

# TaxoNN: A Light-Weight Accelerator for Deep Neural Network Training



Reza Hojabr<sup>1</sup>, Kamyar Givaki<sup>1\*</sup>, Kossar Pourahmadi<sup>1\*</sup>, Parsa Nooralinejad<sup>1\*</sup>,  
Ahmad Khonsari<sup>12</sup>, Dara Rahmati<sup>2</sup>, M. Hassan Najafi<sup>3</sup>

<sup>1</sup>School of Electrical and Computer Engineering, University of Tehran, Tehran, Iran

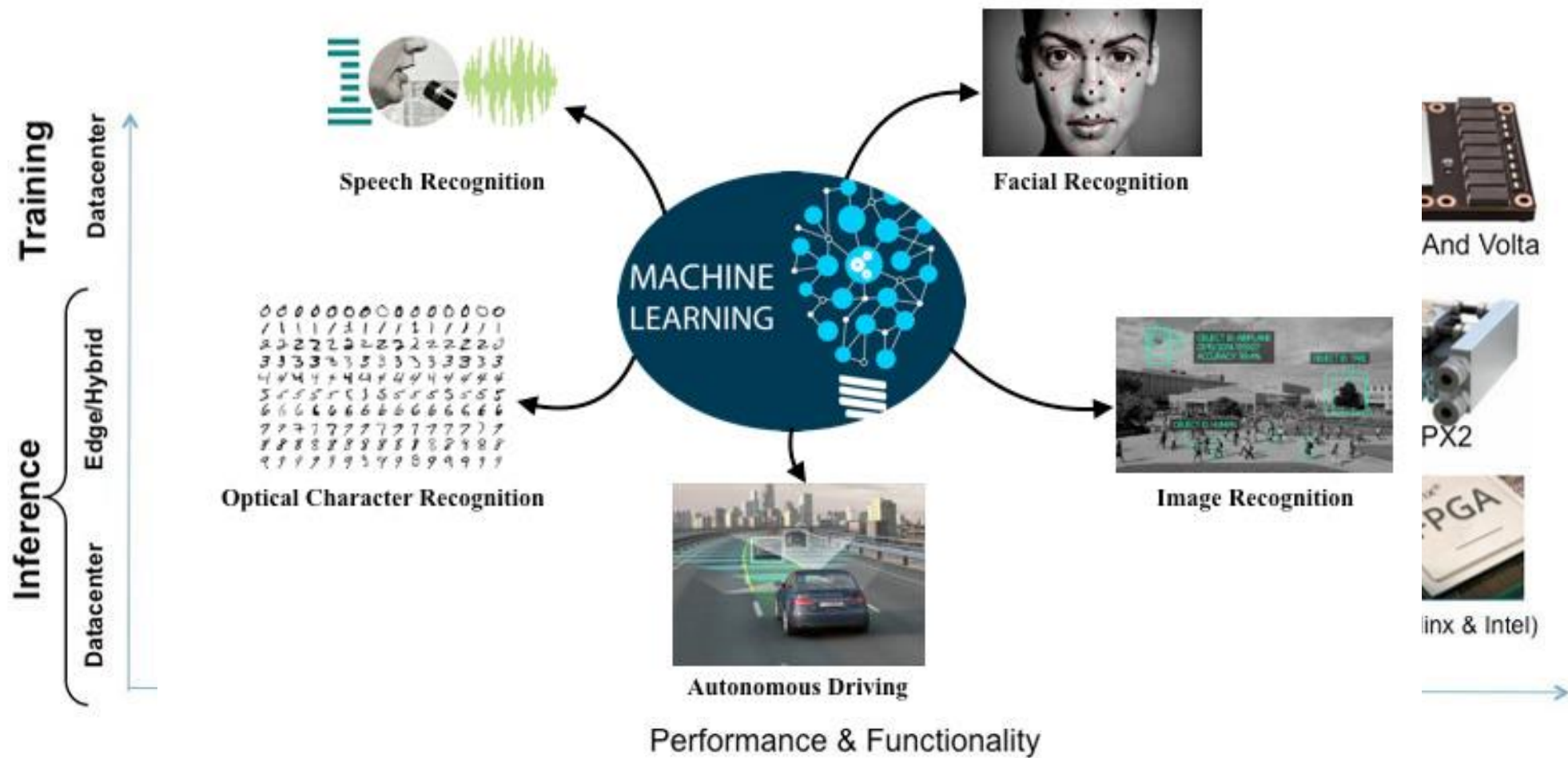
<sup>2</sup>School of Computer Sciences, Institute for Research in Fundamental Sciences (IPM), Tehran, Iran

