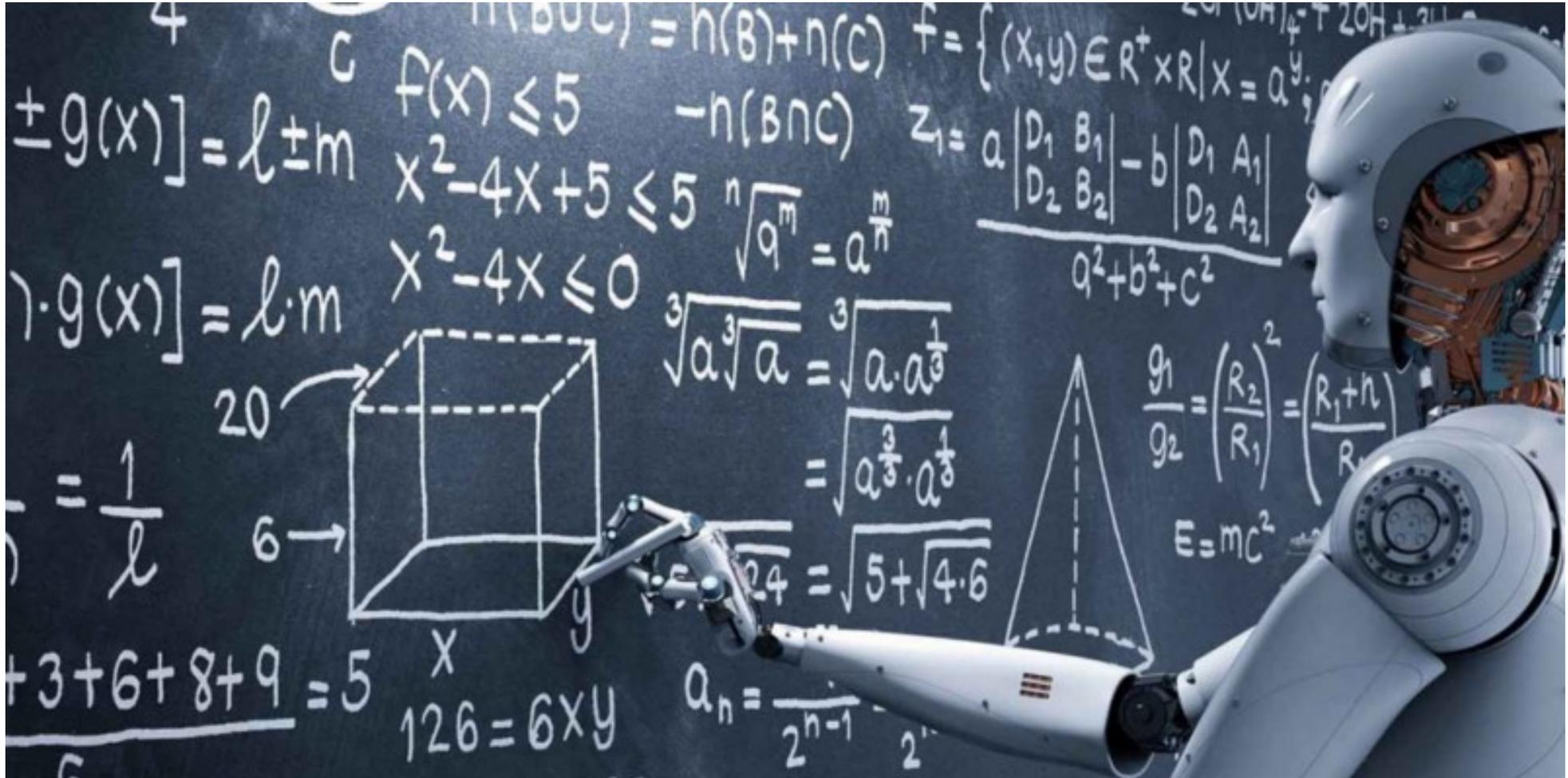


Enabling Continuous Learning through Synaptic Plasticity in Hardware

Tushar Krishna
Georgia Tech

EMC² Workshop
June 23 2019

The Dream!



Deep Learning Applications

“AI is the new electricity” – Andrew Ng

Object Detection

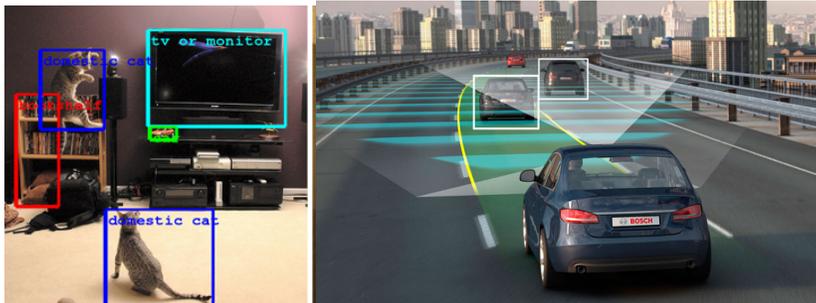
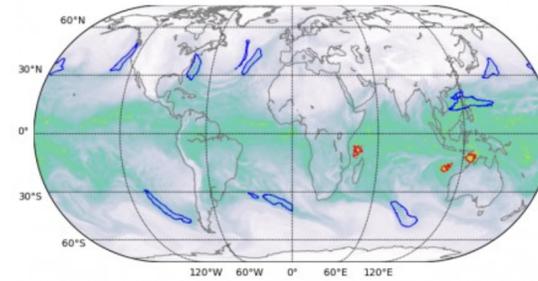


Image Segmentation



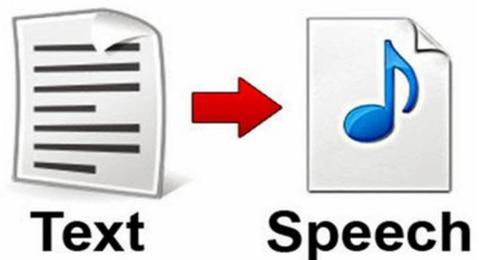
Medical Imaging



Speech Recognition



Text to Speech



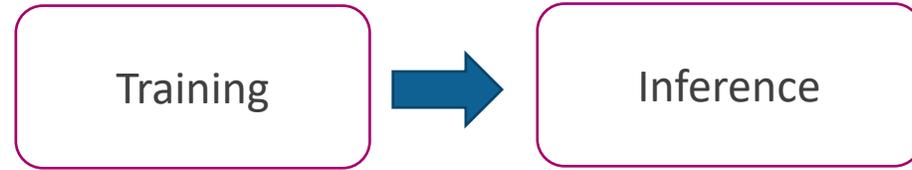
Recommendations



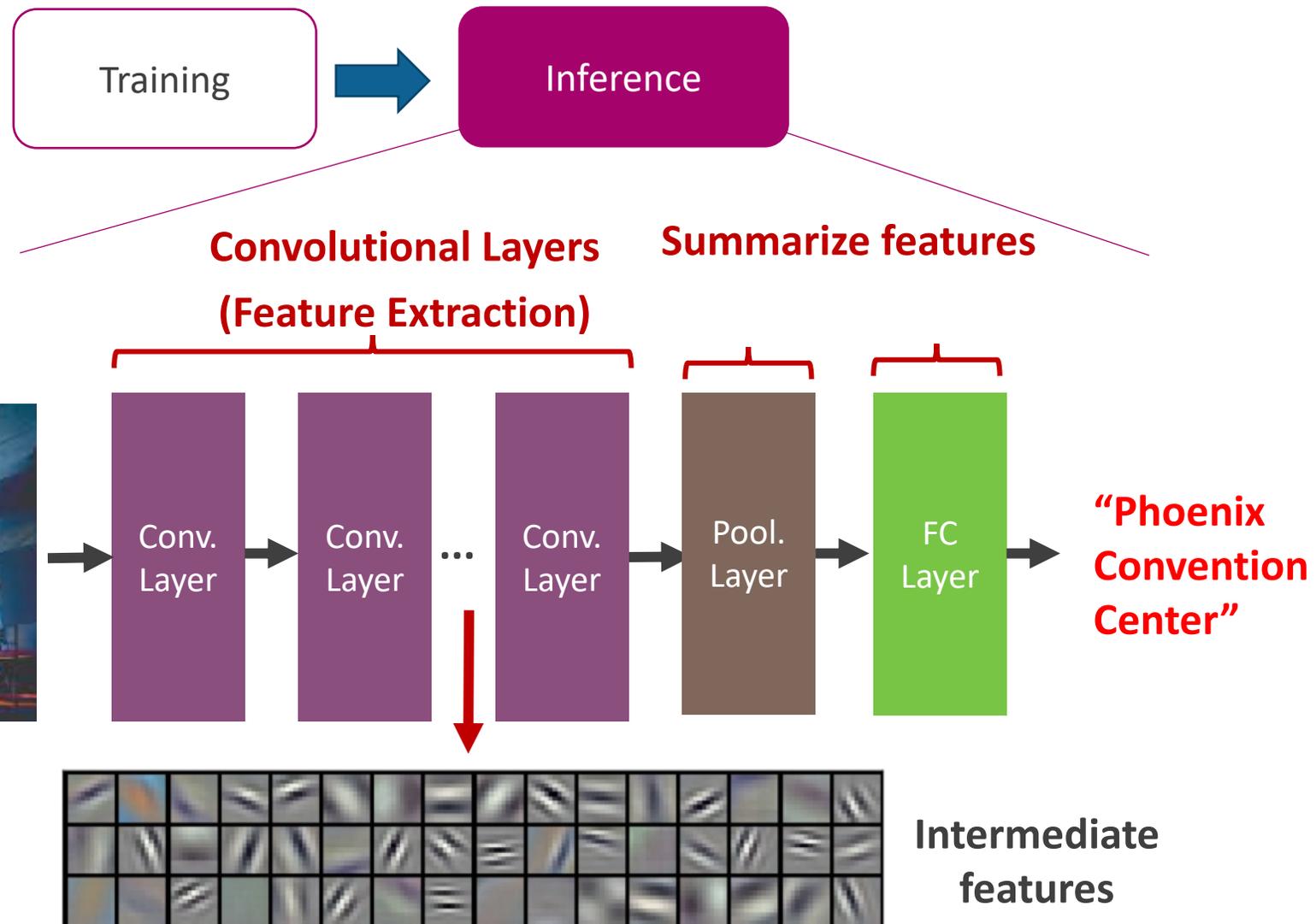
Games



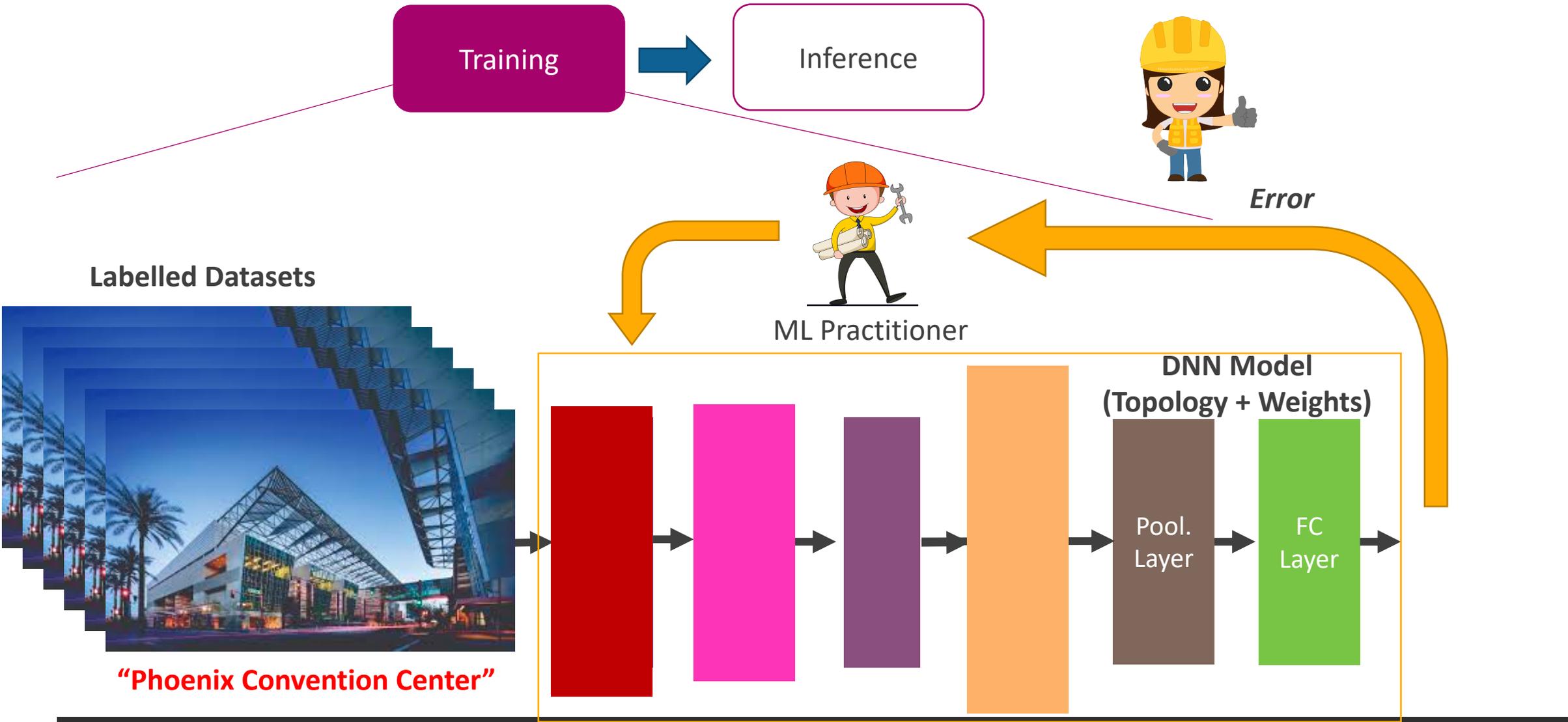
Deep Learning Landscape



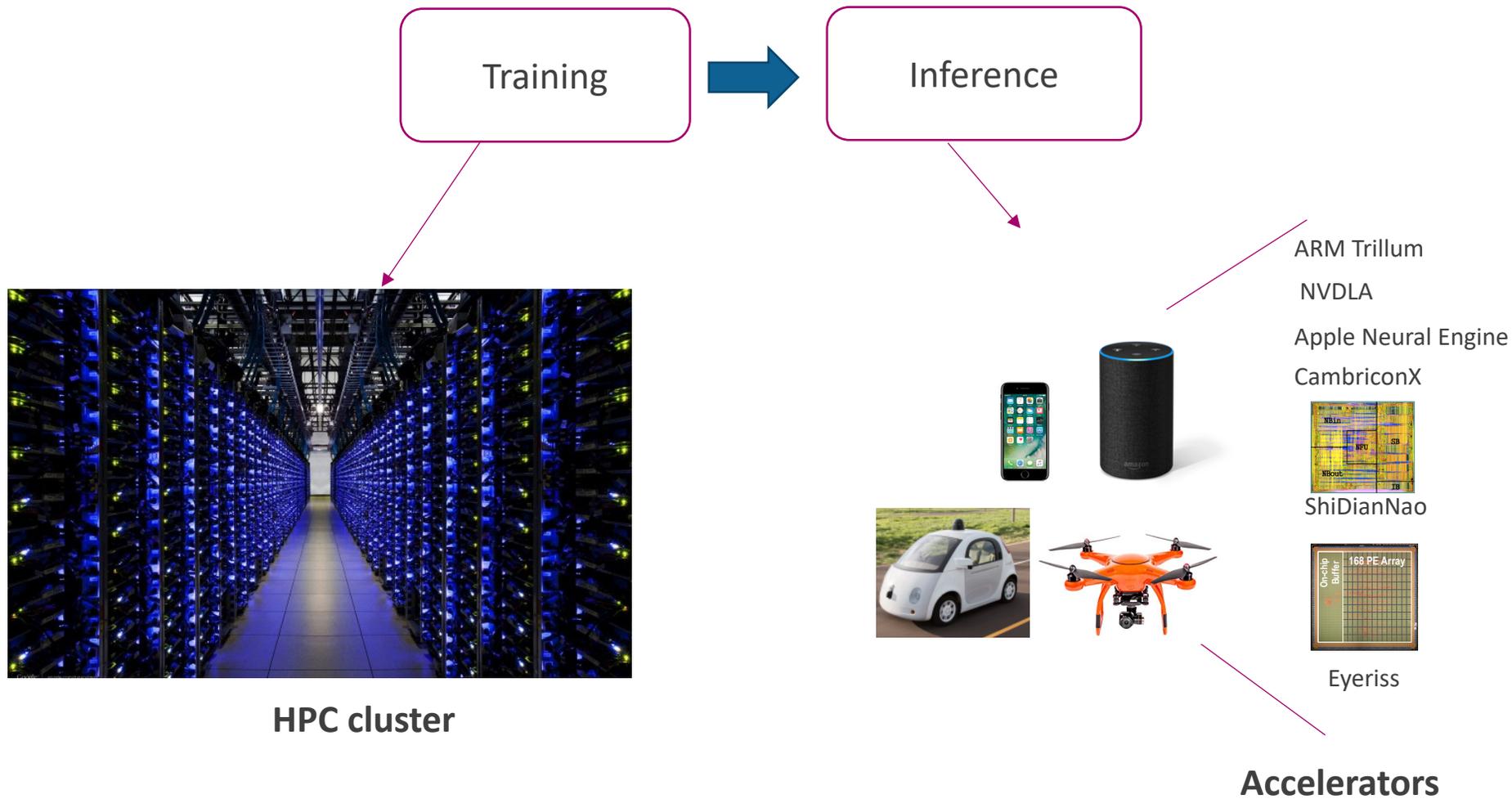
Deep Learning Landscape



Deep Learning Landscape



Computation Platforms

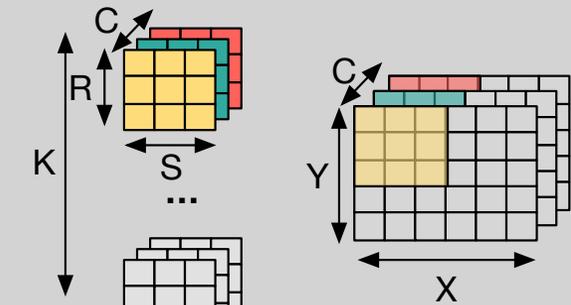


Efficiency of Deep Learning Systems

Training



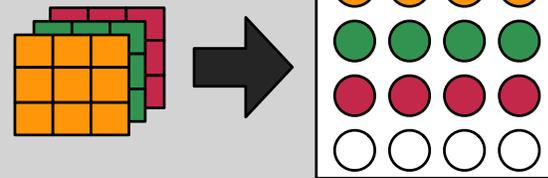
Inference



DNN Architecture

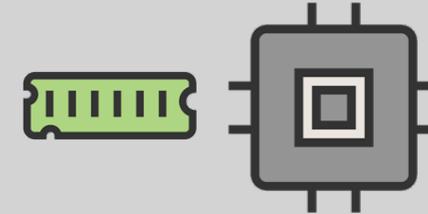
How to design a DNN for the target task?

Mapping (Dataflow)



How to map billions of computations over *limited* compute/memory resources?

Microarchitecture



How to design an efficient accelerator for the DNN model



Energy



Runtime

What is Continuous Learning?

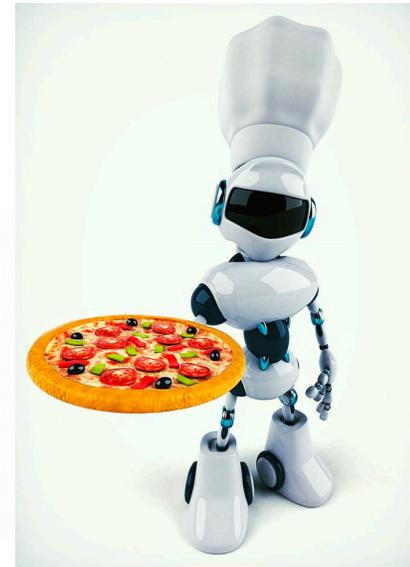
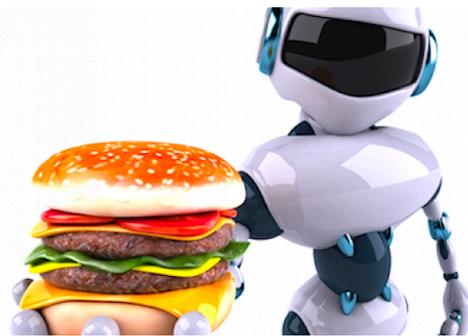


Become better and faster
with experience

Learn new tasks

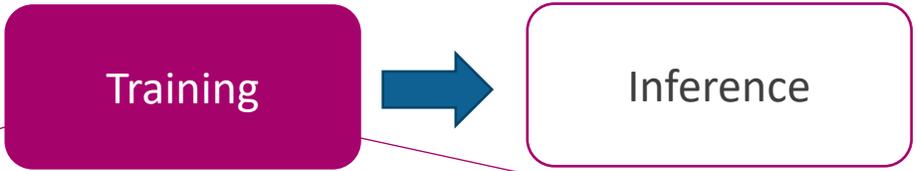
Compute and
energy-efficiency

Can we leverage
Supervised Deep
Learning?



Deep Learning Landscape

Deep Learning not viable for continuous learning



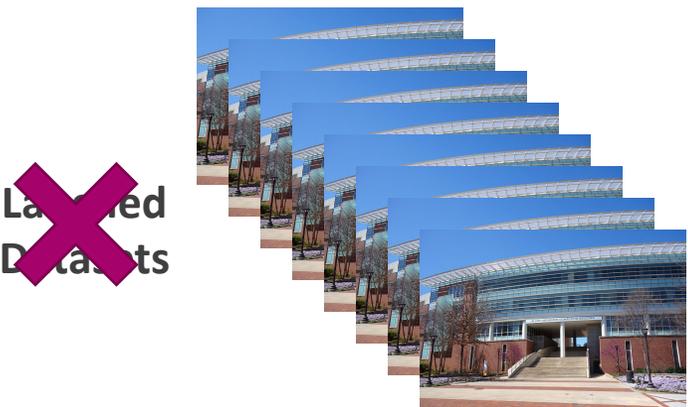
HPC cluster



ML Practitioner

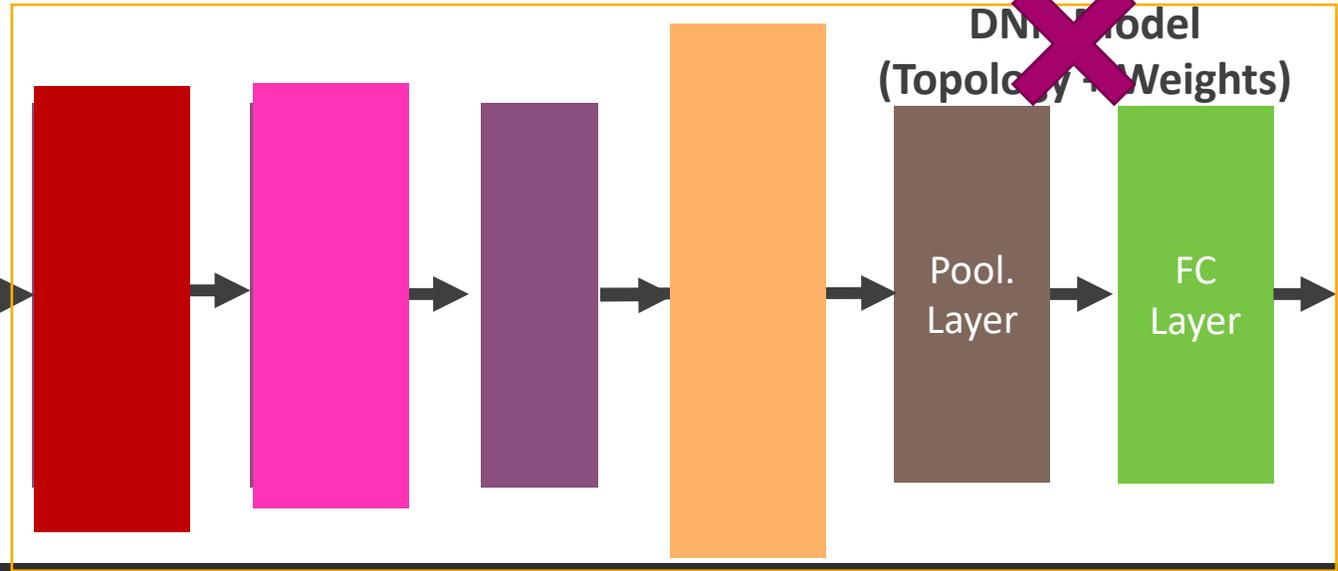


Error



Labeled Datasets

"Klaus Advanced Computing Building"

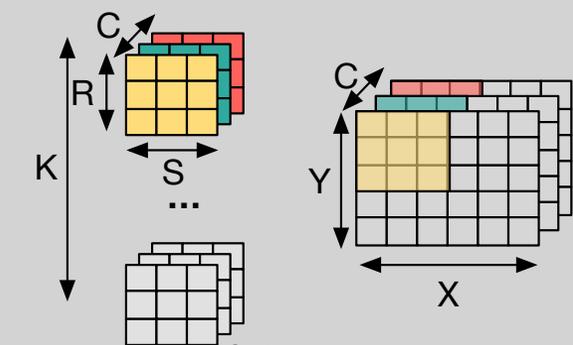


Efficiency of Continuous Learning Systems

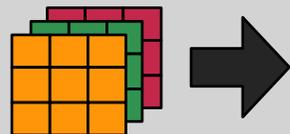
Training



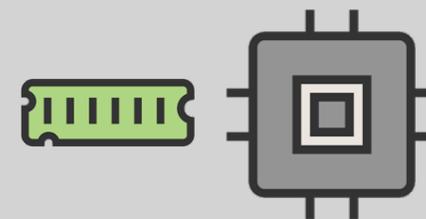
Inference



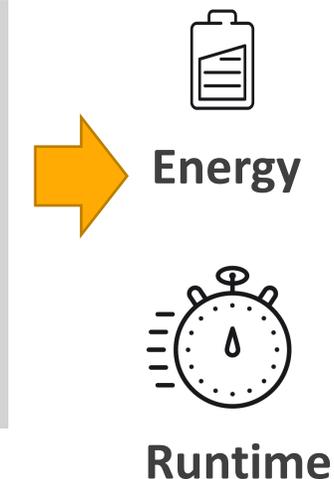
DNN Architecture



Mapping (Dataflow)



Microarchitecture



How to autonomously update DNN models for continuous learning?

