

# Hardware for Machine Learning

Subtitle if any



MS/PhD Thesis Proposal

*Submitted in partial fulfillment of the requirements for the award of the degree of*

Master/DOCTOR OF PHILOSOPHY

IN

ELECTRICAL ENGINEERING

By

**Hazoor Ahmad**

PhDEE17004

Supervisor Name

**Dr. Rehan Hafiz**

Co-Supervisor Name

**Dr. Mohsin Ali**

Department of Electrical Engineering, Faculty of Engineering

INFORMATION TECHNOLOGY UNIVERSITY

Lahore, Punjab, Pakistan

August 2019

## **Abstract**

This section consists of the topic/research problem, theoretical approach, research methodology and significance of the study. ddsdad sggfgs fgfdsgsg. ddsdad sggfgs fgfdsgsg. ddsdad sggfgs fgfdsgsg.

# 1 Introduction

This brief paragraph states the objective and/or goals of the PhD research. The information provided in this area gives the reader a quick overview of what is included in the proposal. reseacha adaadad asdasdasd [Krizhevsky et al., 2012] fsdafh dsafas fast safsa fastfas fas f. This research ais shown in Fig. 1. Our focus is design strategies.

	i24	i23	i22	i19	i18	i17	i14	o2	i1201+i1304	c2	i11c1+i12c4
	i25	i24	i23	i20	i19	i18	i15	a3	i11a1+i12a2+i13a3	c3	i18c1+i19c2+i20c3
	j13	j12	j11	j08	j07	j06	j03	a1	i18a1+i19a2+i20a3+j23a1	c1	i17c1+i18c2+i19c3+j22c1
	j14	j13	j12	j09	j08	j07	j04	a2	i17a1+i18a2+i19a3+j22a1+j23a2	c2	i16c1+i17c2+i18c3+j21c1+j22c2
	j15	j14	j13	j10	j09	j08	j05	a3	i16a1+i17a2+i18a3+j21a1+j22a2+j23a3	c3	i13c1+i14c2+i15c3+i18c1+j19c2+j20c3
k18	k17	k16	k13	k12	k11	k08	k07	a1	i13a1+i14a2+i15a3+j18a1+j19a2+j20a3+k08a1	c1	i12c1+i12c2+i14c3+j17c1+j18c2+j19c3+k07c1
k19	k18	k17	k14	k13	k12	k09	k08	a2	i12a1+i13a2+i14a3+j17a1+j18a2+j19a3+k07a1+k08a2	c2	i11c1+i12c2+i13c3+j16c1+j17c2+j18c3+k06c1+k07c2
0	k19	k18	k15	k14	k13	k10	k09	a3	i11a1+i12a2+i13a3+j16a1+j17a2+j18a3+k06a1+k07a2+k08a3	c3	i08c1+i09c2+i10c3+j13c1+j14c2+j17c3+k03c1+k04c2+k05c3
$(i03a1+i04a2+i05a3+j08a1+j09a2+j10a3+k13a1+k14a2+k15a3)+(i08a1+i09a2+i10a3+j13a1+j14a2+j15a3+k03a1+k04a2+k05a3)$									$(i02c1+i03c2+i04c3+j07c1+j08c2+j09c3+k12c1+k13c2+k14c3)+(i07c1+i08c2+i09c3+j12c1+j13c2+j14c3+k02c1+k03c2+k04c3)$		
$(i02a1+i03a2+i04a3+j07a1+j08a2+j09a3+k12a1+k13a2+k14a3)+(i07a1+i08a2+i09a3+j12a1+j13a2+j14a3+k02a1+k03a2+k04a3)$									$(i01c1+i02c2+i03c3+j06c1+j07c2+j08c3+k11c1+k12c2+k13c3)+(i06c1+i07c2+i08c3+j11c1+j12c2+j13c3+k01c1+k02c2+k03c3)$		
$(i01a1+i02a2+i03a3+j06a1+j07a2+j08a3+k11a1+k12a2+k13a3)+(i06a1+i07a2+i08a3+j11a1+j12a2+j13a3+k01a1+k02a2+k03a3)$									$i13c1+i14c2+i15c3+j18c1+j19c2+j20c3+k23c1+k24c2+k25c3$		
$i13a1+i14a2+i15a3+j18a1+j19a2+j20a3+k23a1+k24a2+k25a3$									$i12c1+i13c2+i14c3+j17c1+j18c2+j19c3+k22c1+k23c2+k24c3$		
$i12a1+i13a2+i14a3+j17a1+j18a2+j19a3+k22a1+k23a2+k24a3$									$i11c1+i12c2+i13c3+j16c1+j17c2+j18c3+k21c1+k22c2+k23c3$		
$i11a1+i12a2+i13a3+j16a1+j17a2+j18a3+k21a1+k22a2+k23a3$									$i08c1+i09c2+i10c3+j13c1+j14c2+j15c3+k18c1+k19c2+k20c3$		
$i08a1+i09a2+i10a3+j13a1+j14a2+j15a3+k18a1+k19a2+k20a3$									$i07c1+i08c2+i09c3+j12c1+j13c2+j14c3+k17c1+k18c2+k19c3$		
$i07a1+i08a2+i09a3+j12a1+j13a2+j14a3+k17a1+k18a2+k19a3$									$i06c1+i07c2+i08c3+j11c1+j12c2+j13c3+k16c1+k17c2+k18c3$		
$i06a1+i07a2+i08a3+j11a1+j12a2+j13a3+k16a1+k17a2+k18a3$									$i03c1+i04c2+i05c3+j08c1+j09c2+j10c3+k13c1+k14c2+k15c3$		

