2 ANN Layers Problem	12
2.2.1 Deep Neural Networks	12
2.2.2 Feature Layers	14
2.3 The Problem	16
2.3.1 ANN Processing	17
Chapter 3 Three-Dimensional Integrated Circuits (3DIC)	23
Chapter 4 DRAM Customizations	30
4.1 Customization One: Very-Wide Bus	33
4.2 Customization Two: Write Mask	34
Chapter 5 System Description	36
5.1 The Solution	37
5.1.1 Novelty	38
5.2 System Overview	38
5.2.1 3D-DRAM	40
5.2.2 Manager Layer	40
5.2.3 Processing Layer	41
5.2.4 Layer Interconnect	42
5.2.5 Inter-Manager Communication	42
5.3 System Operations	42
5.3.1 Manager Operations	44
5.3.1.1 Instructions	44
5.3.1.2 Accessing of Pre-synaptic AN states and connection weights	46
5.3.1.3 Writing AN state results to memory	48
5.3.2 PE Operations	50
5.3.2.1 Streaming Operations (Streaming Operation (StOp))	50
5.3.2.2 SIMD	50
5.3.2.3 Configuration	50
5.4 Detailed System Description	51
5.4.1 Manager	51

5.4.1.1	Operation Decode	51
5.4.1.2	Argument Decode	52
5.4.1.3	Result data Processing	52
5.4.1.4	Memory Write Controller	53
5.4.2	Processing Engine	54
5.4.2.1	Configuration	54
5.4.2.2	Streaming Operations	54
5.4.2.3	SIMD	55
5.4.2.4	Result Data	55
Chapter 6	Results	57
BIBLIOGRAPHY	61
APPENDICES	63

LIST OF TABLES

Table 3.1	TSV Design Characteristics	29
Table 6.1	Simulation-based estimates	59

LIST OF FIGURES

Figure 2.1	Artists impression of a Mammalian Neuron	6
Figure 2.2	Artificial Neural Network	9
Figure 2.3	Example Rated-Based Model Activation functions	9
Figure 2.4	Example Spiking Activation Function Model	10
Figure 2.5	Layered Artificial Neural Networks	11
Figure 2.6	Classification using ANNs layers	13
Figure 2.7	Deep network showing feature layers	13
Figure 2.8	Layer Connection Types	15
Figure 2.9	Layer Feature Planes	21
Figure 2.10	Typical DNN[Kri12]	22
Figure 3.1	3DIC Stack of Die	24
Figure 3.2	Die Stack profile [ITR15] with Through-Silicon Vias (TSVs)	27
Figure 4.1	Typical Memory Block Diagram [Jac07]	31
Figure 4.2	RAM Storage Cell Types	32
Figure 4.3	Typical DRAM Block Diagram	34
Figure 4.4	Exposing more of the DRAM page	35
Figure 5.1	3DIC System Stack	39
Figure 5.2	DRAM Physical Interface Layout showing area for Sub-System Column (SSC)	41
Figure 5.3	System Flow Diagram	43
Figure 5.4	Instruction 4-tuple	45
Figure 5.5	Descriptor 6-tuple	46
Figure 5.6	ROI Storage	47
Figure 5.7	Storage Descriptor	49
Figure 5.8	Sub-System Column (SSC) Flow Diagram	51
Figure 5.9	Sub-System Column (SSC) Block Diagram	56
Figure 6.1	Manager and PE Die layouts	60

CHAPTER

1

INTRODUCTION

1.1 Overview

Machine Learning in the form of Deep Neural Networks (DNNs) have gained traction over the last few years. They get good press in applications such as image recognition and speech recognition. DNNs are constructed from a basic building block, the Artificial Neuron (ANe). With popular DNNs, the Artificial Neural Network (ANN) is often formed from tens of layers with each layer containing many ANes. In most cases, these layers are processed in a feed-forward manner with one layer being the inputs to the next layer. Therefore, useful DNNs often require hundreds of thousands of ANes and within the network, each ANe can have hundreds, even thousands of feeder or pre-synaptic ANes.

There have been implementations that use different number formats from double precision

floating point to eight bit integers, but in all cases these useful ANNs have significant memory requirements to store the connection weights (parameters) therefore requiring Dynamic Random Access memory (DRAM) to store the ANe parameters.

There have been many successful attempts to accelerate ANNs, but in most cases the focus is on a subset of the DNN known as the Convolutional Neural network (CNN). CNNs assume a significant amount of reuse of the weights connecting ANes and thus they can take advantage of local memory (SRAM).

Much of the Application-Specific Integrated Circuit (ASIC) and Application-Specific Instruction-set Processor (ASIP) ANN research has focused on taking advantage of the performance and ease of use of SRAM. These implementations can be shown to be effective with specific ANN architectures, such as CNNs where the ANN parameters can be stored in SRAM in a cache-like architecture avoiding constant accessing of the "slower" DRAM. In addition, to achieve a high performance, these rely on processing a batch of inputs, such as processing a batch of images or voice recordings using the same ANN.

The work in this paper considers "edge" applications that require the processing of a disparate set of useful sized ANNs. The work assumes that the application system is utilizing ANNs for the processing of various sub-systems, such as navigation, engine monitoring etc.. This work also does not assume the ANN is specifically a CNN but a DNN where there may not be opportunities to store and reuse portions of the ANN in SRAM. A further assumption is that the target edge devices do not include opportunities to perform batch processing. Under these circumstance, when these implementations need to constantly load ANN parameters directly from main memory, the performance is constrained to the DRAM interface bandwidth. Therefore, the performance of SRAM-based ASIC/ASIP implementations are severely degraded to the point of being unacceptable.

This work uses the DRAM as the primary processing storage and employs minimal SRAM for the processing of the ANe. In addition, the work considers 3D integrated circuit technology and a custom 3D-DRAM. By employing 3DIC technology, this work takes advantage of the reduced energy and area

and increased connectivity and bandwidth to allow the DRAM to be employed efficiently without the need for local SRAM. This work demonstrates that a 3DIC system based on a customized 3D-DRAM could be used in edge applications requiring at or near real-time performance for systems running multiple ANNs.

It should be noted that this work does not design a custom 3D-DRAM but answers the question "if such a device were available, can we employ it within a useful ANN system".

1.2 Abbreviations

Acronyms

3DIC Three-Dimensional Integrated Circuit

ANe Artificial Neuron

ANN Artificial Neural Network

ASIC Application-Specific Integrated Circuit

ASIP Application-Specific Instruction-set Processor

CNN Convolutional Neural Network

DDR Double Data Rate

DiRAM4 Dis-Integrated 3D DRAM

DNN Deep Neural Network

DRAM Dynamic Random Access Memory

ESD Electrostatic discharge

HBM High Bandwidth Memory

HMC Hybrid Memory Cube

IC Integrated Circuit

IP Intellectual property

KGD Known Good Die

LSTM Long Short-term memory

NoC Network-on-Chip

PC Program Counter

PE Processing Engine

SIMD Single-Instruction Multiple-Data

SoC System-on-Chip

SRAM Static Random Access Memory

SSC Sub-System Column

StOp Streaming Operation

TSV Through-Silicon Via

CHAPTER

2

ARTIFICIAL NEURAL NETWORKS

Recently, there has been much interest in the use of artificial neural networks in systems that employ tasks such as image recognition[Kri12], text recognition[Qiu13] and game playing[Mad14]. In particular, in the field of image recognition these artificial neural network models have demonstrated superior performance over other state-of-the-art technology[Kri12]. These artificial neural networks will continue to be applied to numerous other areas such as voice recognition, text recognition, face recognition and autonomous control.

Artificial neural networks (ANN) take their inspiration from neuron behavior observed in the mammalian brain, although implementations are simplifications of what actually exists in the brain.

The mammalian neuron is a cell that receives input and generates output in the form of electrical and chemical processes. The neuron has a cell body (or soma), a group of dendrites which provide

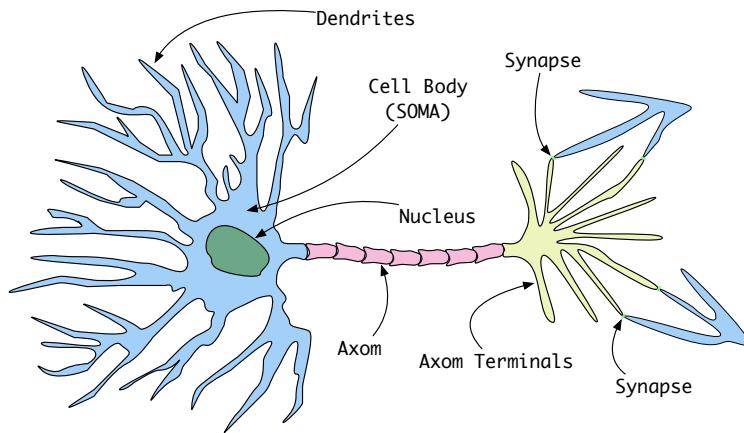


Figure 2.1 Artists impression of a Mammalian Neuron

the inputs from other cells, a cell body, an axon which generates the output signals, and the axon terminals which are the outputs of the cell. The connection from a cell's output, or axon terminal to another cell's input, or dendrite is known as a synapse. The connection in the synapse is a chemical process stimulated by electrical impulses. The neuron can be seen in figure 2.1.

The connection from one cell to another has both an associated delay and a strength. The strength of the connection can be influenced by the size of the pre-synaptic neuron spike or by the pre-synaptic neuron generating a series of spikes rather than a single spike.

So it is known that mammalian neurons generate "spikes" in response to inputs which for humans include sight, touch, sound etc.. This spiking behavior is often referred to as the neuron being activated. When these neurons are activated, their spikes propagate to other neurons. Under certain conditions, the combination of the various inputs to a neuron cause it to activate. A particular neuron may have many hundreds, perhaps thousands of other neurons connected to its "input". These input neurons are referred to as pre-synaptic neurons. These pre-synaptic neurons may provide input to many neurons which are referred to as post-synaptic neurons. A particular neuron can get activated by a particular arrival pattern of pre-synaptic neuron spikes or simply by the intensity of the pre-synaptic spikes.

The spiking behavior of a neuron also varies and many spiking profiles have been observed, including single spikes, groups of spikes and repetitive spiking. It is believed that information is carried in the delay and strength of the connections and how pre-synaptic neurons combine to cause a neuron to activate. In simple terms, if a neuron is activated by its pre-synaptic neurons, then the activation of the neuron means a pattern has been detected which will influence a reaction. In mammalian terms, that might be the detection of a threat from both smell and sight neurons and the reaction is to control muscles resulting in flight.

The various chemical and electrical processes that result in the generation and propagation of these neuron spikes is beyond the scope of this dissertation, but how neurons and networks of neurons are artificially emulated is what we will discuss next.

2.1 Artificial Neural Networks

When modeling these neurons in artificial neural networks, the neuron models either generate actual spikes similar to actual neurons or produce a value which is proportional to the rate at which spikes occur. These artificial neural networks can be categorized as rate-based coded or spike time coded neurons.

When used in networks of neurons, both model types employ a connection weight between the pre and post-synaptic neuron, however, the spiking neuron network also introduces a time delay associated with the connection.

The spiking neuron model is characterized by:

- Connections between neurons have both a strength and a delay
 - The pre-synaptic neuron output is multiplied by the connection weight and delayed
- The weighted inputs from all pre-synaptic neurons are accumulated
- The accumulated inputs drives an activation function

- the activation function $f(x)$ is a spiking model is based on differential equations
- many models have been proposed with varying levels of complexity

examples are:

- Leaky integrate and fire
- Izhikevich [Izh04] (see Fig. 2.4a)

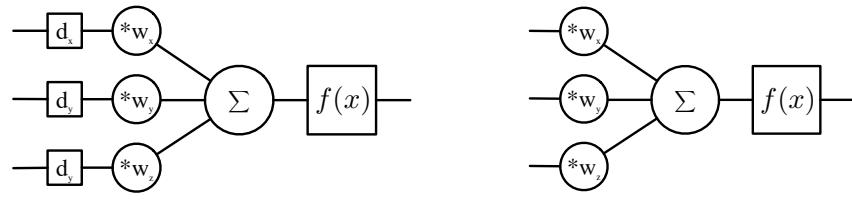
The Rate-based neuron model is characterized by:

- Connections between neurons have only a strength
 - The pre-synaptic neuron output is multiplied by the connection weight
- The weighted inputs from all pre-synaptic neurons are accumulated
- The accumulated inputs drives an activation function
 - the activation function $f(x)$ is a non-linear function
 - early models used binary functions although in practice the function needs to be differentiable

examples are (see Fig. 2.3):

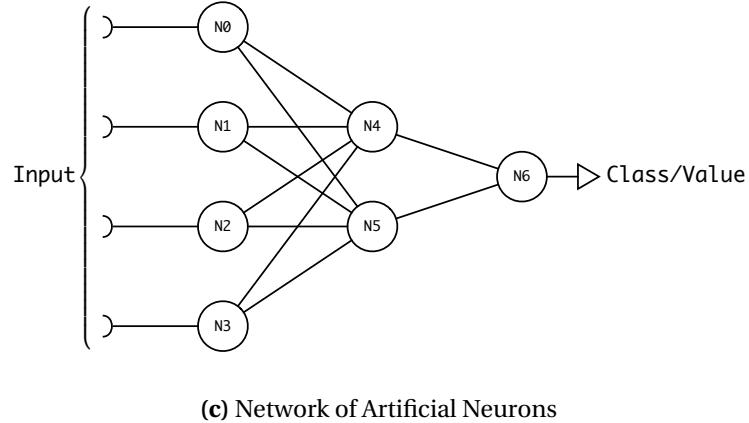
- sigmoid
- rectified linear unit

To emulate complex behavior, the artificial neurons are connected in networks, typically with layers of sub-networks which are in effect separated by the non-linear activation function. Examples of both rate-based and spiking artificial neural networks can be seen in Fig. 2.5a and Fig. 2.5b respectively. Typically neural networks process in a feed-forward fashion. Considering Fig. 2.2, this means the input arrives on the left, the inputs propagate to neurons N0 through N3. When N0 through N3 are processed, their values propagate forward to neurons N4 and N5 etc.. Sometimes ANNs also include