How to efficiently map changing DNNs over accelerator?

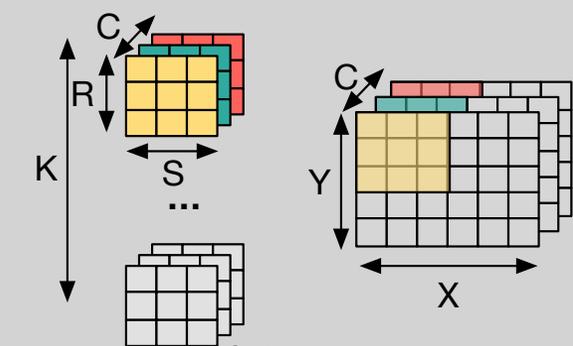
How to design an efficient accelerator for changing DNN models

Outline of Talk

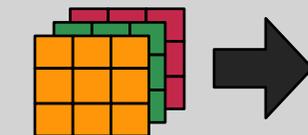
Training



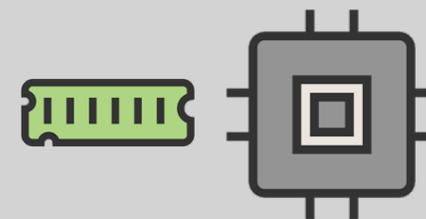
Inference



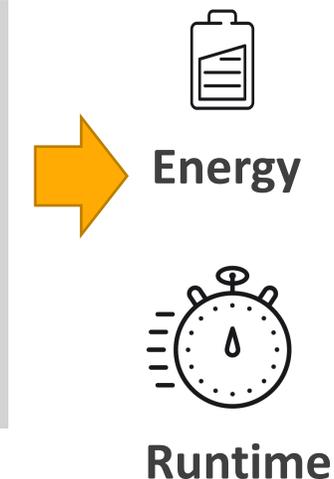
DNN Architecture



Mapping (Dataflow)



Microarchitecture



How to autonomously update DNN models for continuous learning?

GeneSys

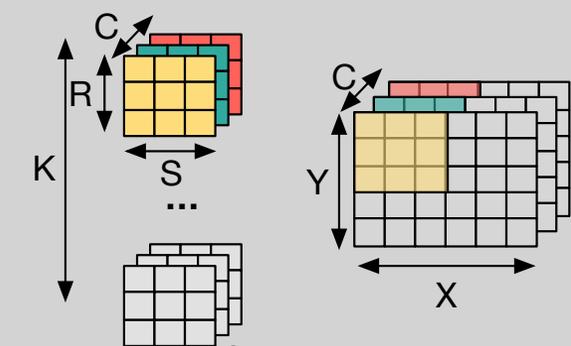
How to efficiently map changing DNNs over accelerator?

How to design an efficient accelerator for changing DNN models

MAERI

Outline of Talk

Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna, *GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware*, *MICRO 2018*



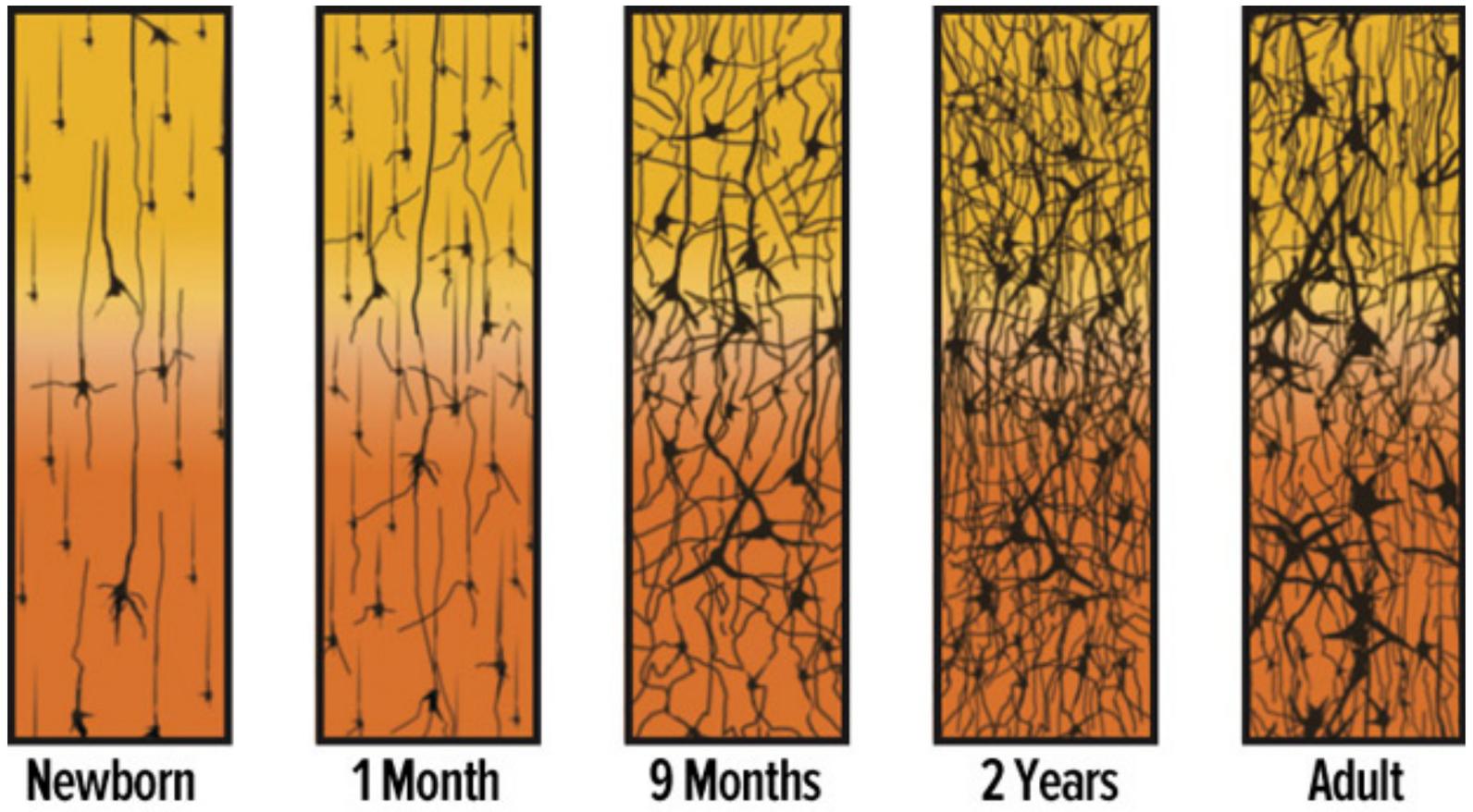
DNN Architecture

How to autonomously update DNN models for continuous learning?

GeneSys

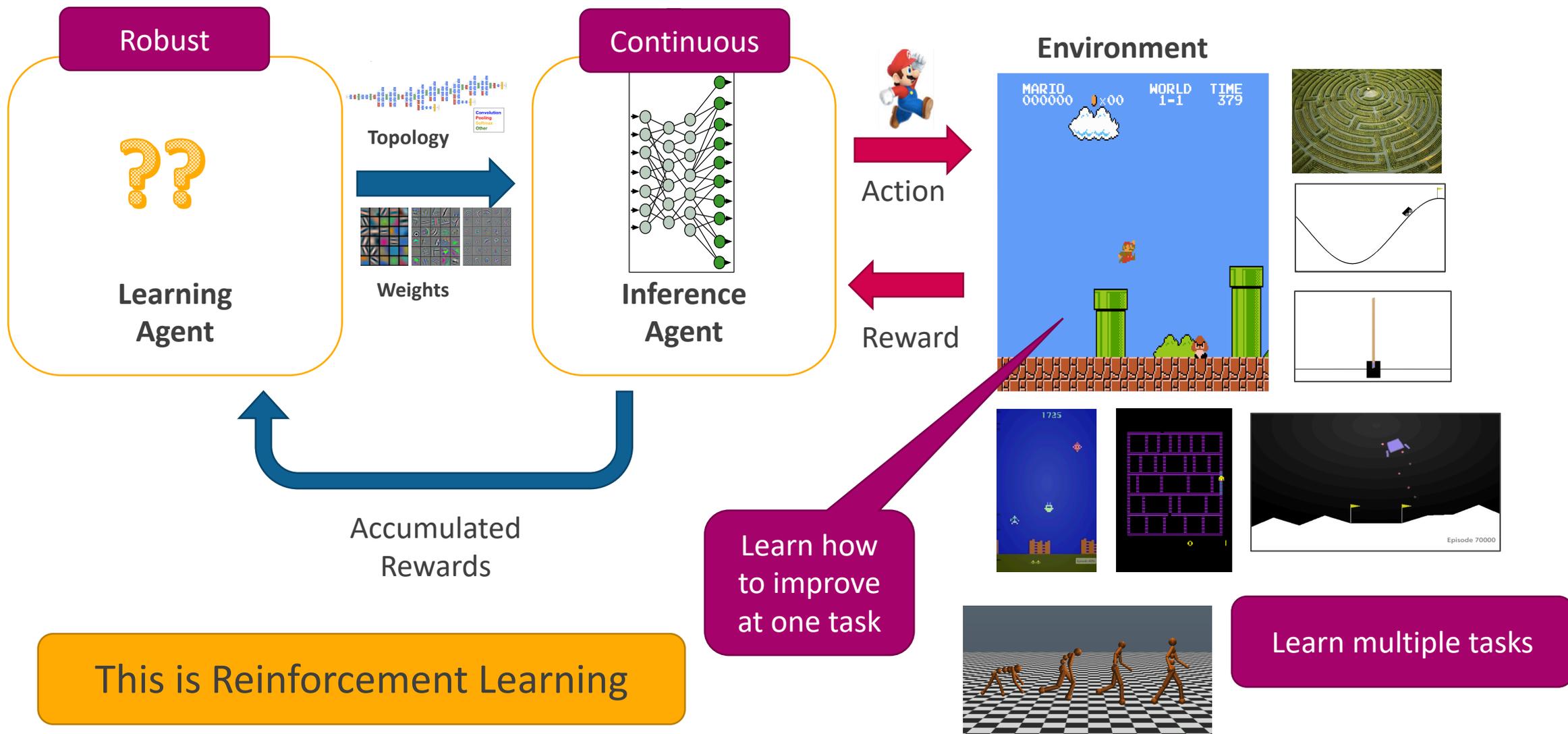
- Continuous Learning Template
- Neuro-Evolutionary Algorithms
 - Algorithm Description
 - Characterizing NEAT
- Microarchitecture
- Evaluations

Continuous Learning in Brains



Constant synapse formation and pruning

Template for Continuous Learning



Conventional RL: Challenges

Deep NNs used internally

! Manual hyperparameter tuning

Each update results in **Backpropagation**

! High compute requirement at every update

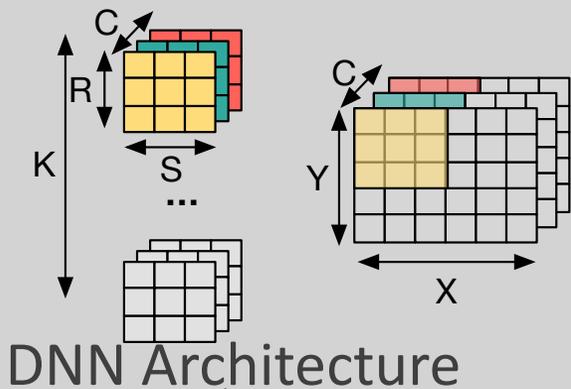
! High memory overhead

! Not scalable

**Not viable for continuous
learning on the edge**

Outline of Talk

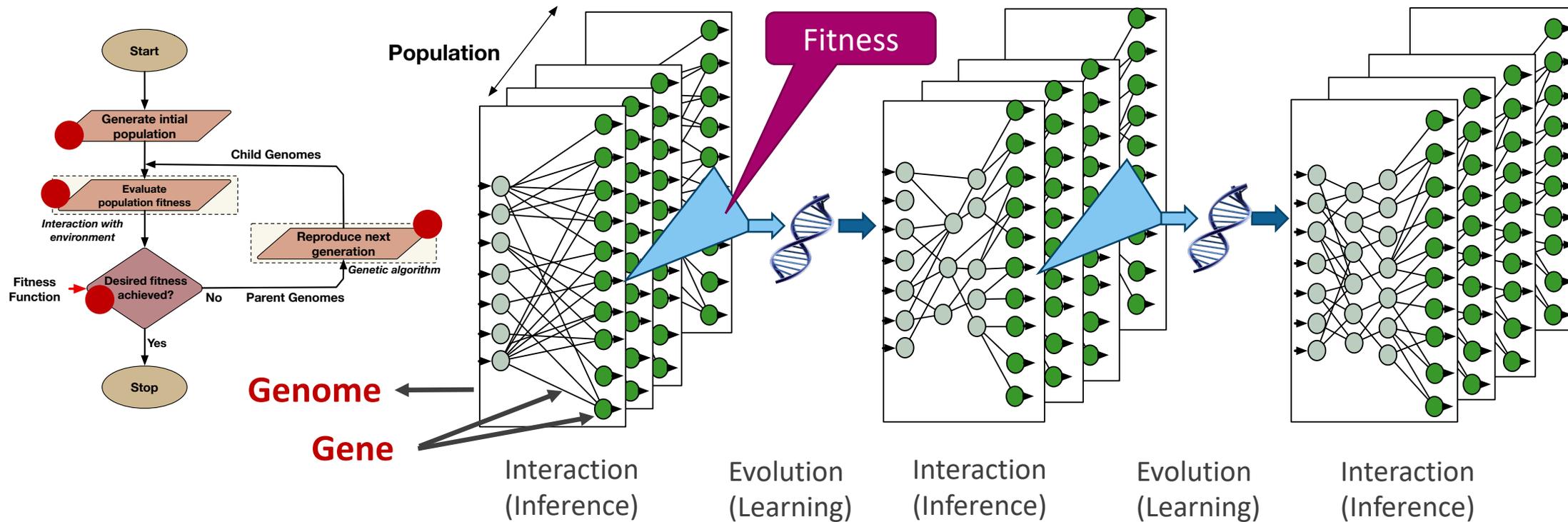
Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna, *GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware*, *MICRO 2018*



GeneSys

- Continuous Learning` Template
- **Neuro-Evolutionary Algorithms**
 - Algorithm Description
 - Characterizing NEAT
- Microarchitecture
- Evaluations

Neuro-Evolutionary (NE) Algorithm



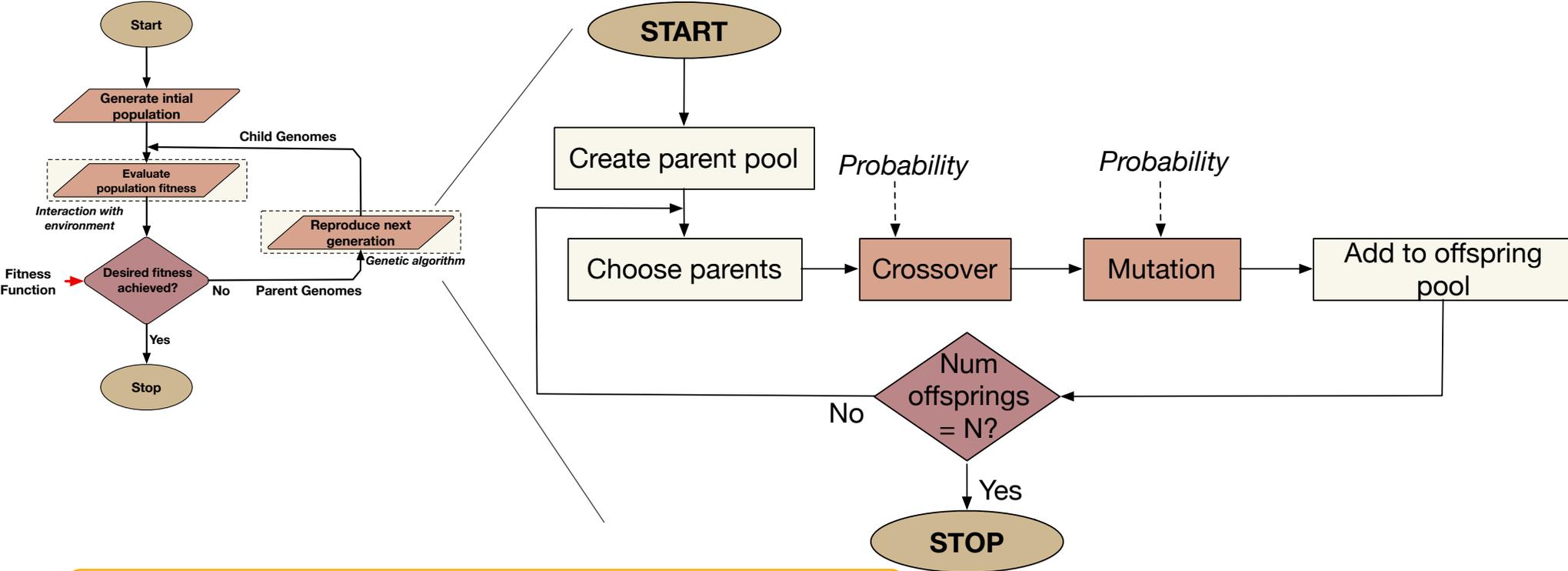
Neural Network (NN) expressed as a graph

Gene: Vertex or Edge in the graph

Genome: Collection of all genes (i.e., a NN)

[1] Stanley, K. O., & Miikkulainen, R. (2002). Evolving neural networks through augmenting topologies. *Evolutionary computation*, 10(2), 99-127.

Neuro-Evolutionary (NE) Algorithm



Neural Network (NN) expressed as a graph

Gene: Vertex or Edge in the graph

Genome: Collection of all genes (i.e., a NN)

NeuroEvolution of Augmented Topologies (NEAT) [1]

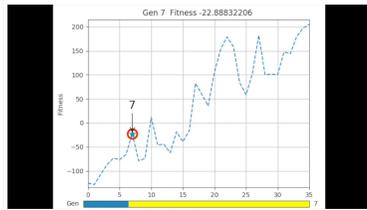
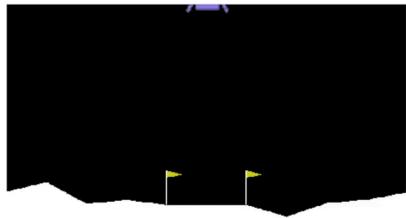
[1] Stanley, K. O., & Miikkulainen, R. (2002). Evolving neural networks through augmenting topologies. *Evolutionary computation*, 10(2), 99-127.

Properties of NE algorithms

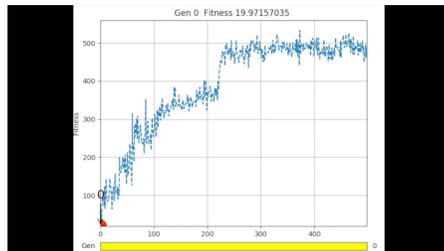
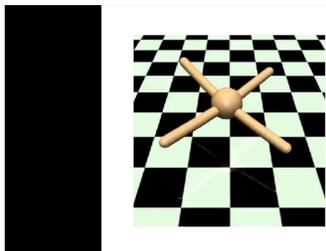
Algorithmic

Robustness

No Training



Change fitness function



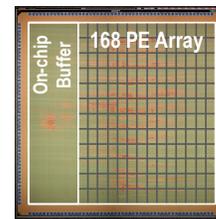
Accuracy?

Systems

Too much compute!

Convergence time?

déjà vu! Looks like Deep Neural Networks in the 90s



Eyeriss



GPU



FPGA

HW solutions enabled Deep Learning

Can we do the same with EA?

Outline of Talk

Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna, *GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware*, *MICRO 2018*

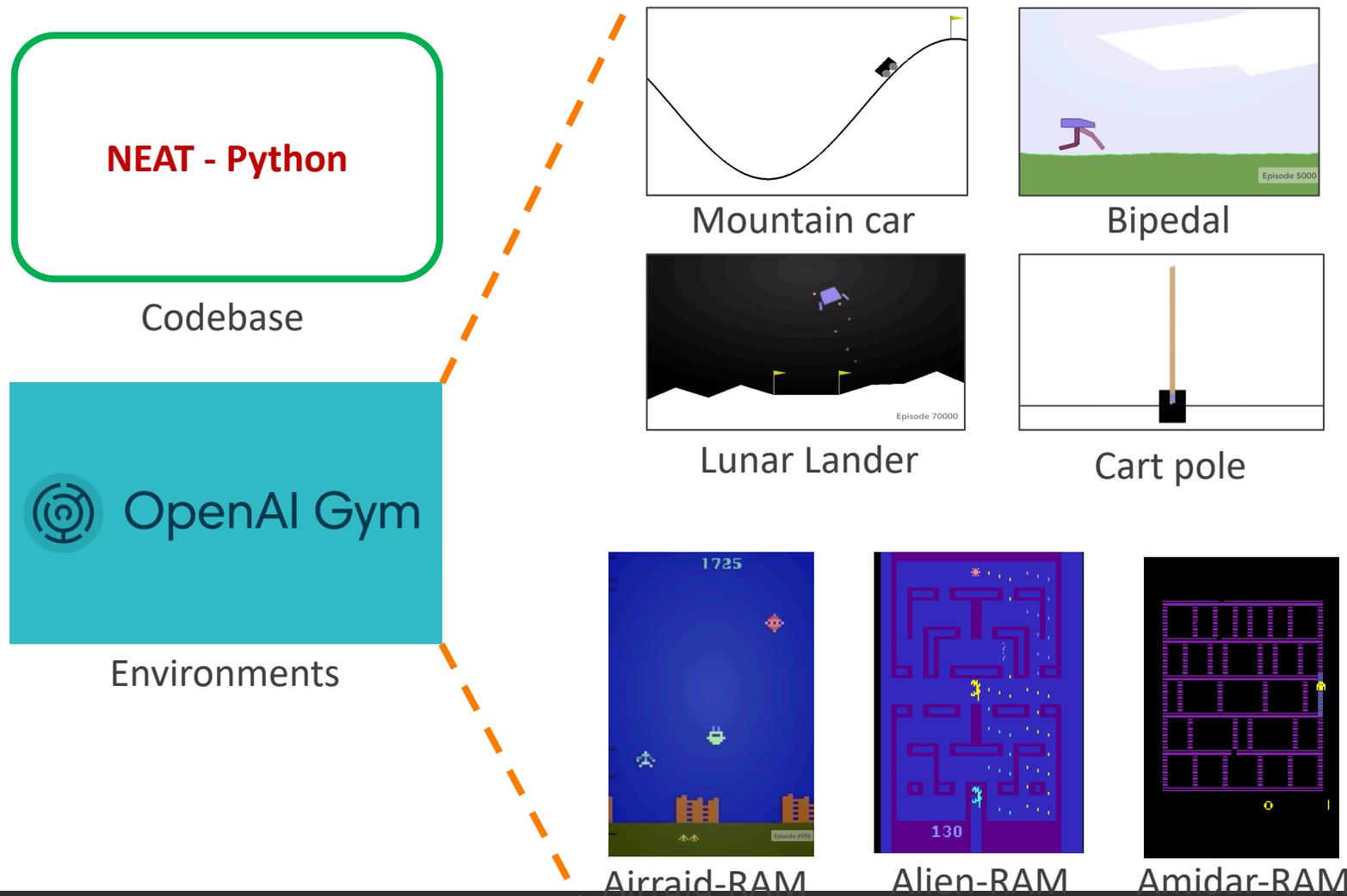
DNN Architecture

How to autonomously design DNN models for continuous learning?