<sup>3</sup>School of Computing and Informatics, University of Louisiana at Lafayette, Lafayette, LA, USA

**2020 IEEE International Symposium on Circuits and Systems**  
**Virtual, October 10-21, 2020**

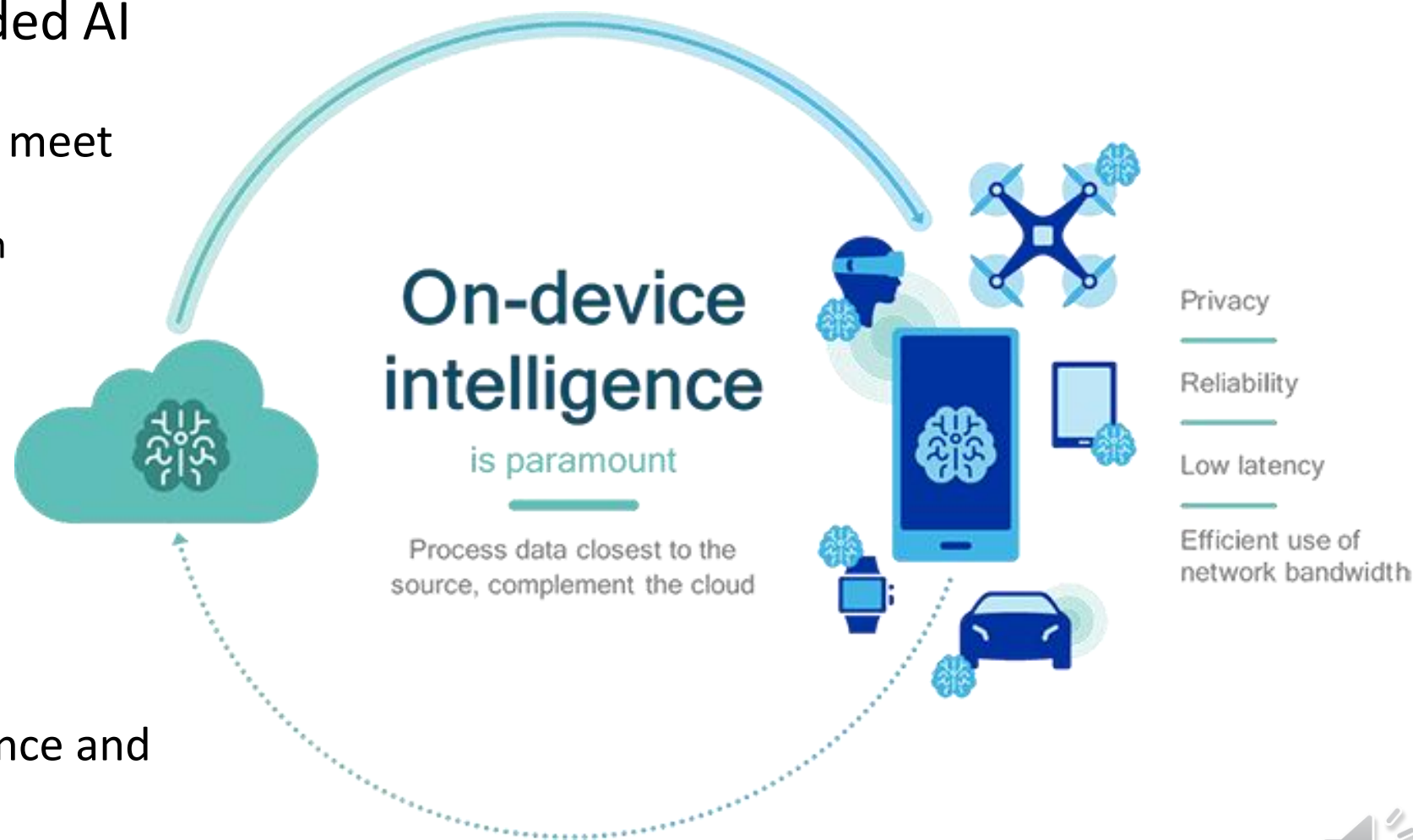


# Deep Learning Accelerators



# Edge AI – The Next Wave of AI

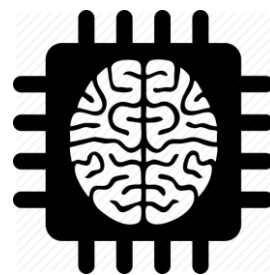
- The need for embedded AI accelerators
  - GPU and TPU cannot meet the requirements:
    - Power consumption
    - Price \$
- On-device training
  - Efficient use of BW
  - Privacy
  - Reliability
  - Low latency
- TaxoNN
  - Supports both inference and training



# Main Contributions

Split the SGD  
algorithm into  
smaller  
computation  
units

**TaxoNN**



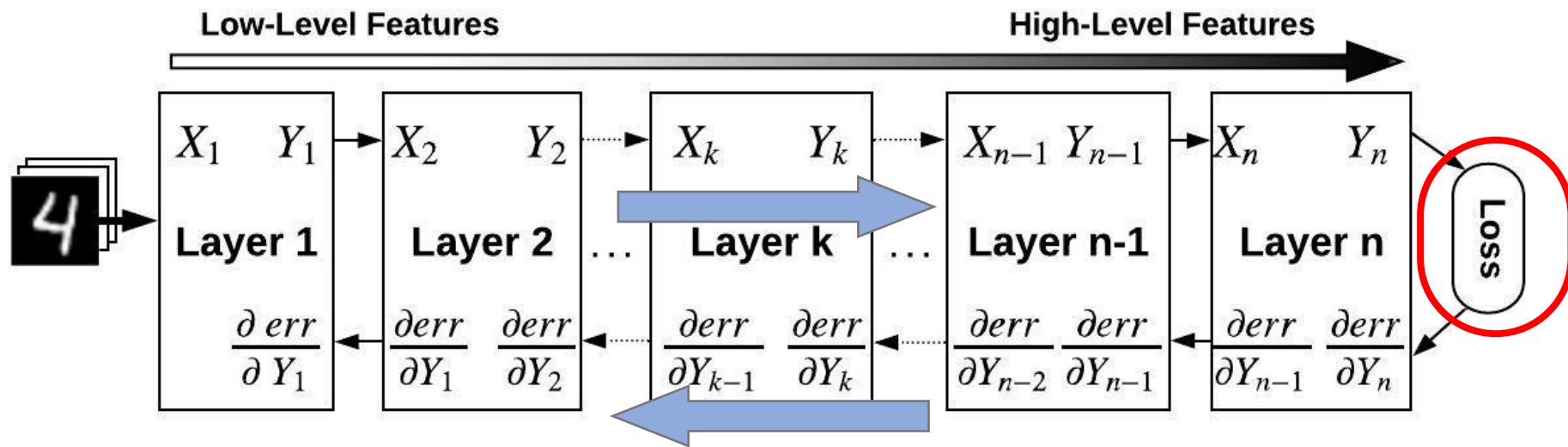
Adding the  
training ability  
to the baseline  
inference only  
Processing  
Elements