(a) Spiking Model

(b) Rate-Based Model



(c) Network of Artificial Neurons

Figure 2.2 Artificial Neurons and Network

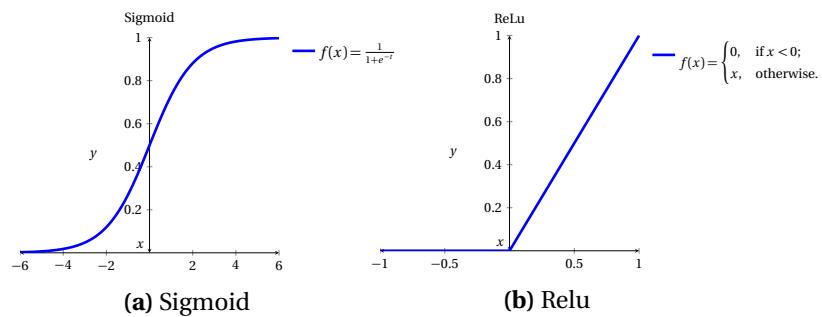


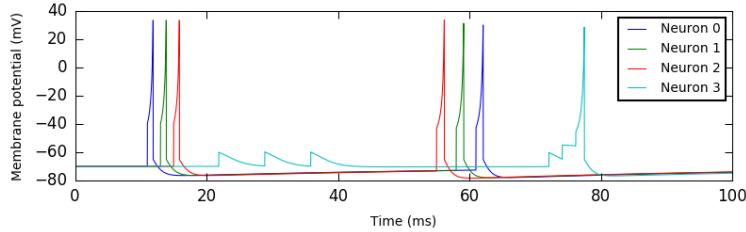
Figure 2.3 Example Rated-Based Model Activation functions

$$v' = 0.04v^2 + 5v + 140 - u - I$$

$$u' = a(bv - u)$$

$$\text{if } v \geq 30 \text{ mV, then } \begin{cases} v \leftarrow c \\ u \leftarrow u + d \end{cases}$$

(a) Izhikevich Model[Izh04]



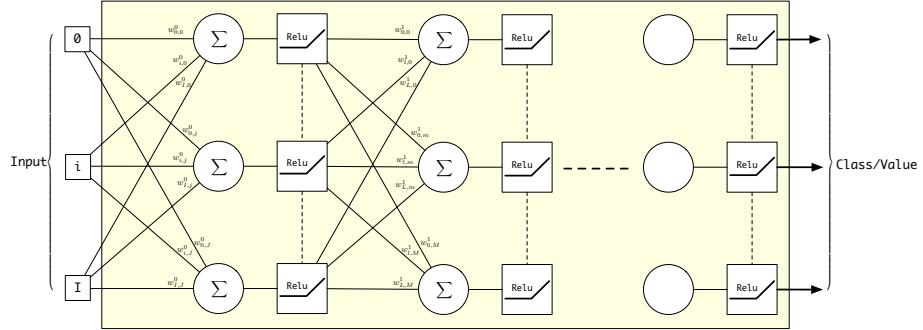
(b) Izhikevich Model Simulation [Izh04][CH06]

Figure 2.4 Example Spiking Activation Function Model

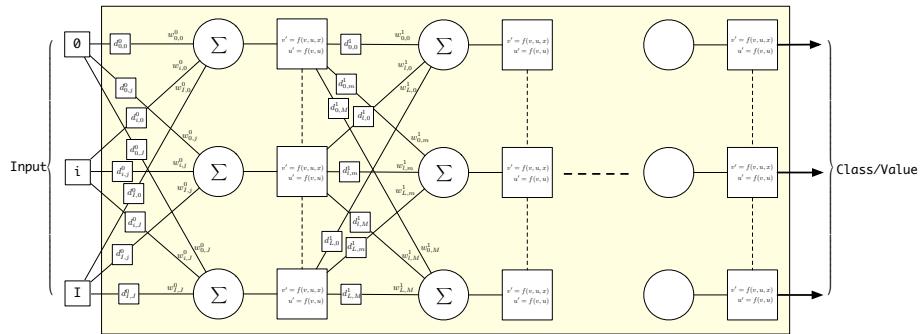
recursion where for example neurons N0 through N4 are not only influenced by the input, but also by themselves. Many ANNs operate only in feed-forward fashion but some popular ANNs, such as Long short-term memory (LSTM), employ recursion.

Another popular ANN known as Deep Neural Networks (DNN) have proved very popular over the last few years. They get good press in applications such as image recognition and speech recognition. Deep Neural Networks are often formed from tens of layers of ANes with each layer containing many ANes. DNNs are also processed in a feed-forward manner with one layer being the inputs to the next layer. As mentioned [Kri12], these useful DNNs often require hundreds of thousands of ANes and within the network, each ANe can have hundreds, even thousands of feeder or pre-synaptic ANes. There have been implementations that use different number formats from double precision floating point to eight bit integers, but in all cases these useful ANNs require a significant amount of memory to store the connection weights (parameters).

Although the spiking neural network more closely models the behavior of real neurons, over the



(a) Rate-based Model Artificial Neural Network (with ReLu activation function)



(b) Spiking-based Model Artificial Neural Network

Figure 2.5 Layered Artificial Neural Networks

last 20 years there have been breakthroughs in the configuring of rate-based models especially with the introduction of the back-propagation algorithm and stochastic gradient descent. Along with the abundance of data now available in the form of voice, images etc. to "teach" these networks using back-propagation, most of the effective applications of artificial neural networks have employed these rate-based models.

This work does not address the training of these rate-based ANNs, the training is mostly performed offline. This work is addressing the inference of ANNs. During inference, the most computationally intensive operation is the multiply accumulate associated with the ANe activation, which can involve hundreds or thousands of multiply-accumulates. The ANe activation calculation for the rate-based ANe in figure 2.2b is shown in equation (2.1).

$$\text{ANe Activation} \quad A = f \left(\sum_{n=0}^{C_p} W_n \cdot A_n \right) \quad (2.1)$$

C_p is the number of pre-synaptic connections

W_p is the weight of a connection

A_p is the activation value of the pre-synaptic neuron

and $f(x)$ is the activation function such as ReLu

2.2 ANN Layers

In figures 2.2c and 2.5a, the ANN is shown to be constructed using layers of ANes. It has long been known that a single layer of ANes can be used linearly partition an n-dimensional input, as shown in figure 2.6a. However, if a more complex partition is required, this cannot be achieved using a single layer of ANes. A higher order classification, as shown in figure 2.6a can only be achieved using multiple layers of ANes. In addition, to ensure the multiple layers cannot be mathematically collapsed into a single layer, the activation function $f(x)$, as shown in figure 2.2b must be a non-linear function.

2.2.1 Deep Neural Networks

As mentioned, a single layer of neurons can be used as a linear classifier as long as the classes can be separated using a linear function. Even some simple cases cannot be linearly separated, an example often used is an exclusive-OR gate.

Even with a layered ANN, the final output comes from a single layer. To allow this final layer to linearly separate classes, the original input needs to be transformed into a space where the classes can be linearly separated. Deep Neural Network (DNN) are ANNs that incorporate many layers of ANes, often which are often up to tens of layers deep. The additional layers are incorporated to translate the

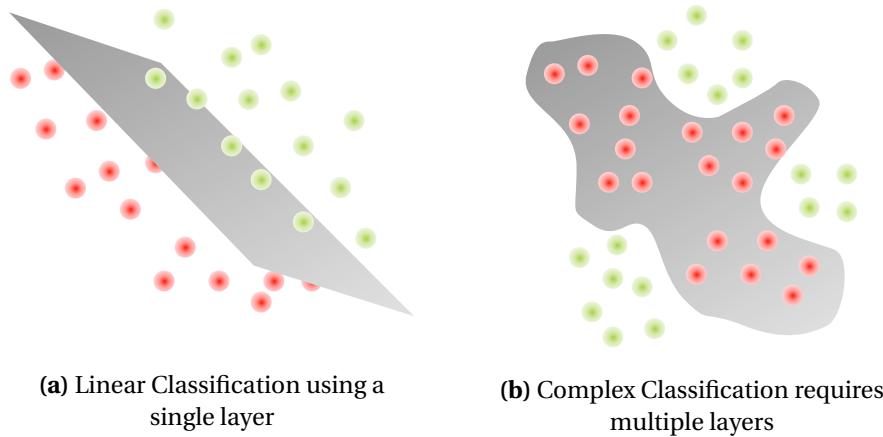


Figure 2.6 Classification using ANNs layers

space of the input so the various classes being identified can be separated using linear classifiers in the later layers.

A recent very successful example of a DNN, known as a Convolutional Neural Network (CNN) was used to classify objects in images [Kri12]. These CNNs use the early layers to identify low-level features and later layers are used to combine these features into yet more higher-level features. Finally, the combination of high-level features are used to identify the required classes. This layering is shown in Fig. 2.7.

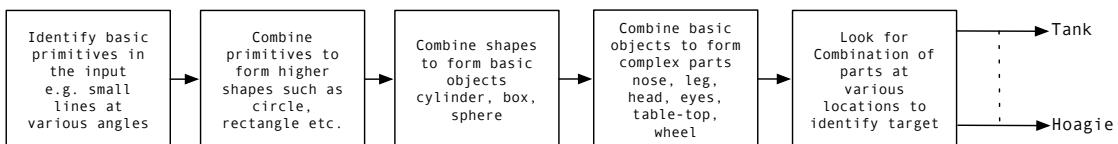


Figure 2.7 Deep network showing feature layers

In figure 2.7, the final layer is often a fully connected linear classifier with the output representing the probability of a particular class being present in the image. In practice, these DNNs can be used as classifiers or as function approximators.

2.2.2 Feature Layers

For the most part, different ANNs are characterized by how the ANEs are interconnected and the activation function employed.

The typical DNN layers are processed in a feed-forward fashion where the pre-synaptic ANEs are formed from ANEs in the previous layer. There are some types of DNN which also include recursive connections where the pre-synaptic ANEs include ANEs from the current layer. A popular recursive DNN is Long Short-term memory (LSTM). Although this work does not preclude supporting LSTM in the future, the focus of this work is on the feed-forward type DNN.

As described in 2.2.1, a DNN layer transforms the previous layer with each higher layer providing a coarser grained transformation. This is best seen in image recognition application were the early layers identify low level shapes or features, such as angled lines with higher layers identifying higher order shapes such as circles, blocks etc. Although the image recognition application is somewhat intuitive, it is believed that in less intuitive applications the DNN performs a similar fine to coarse feature extraction.

The connections between layers can be locally-connected or fully-connected. With locally-connected layers (see figure 2.8a), a particular layer ANEs pre-synaptic ANEs are formed from regions of the previous layer. With fully-connected layers (see figure 2.8b), a particular layer ANEs pre-synaptic ANEs are formed from all ANEs in the previous layer. In many cases, a DNN is constructed with lower layers being locally-connected and higher layers being fully connected.

With locally-connected layers, the pre-synaptic ANEs connection weights are often formed from identifying particular features. In early uses of locally-connected ANNs, the first layers weights were often hand-generated. With automatically trained ANNs, the feature detectors at each layer are often created during training. With image recognition ANNs, the lower level features generated during automatic training are often intuitive and are trained to detect small features such as lines at various angles, different curves etc.. In the general case, the trained feature detectors may not be as intuitive.

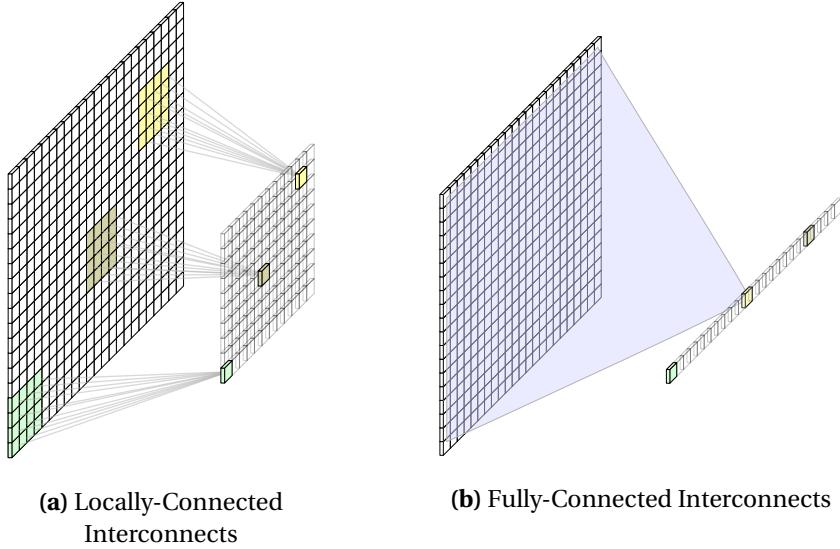


Figure 2.8 Layer Connection Types

Some contrived examples of locally-connected feature detectors are shown in figure 2.9a.