GeneSys

- Continuous Learning
- **Neuro-Evolutionary Algorithms**
 - Algorithm Description
 - **Characterizing NEAT**
- Microarchitecture
- Evaluations

Characterization of NEAT

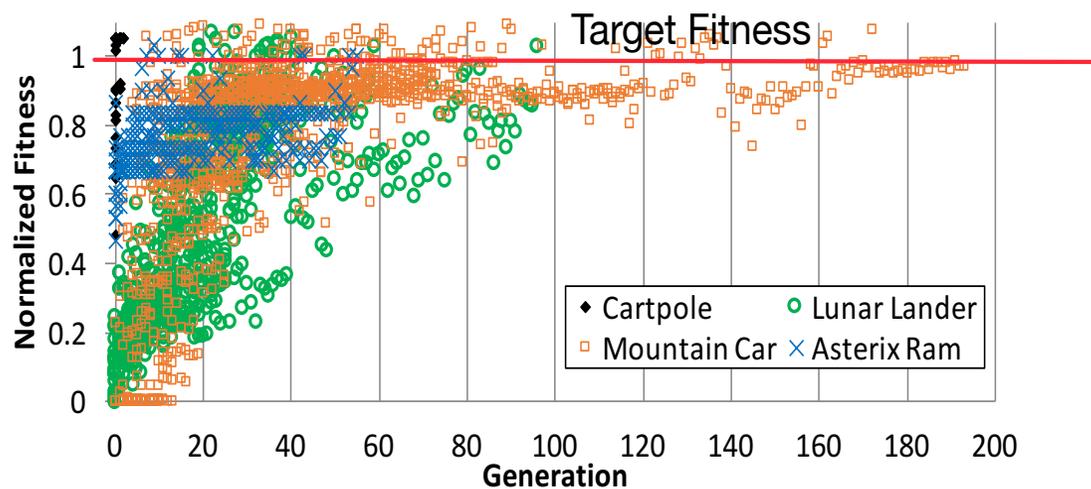


Ran each environment till convergence, multiple times

Only changed fitness function between workloads

Characterization of NEAT

Computations



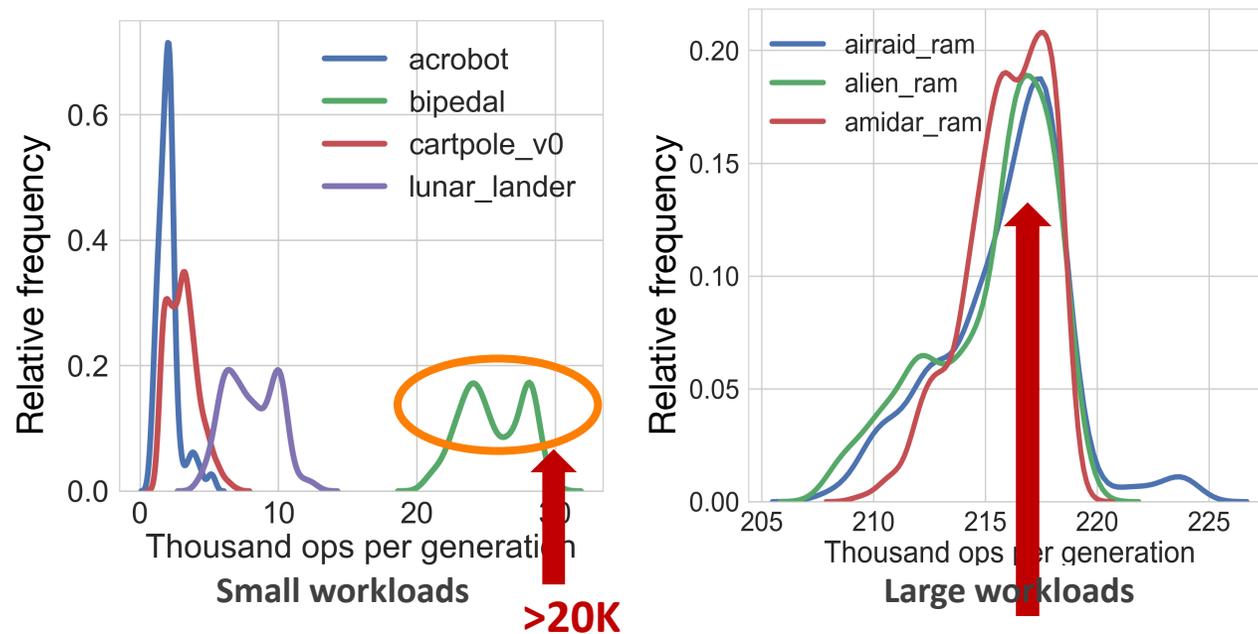
Inference:

Population level parallelism (PLP)

Evolution:

Gene level parallelism (GLP)

Distribution of Operations/Generation



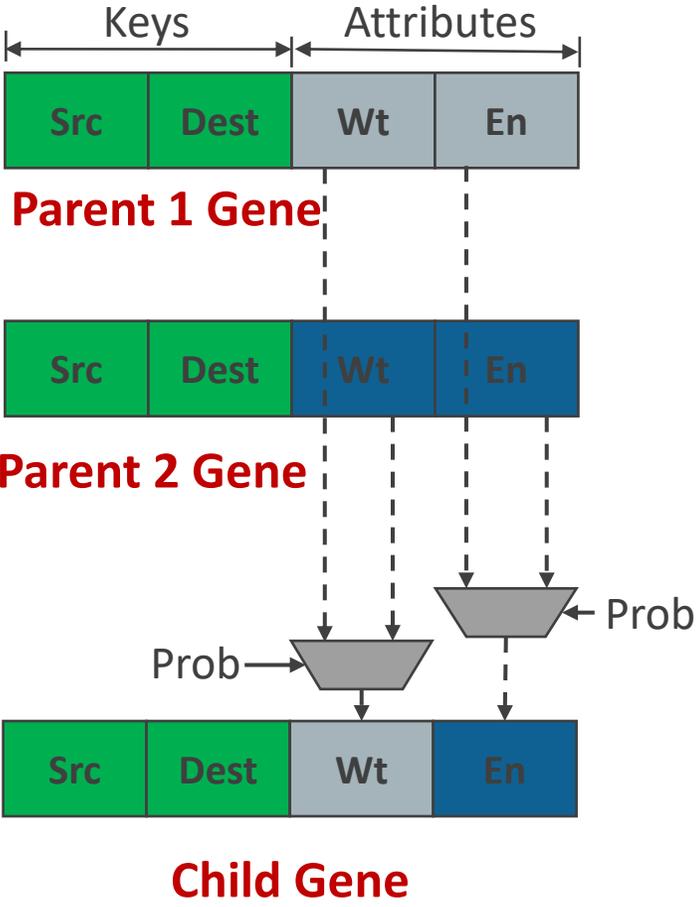
All operations are independent

Large operation level Parallelism

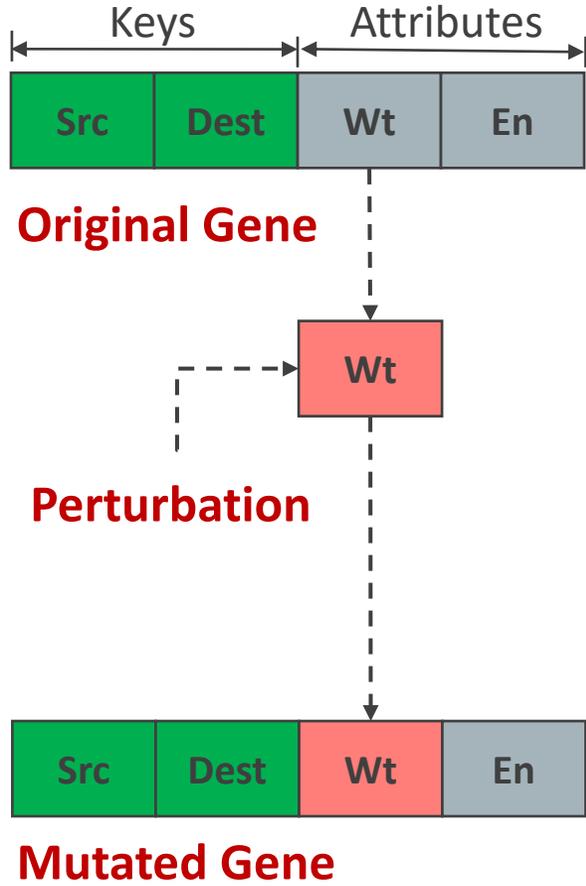
Operations in NEAT

Evolution

Crossover



Mutation



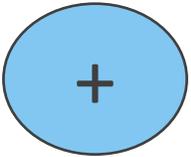
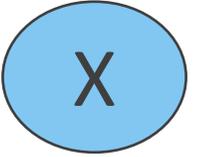
Addition mutation

- Add new node
- Add new connection

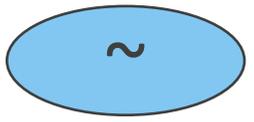
Deletion mutation

- Delete connection
- Delete node

Inference



MAC

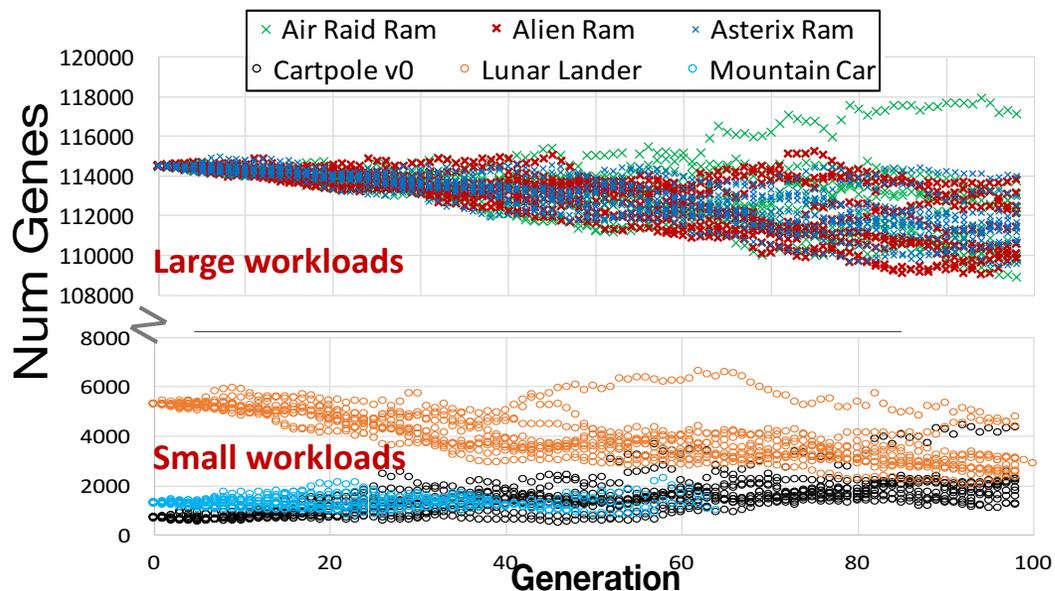


Activation

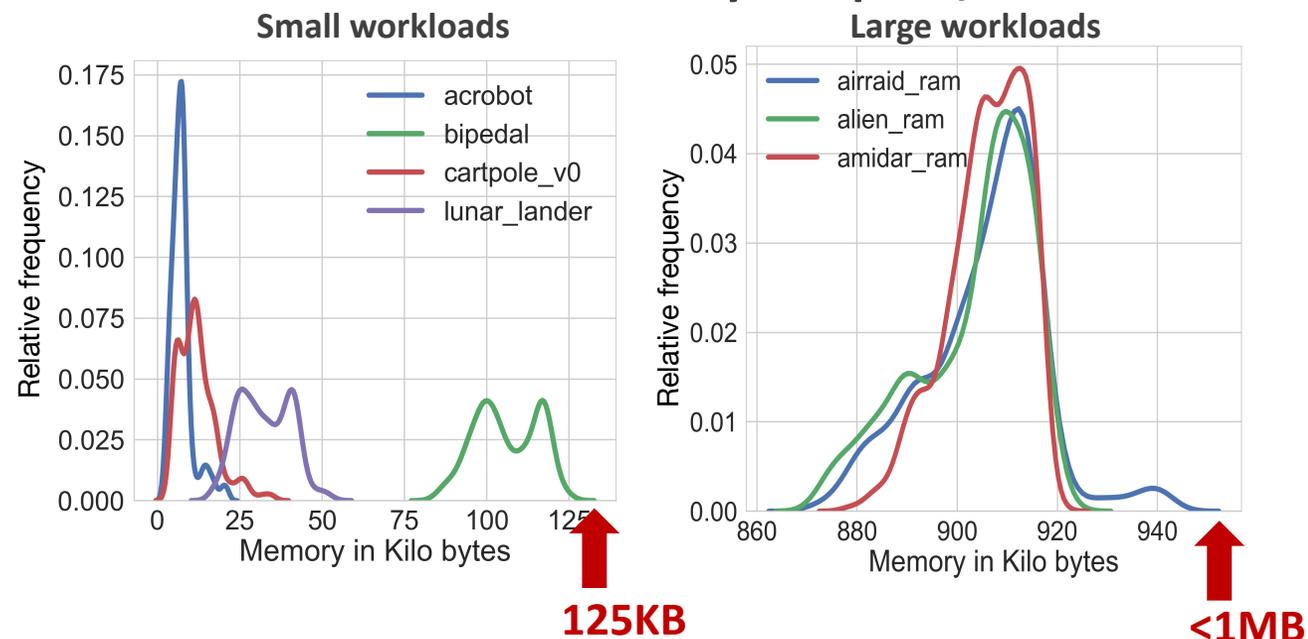
Simple operations

Characterization of NEAT

Memory



Distribution of Memory footprint/Generation



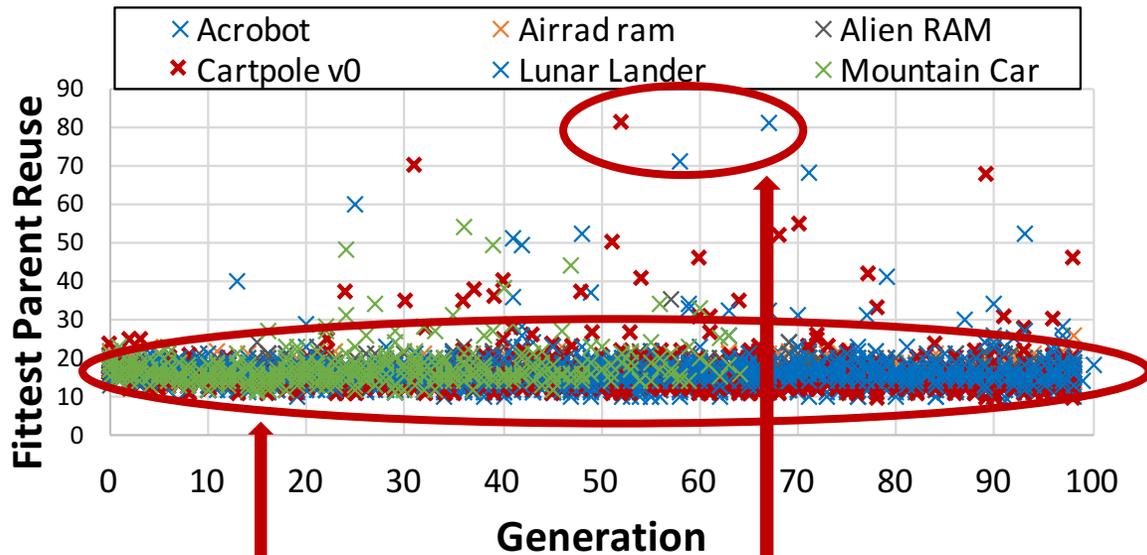
Entire population can fit on-chip

Only need to store the weights and node info

Characterization of NEAT

Memory

Opportunity for Reuse

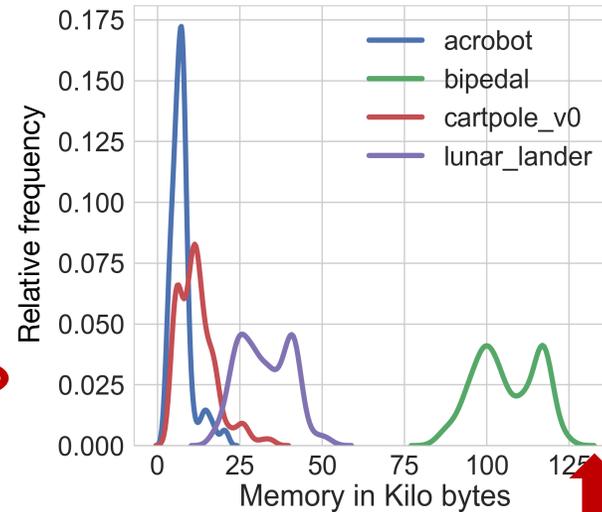


Fittest parent genome is used about ~10-20 times each generation

Even higher in certain cases

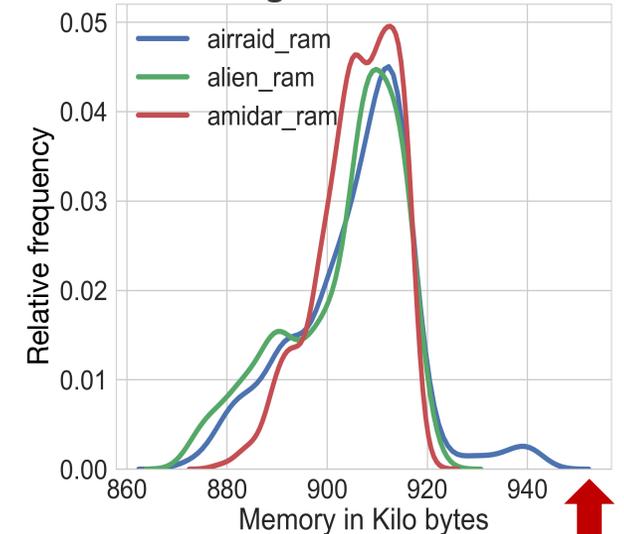
Distribution of Memory footprint/Generation

Small workloads



125KB

Large workloads



<1MB

Entire population can fit on-chip

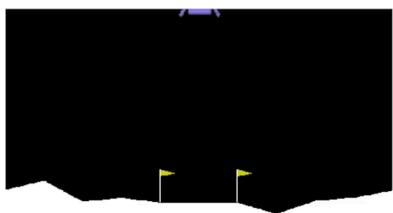
Only need to store the weights and node info

Properties of NE algorithms

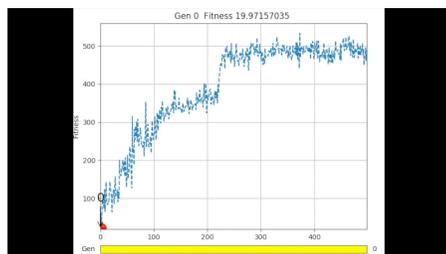
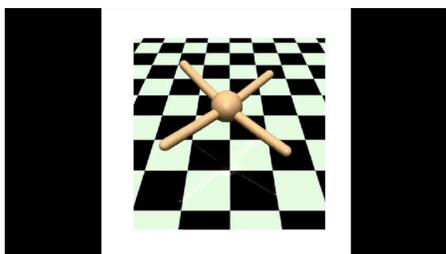
Algorithmic

Robustness

No Training



Change fitness function



Systems

Massive Parallelism

Low Memory Footprint

Genomes within Population

Only store genomes in current generation

Genes within a Genome

No backprop

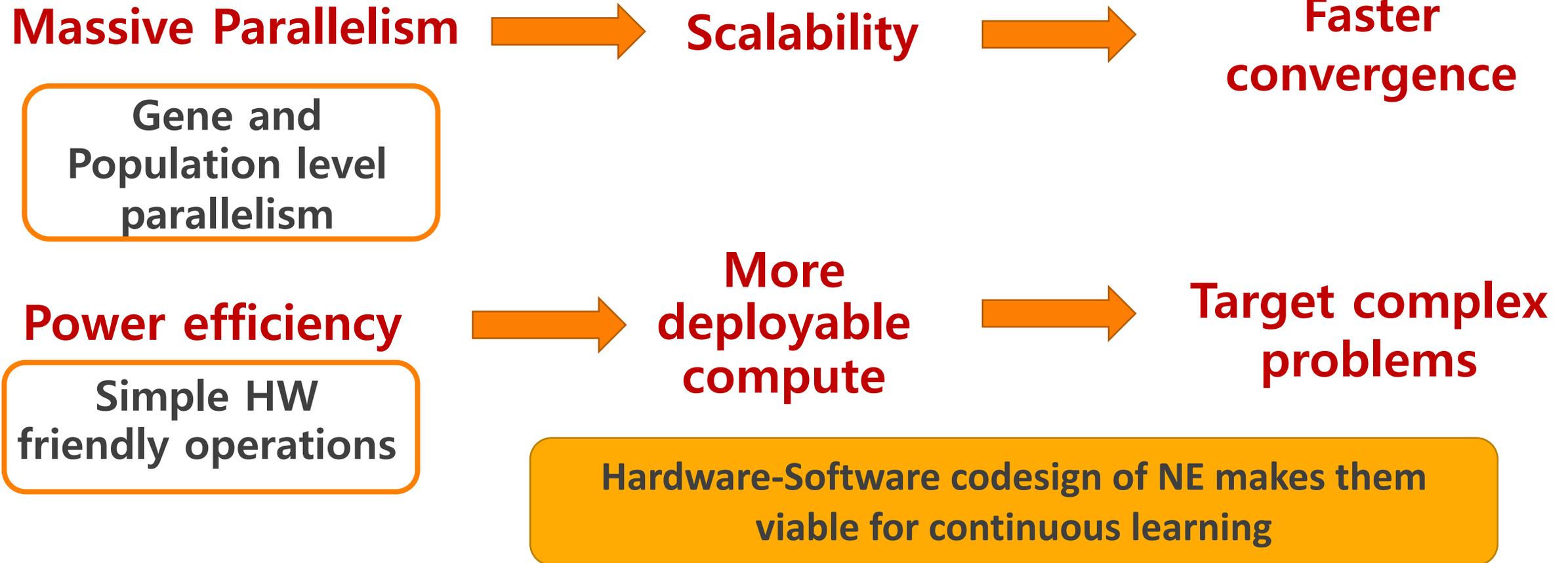
Simple HW-friendly Ops

No gradient calculations or storage

MACs in Inference
Crossover and Mutation in Evolution

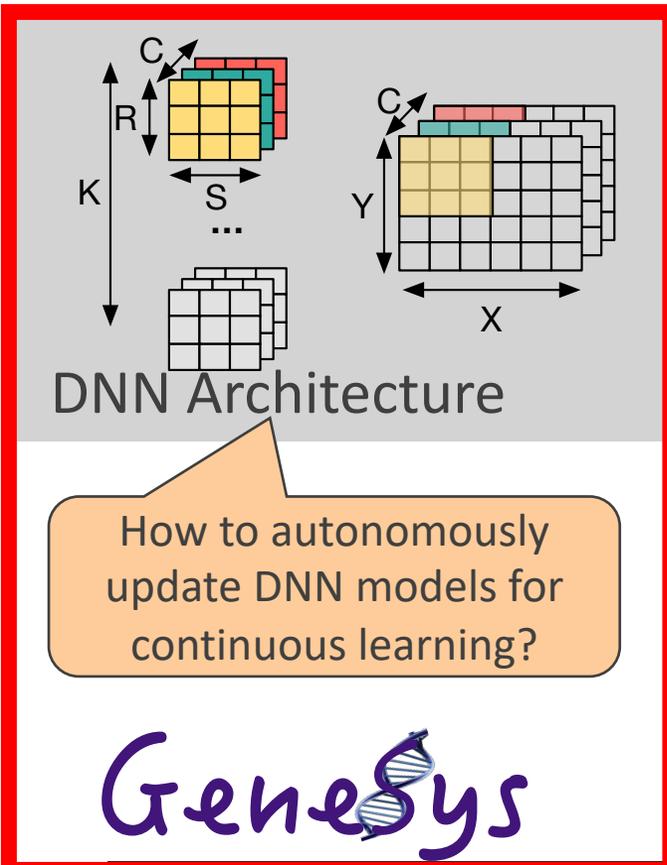
HW-SW Co-Design of NE makes them viable for continuous learning on edge