Figure 1: Overview of Hardware for Machine Learning

Organized as follows: section 2 contains a comprehensive literature survey on deep learning, section 3 limitations of state-of-the-art , section 4 includes problem statement, its scope, novelty and justification, section 5 explains significance and impact of proposed approach, section 6 provides a detailed design, modelling, analysis, simulation and evaluation of preliminary research results, section 7 enlists a clear set tasks or major milestones and evaluation strategies with their time-line in the form of Gantt chart.

## 2 Literature Survey

A detailed review of literature establishes your command in your area of research. This chapter should provide a complete and critical review of the state of the art work and not just summaries of books/articles.

This survey is focused on NeuFlow [Farabet et al., 2011]. All these implementations include Zhang [Zhang et al., 2015], Ayat [Ayat et al., 2018], and Caffeine [Zhang et al., 2018].

gsdfgsdfgg gfsdgsdf fdsgdgsfdgsdf . SysArrAccel [Wei et al., 2017] pointed out exploration. SmartShuttle [Li et al., 2018] provides configuration. Systimator [Hazaar et al., 2018] have tested networks .

Fused-Layer [Alwani et al., 2016] performs t optimize fo a ngle FPGA.

Escher [Shen et al., 2017] utputs. .

## 3 Limitations of the State of the Art

In this section you need to clearly report the current Limitations of the State of the Art. The purpose is to clearly identify the limitations and weakness of the state of the art so that your chosen problem statement is justified.

State-of-the-art techniques lack [Wei et al., 2017, Li et al., 2018] or do not utilize [Shen et al., 2017, Alwani et al., 2016, Kwon et al., 2018]. In addition to above, 3D systolic array has not yet been fully deployed for DNNs [Huan et al., 2018].

## 4 Problem Statement & Proposed Research

In this section you need to report your identified Problem Statement. Furthermore, also provide a clear scope of your proposed research. The novelty of the problem being addressed and the particular approach being explored should also be commented upon/justified.

This research is focused on . We will . The sessgdfg fgsgdfgdf fdgsdg gsdg . The sessgdfg fgsgdfgdf fdgsdg gsdg . The sessgdfg fgsgdfgdf fdgsdg gsdg .