These locally-connected layers have multiple filters applied to the same region-of-interest in the previous layer. The next layer therefore becomes a 3D array with the Z-axis representing the features. The number of feature filters applied at each layer can be tens to hundreds of filters. The filters operating on these 3D layers are themselves 3D and therefore the number of weights associated with each filter is usually hundreds to thousands of weights.

The feature kernels employed in the locally-connected layers can be unique to the regions of the previous layer or the same filter kernels can be employed across the entire previous layer. In the case of employing the same kernels across the entire input layer the ANN is known as a Convolutional Neural Network (CNN). The CNNs are examples of ANNs which can take advantage of reuse described in 1. These CNNs can store the filter kernels in local SRAM and construct an entire feature plane. These CNNs can be considered a subset of the generic case of DNN where the kernels can be unique across the entire input layer. This work considers the more generic DNN case and supports acceleration of both DNNs and the subset CNNs.

An example of a CNN can be seen in figure 2.10b with the layer configurations shown in table 2.10a. An CNN similar to this has demonstrated high levels of efficacy in image recognition applications. **Therefore, this work considers this ANN to be a template for estimating the storage and processing requirements and the range of pre-synaptic fanins for a .**

2.3 The Problem

To approach the capabilities observed in human behavior, such as object recognition these ANNs become very large. They often utilize hundreds of thousands of neurons to implement what a human would consider a relatively straightforward task. For example, a "useful" ANN similar to that described in [Kri12] that is used to recognize up to 1000 different object classes, has a network size of approximately 650,000 neurons and 630 million synaptic connections [Kri].

The increased performance of ANNs over classical methods in image recognition and voice recognition might suggest that ANNs might out-perform techniques in other applications. There is reason to believe that ANNs will replace existing control and monitoring functions in existing systems.

If ANNs fulfill their potential, systems employing ANNs will utilize them for various functions, such as engine monitoring, anomaly detection, navigation etc. all within the same system. Considering the various functions a complex customer facing or edge application system performs, it is likely that many real-world applications will employ multiple disparate instances of these useful sized ANNs. Assuming these complex functions will require ANNs similar in size to [Kri12], these implementations will be processing multiple large ANNs at or near real-time.

Considering the storage required for the input, the ANe states and most significantly the weights for each of the ANes, the storage requirements results in gigabytes of memory. When these ANNs are required to be solved in fractions of a second, the processing and memory bandwidth becomes prohibitive.

2.3.1 ANN Processing

As a metric, this work assumes that any useful ANN will employ 100s of thousands of ANes. Although there is a lot of debate regarding number formats for ANNs, this work also assumes single-precision floating point. Assuming an ANN with 250K neurons and an average fanin to each ANe of 2000, a system employing 10 ANNs for various disparate functions and an average processing time of 10 ms suggests a average bandwidth of 16 Tbit/s (see equation 2.2).

$$\begin{aligned}
 \text{Average Bandwidth} &= \sum_{n=0}^{\bar{N}_n} \left(\frac{\bar{N}_a \cdot \bar{C}_p \cdot \bar{b}_w}{\bar{T}_p} \right) \\
 &= \sum_{n=0}^9 \left(\frac{250 \times 10^4 \cdot 2 \times 10^3 \cdot 32}{10 \times 10^{-3}} \right) \\
 &= 16 \text{ Tbit/s}
 \end{aligned} \tag{2.2}$$

where N_n is the number of ANNs

N_a is the average number of ANEs

C_p is the average number of connections

and T_p is the processing time

Therefore, when implementing ANNs, the memory requirements are significant. The storage is required for the input, the ANe state and most significantly the weights for each of the ANes pre-synaptic connections. This storage requirement often results in gigabytes of memory.

In addition, in edge applications, it is anticipated that these ANNs will be required to be solved in fractions of a second, in which case the processing and memory bandwidth becomes prohibitive.

The problem becomes “**to provide deterministic at or near real-time performance within tolerable power and space constraints for edge systems employing inference on multiple disparate useful-sized neural networks.**”

Considering that DRAM is required to store the ANN parameters, why is it that much of the ASIC and ASIP ANN research employs SRAM as an intermediate store? Well, in practice there are benefits if you can operate solely out of SRAM. Certainly good performance and potentially low power. But use of SRAM makes assumptions on the type of ANNs that can be supported and the application in which the ANN is being deployed. The primary requirement of the type of ANN and the deployed application to allow effective use of SRAM is "reuse". Reuse means that once parameters are transferred and stored in SRAM, these parameters can be reused such that the SRAM isn't simply an intermediate memory but something akin to a cache.

In some ANNs there are reuse opportunities. A prime example is CNNs, where the connection weights are reused. In CNNs, a common filter is passed across an input to form the next layer. These filter "kernels" can be held in memory and the input is read from DRAM thus reducing the DRAM bandwidth. Even with DNNs where weights may not be reused, when implementing multiple DNNs, there is opportunity to hold the input in memory. If the system is being employed in cloud applications or in training, there is opportunity to reuse inputs whilst performing batch processing.

But SRAM comes at a price, its big. Often when we see physical layouts of NN processors, they are dominated by the silicon area of the SRAM. The area required for SRAM has been understood for quite some time and companies attempt to create custom SRAMs to minimize the area impact.

So the question becomes, can a system employ DRAM with minimal SRAM and still provide a high performance system within acceptable area constraints?

Even in cloud applications, there are limitations on reuse. We paraphrase a quote from a Google paper [Aba15] on their Tensor Processing Unit ASIC (TPU):

"the architecture research community is paying attention to NNs, but of all the papers at ISCA 2016 on hardware accelerators for NNs, alas, all nine papers looked at CNNs, and only two mentioned other NNs. Unfortunately CNNs represent only about 5% of our datacenter NN workload"

The applications targeted by the google TPU [Aba15] assume multiple requests, so reuse in the form of batch processing is still of great benefit, but the bulk of the requests in [Aba15] are fully-

connected DNNs and in these cases weight reuse is not as beneficial and the performance of the TPU is degraded when implementing these fully-connected DNNs.

Therefore, implementations that focus on CNNs can suffer from severe degradation in performance when targeting generic types of ANN, such as locally and fully connected DNNs and LSTMs.

This work focuses on edge applications employing disparate ANNs and assumes both weight reuse and batch processing do not apply. Considering systems will want to perform multiple DNNs simultaneously suggests that these edge systems will require usable memory bandwidth of the order of 10s of Tbit/s.

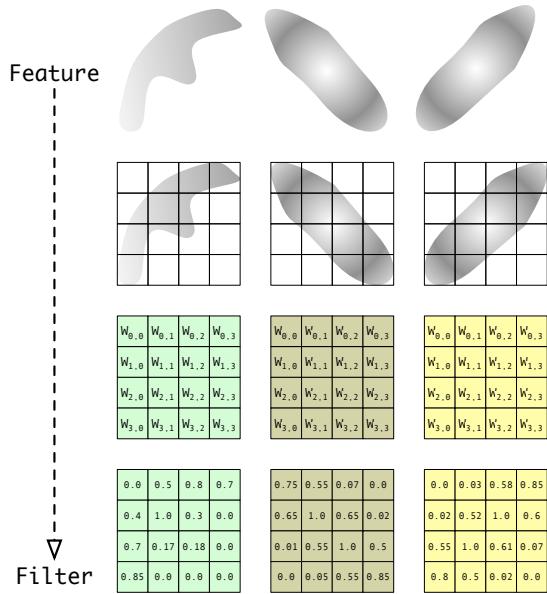
In these cases, **DRAM bandwidth is the bottleneck**.

Some might suggest the requirements of these applications would be satisfied by employing multiple graphics processor units(GPU). In fact, Graphics processing Units (GPU) are used to implement large ANNs and in some ANN architectures, such as Convolutional NNs (CNN), they are quite effective. However, we should not forget they are not optimized purely for ANN processing and are restricted by available SRAM and they are power hungry. These limitations will limit the effectiveness of GPUs regardless of what we might hear from the GPU community. Even in the case of newer GPUs which are employing 2.5DIC technology, the memory bandwidth will still be limited by available DRAM tecnology. For example, a 2.5D solution employing High bandwidth Memory (HBM) would be limited to a maximum raw bandwith of the order of 4 Tbit/s. Also, its has proven very difficult, if not impossible to take advantage of the available memory bandwidth [Far11] [Aba15]. Given these multiple GPU systems have high real-estate and power requirements and given each instance consumes of the order of 100 W to 200 W. Overall GPUs have limited suitability to meet edge application requirements.

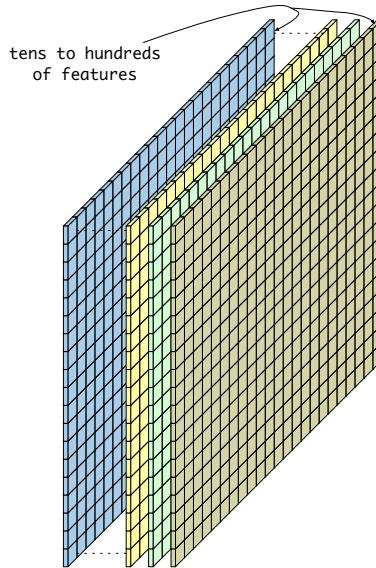
Much of the ANN application specific (ASIC/ASIP) research has focused on taking advantage of the performance and ease of use of Static Random Access Memory or SRAM. These implementations can be shown to be effective with specific ANN architectures (CNN), server applications or the "toy examples" but when a system requires multiple disparate ANNs in an edge application,

these implementations do not provide the required flexibility, storage capacity and deterministic performance.

How this work addresses the problem are outlined in section 5.



(a) Features and locally-connected filters (kernels)

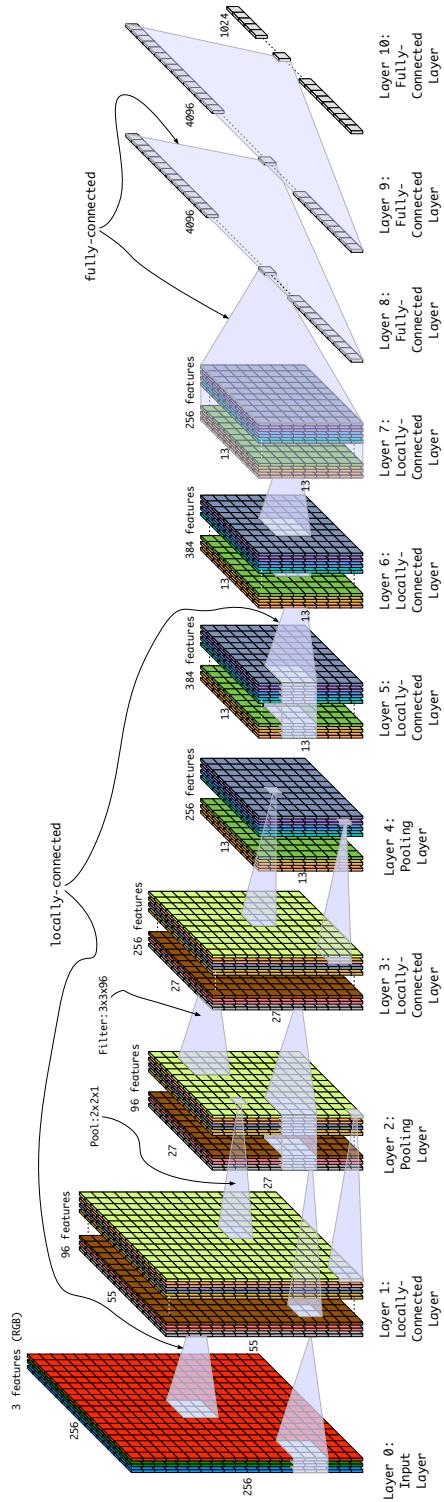


(b) 3D layers of Features

Figure 2.9 Layer Feature Planes

	Type	Input	2	3	4	5	6	7	8	9	10	11	Layers	
Dimensions	X	256	55	27	27	13	13	13	13	4096	4096	4096	Fully	
	Y	256	55	27	256	256	384	384	256	1	1	1	Fully	
	Z	3	96	96	2	5	2	3	3	13	13	13	Fully	
Filter Dimensions	X	na	11	2	5	2	3	3	3	13	13	13	Fully	
	Y	na	11	2	5	2	3	3	3	13	13	13	Fully	
	Z	na	3	1	96	1	256	384	384	256	1	1	Fully	
Stride	na	na	4	2	2	2	1	1	1	na	na	na	Fully	
Pre-synaptic Fanin	na	363	4	2400	4	2304	3456	3456	43264	4096	4096	4096	Total	
Number of ANe	196608	290400	69984	186624	43264	64896	64896	43254	4096	1024	1024	969152		
Number of Weights	na	34848	na	614400	na	884736	1327104	884736	177209344	16777216	4194304	2.02 × 10 ⁸		

(a) Layer Configuration [Kri12]



(b) Example DNN [Kri12]

Figure 2.10 Typical DNN [Kri12]

CHAPTER

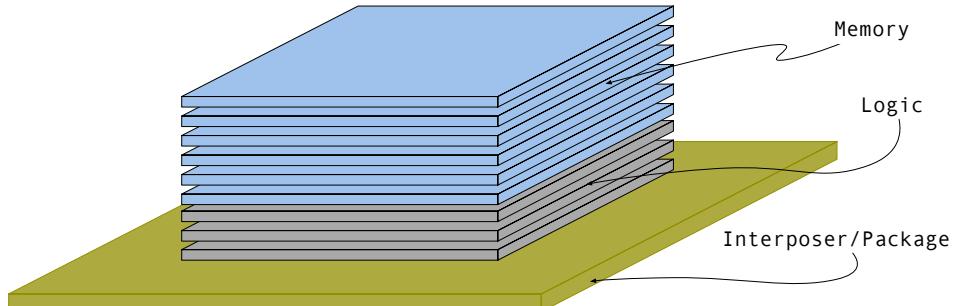
3

THREE-DIMENSIONAL INTEGRATED CIRCUITS (3DIC)

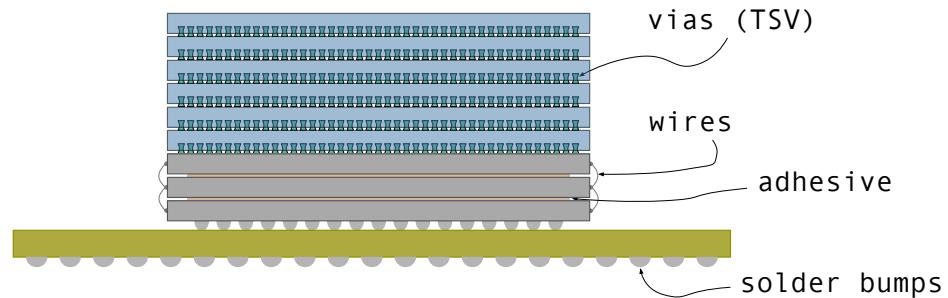
Over the last couple of decades, the ever shrinking world of Integrated Circuits (ICs) has enabled the introduction of devices society takes for granted such as personal computing and cell phones.