Motivating Hardware Solution



Outline of Talk

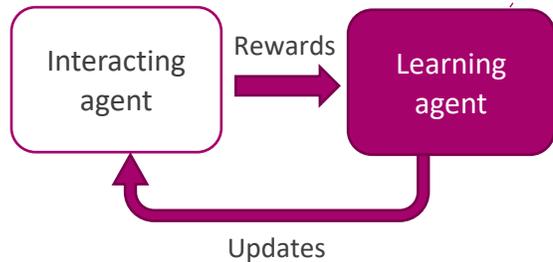
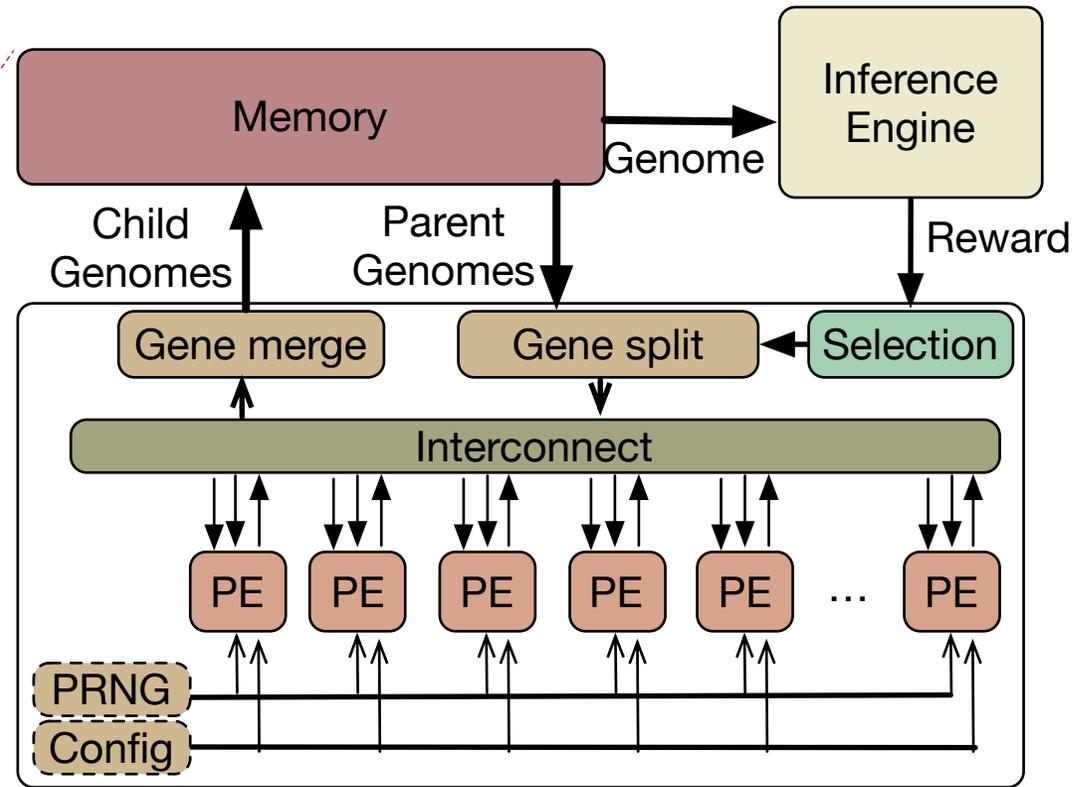
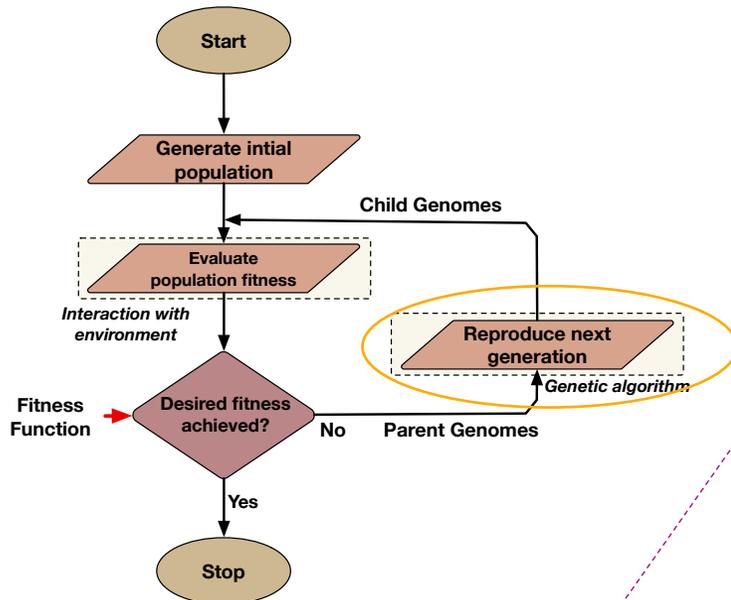
Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna, *GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware*, *MICRO 2018*



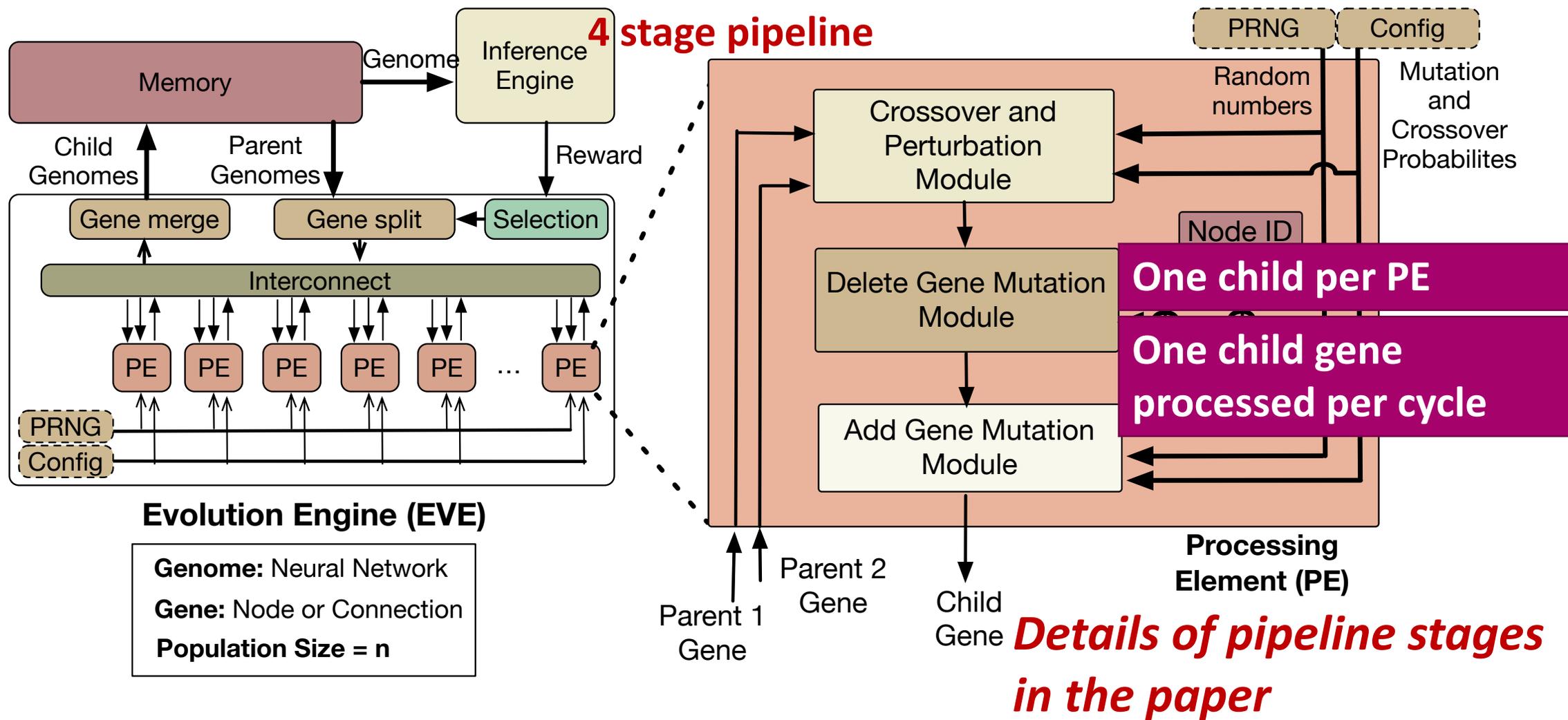
- Continuous Learning
- Neuro-Evolutionary Algorithms
 - Algorithm Description
 - Characterizing NEAT
- **Microarchitecture**
- Evaluations

Evolution Engine: EvE Microarchitecture

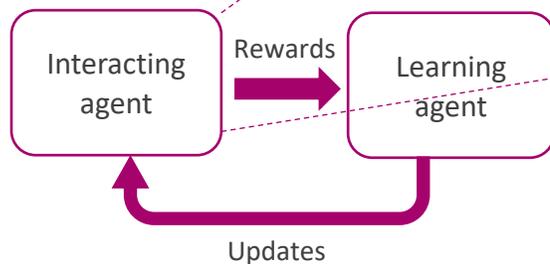
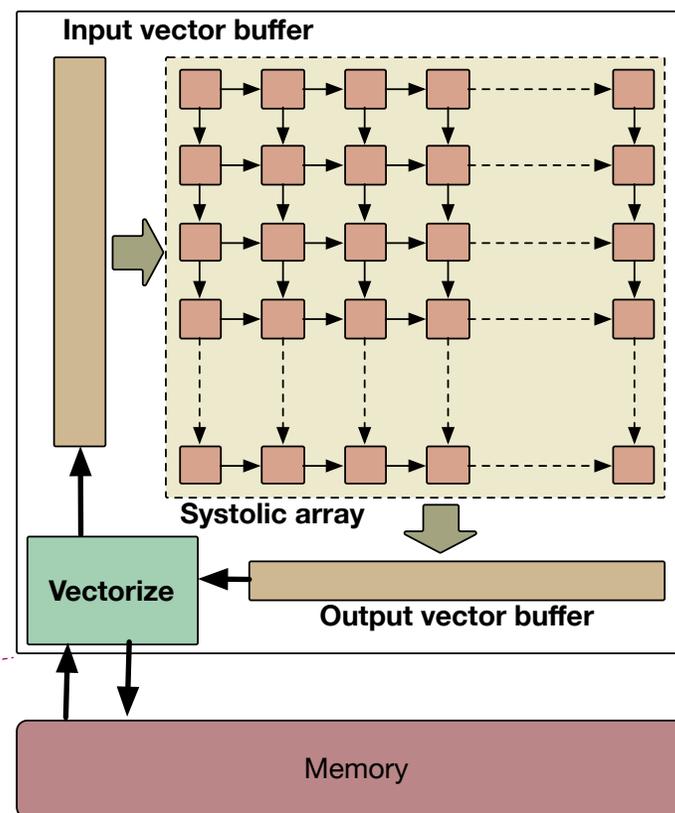
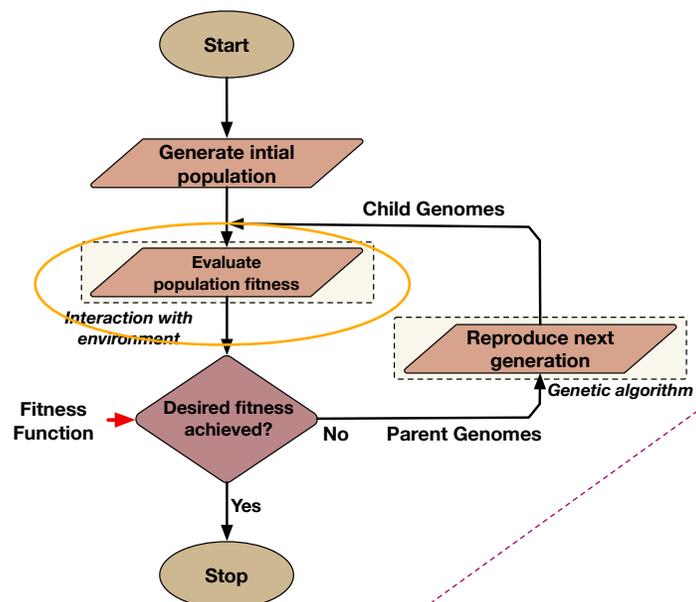
Large number of PE to exploit parallelism



PE Microarchitecture



Inference Engine: ADAM Microarchitecture



Conventional DNN
Inference Accelerator

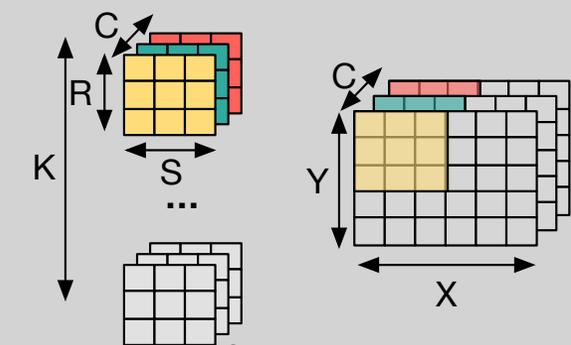
Exploit Population Level
Parallelism

Networks generated by
NEAT are irregular (thus
sparse)

Details later in
talk!

Outline of Talk

Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna, *GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware*, *MICRO 2018*



DNN Architecture

How to autonomously update DNN models for continuous learning?

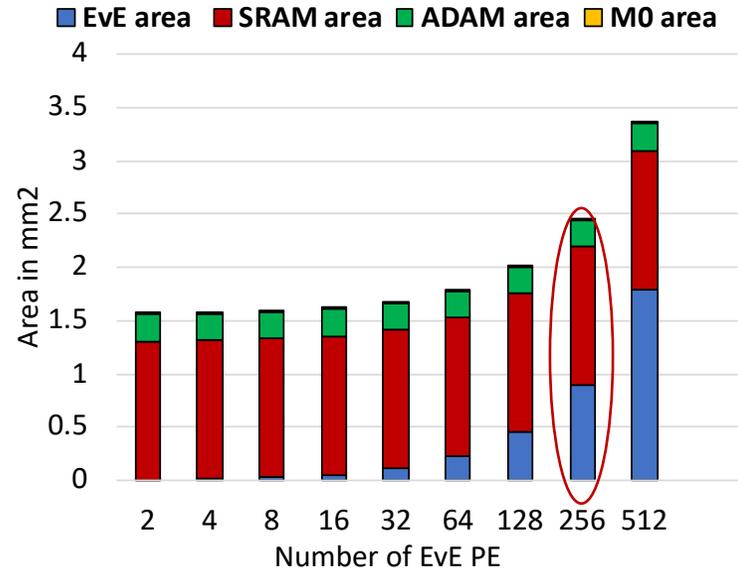
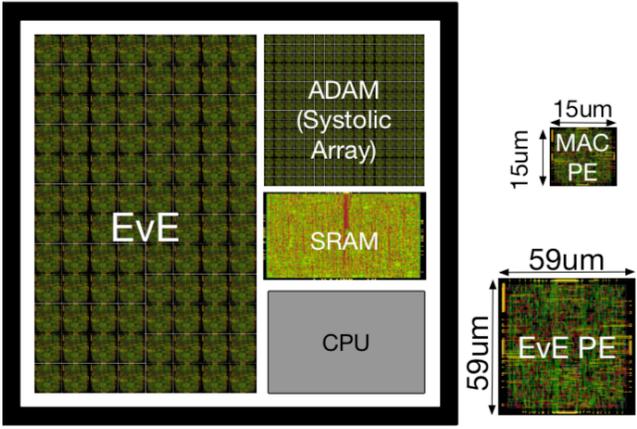
GeneSys

- Continuous Learning
- Neuro-Evolutionary Algorithms
 - Algorithm Description
 - Characterizing NEAT
- Microarchitecture
- Evaluations

Implementation

GeneSys Parameters

Tech node	15nm
Num EvE PE	256
Num ADAM PE	1024
EvE Area	0.89 mm ²
ADAM Area	0.25 mm ²
GeneSys Area	2.45 mm ²
Power	947.5 mW
Frequency	200 MHz
Voltage	1.0 V
SRAM banks	48
SRAM depth	4096



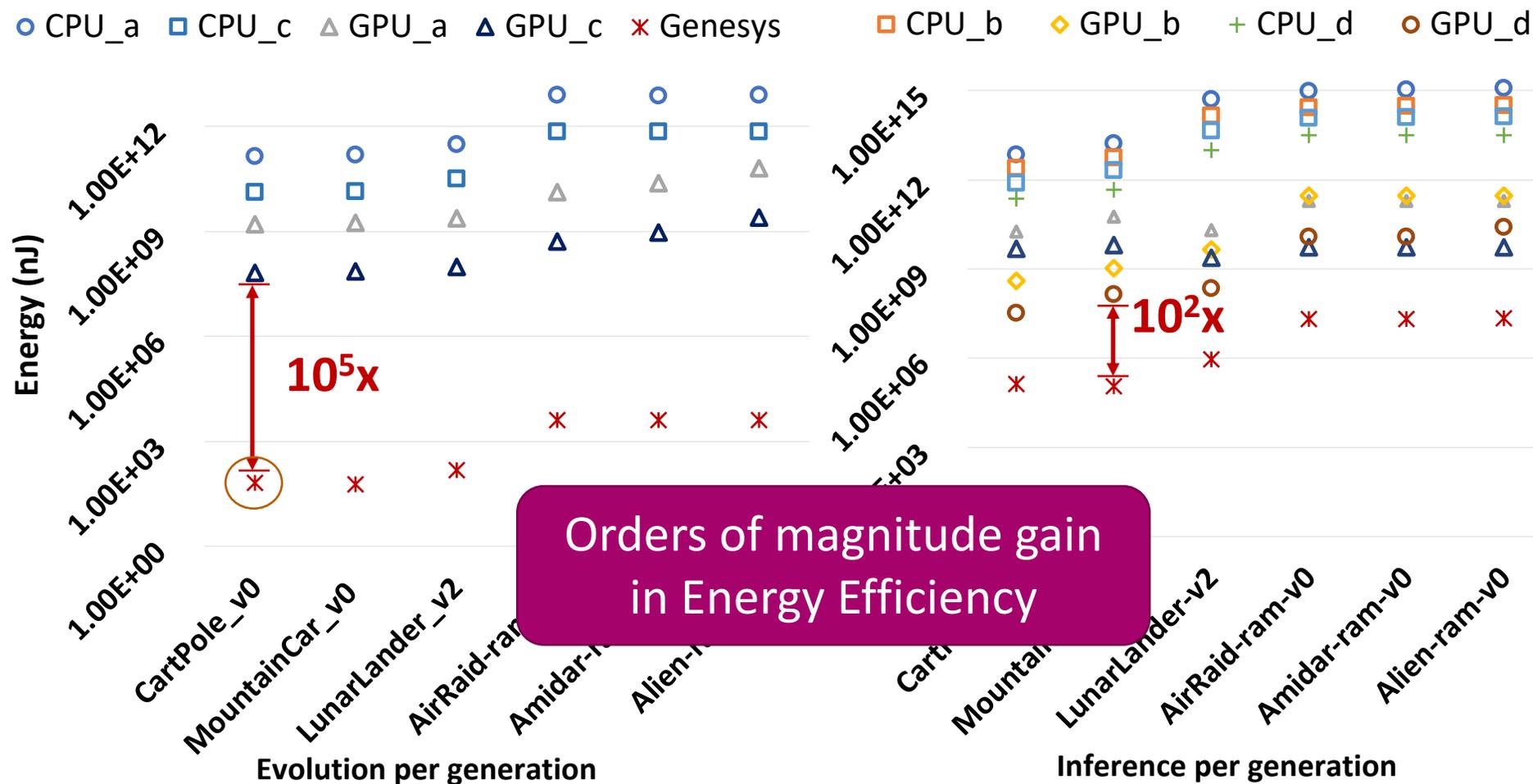
Evaluations

Legend	Inference	Evolution	Platform
CPU_a	Serial	Serial	6th gen i7
CPU_b	PLP	Serial	6th gen i7
GPU_a	BSP	PLP	Nvidia GTX 1080
GPU_b	BSP + PLP	PLP	Nvidia GTX 1080
CPU_c	Serial	Serial	ARM Cortex A57
CPU_d	PLP	Serial	ARM Cortex A57
GPU_c	BSP	PLP	Nvidia Tegra
GPU_d	BSP + PLP	PLP	Nvidia Tegra
GENESYS	PLP	PLP + GLP	GENESYS

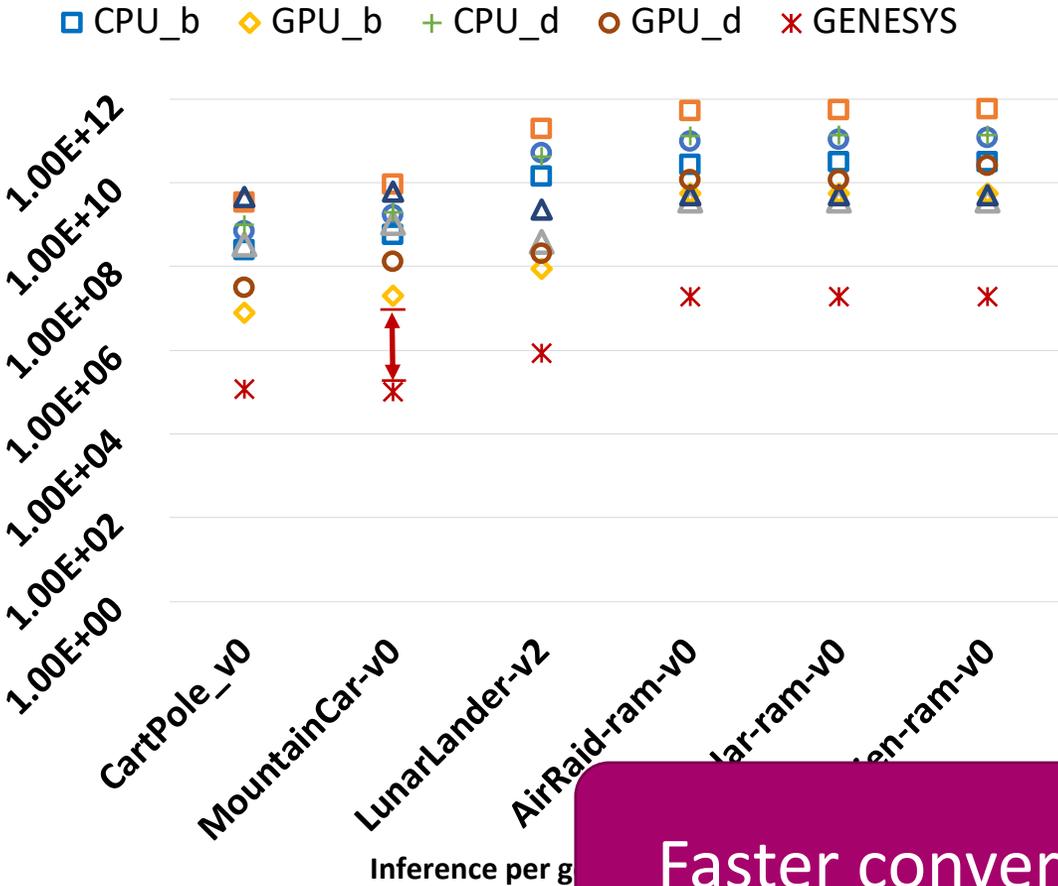
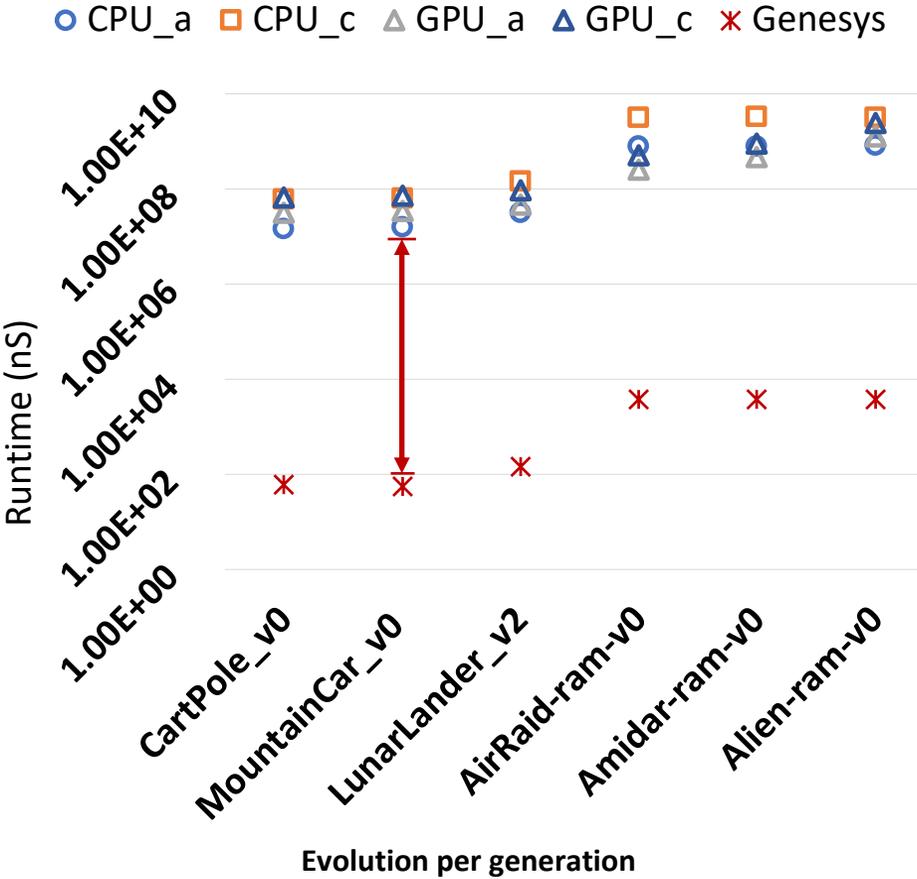
PLP (GLP) - Population (Gene) Level Parallelism

BSP - Bulk Synchronous Parallelism (GPU)

Evaluations: Energy



Evaluations: Runtime



Faster convergence

Summary for GeneSys

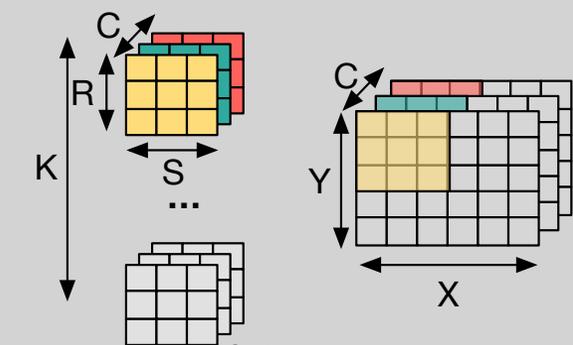
- **Robust, Scalable** and **Energy** efficient solutions needed for continuous learning
 - Look **beyond DL and RL**
- NEs offer promise
 - Parallelism
 - Low-memory Footprint
 - HW friendly
- GeneSys: *100x – 100000x energy efficiency and performance*
 - **More deployable compute**
 - **Enables** AI solutions for a large gamut of problems

Outline of Talk

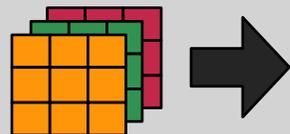
Training



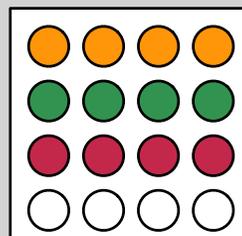
Inference



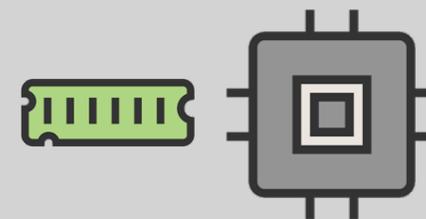
DNN Architecture



Mapping (Dataflow)



Microarchitecture



Energy



Runtime

How to autonomously design DNN models for continuous learning?