Low-bitwidth  
operations for  
different layers



# Stochastic Gradient Descent

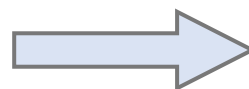
$$W_i = W_i - \alpha \frac{\partial error}{\partial W_i}$$



# Simplifying SGD Algorithm

$$W_i = W_i - \alpha \frac{\partial error}{\partial W_i}$$

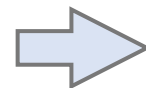
Chain rule



$$\frac{\partial error}{\partial W_i} = \frac{\partial error}{\partial Y_{i+1}} \times \frac{\partial Y_{i+1}}{\partial Y_i} \times \frac{\partial Y_i}{\partial W_i}$$

$Y_i$  is the output of  $i^{th}$  layer

$$\frac{\partial Y_{i+1}}{\partial Y_i} = \frac{\partial f_{i+1}(W_{i+1}Y_i)}{\partial Y_i} = W_{i+1} f'_{i+1}(W_{i+1}Y_i)$$



$$\frac{\partial Y_i}{\partial W_i} = \frac{\partial f_i(W_iY_{i-1})}{\partial W_i} = Y_{i-1} f'_i(W_iY_{i-1})$$

$$\frac{\partial error}{\partial W_i} = \frac{\partial error}{\partial Y_{i+1}} \times f'_{i+1} \times W_{i+1} \times f'_i \times Y_{i-1}$$

$$G_i = G_{i+1} \times W_{i+1} \times f'_i$$



$$\frac{\partial error}{\partial W_i} = G_i \times X_i$$

$f_i()$  is the activation function of the  $i^{th}$  layer

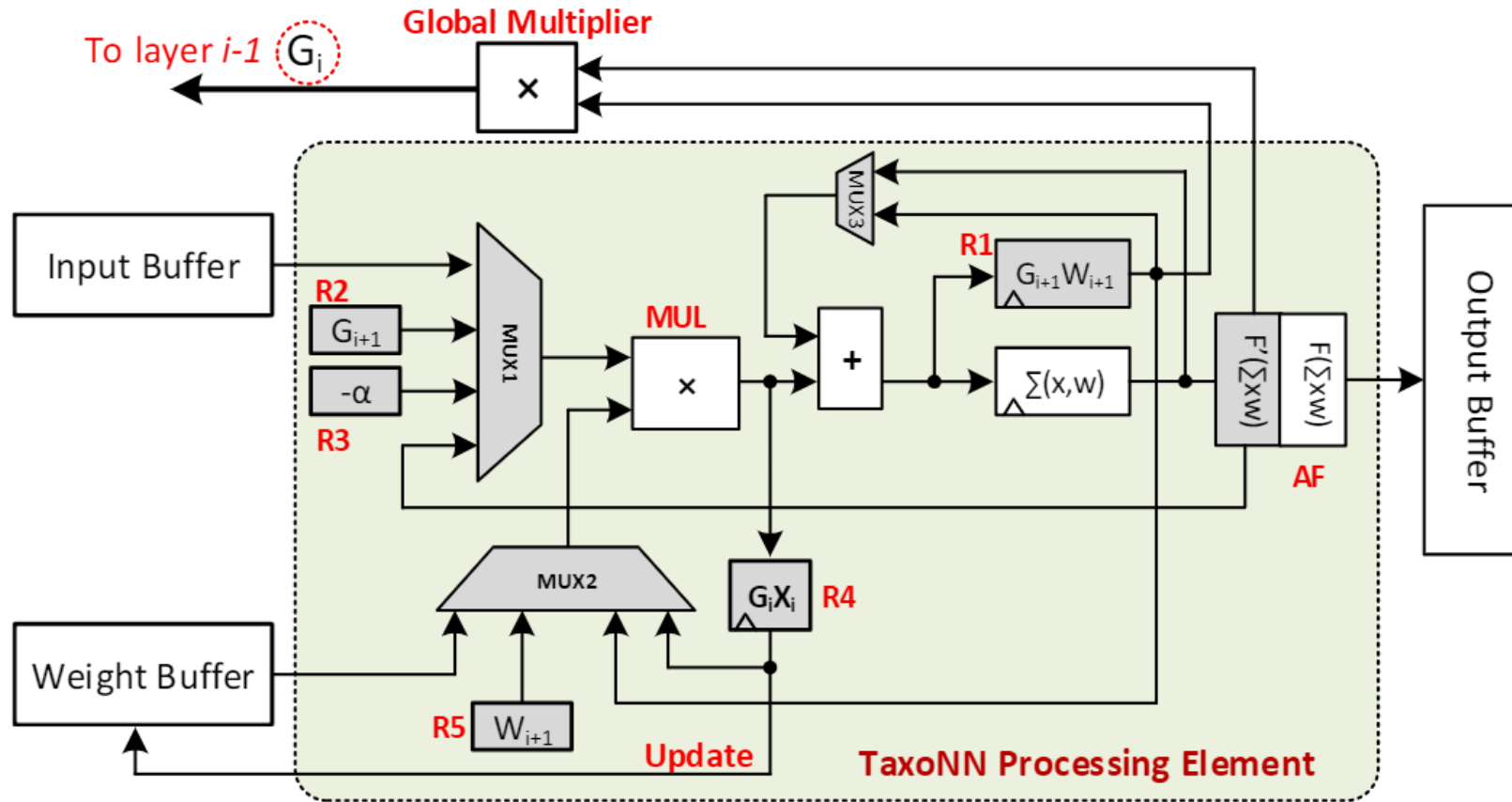
$f'()$  is the derivation of the activation function