As IC technology has shrunk, design complexity has grown to take advantage where hundreds of millions of transistors are placed on a typical IC. These ICs have evolved from performing small functions to becoming systems on a chip.

However, at some point, an IC has to interface with another function and often that communication involves the moving of data to and from memory and the increase in complexity often drives a need for higher memory data bandwidth.



(a) Different Die types stacked mounted on an interposer/package substrate



(b) Connection Types

Figure 3.1 3DIC Stack of Die

The IC complexity has been doubling approximately every two years but the external interfaces are restricted by physical limitations. Within the System-on-Chip (SoC), the designer can take advantage of very wide interfaces, often thousands of bits wide to increase bandwidth, but when data is moved to off-chip memory, the widest buses are usually hundreds of bits wide.

One way to avoid this limitation is to employ 3DIC (see figure 3.1). The advantages of 3DIC are well understood. Reducing the amount of off-chip communication increases bandwidth and reduces power. The power reduction comes from not having to drive the relatively high capacitance inputs and outputs.

Taking advantage of 3DIC means stacking die on top of one another and making connections directly between the die. These connections can be in the form of wire made at the edges of the die or

using vias buried in the die itself.

Below is a summary the benefits of 3DIC :

- Reduced Power
 - mainly from not having to drive external outputs and receiving external inputs
- Increased Connectivity
 - maintaining very wide buses through the SoC increases bandwidth
- Ability to mix heterogeneous technology
 - Mixed Analog/Digital
 - Mixing memory technology and logic technology
- Increase density and mitigation against the slowing of Moore's Law
 - using the vertical domain to increase perceived transistor per mm^2
- Potentially lower costs by combining simpler die rather than build a large die
 - yield benefits from combining higher yield die
- Possibility of novel architectures [Kim16]

Some disadvantages of 3DIC are:

- Reliability
- Cost
 - being a relatively new technology it is still expensive
 - TSV technology is still unreliable

There is still some reluctance to fully embrace 3DIC but undoubtly the various barriers will be broken down.

The ability to mix heterogenous technology is of particular interest to this works target application because mixing technology targeted toward DRAM with cmos logic technology is a heterogenous mix of which this work takes advantage.

In [ITR15] there are four definitions of 3DIC interconnects:

- 3D-Wafer level package [ITR15]
 - In this case, different die are stacked and then connected using traditional bond bumps and/or bond wires at the periphery of the chip.
 - This technique provides better transistor density compared to traditional 2D-IC with improvements in interconnect density.
- 3D-Stacked SoC [ITR15]
 - In this case, different die are stacked and then connected using TSVs. The TSVs connect the dies to intermediate metal layers known as global metal layers. This allows the individual die to maintain a high level of functionality and thus is similar to connecting functional building blocks meaning the individual die are likely to be significant functional pieces of intellectual property (IP).
 - using TSVs provides a medium level of interconnect
- 3D-Stack IC [ITR15]
 - In this case, different die are stacked and then connected using TSVs. The TSVs connect the dies to intermediate higher metal layers known as global metal layers. This infers the individual die are not large functioning pieces of IP.
 - using TSVs provides a high level of interconnect

- 3D-Integrated Circuit [ITR15]
 - In this case there are not multiple dies. Instead the additional silicon layers are deposited on top of each other with the final 3DIC device having multiple layers of transistors
 - Local metal layers are used which along with TSVs provides a very high level of interconnect

A die stack with TSVs can be seen in Fig. 3.2.

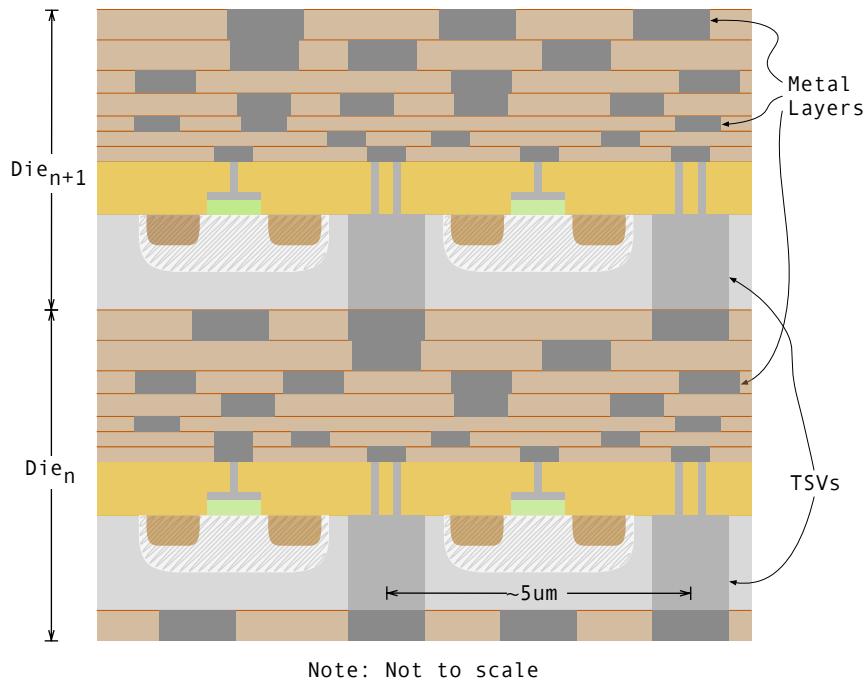


Figure 3.2 Die Stack profile [ITR15] with TSVs

There are other definitions on how the dies are bonded together:

- Wafer-to-Wafer
 - current Electrostatic discharge (ESD) mitigation allows implementation of unbuffered IO

- potential low yield because of lack of knowledge regarding Known Good Die (KGD)
- Die-to-Wafer
 - will need additional ESD mitigation support
 - higher yield because of knowledge of KGD
- Die-to-Die
 - will need additional ESD mitigation support
 - higher yield because of knowledge of KGD

This work is targeting 3DIC technology that supports 3D-Stacked SoC or 3D-Stack IC with high levels of interconnect. To avoid using large IO buffers for the TSV interconenct, this work assumes that the 3DIC technology supports unbuffered interconnects. This would suggest wafer-to-wafer bonding because of the existing ESD mitigation during wafer handling although it is anticipated that improved ESD mitigation will be introduced in future manufacturing steps.

The technology roadmap in [ITR15] and the information in [Pat14] suggests 5 μm pitch TSVs is a reasonable design goal. This work assumes a one-to-one ration of signal TSVs to power/grod TSV so when accounting for area associated with TSVs, the number of signal TSVs are doubled.

As a large amount of TSVs are employed, TSV energy cannot be ignored. Most of the energy dissipated in the TSV is associated with the charging and discharging of the TSVs capacitance. For a 5 μm pitch and 2 μm radius TSVs, [BGO17] table I suggests an average capacitance of 4.2 fF¹.

Based on (3.1) and assuming a supply voltage of 1.0 V, the power associated with a TSV is shown in equation (3.2).

$$\text{Energy to charge a TSV, } E_{tsv} = \frac{1}{2} \cdot C_{tsv} \cdot V^2 \quad (3.1)$$

¹[Tezb] suggests a lower capacitance

$$\text{Energy to charge a TSV, } E_{tsv} = \frac{1}{2} \cdot C_{tsv} \cdot V^2 = \frac{1}{2} \cdot 4.2 \times 10^{-15} \cdot 1.0 = 2.1 \text{ fJ}$$

$$\text{Power per TSV, } P_{tsv} = E_{tsv} \cdot \text{bit rate}$$

normalizing to a clock of 1.0 GHz

$$\text{Power per TSV per Hz} = 2.1 \mu\text{W}/\text{Gbit/s/TSV} \quad (3.2)$$

The TSV design guidelines used by this work are summarized in table 3.1.

Table 3.1 TSV Design Characteristics

Parameter	Dimensions		Power
	Pitch	Radius	
Value	5 μm	2 μm	2.1 μW/Gbit/s/TSV[BGO17]

CHAPTER

4

DRAM CUSTOMIZATIONS

There are two types of memory employed in ASICs and ASIPs, Static Random Access Memory or SRAM and Dynamic Random Access Memory, or DRAM. Both of these technologies have a similar top level block diagram which contains an array of storage elements, a means to address into a particular row of memory cells and a means to read and write a column of those cells. A basic block diagram is shown in figure 4.1.

Accessing a typical SRAM involves providing an address and either reading or writing the contents of that location. The read or write often completes in one or two clock cycles depending on whether the SRAM employs internal registers which are used run the SRAM with a faster clock.

The storage cell inside the SRAM is formed from cross-coupled transistors (see figure 4.2a) which latch the contents and hold the contents indefinitely or until power is removed from the device. The

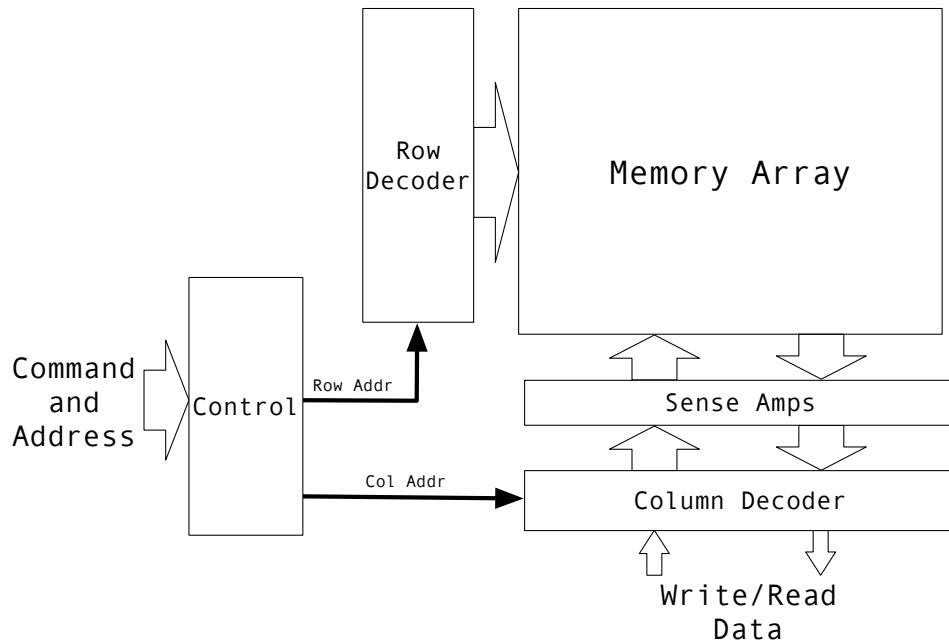


Figure 4.1 Typical Memory Block Diagram [Jac07]

storage structure employs six transistors and allows the access logic to be relatively simple and fast but has a relatively low XXXXXXXXXXXX capacity.

Accessing a "typical" DRAM is much more involved, it involves opening a page in a bank, reading or writing a portion of the contents of the page then closing the page.

The reasons behind this added complexity is the memory cell inside the DRAM which is formed from a capacitor (see figure 4.2b) which holds a charge reflecting a logic zero or one. The major difference between the SRAM cell is the number of elements it takes to implement. The size of the DRAM cell means DRAM arrays provide five to six times more storage when compared to a similar sized SRAM array. The major disadvantage is the capacitor cannot hold a indefinitely because the leakage currents in ICs cause the charge to leak away. If kept unchecked, the stored value will dissipate and it is this behavior that makes accessing a DRAM array more complicated than a similar SRAM array.