How to efficiently map changing DNNs over accelerator?

How to design an efficient accelerator for changing DNN models

GeneSys

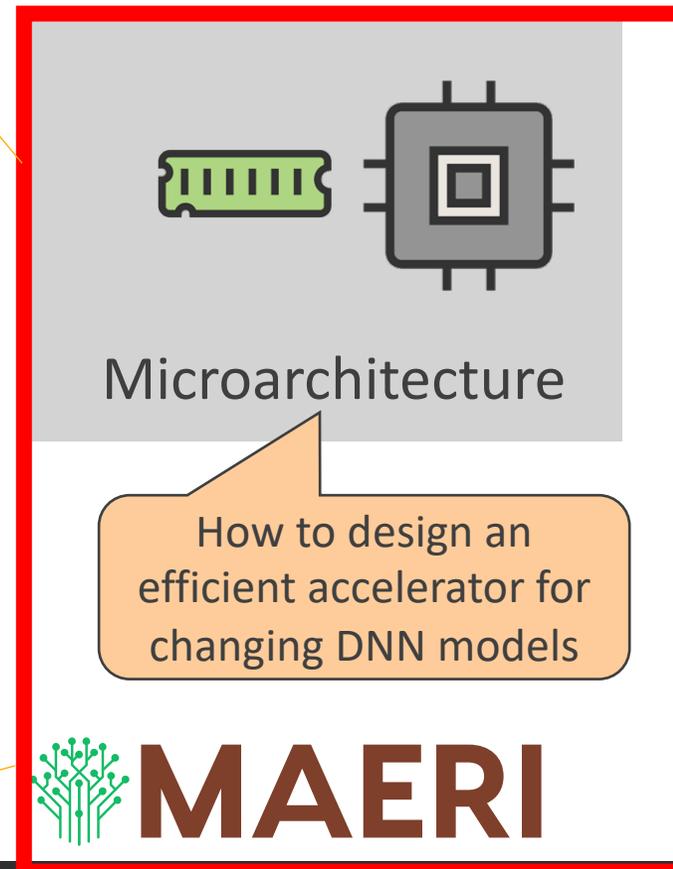
MAERI

Outline of Talk

- Motivation
 - Irregular Dataflows
 - DNN Computation
- MAERI
 - Abstraction
 - Implementation
 - Operation Example
 - Mapping Strategies
- Evaluations

Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:

ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention



Myriad Dataflows in DNN Accelerators

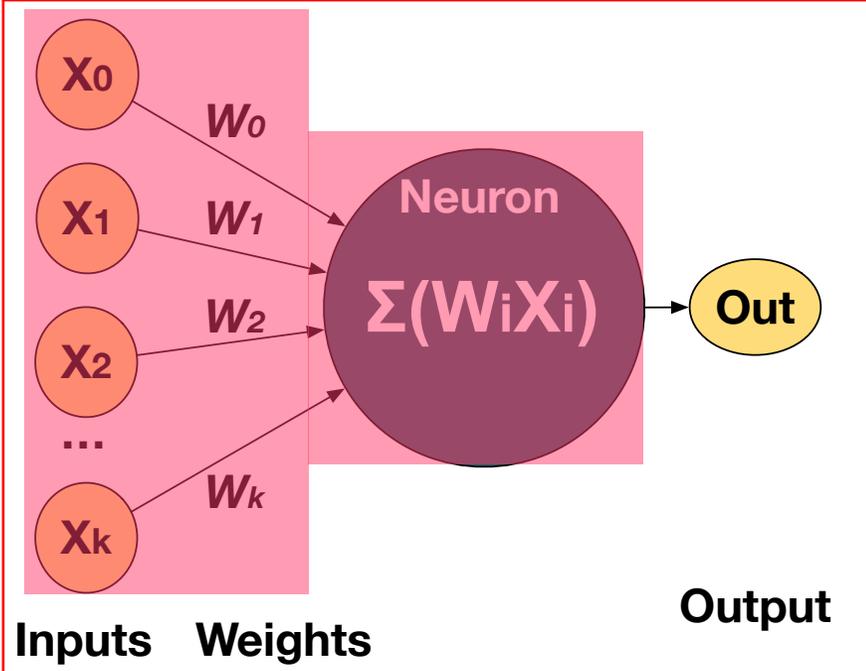
- **DNN Topologies**
 - Layer **size / shape**
 - Layer **types**: Convolution / Pool / FC / LSTM
 - New **sub-structure**: e.g., Inception in Googlenet
- **Compiler/Mapper**
 - Loop Scheduling
 - **Reordering and Tiling**
 - Mapping
 - **Output/Weight/Input/Row-stationary**
- **Algorithmic Optimization (e.g., Sparsity)**
 - Weight pruning
 - GeneSys



Can we have **one architectural solution** that can handle arbitrary dataflows and provides ~100% utilization?

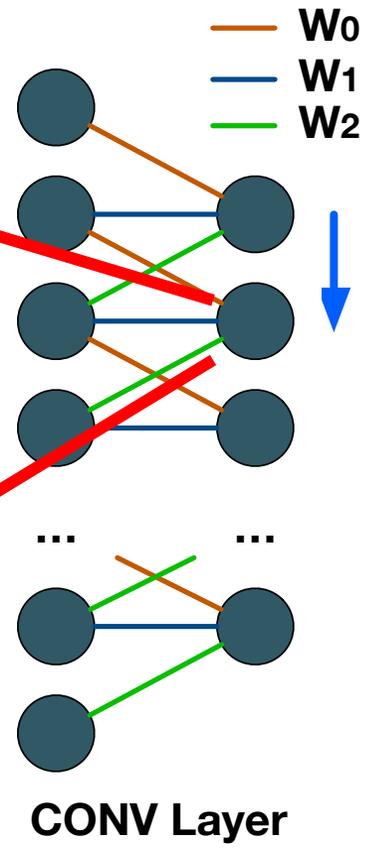
What is the computation in a DNN?

Independent multiplication



Compute weighted sum

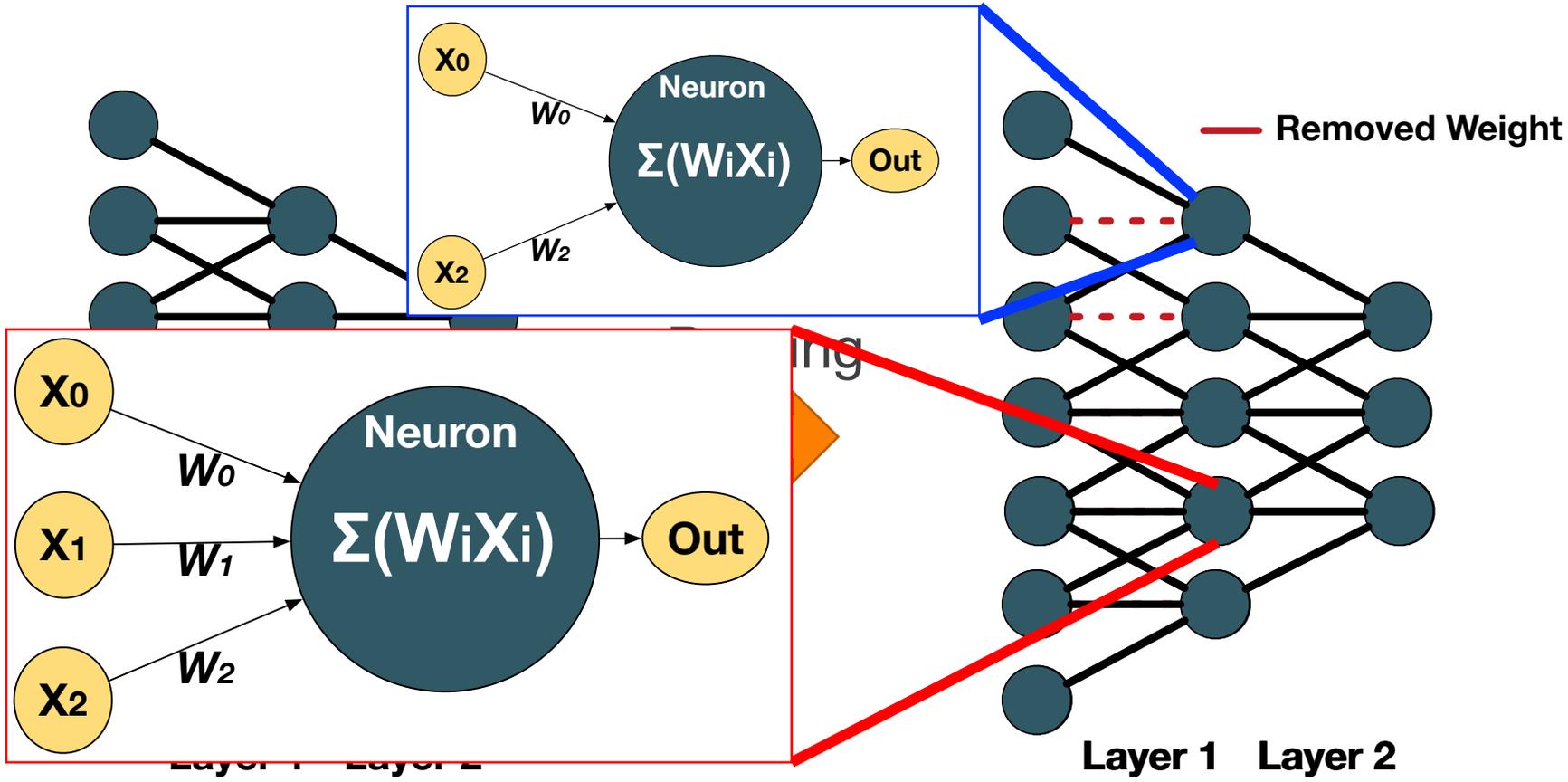
Accumulation of partial products



Our Key insight: Each DNN/dataflow translates into neurons of different sizes

Irregular Dataflow: Pruning

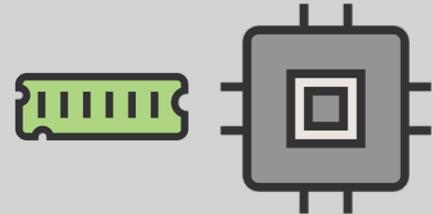
Example: Weight Pruning (Sparse Workload)



Our Key insight: Each DNN/dataflow translates into neurons of different sizes

Outline of Talk

- Motivation
 - Irregular Dataflows
 - DNN Computation
- **MAERI**
 - **Abstraction**
 - Implementation
 - Operation Example
 - Mapping Strategies
- Evaluations

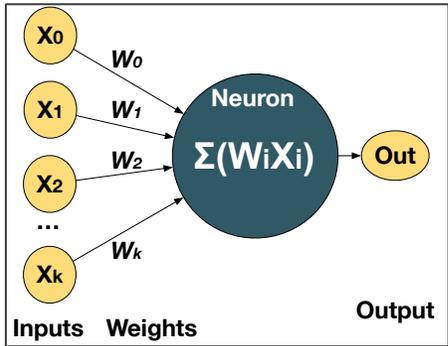


Microarchitecture

How to design an efficient accelerator for changing DNN models



The MAERI Abstraction



Multiplier Pool



Adder Pool



Virtual Neuron (VN): Temporary grouping of compute units for an output

How to enable flexible grouping?

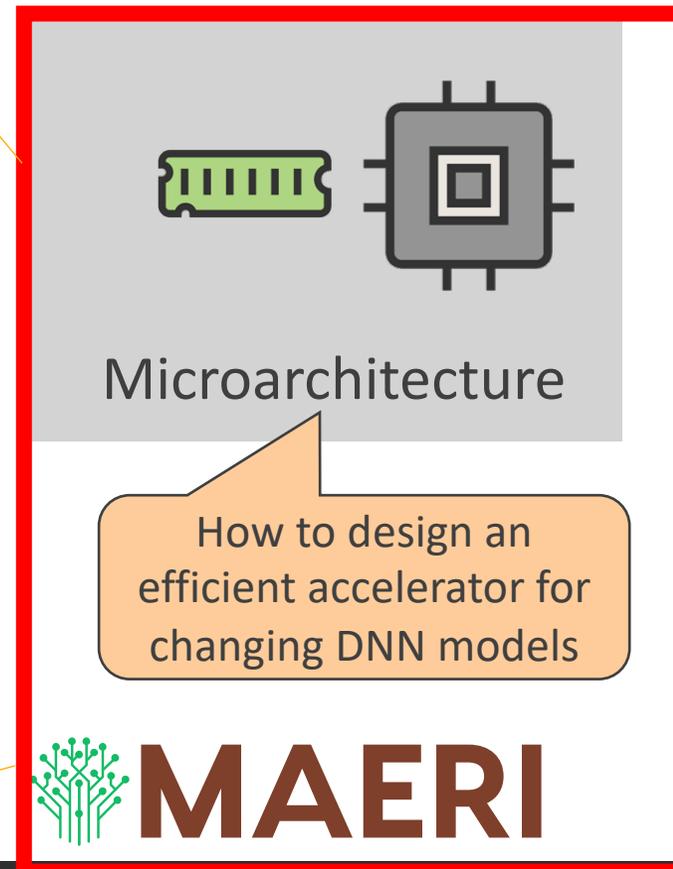
Need flexible connectivity!

Outline of Talk

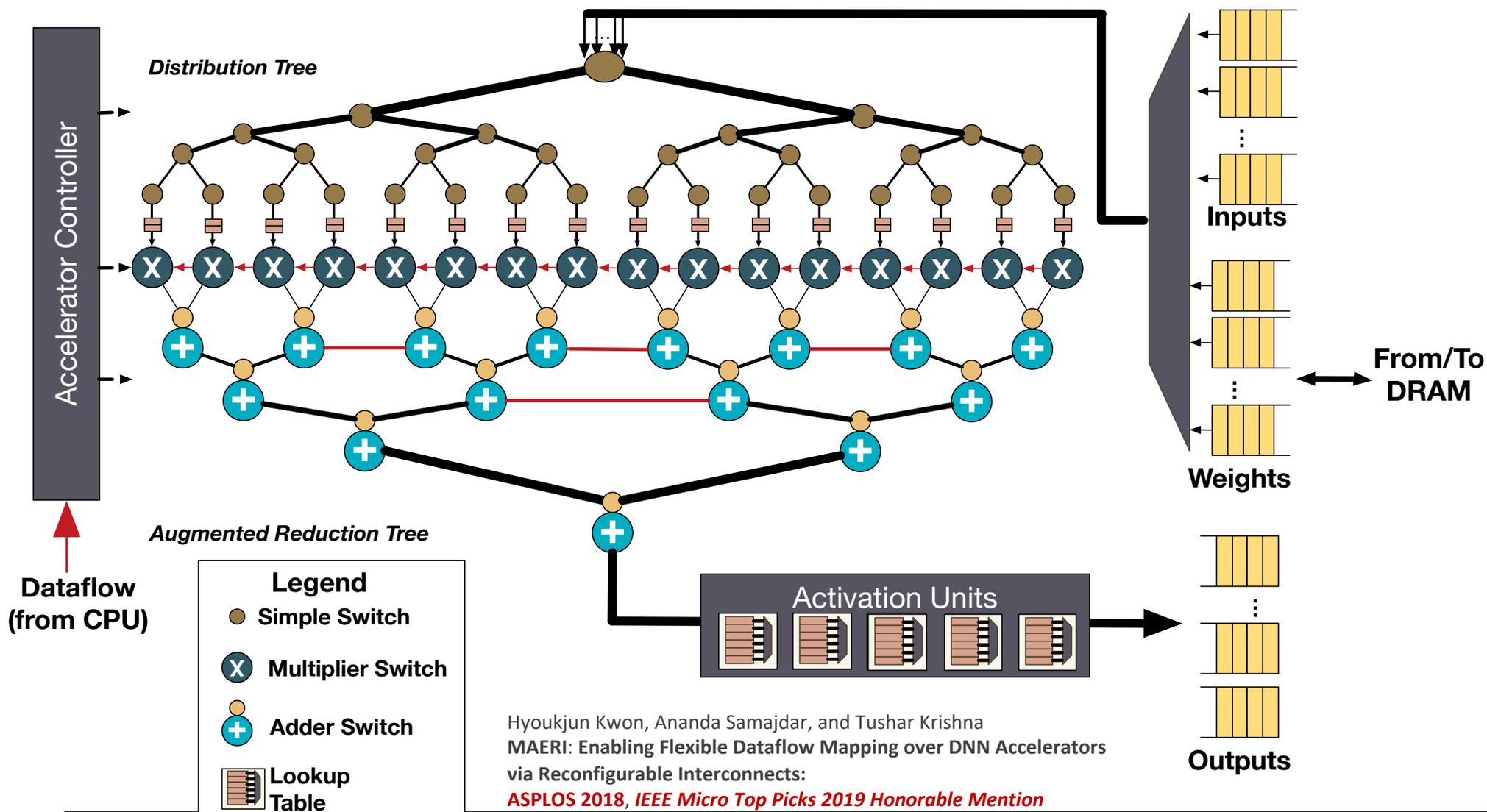
- Motivation
 - Irregular Dataflows
 - DNN Computation
- **MAERI**
 - Abstraction
 - **Implementation**
 - Operation Example
 - Mapping Strategies
- Evaluations

Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:

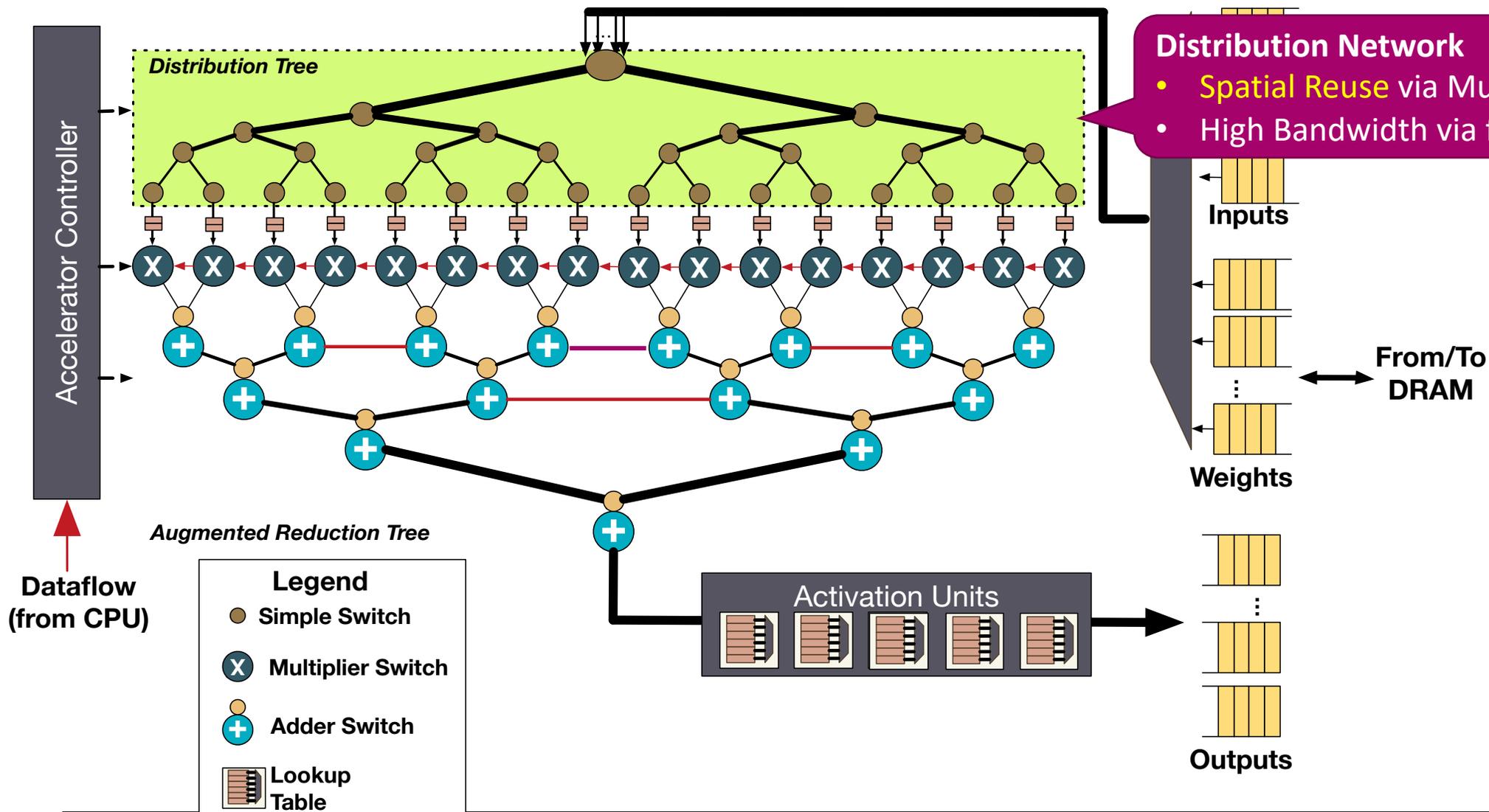
ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention



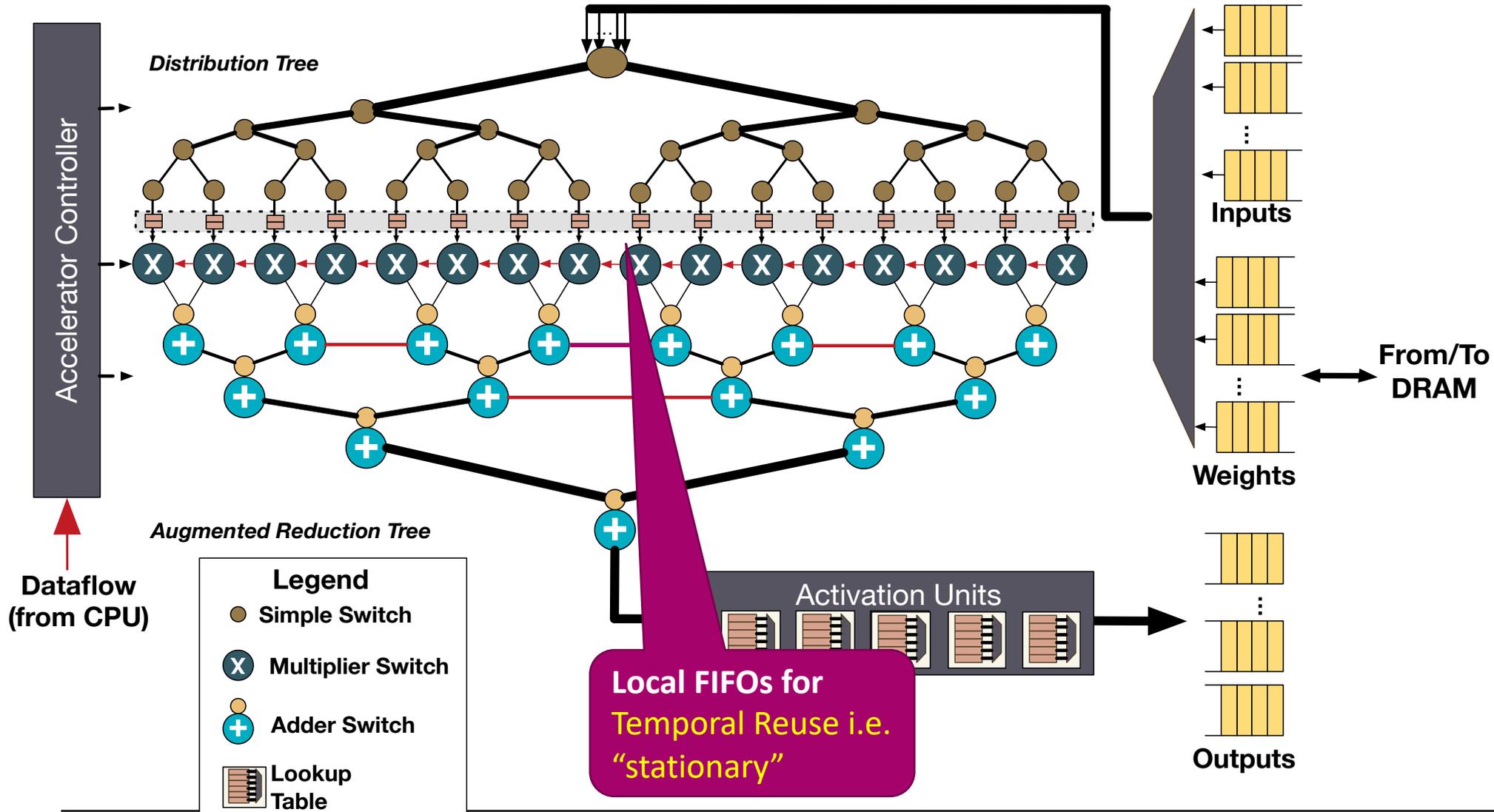
The MAERI Implementation



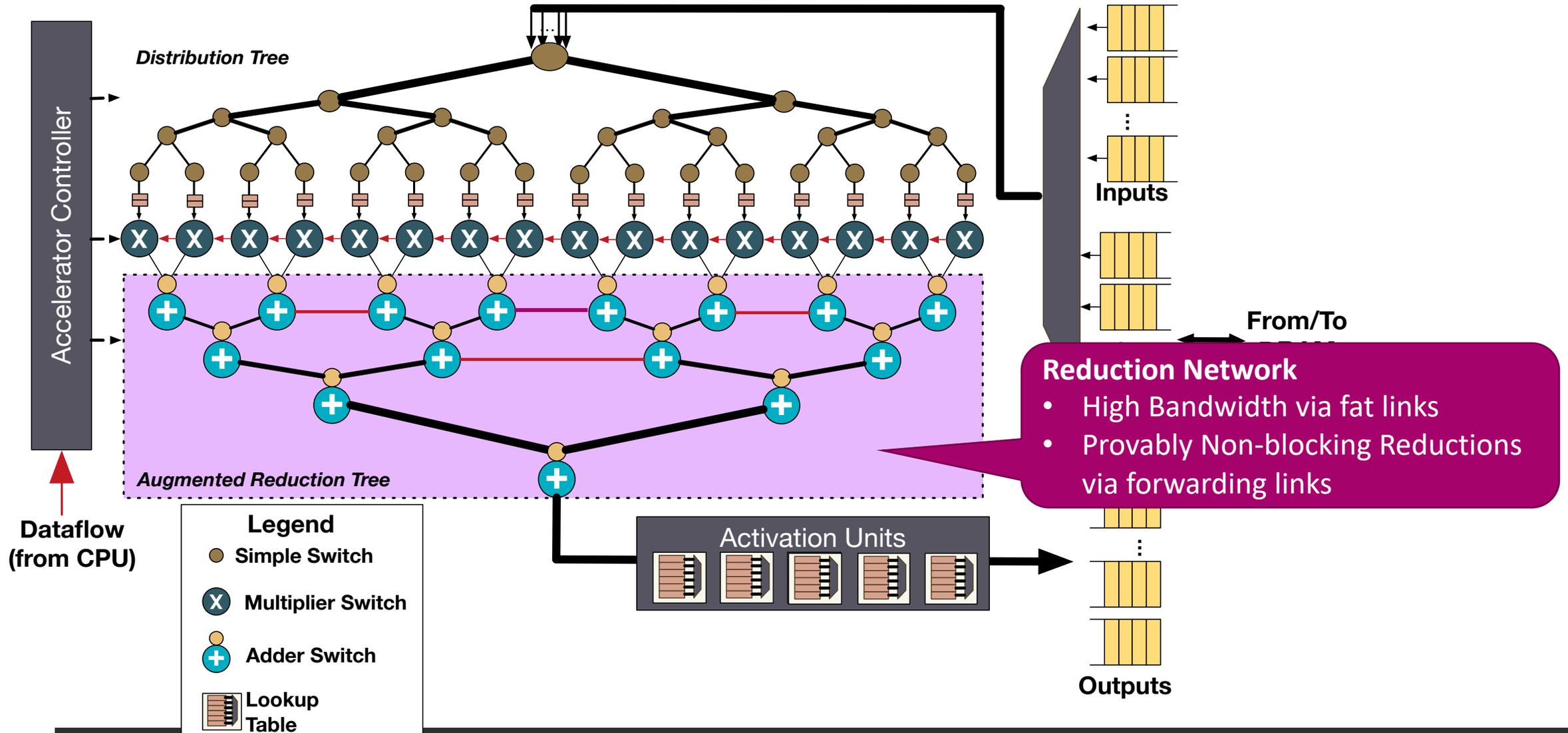
The MAERI Implementation



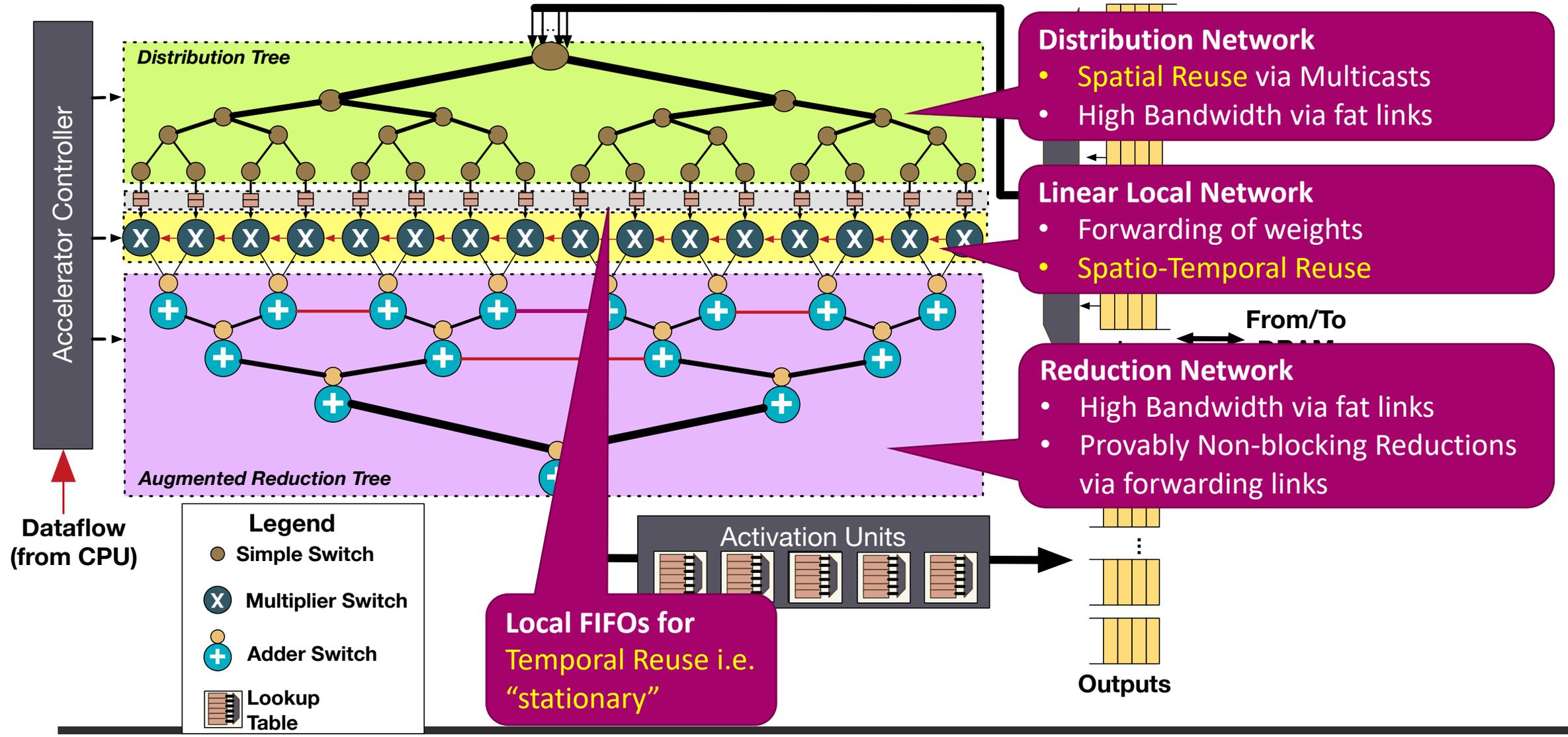
The MAERI Implementation

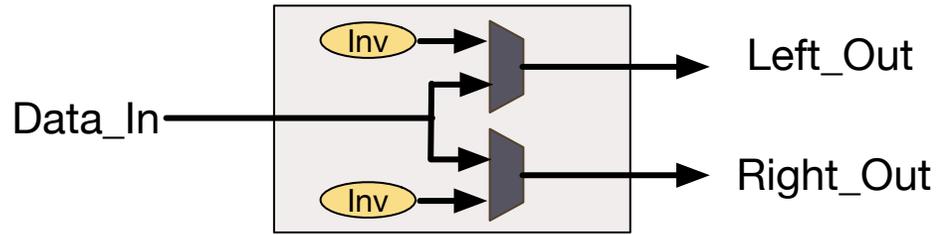


The MAERI Implementation

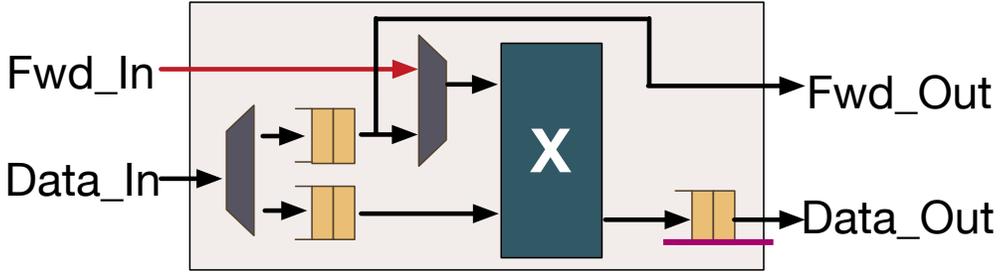


The MAERI Implementation

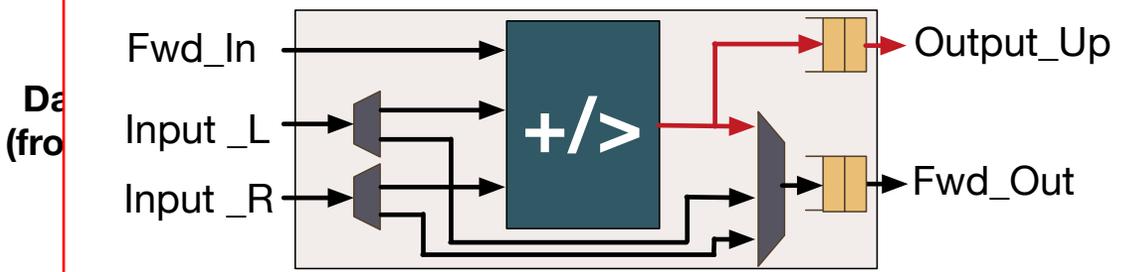




Distribute Switch (1x2 Switch)

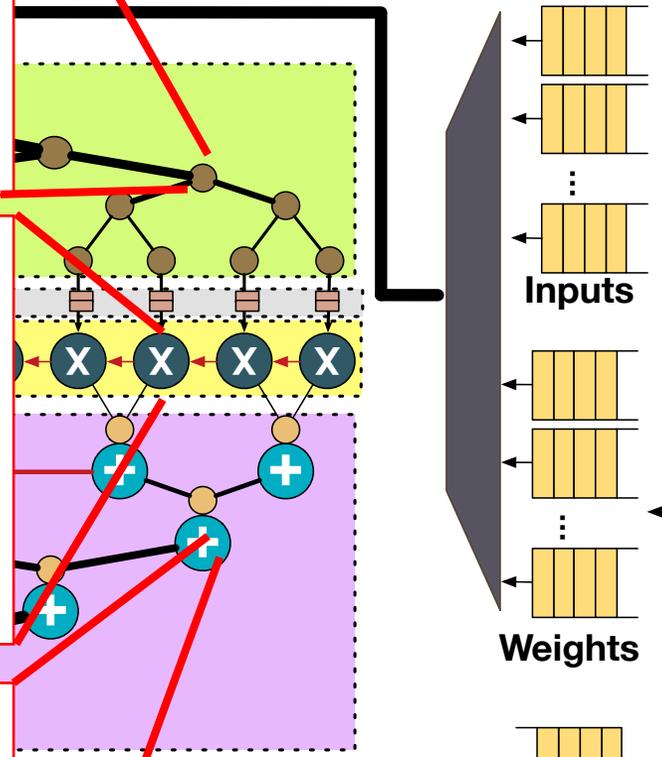


Multiplier Switch (multiplier+ 2x2 switch)

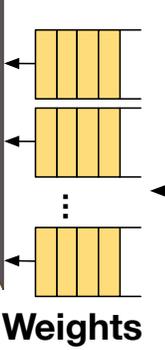
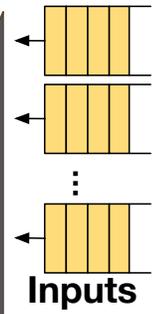


Adder Switch (adder+ 3x2 switch)

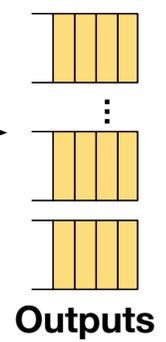
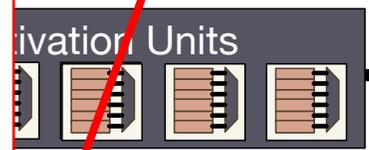
Representation



Micro-Switches



From/To DRAM

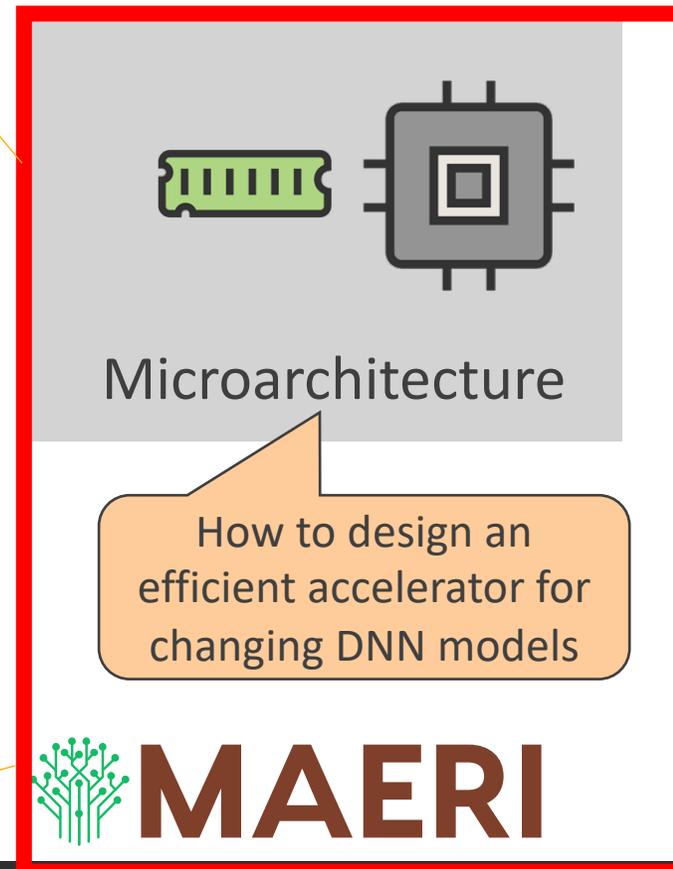


Outline of Talk

- Motivation
 - Irregular Dataflows
 - DNN Computation
- **MAERI**
 - Abstraction
 - Implementation
 - **Operation Example**
 - ▶ • Mapping Strategies
- Evaluations

Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:

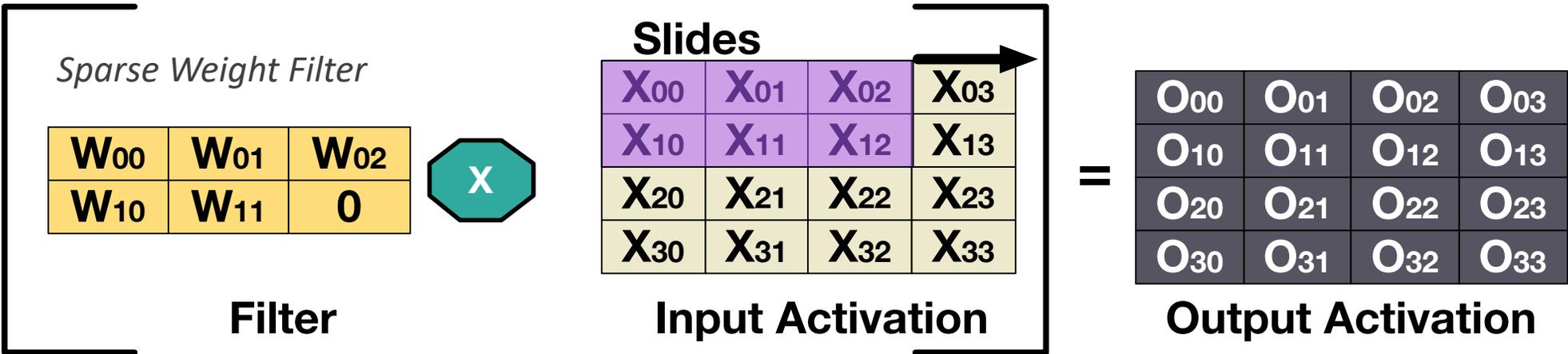
ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention



Example: Computing a CONV layer

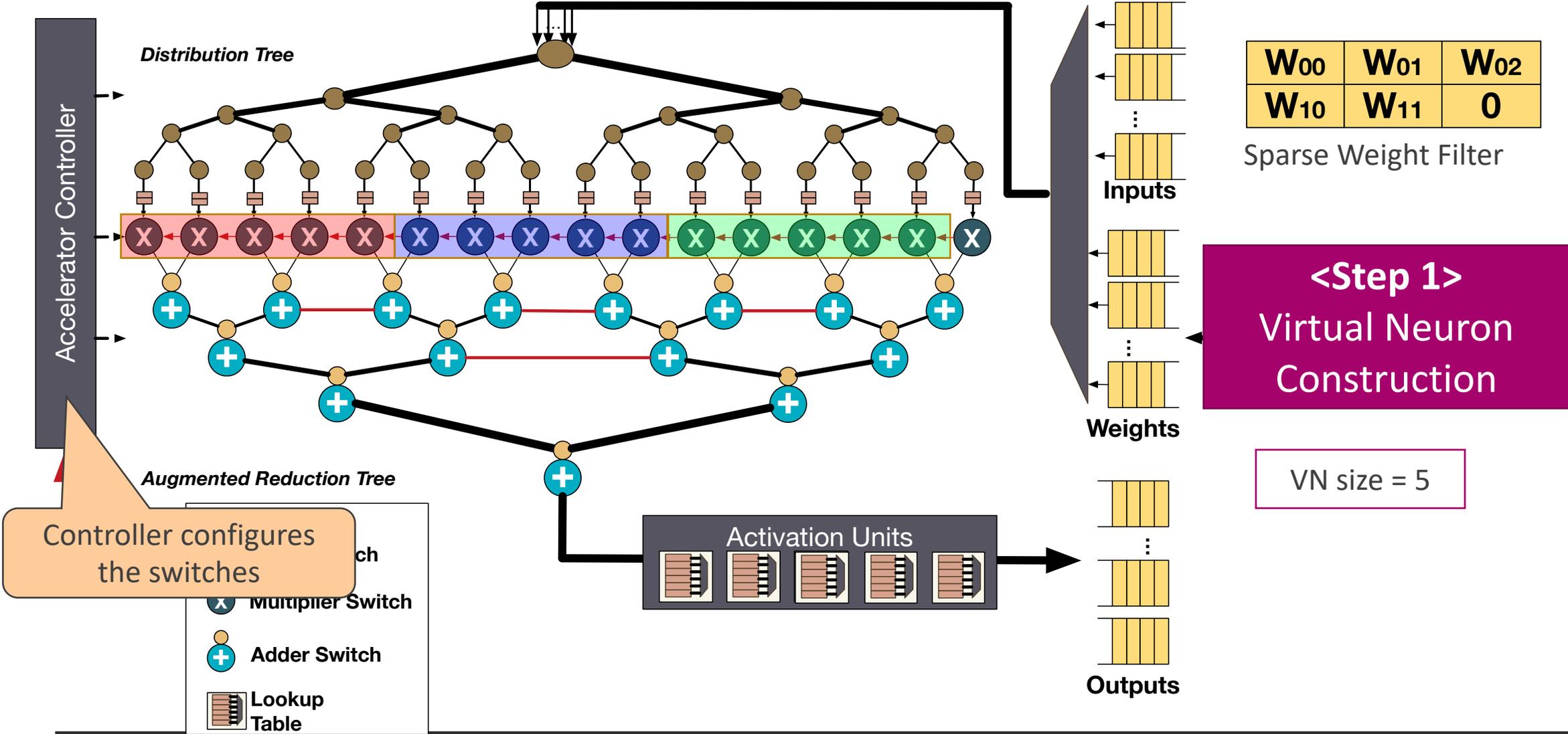
- **[Communication]** Distribute weights and inputs (image pixels) to multiplier switches
 - *Assume: weight stationary, conv reuse of inputs via local links*
- **[Computation]** Compute partial sums
- **[Computation]** Reduce partial sums
- **[Communication]** Collect outputs to buffer