## 5 Significance of the Study

This section is to provide the significance and impact of your proposed research. It should answer the following questions: Why you believe the study is significant? What implications your findings may have? Who will benefit from it? What will it contribute to the existing body of knowledge?

Systolic arrays are most useful .

## 6 Preliminary Research and Results

A written summary of some or all of the research performed is presented in a coherent manner. This section would include the approach taken and some preliminary results.

In this section we will demonstrate .

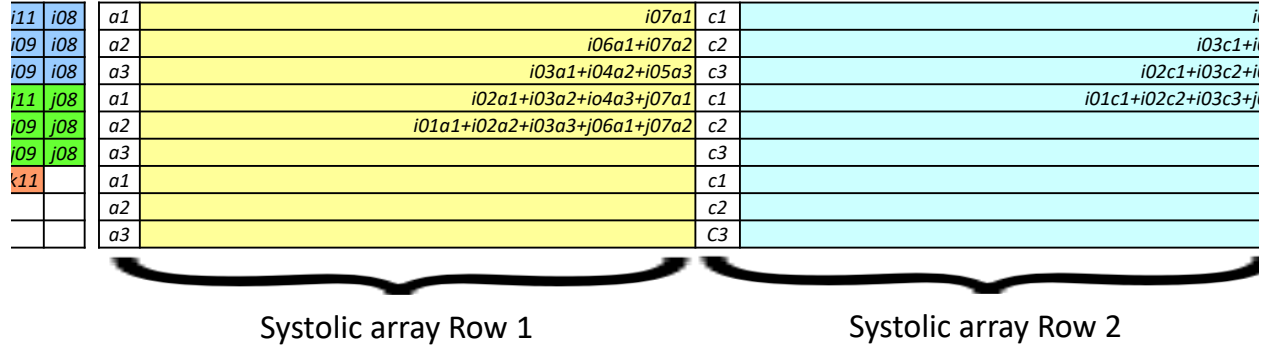


Figure 2: Dataflow engine.

### 6.1 E-Syzer

In this section we propose E-Syzer: An Efficient Systolic array sizer for DNN . Our proposed Fig. 2. We consider three different types of data traversals:

1. Input:
2. Weight:
3. Output:

systolic array in Fig. 2.

#### 6.1.1 Transr Enges (WTE, ITE, OTE)

Weiht, and oput tansfr gines.