# TaxoNN Architecture

$$W_i = W_i - \alpha \frac{\partial error}{\partial W_i}$$

$$\frac{\partial error}{\partial W_i} = G_i \times X_i$$



# Timing and Pipelining

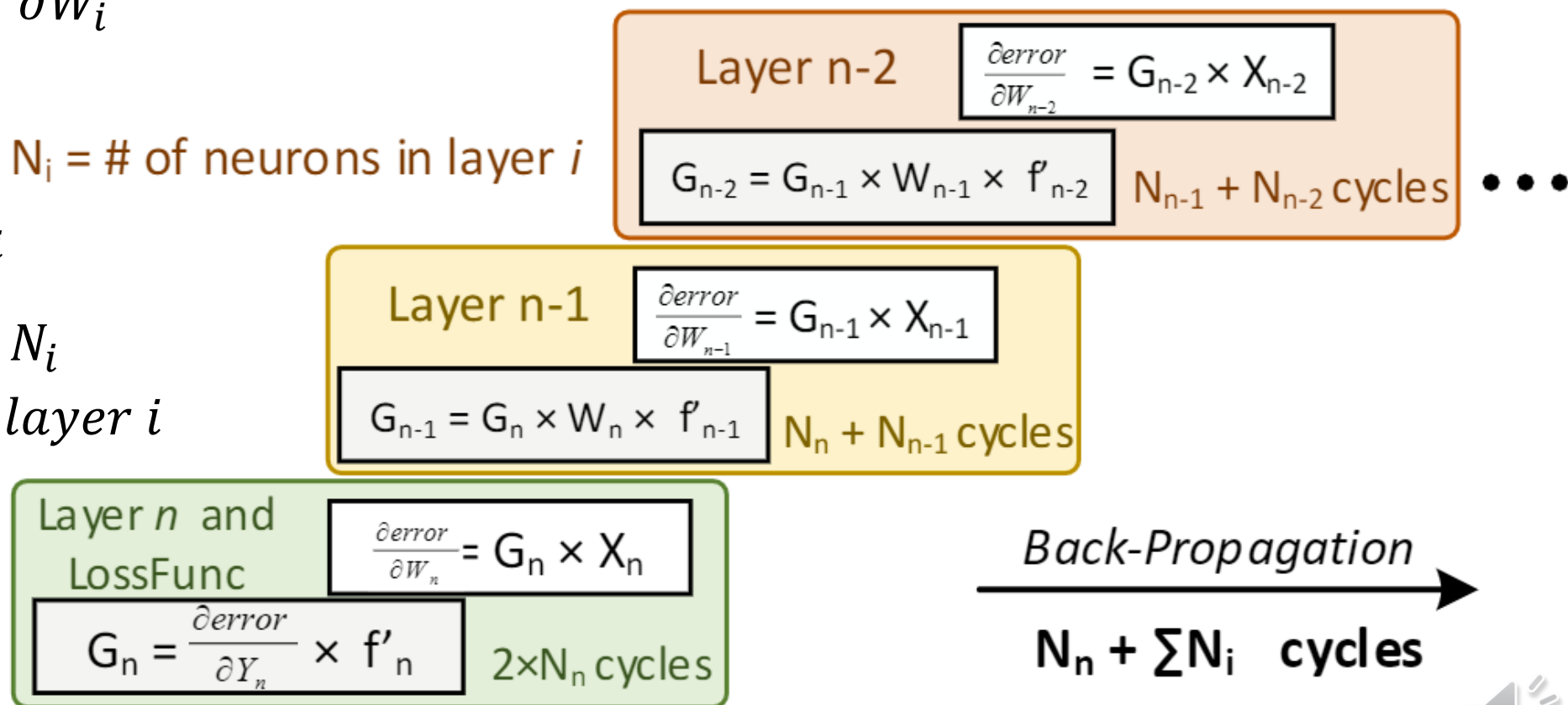
$$W_i = W_i - \alpha \frac{\partial error}{\partial W_i}$$

