

ARM Research Summit
September 16, 2019

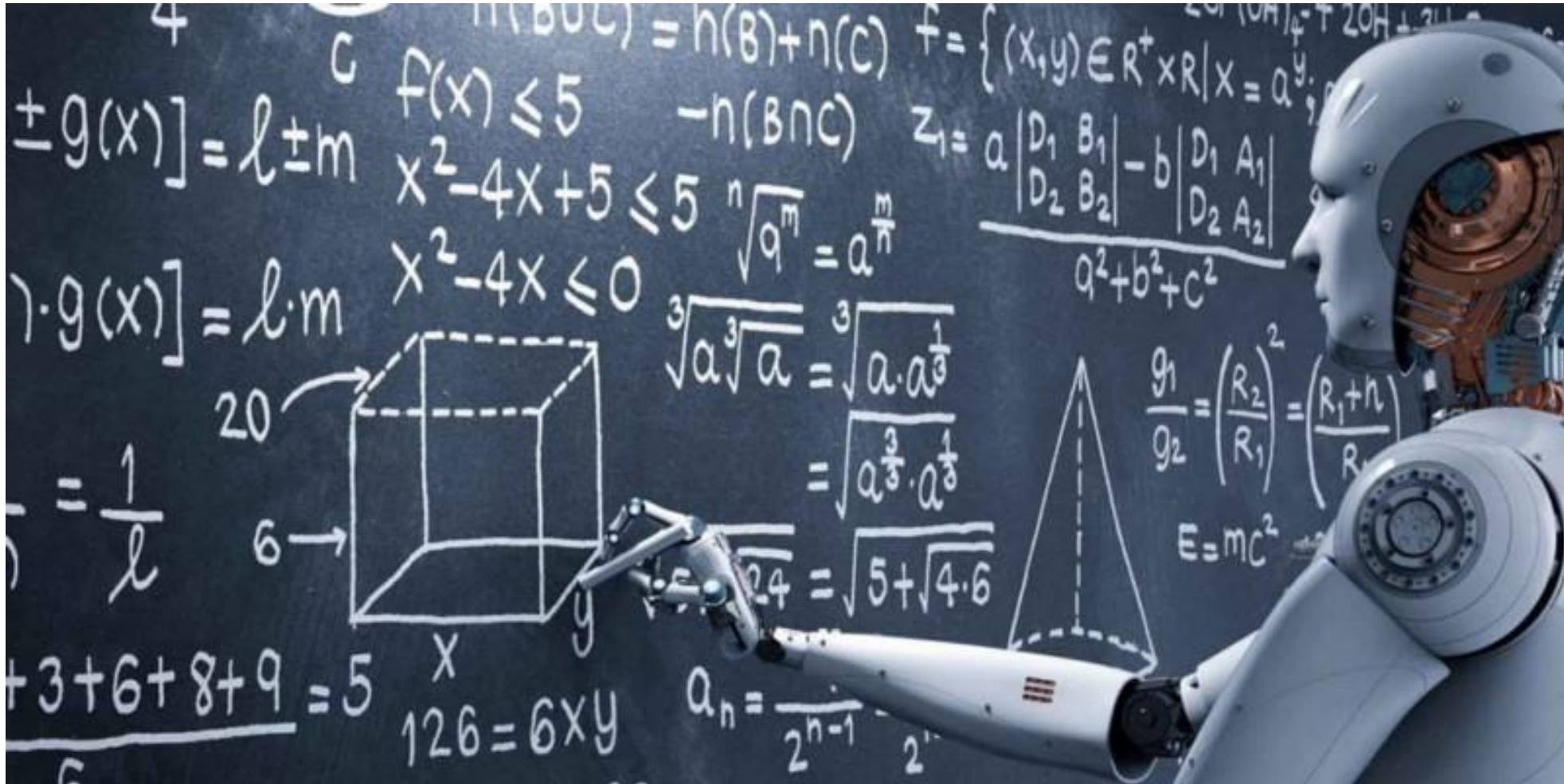
Georgia Tech School of Electrical and
Computer Engineering
College of Engineering