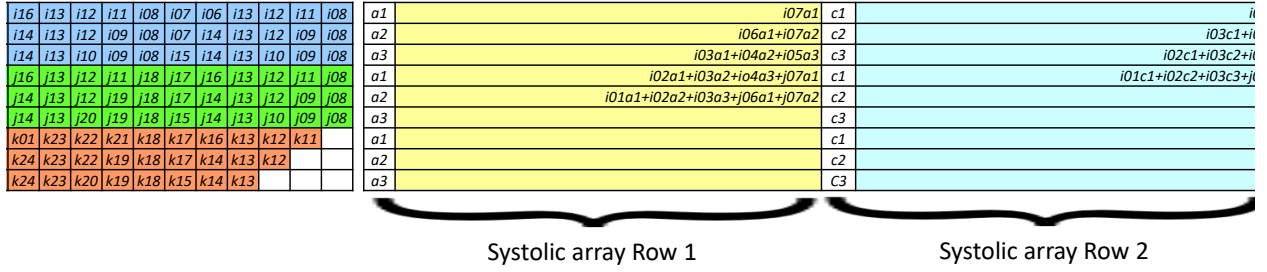


Figure 3: Systolic array

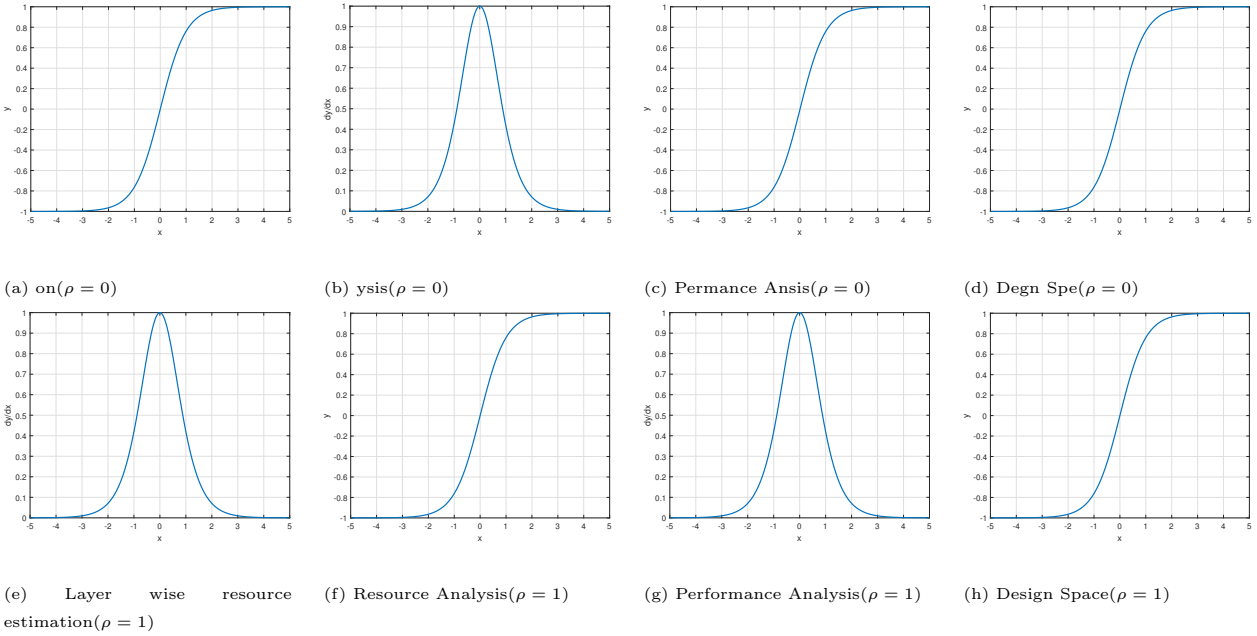


Figure 4: Layer rders.

### 6.1.2 Systolic Aay (SA)

Onevevy klok cycle in Fig. 2.

Table 1: Syzer Parameters

CNN Parameters			
Parameter Name	Symbol	Tiling	indexing
output	$M$	$T_m$	$m$
input	$N$	$T_m$	$n$
rows	$R$	$T_m$	$r$
columns	$C$	$T_m$	$c$
FPGA/Hardware Resources			
No. of available DSP blocks	$\mathbb{D}$		
No. of available slice LUTs	$\mathbb{G}$		

Before we move on,  $\mathcal{O}$

$$T_m = \min(\lceil \frac{T_m}{T_m} \rceil, R) \quad (1)$$

For simplification we take all columns of a layer to be the tiling factor for ofm columns.

$$T_c = \min(\lceil \frac{C_{TH}}{\mathcal{O}} \rceil, C) \quad (2)$$

To minimize optimization problem

$$\text{minimize}_{\langle \mathcal{A}, \mathcal{B} \rangle} \sum_{i=1}^c \mathcal{D}_i$$

subject to

$$T_m \leq R$$

$$T_m \leq C$$

$$\mathfrak{d} \leq \mathbb{D}$$

### 6.1.3 Evaluation



## 7 Proposed Research Plan

This objective of this section is to provide a clear set of tasks and intended approaches that shall be executed and evaluated to complete the PhD research work. All major milestones should be clearly mentioned. Evaluation plan can also be provided as to how the results shall be evaluated. A Gantt chart should be provided highlighting the outline of the work needed to complete the PhD research, and the time required for completion

There are following milestones to study

1. First Goal
2. Second Goal
3. Third Goal
4. Froth GOal

MATLAB<sup>®</sup> will be used for all simulations. Keras<sup>®</sup>, TensorFlow<sup>®</sup> or Caffe<sup>®</sup> for training. Hardware implementation will be using Vivado Design Suite - HLLX Editions by Xilinx<sup>®</sup>.