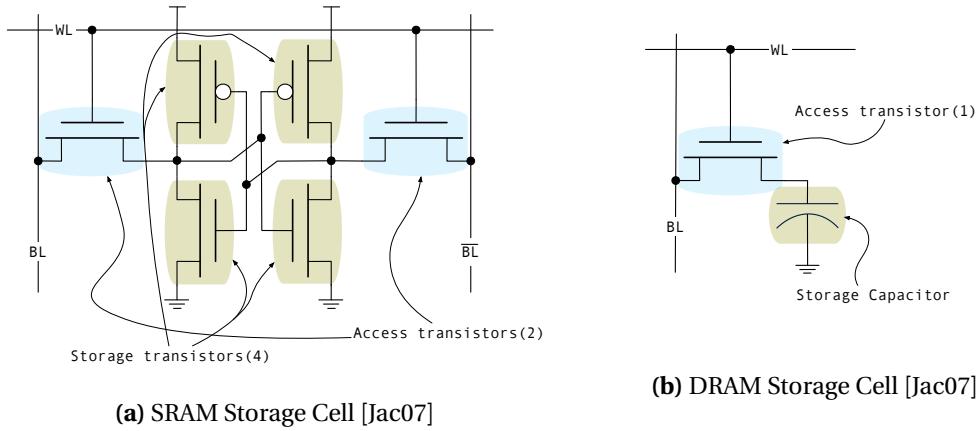


Figure 4.2 RAM Storage Cell Types

When an SRAM cell is read, the cross-coupled transistors retain the stored value. The issue with DRAMs is when a storage cell is read, the capacitor is discharged. In addition, sensing this discharging capacitor takes long time to sense when compared to the SRAM cell. To alleviate this problem, DRAM arrays are formed from a column of storage elements known as a page. When the read occurs, the entire contents of the page are transferred to registers. Once this transfer is complete, portions of the page can be read similar to reading an SRAM. The problem is that if another read wants to access data that is not in the page, the page has to be closed and another page opened. This involves transferring the previously registered page back to the array to recharge the capacitors in the pages storage elements. The next page can then be read and transferred to the page registers. The process of opening and closing pages is a relatively long time, typically 10-20ns. So if the accesses are somewhat random, accesses are very slow when compared to SRAM. To get around this issue, the DRAM is formed from more than one array of storage elements, between 8 and 32, known as banks. The idea is that while a page from one bank is being accessed, another bank can be opened in preparation for a future read (or write). This access protocol is rather complicated and if consecutive accesses are not sequenced carefully, the performance of DRAM can be poor.

The sequence of accesses cannot always be controlled, especially in general purpose computers,

so SRAM is typically used as the first level of memory with DRAM used as the primary storage. Using SRAM as this first level memory is called a cache and these have been used for decades to isolate the computing system from unpredictable access behavior of the DRAM.

As mentioned in section 1.1 much of the ASIC and ASIP ANN research has focused on taking advantage of the performance and ease of use of SRAM, but with this work's target application an ASIC or ASIP employing SRAM as its primary memory cannot be implemented with adequate storage capacity. Under these circumstances, DRAM bandwidth will be the bottleneck.

This work focuses on using DRAM as the primary storage and managing the accesses to ensure the DRAM is used effectively. With the additional DRAM customizations discussed in section 4.1 and 4.2, this work demonstrates DRAM bandwidths 10X faster than what is available with 2 or 2.5D solutions. This high level of DRAM bandwidth provides this work the ability to process multiple disparate ANNs at or near real-time whilst being **10-100X faster than state-of-the-art solutions**.

4.1 Customization One: Very-Wide Bus

Typically a bank may contain of the order of a few thousand pages and a page may contain of the order of a few thousand bits.

Once the page is open, the user accesses a portion of the requested page over a bus. With PCB based DRAMs the bus might vary from four to 16 bits wide, but with 3D DRAMs, such as HBM the bus might be up to 128 bits wide.

Figure 4.3 shows a block diagram of a typical DRAM.

This work achieves the increase in bandwidth by proposing that the DRAM expose more of its currently open page.

Without the limitations of having to transfer data beyond the chip stack, this work suggests exposing a larger portion of the page over a very wide bus. By staying within the 3D footprint, this bus can be implemented using fine pitch through-silicon-vias. (see figure 4.4).

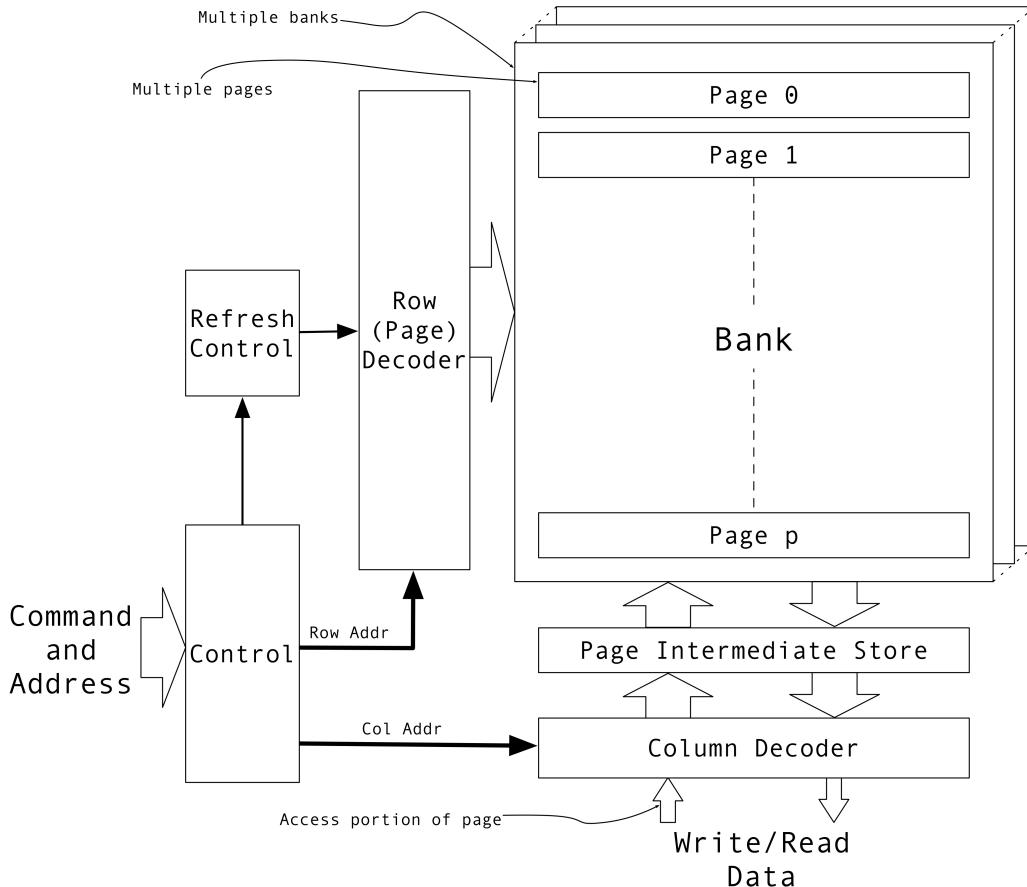


Figure 4.3 Typical DRAM Block Diagram

This work assumes the DRAM interface protocol uses Double Data Rate (DDR) with a bus width of 2048. Given the Dis-Integrated 3D DRAM (DiRAM4) employs a burst of two for read and write cycles, an entire DiRAM4 page of 4096-bits is accessed during each read or write.

4.2 Customization Two: Write Mask

When processing an ANN, to compute the activation of an individual AN involves reading the pre-synaptic AN activations and the weights of the connections between the pre-synaptic ANs and the AN

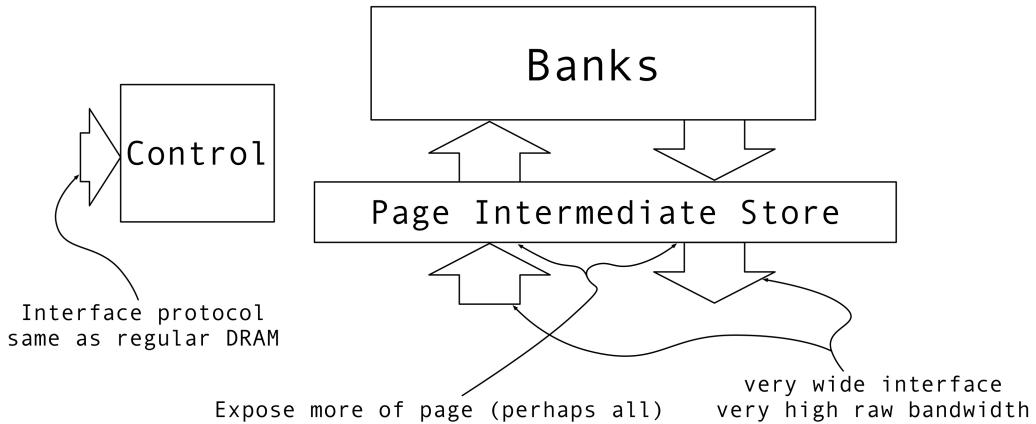


Figure 4.4 Exposing more of the DRAM page

being processed. The activation of the processed AN is written back to memory. The ratio of reads to writes is high, 100s or 1000s to one. Therefore, the system often needs to write a portion of the page back to memory. To avoid a read/modify/write, a customization to the DRAM is the addition of a write data mask to the DRAM write path.

This work assumes single precision floating point for ANN weights and activation, so a mask bit will be provided on a word basis or every 32-bits.

CHAPTER

5

SYSTEM DESCRIPTION

The primary design considerations that drove the architecture of this work are :

- Target applications are unable to take advantage of memory reuse opportunities and therefore not able to achieve high performance using local SRAM to store ANN parameters or the ANN input
- DRAM is required for storage of ANN parameters
- Target application will likely apply many disparate ANNs to perform various system functions
- Target application will have space and power limitations

When performing inference in ANN, the computational hotspot is the ANe activation calculation shown in figure 2.2b and equation (2.1). This ANe activation involves hundreds or thousands of

multiply-accumulates of ANe activations and corresponding connection weights. In this work, the ANe activations and weights are stored in DRAM. Therefore, because of the complex access protocol associated with DRAM, one of the main objectives is to demonstrate the 3DDRAM can be accessed with relatively high levels of bus efficiency.

The system has to process thousands of ANes concurrently and do this with minimal unused bus cycles. Therefore, the system must decode instructions, configure the various functions, pre-fetch and pipeline DRAM data and perform the actual activation calculation. To achieve high average DRAM bandwidth usage, these functions are all performed in parallel and this work demonstrates a average system bandwidth able to meet the tens of 16 Tbit/s of bandwidth requirement outlined in equation (2.2).

The problem associated with processing ANNs are outlined in section 2.3.

5.1 The Solution

Most researchers acknowledge that realistically, DRAM is required to meet the main storage requirements of useful sized ANNs.

We further believe that to support all types of disparate ANNs, we need to be able to operate directly from the DRAM memory.

This is because SRAM-based solutions assume memory locality when processing a neural network. However, when ANNs do not provide sufficient locality these solutions become DRAM bandwidth bound. If we then ensure the DRAM can feed the SRAM at the necessary bandwidth, why use an SRAM and waste the significant silicon area they require.

This works system operates directly from DRAM, but not just DRAM, 3D-DRAM.

In addition, this work has designed a system that can stay within the physical footprint of the 3D-DRAM.

By ensuring the system stays physically within the 3D stack, we take advantage of high density

connectivity provided by TSVs. Therefore, this work is able to propose a custom 3D-DRAM that exposes more of the DRAMs internal page and thus generates interface bandwidth that is of the order of 64 times that of the standard 3D-DRAM.

5.1.1 Novelty

The novelty of this work includes:

- A custom 3D-DRAM providing a 64X bandwidth benefit compared to standard 3D-DRAM
- A system that benefits from the power and performance benefits of 3DIC technology by remaining within a 3DIC stack
- New data structures that allow use to operate directly out of DRAM whilst ensuring effective use of DRAM bandwidth
- A system that can simultaneously process multiple disparate ANNs at or near real-time

This research explores a 3DIC solution using a custom organized 3DIC memory in conjunction with unique data structures and custom processing modules to significantly reduce the area and power footprint of an application that needs to support the processing associated with multiple ANNs. This works system will provide at or near real-time performance required for systems employing multiple disparate ANNs whilst staying within acceptable area and power limits and will provide greater than an order of magnitude benefit over comparable solutions.

5.2 System Overview

This work employs 3DIC technology along with a customized 3D-DRAM. The objective was to demonstrate that a pure 3DIC system can implement multiple disparate ANNs. By staying within the 3DIC footprint and taking advantage of high density TSVs this work is able to maintain a significantly higher bandwidth over 2D or 2.5D ASIC or ASIP solutions.

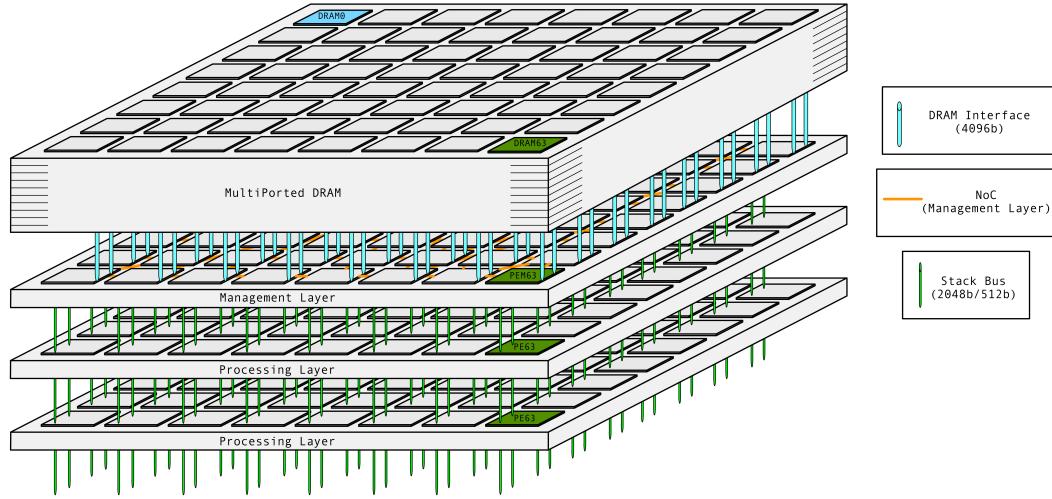


Figure 5.1 3DIC System Stack

The 3DIC system die stack (figure 5.1) includes the 3D-DRAM with a system manager below and one or more processing layers below the manager.