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

Enabling Continuous Learning through Neural Network Evolution in Hardware

Tushar Krishna

Assistant Professor
ON Semiconductor Professor
School of ECE
Georgia Institute of Technology



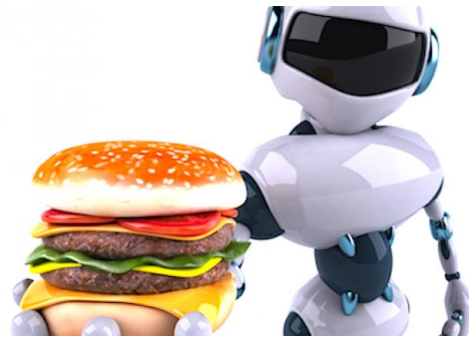
General-Purpose AI capable of emulating Human-like Continuous Learning

What is Continuous Learning?



Become better and faster with experience

Learn new tasks



Autonomous

Deep Learning Applications

Object Detection

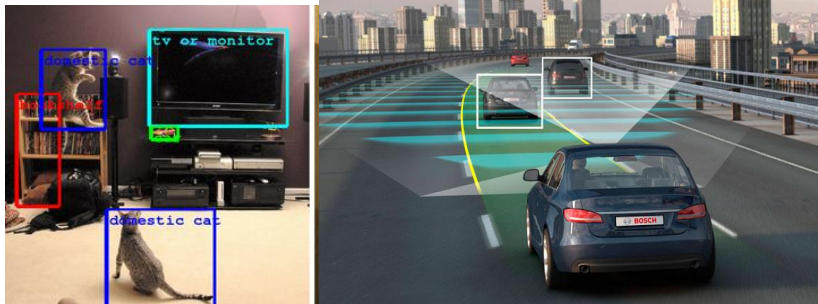
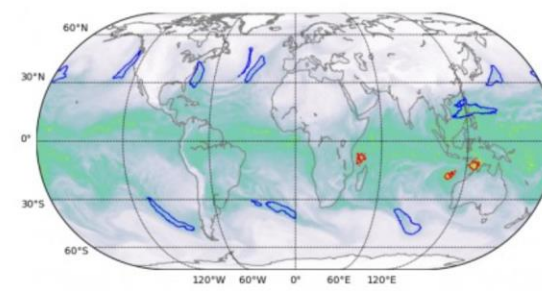
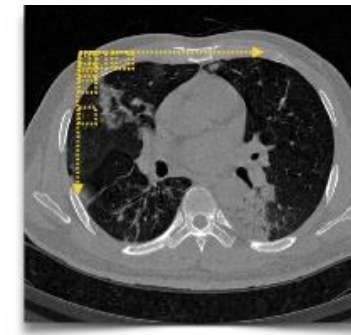