MAERI Operation Example

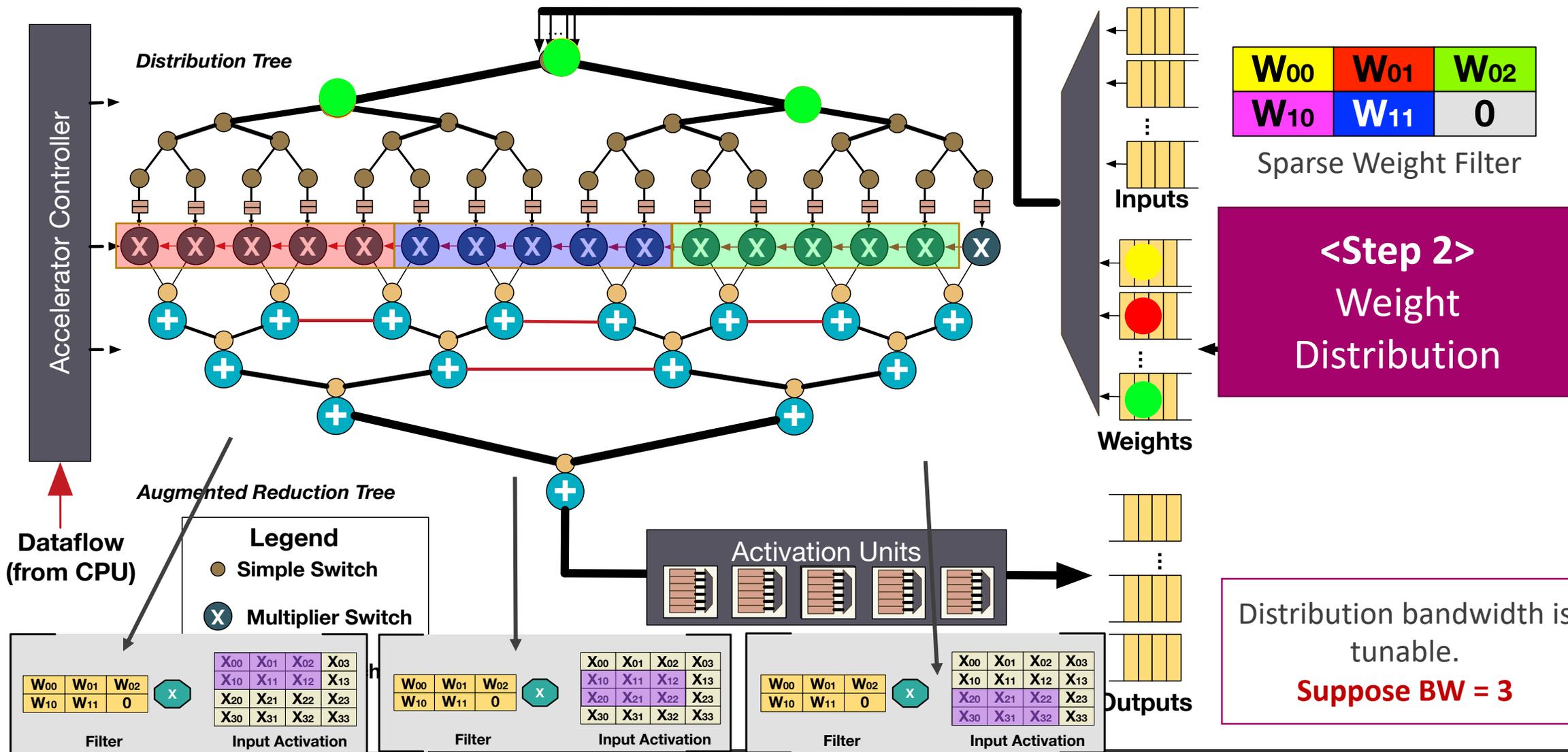


$$O_{00} = [W_{00} \times X_{00}] + [W_{01} \times X_{01}] + [W_{02} \times X_{01}] + [W_{10} \times X_{10}] + [W_{11} \times X_{11}]$$

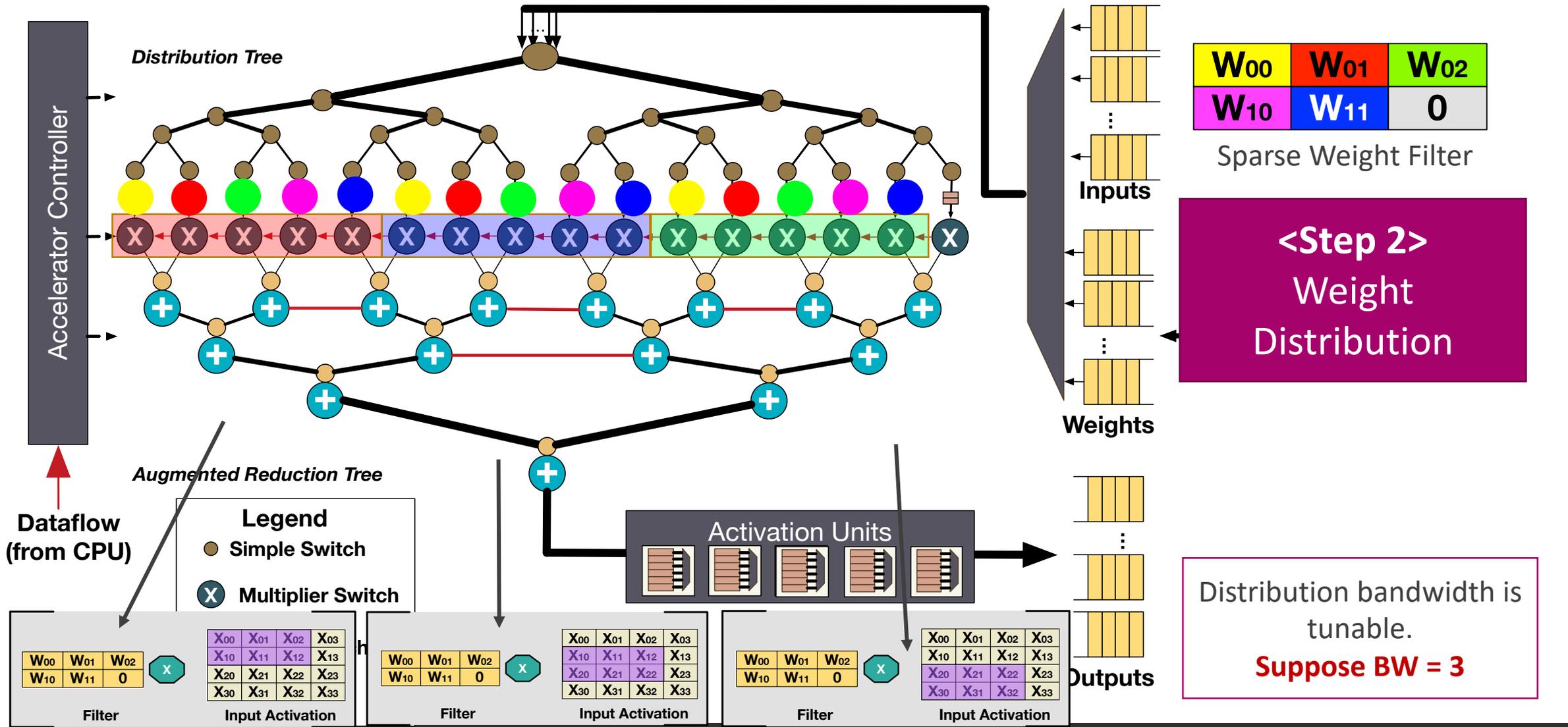
MAERI Operation Example



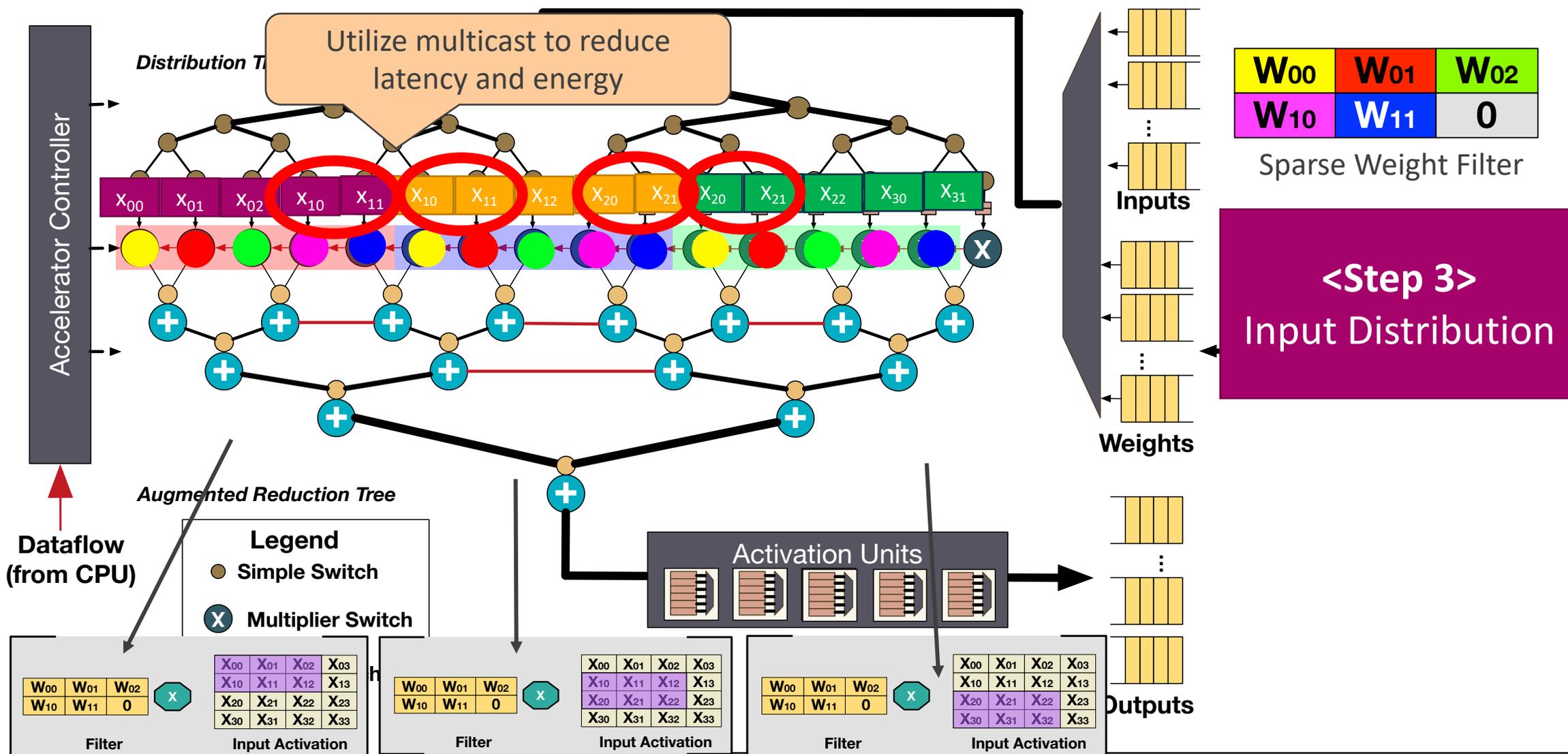
MAERI Operation Example



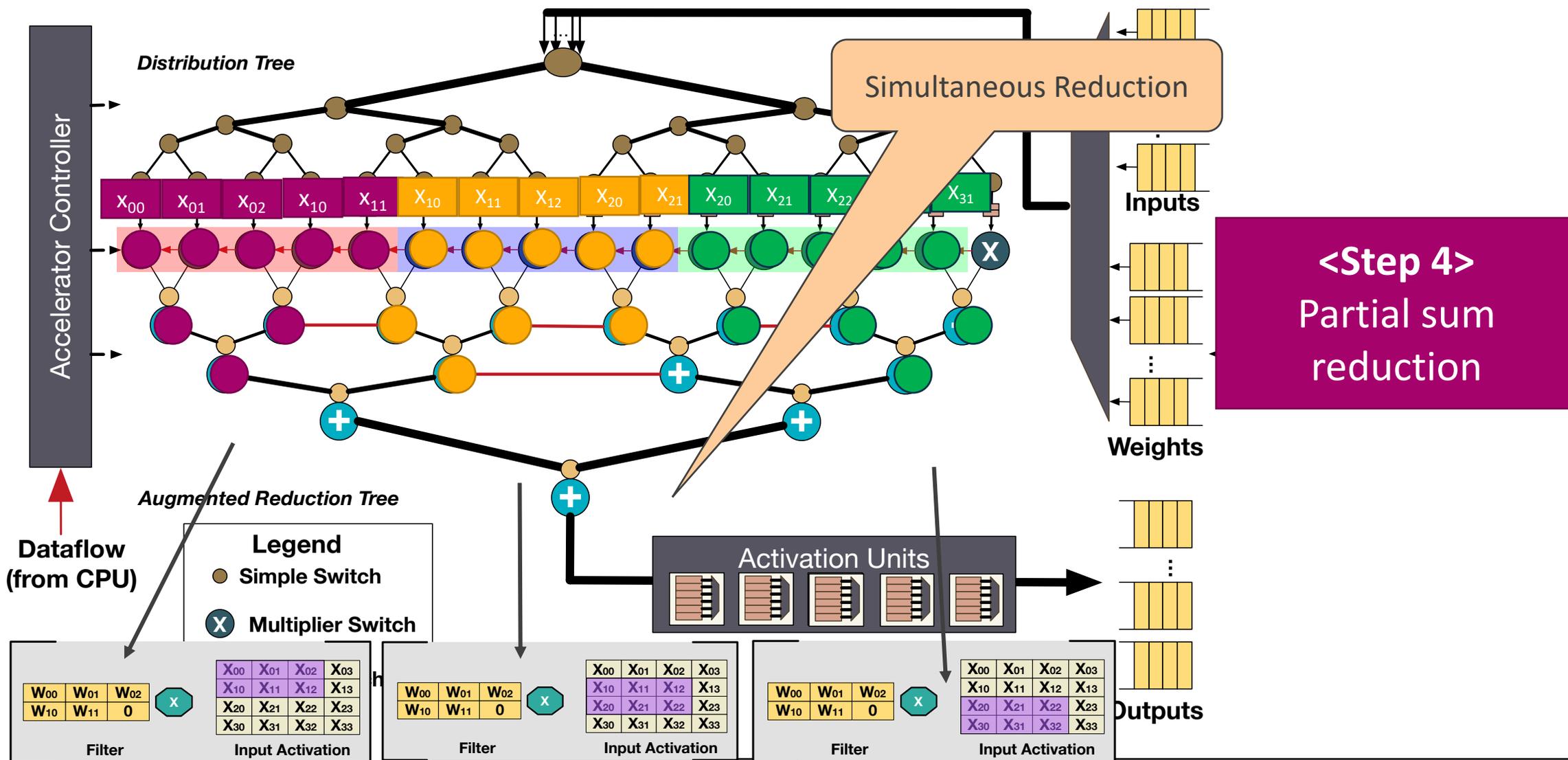
MAERI Operation Example



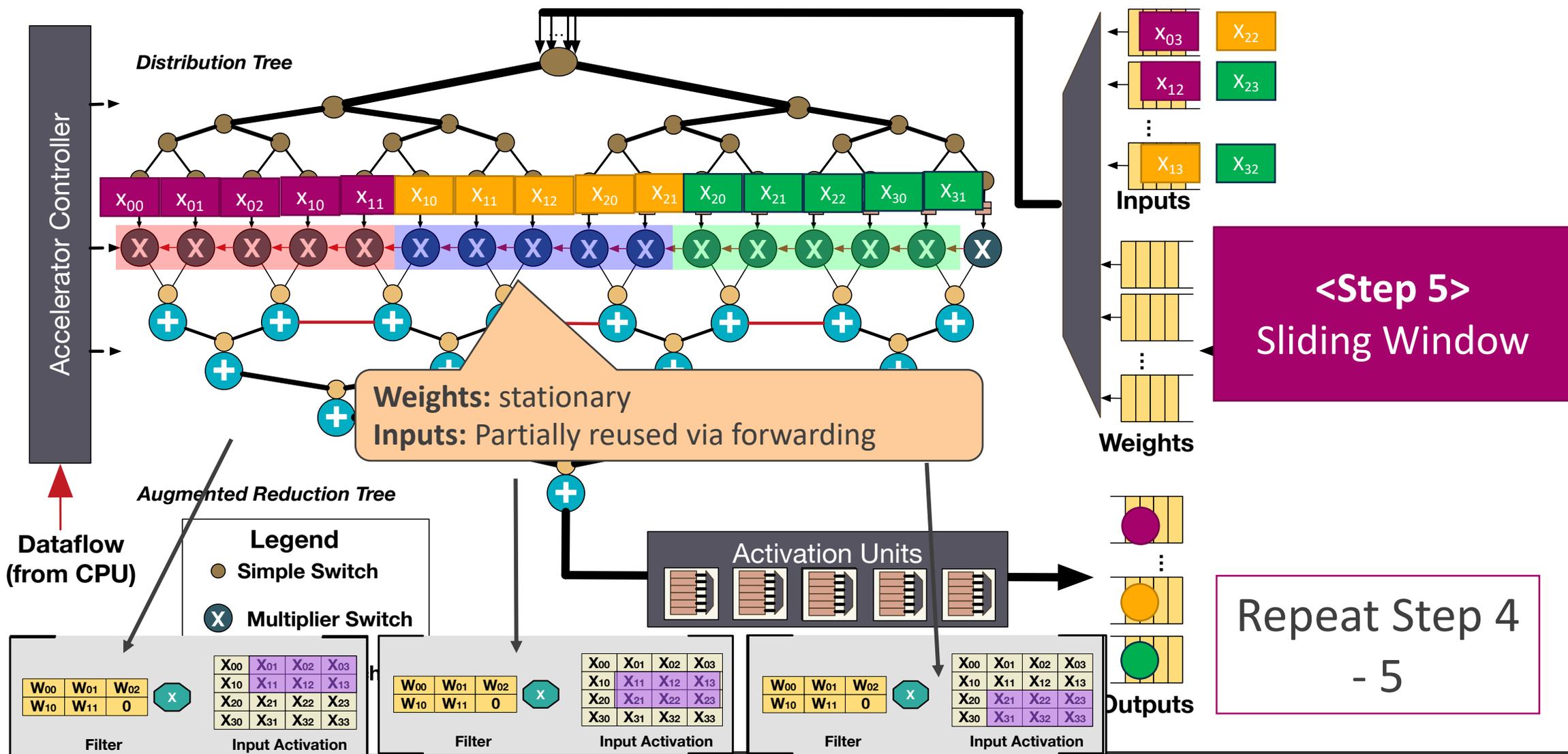
MAERI Operation Example



MAERI Operation Example



MAERI Operation Example



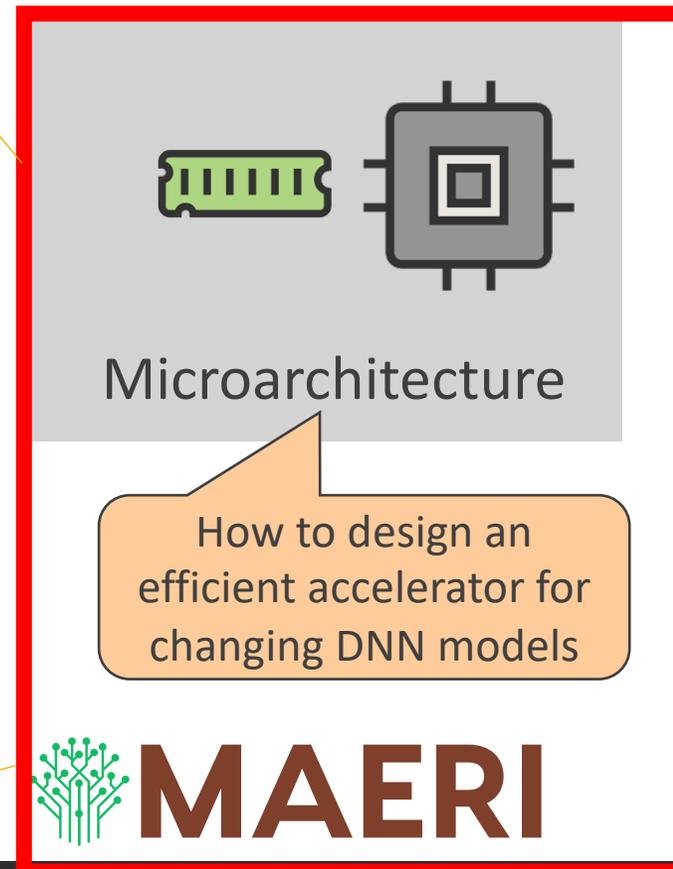
Outline of Talk

- Motivation
 - Irregular Dataflows
 - DNN Computation
- **MAERI**
 - Abstraction
 - Implementation
 - Operation Example
 - **Mapping Strategies**
- Evaluations

Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna

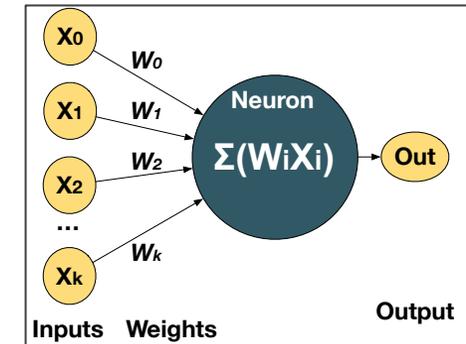
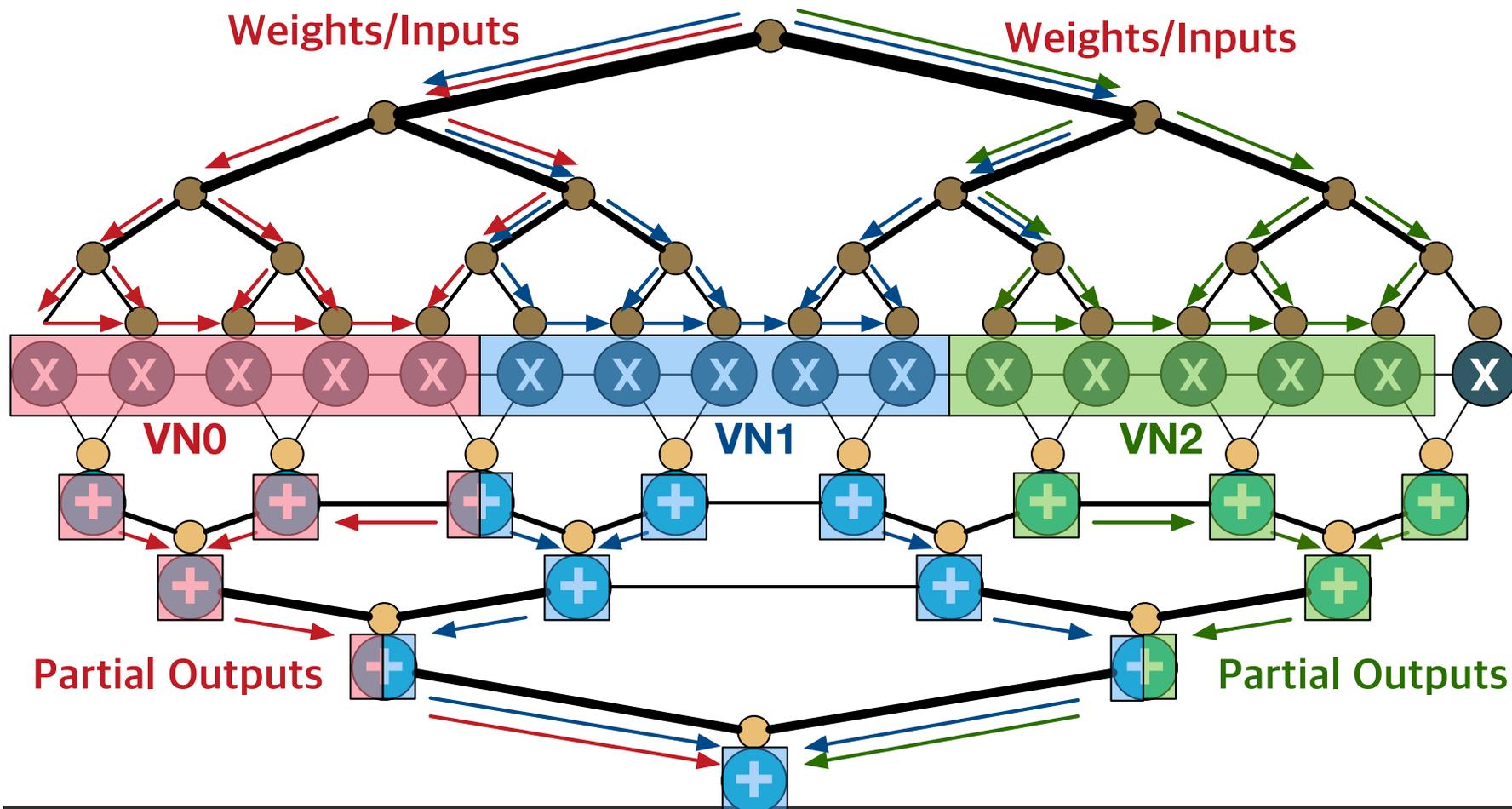
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:

ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention



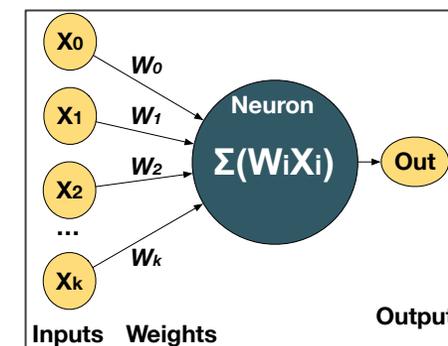
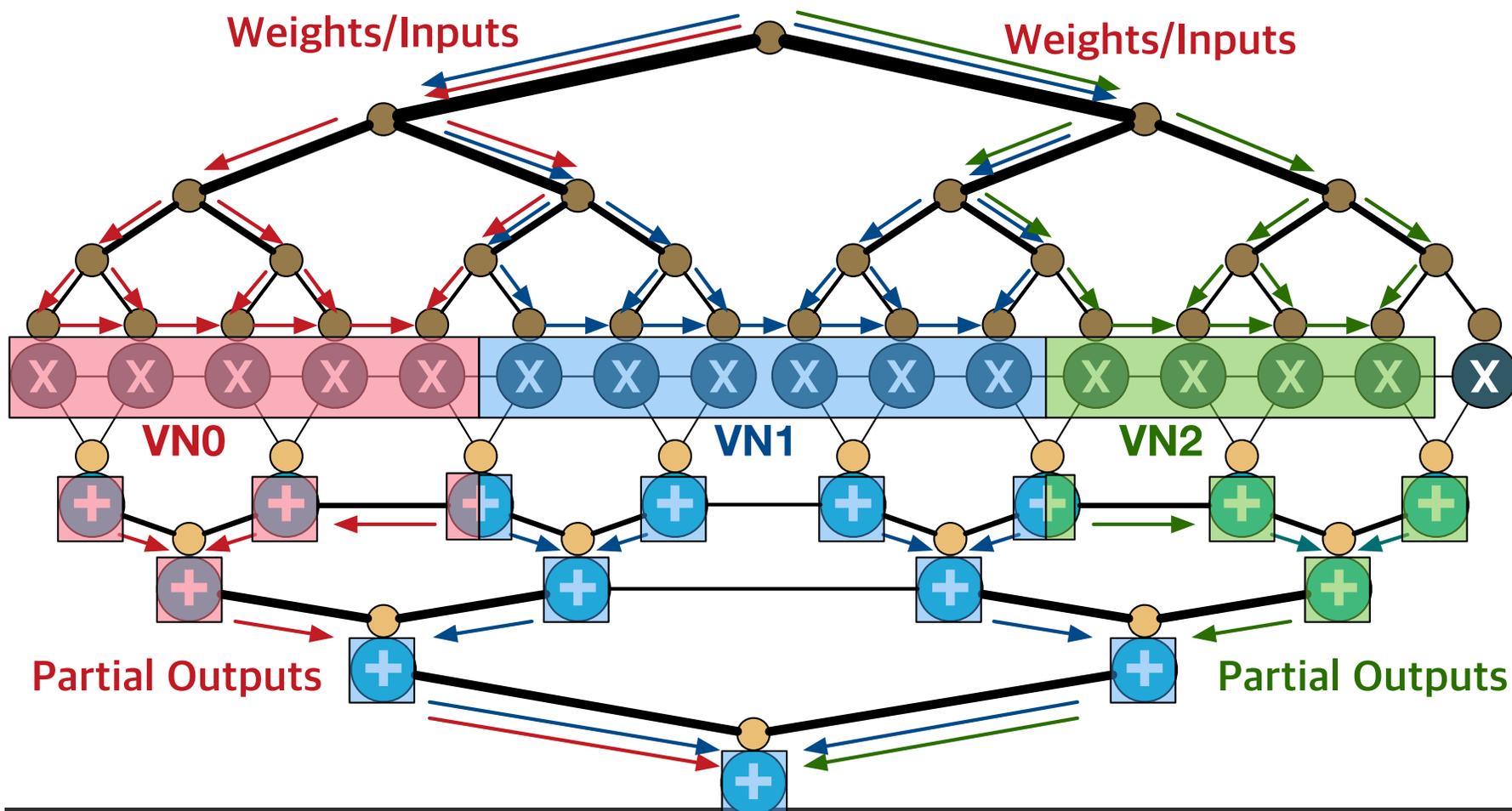
Example Mapping – Dense CNN

Our Key insight: Each **DNN/dataflow** translates into **neurons of different sizes**



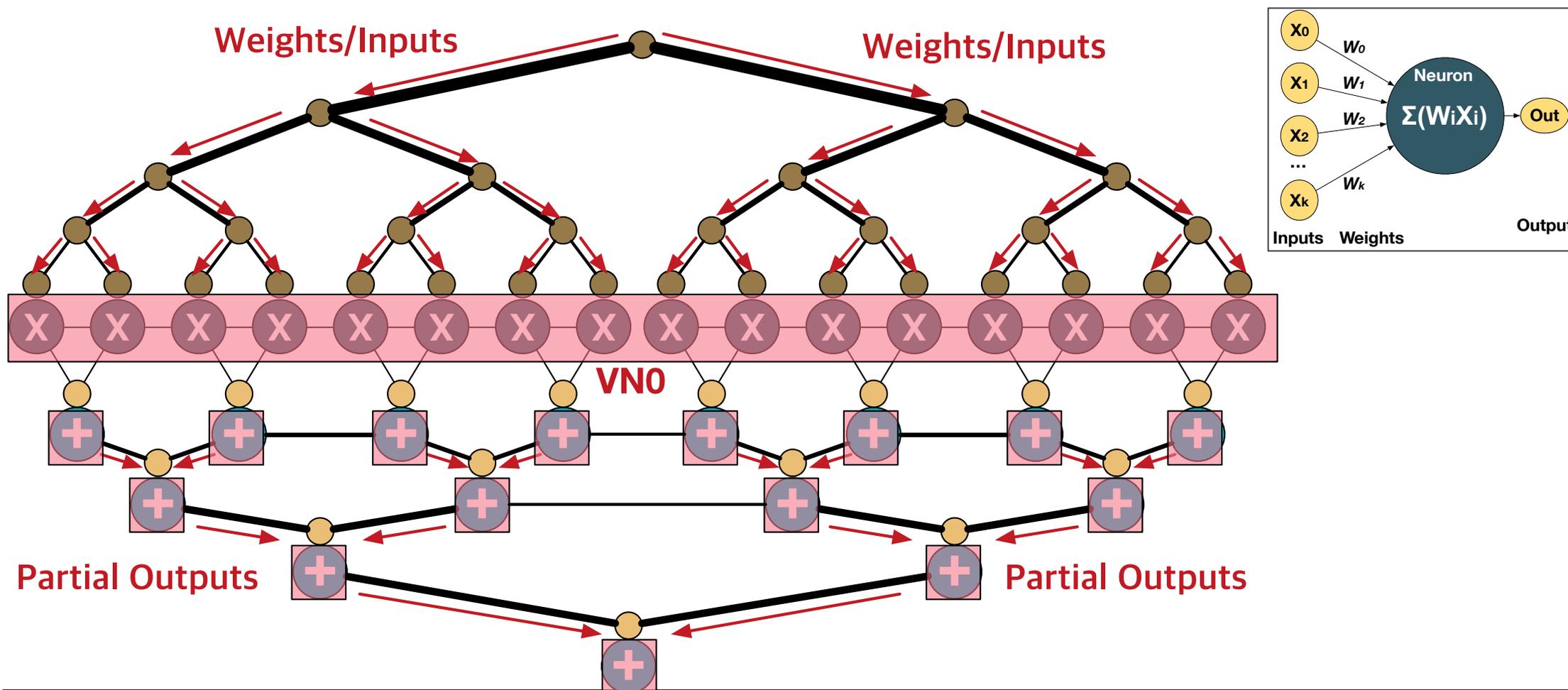
Example Mapping – Sparse DNN

Our Key insight: Each DNN/dataflow translates into neurons of different sizes



Example Mapping – LSTM/FC

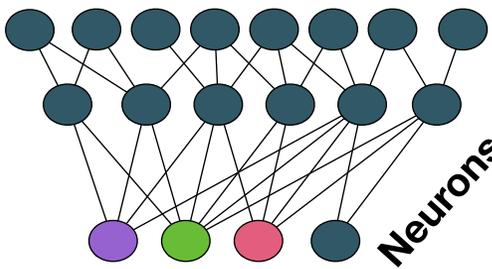
Our Key insight: Each DNN/dataflow translates into neurons of different sizes



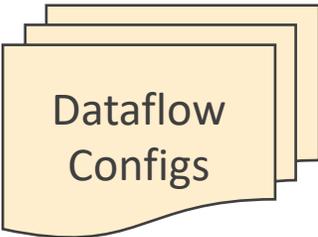
Searching optimal dataflows for MAERI

Find Optimal Mapping

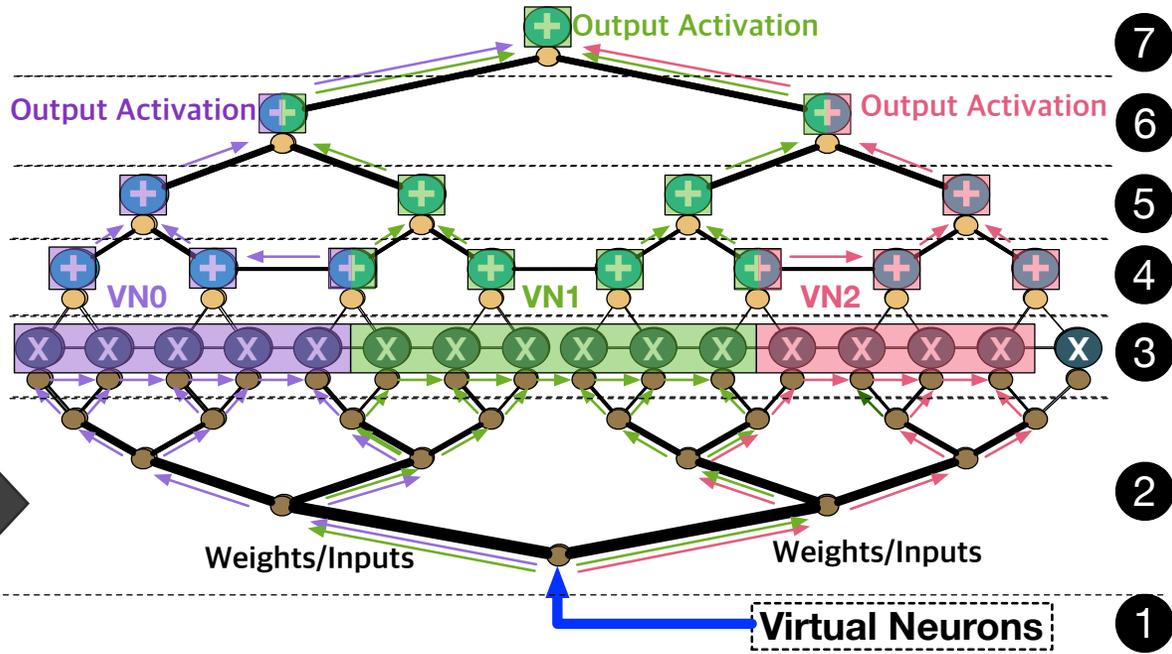
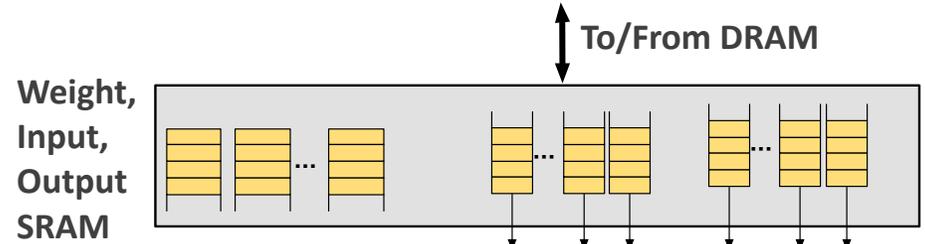
Deep Neural Network



Neurons



Z. Zhao, H. Kwon, S.Kuhar, W. Sheng, Z. Mao, T. Krishna
Efficient Mapping Space Exploration on a Reconfigurable Neural Accelerator
ISPASS 2019



~100% Utilization



Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:
ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention

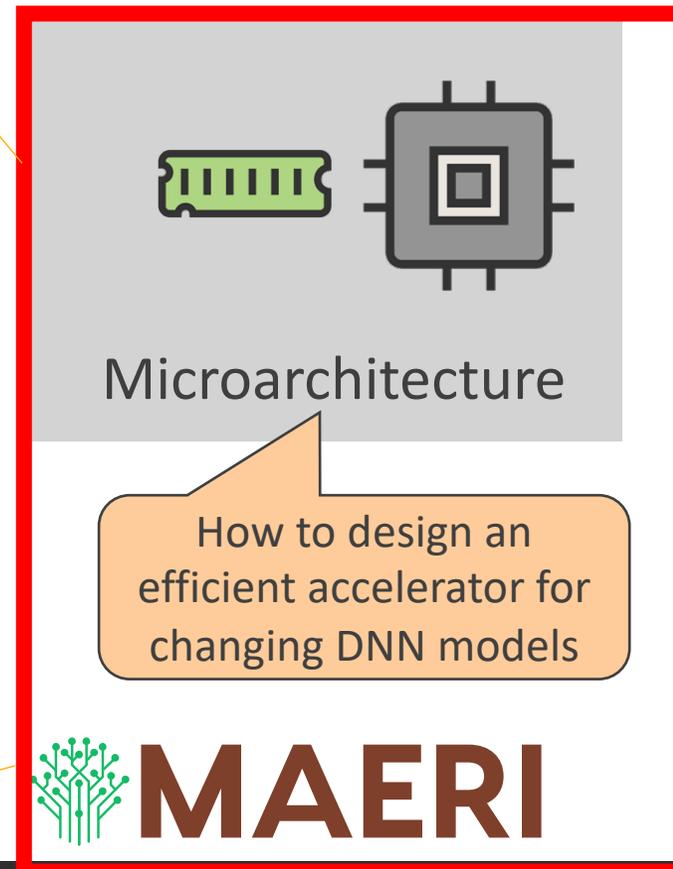
Outline of Talk

- Motivation
 - Irregular Dataflows
 - DNN Computation
- MAERI
 - Abstraction
 - Implementation
 - Operation Example
 - Mapping Strategies
- Evaluations

Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna

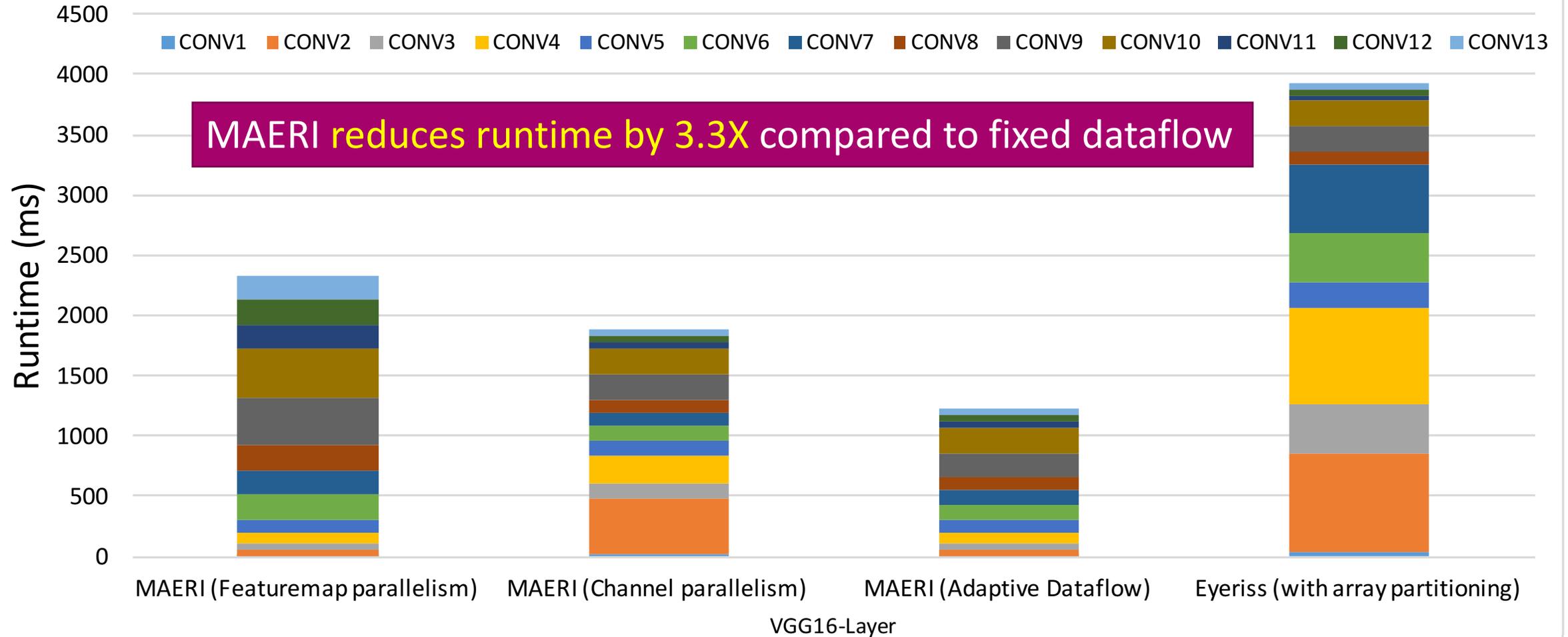
MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:

ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention

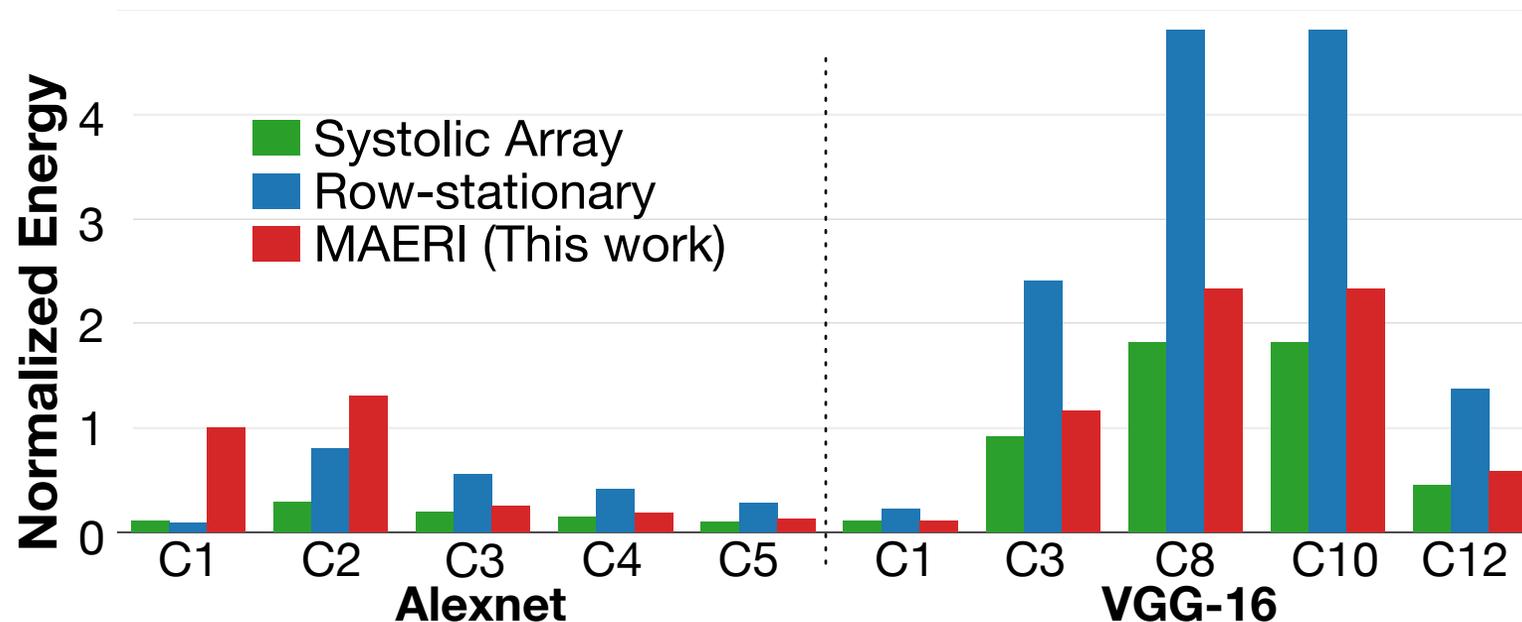


End-to-End Performance

VGG16 End-to-end Runtime (MAERI vs Eyeriss)



Energy with Convolution Layers



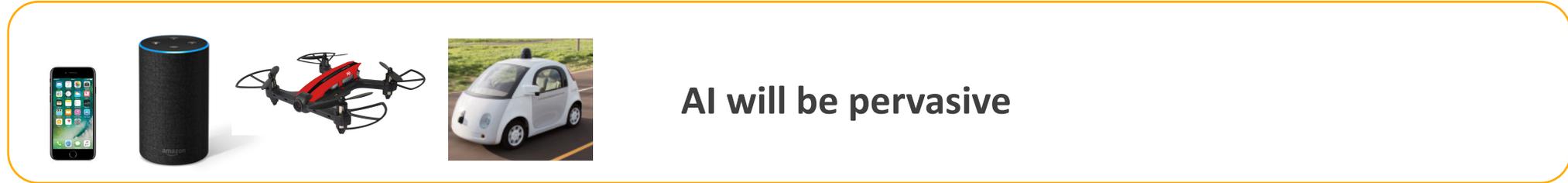
* Normalized to MAERI energy with Alexnet C1

MAERI reduces energy upto 57% and 28% in average compared to Row-Stationary (dense dataflow) and 7.1% in average compared to Systolic Array (sparse dataflow)

Summary of MAERI

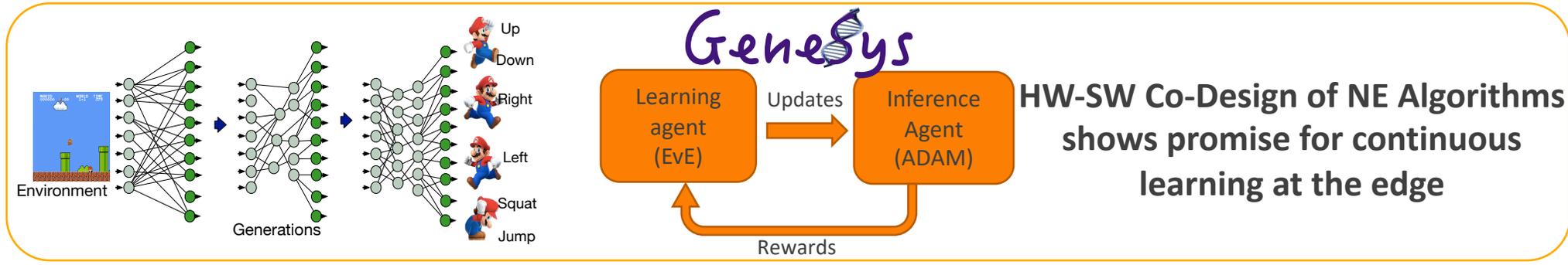
- DNN models evolving rapidly
 - Multiple layer types
 - Sparsity Optimizations
 - Myriad dataflows for scheduling and mapping
- MAERI enables dynamic grouping of arbitrary number of MACCs (“Virtual Neuron”) via reconfigurable, non-blocking interconnects, providing
 - Future proof to DNN models and dataflows
 - Near 100% compute unit utilization

Takeaways



AI will be pervasive

This section features four images: a smartphone, an Amazon Echo smart speaker, a red and black drone, and a white self-driving car (Waymo Firefly).



GeneSys

Learning agent (EvE) → Updates → Inference Agent (ADAM)

Rewards

HW-SW Co-Design of NE Algorithms shows promise for continuous learning at the edge

This section includes a diagram of a neural network training process for Super Mario Bros. It shows an 'Environment' (Mario game screen) feeding into a neural network over 'Generations'. The network outputs actions: Up, Down, Right, Left, Squat, and Jump. Below this is the GeneSys architecture diagram showing a Learning agent (EvE) sending 'Updates' to an Inference Agent (ADAM), which sends 'Rewards' back to the Learning agent.



MAERI

DNN Accelerator with Configurable Interconnects can map Irregular Dataflows

This section contains two diagrams. On the left is a single neuron diagram with inputs $x_0, x_1, x_2, \dots, x_k$ and weights $w_0, w_1, w_2, \dots, w_k$ entering a circle labeled 'Neuron' with the equation $\Sigma(W_i X_i)$ and an 'Output'. On the right is a diagram of a neural network layer with irregular dataflows, showing nodes with '+' and 'x' symbols connected in a non-uniform grid.

Thank you!

<http://synergy.ece.gatech.edu>