3D-DRAM has recently become available in standards such as High Bandwidth Memory (HBM) and Hybrid Memory Cube (HMC) and proprietary devices such as the DiRAM4 available from Tezzaron. These technologies provide high capacity within a small footprint.

In the case of HBM and DiRAM4, the technology can be combined with additional custom layers to provide a system solution.

The question becomes, can a useful system coexist within the same 3D footprint?

This work targeted a baseline system with:

- Computations requiring single precision floating point precision
- Tezzaron Dis-Integrated 3D DRAM (DiRAM4) [Teza] for main memory

The work includes customizing the interface to a 3D-DRAM, researching data structures to describe storage of ANN parameters, designing a memory manager with unique micro-coded instructions

and a Processing Engine (PE) layer. The system is designed such that a sub-system, known as a SSC operates on one of these disjoint memories within the 3D-DRAM (see figure 5.2).

When the SSCs need to share data or neuron activations, the data is passed between SSCs using a Network-on-Chip (NoC).

An overview of the various blocks and interconnects are given below:

5.2.1 3D-DRAM

The targetted 3D-DRAM, the Tezzaron DiRAM4 is a 3D-DRAM employs multiple memory array layers in conjunction with a control and IO layer. The memory is formed from 64 disjoint sub-memories each providing upwards of 1Gigabit with a total capacity of at least 64 gigabit.

5.2.2 Manager Layer

The Manager block is the main controller in the system. The operations required to process an ANN are formed from individual instructions which are decoded by the Manager. These instructions include descriptors to describe memory read operations, processing engine operations and memory write operations. The manager reads these system instructions from an instruction memory, decodes the instruction and configures the various blocks in the system. The configuration includes:

- initiate operand reads from DRAM
- prepare the processing engine (PE) to operate on the operands
- prepare the result processing engine to take the resulting neuron activations from the PE and write those results back to the DRAM
- replicate the resulting neuron activations to neighbor managers for processing of other ANN layers

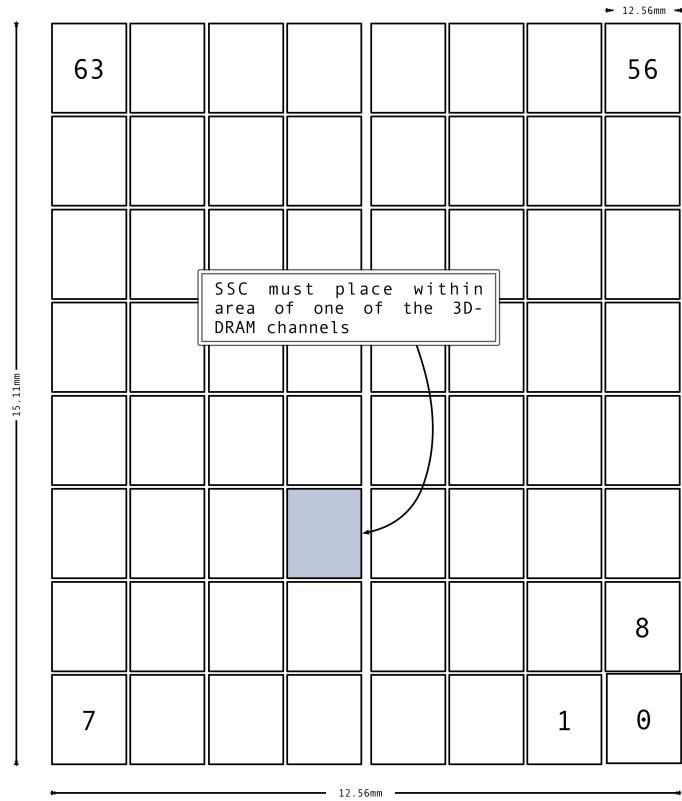


Figure 5.2 DRAM Physical Interface Layout showing area for SSC

5.2.3 Processing Layer

The PE is able to operate on data streamed directly from the DRAM via the Manager layer. The PE is configured by the manager to perform operations on the operand data streamed from the manager. In the baseline system, the main operation is to perform multiply-accumulates on 32 execution lanes of two operands. These operands typically are the pre-synaptic neuron activations and the connection weights. The PE also performs the activation function on the result of the MAC to generate the neuron activation value. These 32 activation values are sent back to the Manager layer.

5.2.4 Layer Interconnect

The layers are connected using through-silicon-vias (TSVs) which provide high connection density, high bandwidth and low energy. By ensuring the system stays within the 3D footprint ensures we can take advantage of the huge benefits provided by TSVs.

5.2.5 Inter-Manager Communication

During configuration and/or computations, data must be transported between managers. This inter-manager communication is provided by an NoC. When computing an ANN across multiple processing sub-systems, often neuron activation data must be shared between these SSCs. The SSC includes the DRAM port, the manager and the PE. An NoC within each management block communicates with each adjacent manager using a mesh network. This NoC has a forwarding table that can be reconfigured to provide more efficient routing for a given processing step.

A control and data flow diagram of the stack showing the 64 sub-system columns can be seen in figure 5.3.

5.3 System Operations

In the context of this system and AN state calculation, the basic operations to determine the state of a neuron is to :

- Inform the Manager and PE which operations are to be performed
- Tell the manager to access the states of the pre-synaptic neurons
- Tell the manager to access the weights of the connections from the pre-synaptic ANs
- Provide the pre-synaptic neuron weights and states to the processing engine execution lanes
- Tell the manager where to store the resulting AN state back to memory

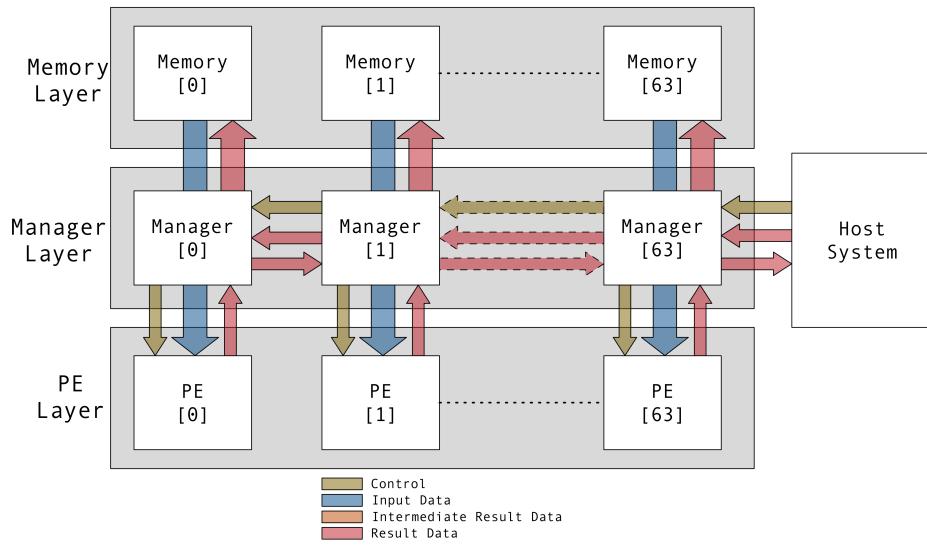


Figure 5.3 System Flow Diagram

This work has researched an instruction architecture to describe the above operations which are interpreted by the manager.

In the baseline system, the manager is not responsible for performing specific algorithm operations but is responsible for coordinating the various data flows and configuration of the modules that make up the system.

The manager's primary responsibility is:

- Instruction decode
- Internal Configuration messages
- Operand read
- Result write

In the baseline system, the PE is responsible for the main algorithm operations.

The PE has three major blocks:

- Streaming operation function (StOp)
 - Processes data from the manager on-the-fly without storing in local SRAM
- Single-Instruction Multiple-Data (SIMD)
 - processes the data from the StOp function, usually neuron activation such as ReLu
- DMA/local memory controller
 - transfer configuration data to PE controller or to store StOp results to a small local SRAM which can be used for access by SIMD or by the StOp function

5.3.1 Manager Operations

5.3.1.1 Instructions

The instructions communicate:

- To the Manager
 - ROI Storage descriptor
 - Parameter/Weight Storage Descriptor
 - Broadcast or Vectored
 - Result write storage descriptor
 - include descriptors for all destination managers
- To the PE
 - StOp operation
 - SIMD operation
 - Number of active lanes

Instruction (4-tuple example)

Operation Descriptor	arg0 Read Descriptor	arg1 Read Descriptor	Result Write Descriptor
----------------------	----------------------	----------------------	-------------------------

Figure 5.4 Instruction 4-tuple

- Operand Vector length

Instructions include information to control the above operations.

Instructions contain sub-instruction called descriptors. These descriptors contain the information to control the various operations associated with the processing of a group of ANs.

The group size is related to the number of execution lanes which for the baseline system is 32. So a group can be anywhere from 1-32. It should be said that unless the group size consistently approaches 32 the system performance will be poor.

An instruction will typically have four descriptors:

1. Operation
2. Memory read for operand stream 0
3. Memory read for operand stream 1
4. Result write

Note: An operand stream will be referred to as an argument.

The instruction is actually an n-tuple where the tuple elements are descriptors and the number of elements can vary based on the operation being performed. In figure 5.4 we see the format of a 4-tuple instruction which we use to perform an activation calculation for a group of neurons.

Now within a descriptor, we need to describe the various options such as storage descriptor pointer, number of operands etc..

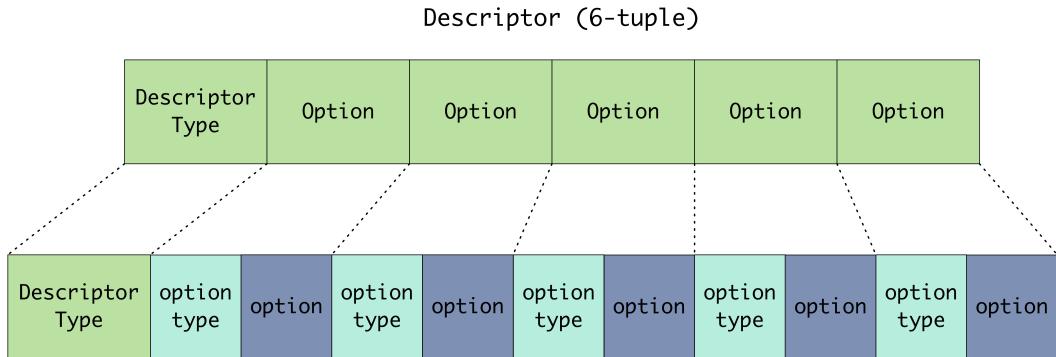


Figure 5.5 Descriptor 6-tuple

Again, we employ a n-tuple format where the first tuple element describes the descriptors operation followed by an m-tuple whose elements contain the options required for the operation.

These option elements are a two-tuple with option and associated value. In figure 5.5 we see the format of a 6-tuple descriptor.

5.3.1.2 Accessing of Pre-synaptic AN states and connection weights

As was discussed previously, the ANN input and configuration is stored in main DRAM memory. A part of the research is determining how to store the ANN input and parameters in such a way to effectively make use of main DRAM bandwidth. To provide parameters for the up to 32 execution lanes within the PE, we store the AN parameters in consecutive address locations. With one read to the DRAM, we access 128 words. This provides four weights for each of the 32 ANs being processed. These weights are sent to each lane of the PE over four cycles. We will discuss memory efficiency later, but by taking advantage of the multiple DRAM banks along with pre-fetching and buffering, we are able to achieve relatively high efficiency of the available maximum bandwidth.

Although AN parameters (weights) are stored in contiguous memory locations, providing the input state to a particular AN presents us with an interesting problem.

