

UNIVERSITY NAME

DOCTORAL THESIS

---

# Thesis Title

---

*Author:*  
John SMITH

*Supervisor:*  
Dr. James SMITH

*A thesis submitted in fulfillment of the requirements  
for the degree of Doctor of Philosophy  
in the*

Research Group Name  
Department or School Name

May 10, 2023



## Declaration of Authorship

I, John SMITH, declare that this thesis titled, "Thesis Title" and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

---

Date:

---



*“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”*

Dave Barry



UNIVERSITY NAME

*Abstract*Faculty Name  
Department or School Name

Doctor of Philosophy

**Thesis Title**

by John SMITH

The landscape of deep learning compiler frameworks has evolved rapidly with the development of various tools, such as TVM, deeptools, TensorFlow, DLVM, nGraph, and Glow. These frameworks offer unique optimizations to address computation and data movement challenges in deep learning accelerators (DLAs). These approaches include graph or IR level optimizations related to intra node memory access optimizations, operator fusion, and various tiling techniques. Despite their unique approaches, these frameworks primarily concentrate on node level optimizations that focus on increasing the performance of executing a scheduled operation in the graph and overlook the potential for inter-node data reuse optimizations within on-chip memory resources.

OnSRAM, a scratchpad management framework build to work with deep learning compilers, addresses this gap by focusing on internode scratchpad management in DLAs. OnSRAM exploits the static graph representations of deep learning models by identifying data structures that can be pinned to on-chip memory based on their reuse rate and cost of transfer from main host memory. OnSRAM has been implemented and evaluated on single DLA that contains a monolithic scratchpad and is integrated as part of a custom deep learning compiler framework.

In this work, we extend the capabilities of OnSRAM by introducing dynamic scratchpad allocation for static graph execution models using any number of scratchpads and allowing the implementation to be interfaced as a compiler extension rather than a core part of a bespoke compiler. This enhancement allows for more fine-grained control over on-chip memory resources, providing increased flexibility and adaptability to better accommodate diverse deep learning accelerators and memory access patterns. By optimizing inter-node data movement and storage across multiple scratchpads, our approach further reduces energy consumption and latency associated with inter-node communication.





## *Acknowledgements*

The acknowledgments and the people to thank go here, don't forget to include your project advisor...



# Contents

<b>Declaration of Authorship</b>	<b>iii</b>
<b>Abstract</b>	<b>vii</b>
<b>Acknowledgements</b>	<b>ix</b>
<b>1 Deep Learning</b>	<b>1</b>
1.1 Deep Learning Architectures . . . . .	1
1.1.1 Perceptron . . . . .	1
1.1.2 CNN . . . . .	1
1.1.3 RNN . . . . .	1
1.1.4 LSTM . . . . .	1
1.2 Main Section 2 . . . . .	1
<b>2 Deep Learning Compilers</b>	<b>3</b>
2.1 Deep Learning Compilers . . . . .	3
2.2 Graph Level Frameworks . . . . .	3
2.2.1 TVM . . . . .	3
2.2.2 TensorFlow . . . . .	3
2.3 IR Level Frameworks . . . . .	3
<b>3 OnSRAM</b>	<b>5</b>
3.1 Introduction . . . . .	5
3.1.1 Subsection 1 . . . . .	5
3.1.2 Subsection 2 . . . . .	5
3.2 Generalizing To Multiple Scratchpads . . . . .	5
<b>Bibliography</b>	<b>7</b>



# List of Figures



# List of Tables





*For/Dedicated to/To my...*



## Chapter 1

# Deep Learning

### 1.1 Deep Learning Architectures

#### 1.1.1 Perceptron

#### 1.1.2 CNN

#### 1.1.3 RNN

#### 1.1.4 LSTM

### 1.2 Main Section 2



## Chapter 2

# Deep Learning Compilers

## 2.1 Deep Learning Compilers

"Few research efforts have proposed end-to-end software stacks that support AI accelerators. These are realized either as standalone custom frameworks or as software extensions to existing Deep Learning frameworks." [1] What are tensors in the context of compilers and why do they matter multi dimensional array that gets accessed a certain way and takes up memory for all we care What operations exist in deep learning - mat mul - convolution - adding? - Training - stuff ig What optimizations exist on deep learning models What is graph level representation What is IR

## 2.2 Graph Level Frameworks

### 2.2.1 TVM

### 2.2.2 TensorFlow

## 2.3 IR Level Frameworks



## Chapter 3

# OnSRAM

### 3.1 Introduction

OnSRAM introduces the notion of two types of optimizations that can be done on graph based compilers: intra-node and inter-node optimizations [1]. Intra node optimizations are those that are focused on optimizing the operation that the node specifies. This means tiling in favorable sizes and across specific dimensions, loop ordering, and DMA pipelining between tiles **Aladdin**. Inter node optimizations are optimizations concerning the overall structure and ordering of the nodes in the graph and their connections. This includes operator fusion and node reordering or rescheduling [1]. OnSRAM focuses on internode optimizations, specifically in the domain of scratchpad memory management techniques for DLAs. OnSRAM exploits the repeating useage of nodes in a graph -> some more detail on how graphs look like and why they're repeateable <- and identifies the nodes that can pinned to minimize memory transfers which increase inference time and decrease energy costs.

#### 3.1.1 Subsection 1

#### 3.1.2 Subsection 2

### 3.2 Generalizing To Multiple Scratchpads





# Bibliography

- [1] Carl E. Wieman and Leo Hollberg. “Using Diode Lasers for Atomic Physics”. In: *Review of Scientific Instruments* 62.1 (Jan. 1991), pp. 1–20. URL: <http://link.aip.org/link/?RSI/62/1/1>.