$$\frac{\partial error}{\partial W_i} = G_i \times X_i$$

$G_i$ : a vec of size  $N_i$

$N_i$ : #of neuron in layer  $i$

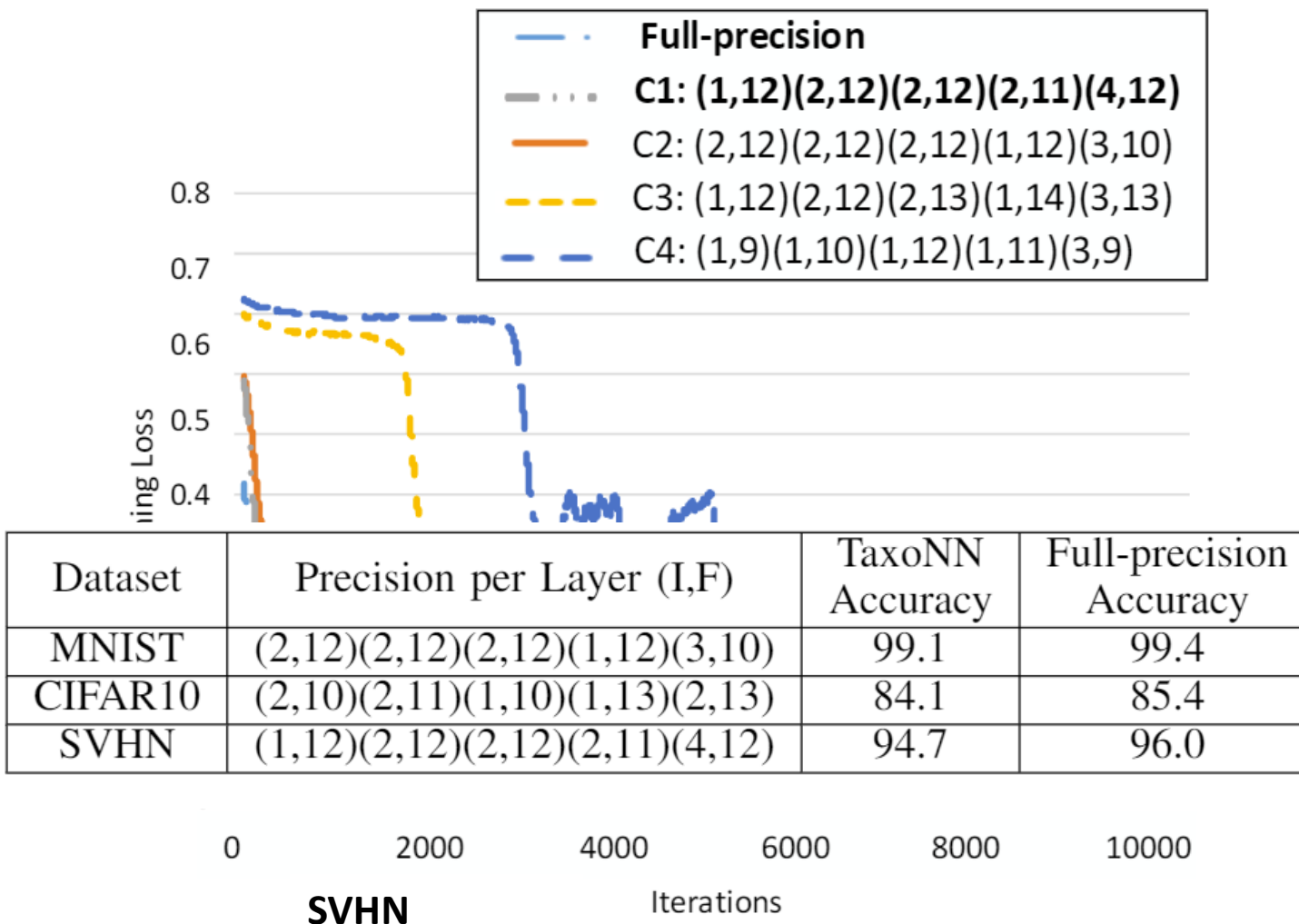
$N_i$  = # of neurons in layer  $i$





# Performance Evaluation

- Network
  - LeNet
- Dataset
  - MNIST
  - CIFAR10
  - SVHN
- Developed a cycle-accurate simulator
- Full-precision results are extracted from TensorFlow



# Hardware Cost

- Implemented TaxoNN in RTL Verilog
- Synthesized using DC with 45nm technology

Bitwidth	21	20	19	18	17	16	15	14	13	Average
Eyeriss	14.3	13.1	11.8	11.1	10.6	10.1	9.7	9.0	8.1	Area
<b>TaxoNN</b>	15.5	14.3	12.9	12.1	11.7	11.2	10.6	9.9	9.0	Overhead
Overhead	8.3%	9.2%	9.1%	8.6%	10.0%	10.8%	8.8%	9.8%	10.5%	9.5%

Bitwidth	21	20	19	18	17	16	15	14	13	Average
Eyeriss	4.54	4.48	4.42	4.31	4.22	4.10	3.98	3.88	3.75	Power
<b>TaxoNN</b>	4.84	4.78	4.70	4.65	4.49	4.31	4.15	4.13	4.04	Overhead
Overhead	6.5%	6.7%	6.2%	7.9%	6.5%	5.2%	4.3%	6.5%	7.7%	6.4%

Dataset	Power Reduction	Area Reduction
MNIST	2.1×	1.7×
CIFAR10	2.3×	1.8×
SVHN	1.9×	1.5×



# TaxoNN Summary

- TaxoNN: a light-weight accelerator for DNN training
- A novel method to unroll and parallelize the SGD algorithm
- A fine-grained and optimized datapath to perform the SGD
- Evaluated TaxoNN using low-bitwidth operations for each layer
- 1.65x area and 2.1x power saving at the cost of 0.97% misclassification rate compared to full-precision





**Reza Hojabr**

[r.hojabr@ut.ac.ir](mailto:r.hojabr@ut.ac.ir)