Most often DNNs are represented by layers of ANs whose pre-synaptic neurons are from the

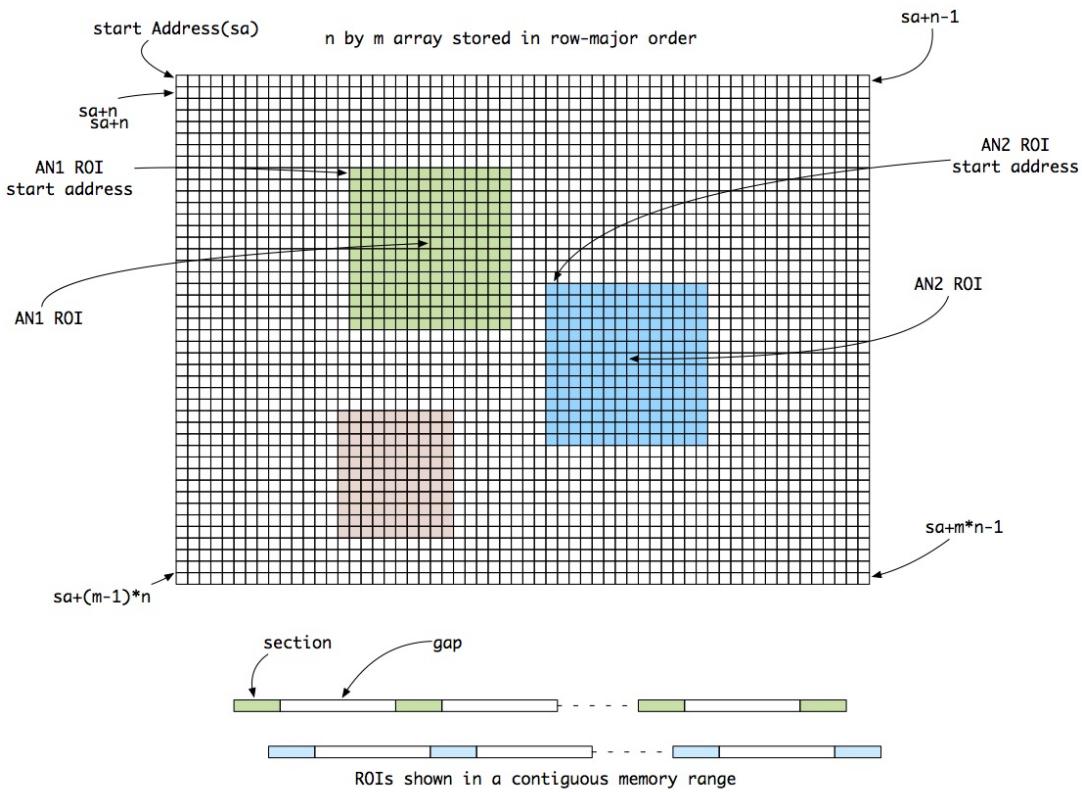


Figure 5.6 ROI Storage

previous layer. These previous layers represent the input to a given layer. The first layers input is the actual input to the ANN.

The input can be represented in the form of a 2-D array of AN states. For the sake of generality, the input array elements are considered as AN states.

Any given AN operates on a region of interest (ROI) within the input array.

In figure 5.6, an input to a ANN layer in the form of a 2-D array along with the ROI of two ANs.

The various connection weights are stored in multiple contiguous sections. However, it's not possible to arrange the input in such a way that each ANs ROI can be stored in contiguous memory locations. The figure above shows a typical ROI arrangement. Assuming the input array is stored in

row-major order, an ROI is drawn from disjoint sections of memory. These disjoint sections contain a number of AN states and the sections are separated by a gap of a number of memory addresses. When the parameters are accessed when performing a particular operation, the memory controller within the manager must be informed of the start address and the lengths of the sections and gaps. Now this looks problematic, and it is, but in practice groups of ANs share a common ROI. So once we solve the problem of efficiently reading an ROI from the DRAM, that ROI can be shared across a group of ANs

The read efficiency problem is solved by again taking advantage of the DRAMs banks and pages. This work proposes a data structure to describe these ROI storage locations.

Although disparate groups of ANs may have a different start addresses for their ROI, a commonality is observed in the ROI section lengths and gaps. So for each AN group, the groups ROI starting address is stored along with a pointer to a common set of section length/gaps. This structure is termed a storage descriptor.

This storage descriptor contains, amongst other things the start address of the ROI and a pointer to a section/gap descriptor. Many storage descriptors point to a common section/gap descriptor. This avoids having to have a unique section/gap descriptors for each AN group.

Figure 5.7 shows the structure of the storage descriptor. The SOD, MOD and EOD are used to delineate each descriptor in memory and stand for start-of-descriptor, middle-of-descriptor and end-of-descriptor.

5.3.1.3 Writing AN state results to memory

When the PE has processed the group of ANs, the new AN states are sent back to the manager. The manager will store these back to DRAM most likely in the array format as described earlier.

A significant difference taken advantage of is that for any given operation, the system is writing far less than is being read. For example, the ROI and parameters are usually vectors that will typically exceed 100 elements and in many cases much higher. When an operation is complete, in almost all cases one word per lane is written back to main memory. Now that sounds like writing back has a very

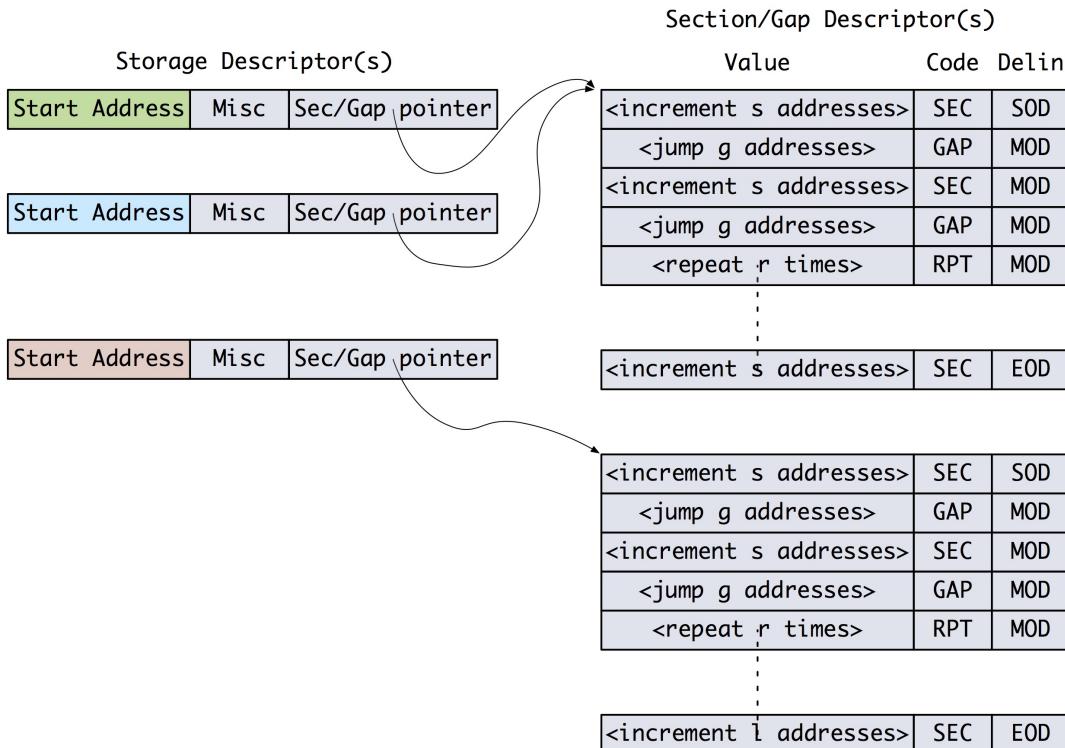


Figure 5.7 Storage Descriptor

small impact on performance but with DRAMs that's not always true.

When the system writes the result of an operation back to memory, it is often writing a small portion of a DRAM page and the nature of the DRAM protocol means this is a very inefficient use of DRAM bandwidth. So although the amount of data written is small the performance impact cannot be ignored.

In addition, in many cases the results from a particular PE has to be provided not only to the PE's local manager but also to other managers. This is handled with a network-on-chip (NoC).

The result storage directives are communicated by using the same storage descriptor mechanism. However, the added complication is because the result will likely have to be replicated to other managers, the storage descriptors must be sent to all destination managers.

5.3.2 PE Operations

5.3.2.1 Streaming Operations (StOp)

The operations performed by the StOp are primarily multiple-accumulate with a transfer to the SIMD or to local memory.

Even though the baseline system focuses on the AN multiply-accumulate followed by a ReLu activation function, the system has built in flexibility into the StOp function to allow other functions to be added

In most cases, the StOp module will operate on the AN state and weights provided by the manager and provide the result to the SIMD.

5.3.2.2 SIMD

The SIMD is a 32-lane processor with some builtin special functions, such as the ReLu operation.

The SIMD will take the result provided by the StOp and perform a ReLu. The result will, in most cases, then transmitted back to the manager.

5.3.2.3 Configuration

To configure these operations, two pointers are sent to the PE. These pointers index into a small local memory which provides a program counter (Program Counter (PC)) to the function to be performed by the SIMD and a configuration entry for the operation to be performed by the StOp.

The PE is able to perform its operation concurrently on 32-lanes. However, there are cases when less than 32-lanes will be employed. This may occur if the number of ANs being processed is not modulo-32. In this case, the manager provides the number of lanes being processed for any given operation. In addition, the length of the vector of operands is also sent by the manager to the PE.

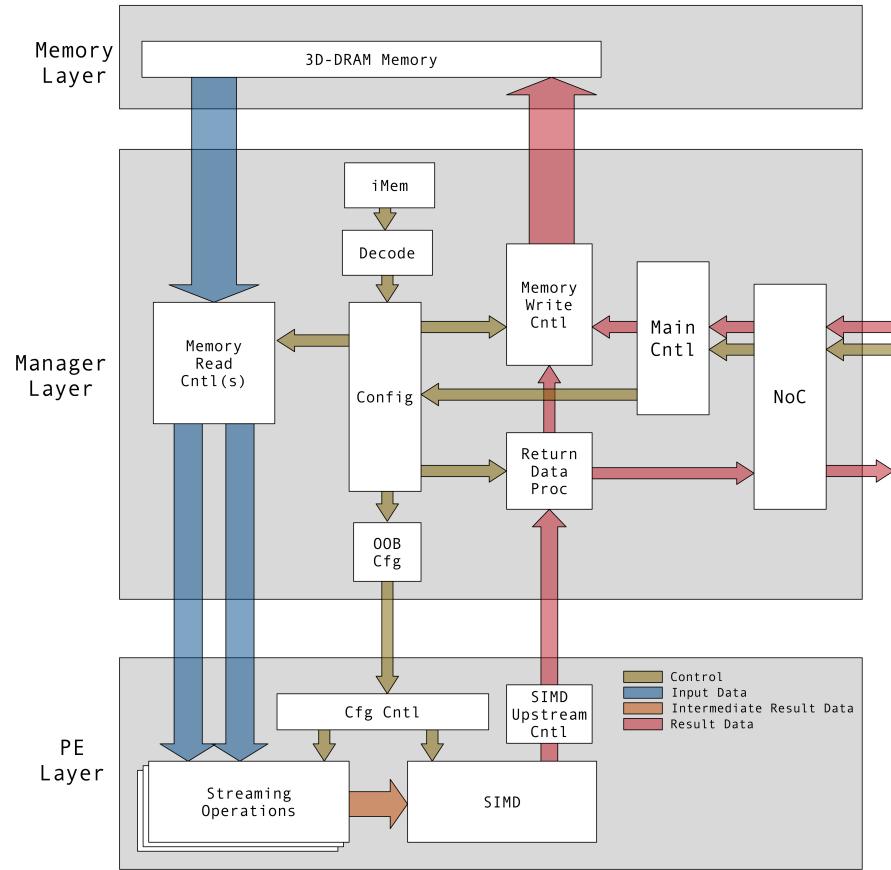


Figure 5.8 Sub-System Column (SSC) Flow Diagram

5.4 Detailed System Description

A detailed flow diagram and block diagram of the sub-system column can be seen in figures 5.8 and 5.9 respectively.

5.4.1 Manager

5.4.1.1 Operation Decode

In figure 5.9, instructions are read from instruction memory and passed to the instruction decoder.

The operation tuple is decoded and a streaming operation (stOp) pointer and a SIMD operation pointer are sent to the PE inside an OOB control packet.

The stOp pointer specifies what streaming operation is to take place on the data directly streamed to the PE. In the baseline system, typically this would be a floating-point multiply accumulate on two arguments, the pre-synaptic neuron states and the pre-synaptic weights.

The SIMD pointer is essentially a program counter that will be invoked when the stOp result is passed to the SIMD.

Note that other types of stOp includes a NOP with a destination of local memory. This allows us to transfer block of instruction or data from the manager to the PE.

5.4.1.2 Argument Decode

The instruction also includes argument descriptors. These descriptors include a storage descriptor pointers that point to a storage descriptor stored in local memory that encodes where data should be read from for the one or two arguments that will be streamed from DRAM to the stOp within the PE. In the case of a AN activation calculation, there are two arguments, the pre-synaptic neuron states and the pre-synaptic weights. The read storage descriptor pointers are passed to the Memory Read Controllers (MRC). The MRCs read the actual storage descriptor from their local memory and immediately start sending read commands to the memory via a Main Memory Controller (MMC). The MMC is not shown in the diagram but essentially takes the memory read requests and converts them into the DRAM read protocol.

As soon as read data is sent back to the MRC via the MMC, that data is aligned with the downstream bus and sent to the 32 Streaming Operations inside the PE.

5.4.1.3 Result data Processing

The instruction also includes argument descriptors. These descriptors include a storage descriptor pointers that point to a storage descriptor stored in local memory that encodes where data should

be read from for the one or two arguments that will be streamed from DRAM to the stOp within the PE. In the case of a AN activation calculation, there are two arguments, the pre-synaptic neuron states and the pre-synaptic weights. The read storage descriptor pointers are passed to the Memory Read Controllers (MRC). The MRCs read the actual storage descriptor from their local memory and immediately start sending read commands to the memory via a Main Memory Controller (MMC). The MMC is not shown in the diagram but essentially takes the memory read requests and converts them into the DRAM read protocol.

As soon as read data is sent back to the MRC via the MMC, that data is aligned with the downstream bus and sent to the 32 Streaming Operations inside the PE.

5.4.1.4 Memory Write Controller

The Memory Write Controller (MWC) receives data from two sources, the NoC via the MCNTL and the RDP.

In both cases, the MWC reads the actual storage descriptor from their local memory and immediately starts forming data that will be written back to main memory.

When the data is formed, a write command is sent to the memory via the MMC. Again, the MMC is not shown in the diagram but takes the memory write requests along with the data and converts them into the DRAM write protocol.

The MWC can only operate on one of the two sources at any one time. However, there are four 4096-bit holding registers where data is formed prior to the write request.