Table 2 provides a summary of major tasks along with their expected time of completion.

Activity \ Time	F 18	S 19	F 19	S 19	F 20	S 20
Literature Review and Study	✓					
Task 1		✓				
Task 2			✓			
Task 3				✓		
Task 4					✓	
Final Write-up & Thesis Submission						✓

Table 2: List of tasks

## Other Considerations

References should be presented in a consistent manner throughout the proposal. Use APA Referencing style for formatting the bibliographic material. It is important that figures, tables, and references in the proposal are presented in a manner consistent with professional publication standards. When placing a figure or table and its identifying description in the proposal, it is important to consider ease of access for the document reader. In most cases the text which introduces a figure or table will precede the placement of the figure or table in the proposal

## References

- [Alwani et al., 2016] Alwani, M., Chen, H., Ferdman, M., and Milder, P. (2016). Fused-layer cnn accelerators. In *The 49th Annual IEEE/ACM International Symposium on Microarchitecture*, page 22. IEEE Press.
- [Ayat et al., 2018] Ayat, S. O., Khalil-Hani, M., and Ab Rahman, A. A.-H. (2018). Optimizing fpga-based cnn accelerator for energy efficiency with an extended roofline model. *Turkish Journal of Electrical Engineering and Computer Science*, 26(2):919–935.
- [Farabet et al., 2011] Farabet, C., Martini, B., Corda, B., Akselrod, P., Culurciello, E., and LeCun, Y. (2011). NeufLOW: A runtime reconfigurable dataflow processor for vision. In *CVPR Workshops*, pages 109–116.
- [Hazoor et al., 2018] Hazoor, A., Tanvir, M., Abdullah, M., Javed, M. U., Hafiz, R., and Shafique, M. (2018). Systimator: A design space exploration methodology for systolic array based cnns acceleration on the fpga-based edge nodes. *arXiv preprint arXiv:1901.04986*.
- [Huan et al., 2018] Huan, Y., Xu, J., Zheng, L., Tenhunen, H., and Zou, Z. (2018). A 3d tiled low power accelerator for convolutional neural network. In *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1–5. IEEE.

- [Krizhevsky et al., 2012] Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- [Kwon et al., 2018] Kwon, H., Samajdar, A., and Krishna, T. (2018). Maeri: Enabling flexible dataflow mapping over dnn accelerators via reconfigurable interconnects. *ACM SIGPLAN Notices*, 53(2):461–475.
- [Li et al., 2018] Li, J., Yan, G., Lu, W., Jiang, S., Gong, S., Wu, J., and Li, X. (2018). Smartshuttle: Optimizing off-chip memory accesses for deep learning accelerators. In *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 343–348. IEEE.
- [Shen et al., 2017] Shen, Y., Ferdman, M., and Milder, P. (2017). Escher: A cnn accelerator with flexible buffering to minimize off-chip transfer. In *2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 93–100. IEEE.
- [Wei et al., 2017] Wei, X., Yu, C. H., Zhang, P., Chen, Y., Wang, Y., Hu, H., Liang, Y., and Cong, J. (2017). Automated systolic array architecture synthesis for high throughput cnn inference on fpgas. In *Proceedings of the 54th Annual Design Automation Conference 2017*, page 29. ACM.
- [Zhang et al., 2015] Zhang, C., Li, P., Sun, G., Guan, Y., Xiao, B., and Cong, J. (2015). Optimizing fpga-based accelerator design for deep convolutional neural networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 161–170. ACM.
- [Zhang et al., 2018] Zhang, C., Sun, G., Fang, Z., Zhou, P., Pan, P., and Cong, J. (2018). Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*.