The holding registers have the potential in future to allow aggregation of data from one or more operations to allow a coalesced write back to main memory.

5.4.2 Processing Engine

5.4.2.1 Configuration

A configuration controller within the PE (PE_CNTL) takes the OOB packet from the Manager and extracts the stOp and SIMD operation pointers.

The stOp pointer is used to point to a local stOp configuration memory. The memory contains the various configuration data required by the streaming operation controller (stOp_CNTL). The stOp_CNTL is not shown.

The stOp_CNTL configures the:

- Operation type
- Number of active execution lanes
- Source of the argument data, which can be downstream data from the manager or from the small local SRAM
- Destination of the result data, which can be the SIMD or the small local SRAM

The SIMD operation pointer is sent to the SIMD.

5.4.2.2 Streaming Operations

The streaming Operations (stOp) are designed to operate on data passed from the Manager at or near line-rate. If line-rate cannot be maintained, a flow-control mechanism is employed to slow the data from the Manager.

Once the stOp has processed the data, it passes the result to the SIMD. Note in some cases the result can be placed in local SRAM or sent to both SIMD and SRAM.

It should also be stated that while the stOp is processing the current data, the SIMD may be operating on the result of the previous operation. It is expected the SIMD will have completed the

previous operation before the stOp completes the current operation, but again, if necessary a flow control mechanism between SIMD and stOP will be engaged if the SIMD is not ready.

5.4.2.3 SIMD

The SIMD takes the result data and performs the operation starting at the program counter (PC) indicated by the SIMD operation pointer provided by the PE_CNTL.

The stOp provides the result to the SIMD via a local register. The result is also written, in most cases to the small local SRAM.

The SIMD performs the specified operation on the data provided by the stOp.

In most cases this will be the AN activation function and in the baseline system is the Rectified Linear function (ReLU).

When the SIMD has completed its operation, it passes the result to the SIMD Upstream controller to be returned to the Manager.

5.4.2.4 Result Data

The SIMD Upstream Controller (SUI) takes the data and encapsulates it in an Upstream packet. Included in the packet is the tag required by the Return Data processor within the Manager.

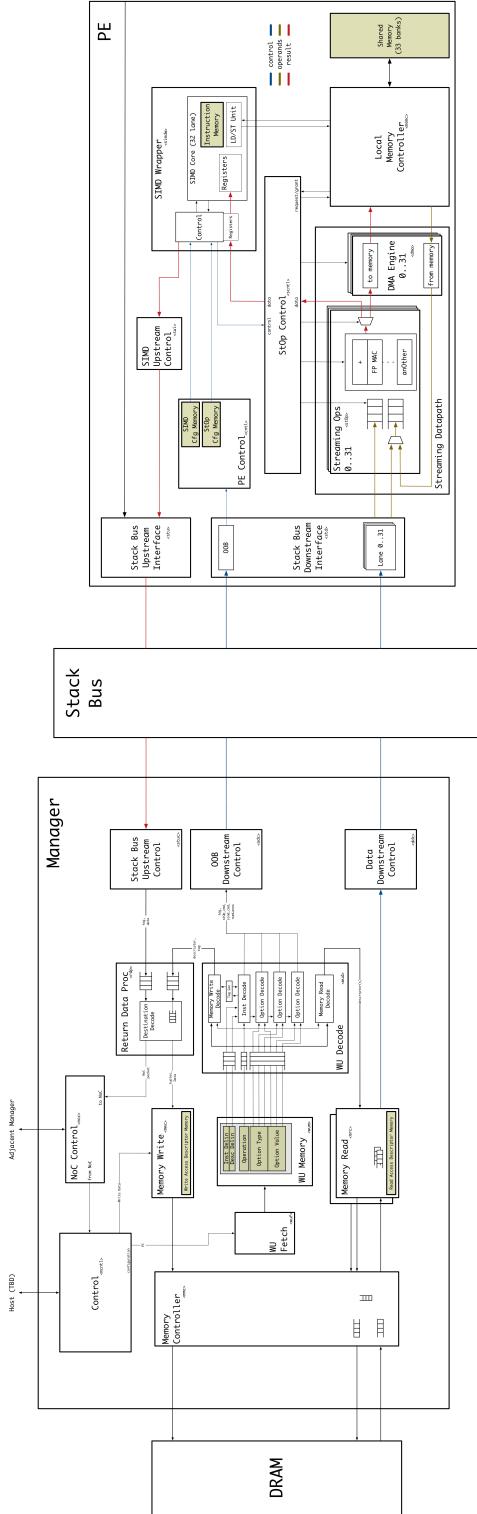


Figure 5.9 Sub-System Column (SSC) Block Diagram

CHAPTER

6

RESULTS

The objectives of this work was to design a system able to accelerate ANNs in customer facing systems implemented at the edge. This means systems that are not designed to process multiple requests of essentially the same operation. One assumption is that these systems implement disparate ANNs performing various functions. These assumptions imply that the system is not able to take advantage of local SRAM when processing the ANN. Another assumption is that the target systems will be space and power constrained. Finally, this work assumes that implementing multiple disparate ANNs cannot be implemented purely with SRAM and that DRAM is required to store the ANN parameters.

To demonstrate such a system, this work targeted 3DIC technology including 3D-DRAM. This work proposes that if a system can be purely 3DIC, the system can take advantage of the benefits of 3DIC which includes reduced energy, area and additional bandwidth due to high levels of connectivity.

In addition, this work proposes that if the system is purely 3DIC, then a customized DRAM would provide a significant bandwidth boost over typical implementations using standard DRAM. This work targeted the Tezzaron DiRAM4 3D-DRAM.

To ensure the system was purely 3DIC, the area of the system Manager and Processing Engine has to stay within the physical footprint of the DRAM.

The target technology node was 28nm. This was chosen because its the technology node employed for some recent GPUs and other ASICs such as [Jou17]. The design was synthesized using an available 65nm technology node and then scaled to 28nm to demonstrate fitting within the 3DIC footprint.

The primary control and datapaths of the system have been simulated in a system verilog environment. It has been synthesized using a 65nm technology node. The design has been coded targeting a frequency of >500 MHz. Timing closure and place and route is ongoing. The system has been designed throughout to meet the timing target so minimal changes are expected to meet timing.

As mention previously (2.2), to process multiple useful sized ANNs requires a sustained bandwidth to the PE of the order of tens of Tbit/s.

The Manager and PE have been place and routed as shown in figure 6.1. Although the design is yet to close timing, the parasitics were extracted from these layouts and simulated against a group of operations. The operations simulated were based on the expected lower and upper limits of pre-synaptic fanin. These were based on layers similar to CONV2 and FC-7 from [Kri12] and represent a pre-synaptic fanin of 225 and 4000 respectively. The simulation generated an activity file which was then used by the Synopsys® Primetime-PX™ power analysis tool to obtain power and bandwidth estimates. The DRAM accesses were captured and energy used calculated from [Teza]. The power dissipated in the TSVs were estimated from data from [BGO17]. These estimates were used to estimate power dissipation for operating frequencies of 500 MHz and 700 MHz which are shown in table 6.1.

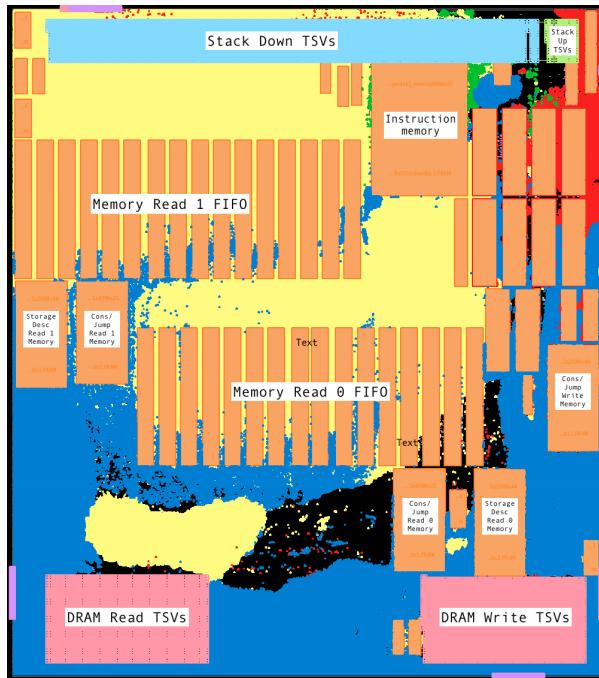
Table 6.1 Simulation-based estimates

Technology Node	Clock Frequency	Expected Power	Bus Efficiency
28nm	500 MHz	53W	~ 70 %
28nm	700 MHz	73W	~ 70 %

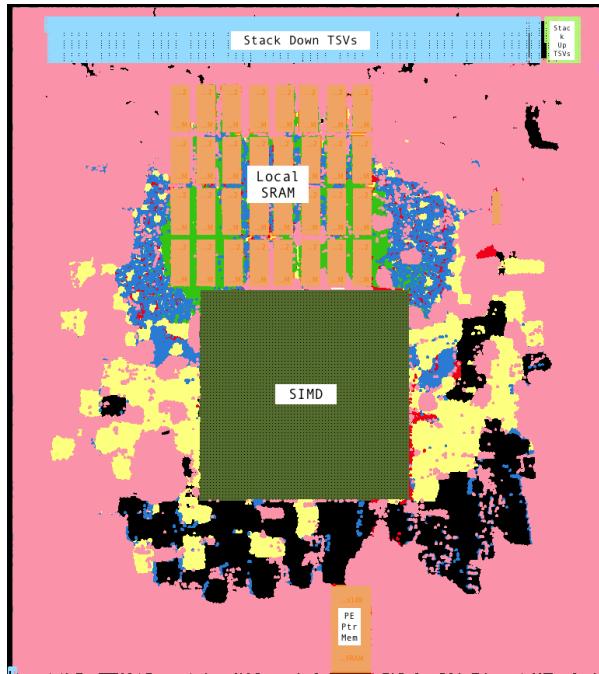
(a) Power Dissipation

Test	Downstream Stack Bus Efficiency	Average Downstream Bandwidth	
		500 MHz	700 MHz
CONV2 [Kri12]	~ 65 %	~ 43 Tbit/s	~ 60 Tbit/s
CONV-294	~ 67 %	~ 44 Tbit/s	~ 61 Tbit/s
CONV-300	~ 73 %	~ 48 Tbit/s	~ 67 Tbit/s
CONV-500	~ 82 %	~ 54 Tbit/s	~ 75 Tbit/s
CONV-1000	~ 89 %	~ 58 Tbit/s	~ 82 Tbit/s
FC-350	~ 78 %	~ 51 Tbit/s	~ 72 Tbit/s
FC-500	~ 84 %	~ 55 Tbit/s	~ 77 Tbit/s
FC-1000	~ 91 %	~ 60 Tbit/s	~ 83 Tbit/s
FC-7 [Kri12]	~ 94 %	~ 62 Tbit/s	~ 86 Tbit/s

(b) Bus Efficiency



(a) Manager



(b) PE

Figure 6.1 Manager and PE Die layouts

BIBLIOGRAPHY

- [Aba15] Abadi, M. et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. 2015.
- [BGO17] Bamberg, L. & Garcia-Ortiz, A. “High-Level Energy Estimation for Submicrometric TSV Arrays”. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* **25**.10 (2017), pp. 2856–2866.
- [CH06] Carnevale, N. & Hines, M. *The NEURON Book*. Cambridge University Press, 2006.
- [Teza] *DiRAM4-64Cxx Cached Memory Subsystem*. Rev. 0.04. Tezzaron Semiconductor. 2015.
- [Tezb] *Evolving 2.5D and 3D Integration*. Tezzaron Semiconductor.
- [Far11] Farabet, C. et al. “Neuflow: A runtime reconfigurable dataflow processor for vision”. *Cvpr 2011 Workshops*. IEEE. 2011, pp. 109–116.
- [ITR15] ITRS. *International Technology Roadmap for Semiconductors 2.0*. 2015.
- [Izh04] Izhikevich, E. M. “Which model to use for cortical spiking neurons?” *IEEE Transactions on Neural Networks* **15**.5 (2004), pp. 1063–1070.
- [Jac07] Jacob, B. et al. *Memory Systems: Cache, DRAM, Disk*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2007.
- [Jou17] Jouppi, N. P. et al. “In-datacenter performance analysis of a tensor processing unit”. *arXiv preprint arXiv:1704.04760* (2017).
- [Kim16] Kim, S. W. et al. “Ultra-Fine Pitch 3D Integration Using Face-to-Face Hybrid Wafer Bonding Combined with a Via-Middle Through-Silicon-Via Process”. *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*. 2016, pp. 1179–1185.
- [Kri] Krizhevsky, A. et al. *ImageNet Classification with Deep Convolutional Neural Networks*. <http://image-net.org/challenges/LSVRC/2012/supvision.pdf>. Accessed: 2016-08-30.
- [Kri12] Krizhevsky, A. et al. “Imagenet classification with deep convolutional neural networks”. *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [Mad14] Maddison, C. J. et al. “Move evaluation in go using deep convolutional neural networks”. *arXiv preprint arXiv:1412.6564* (2014).
- [Pat14] Patti, R. “2.5 D and 3D Integration Technology Update”. *Additional Papers and Presentations 2014*.DPC (2014), pp. 1–35.

- [Qiu13] Qiu, Q. et al. “A parallel neuromorphic text recognition system and its implementation on a heterogeneous high-performance computing cluster”. *IEEE Transactions on Computers* **62**.5 (2013), pp. 886–899.

APPENDICES