Image Segmentation



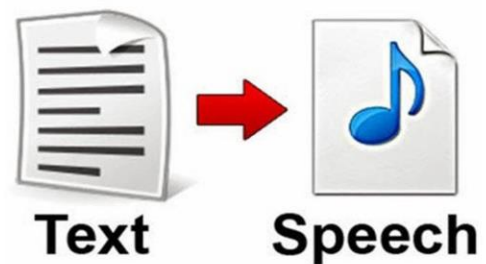
Medical Imaging



Speech Recognition



Text to Speech



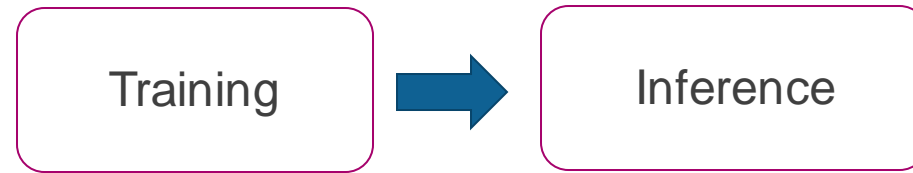
Recommendations



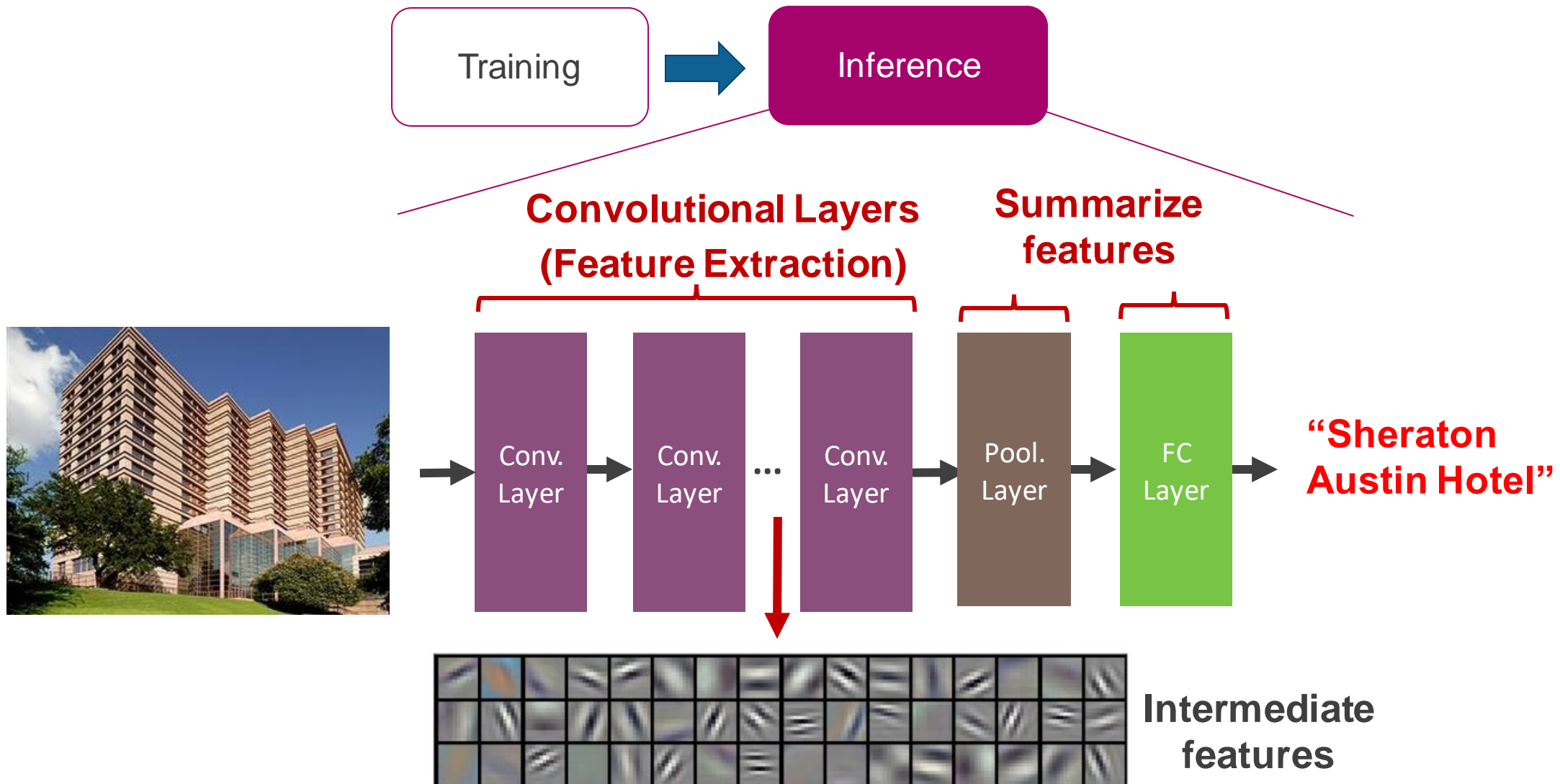
Games



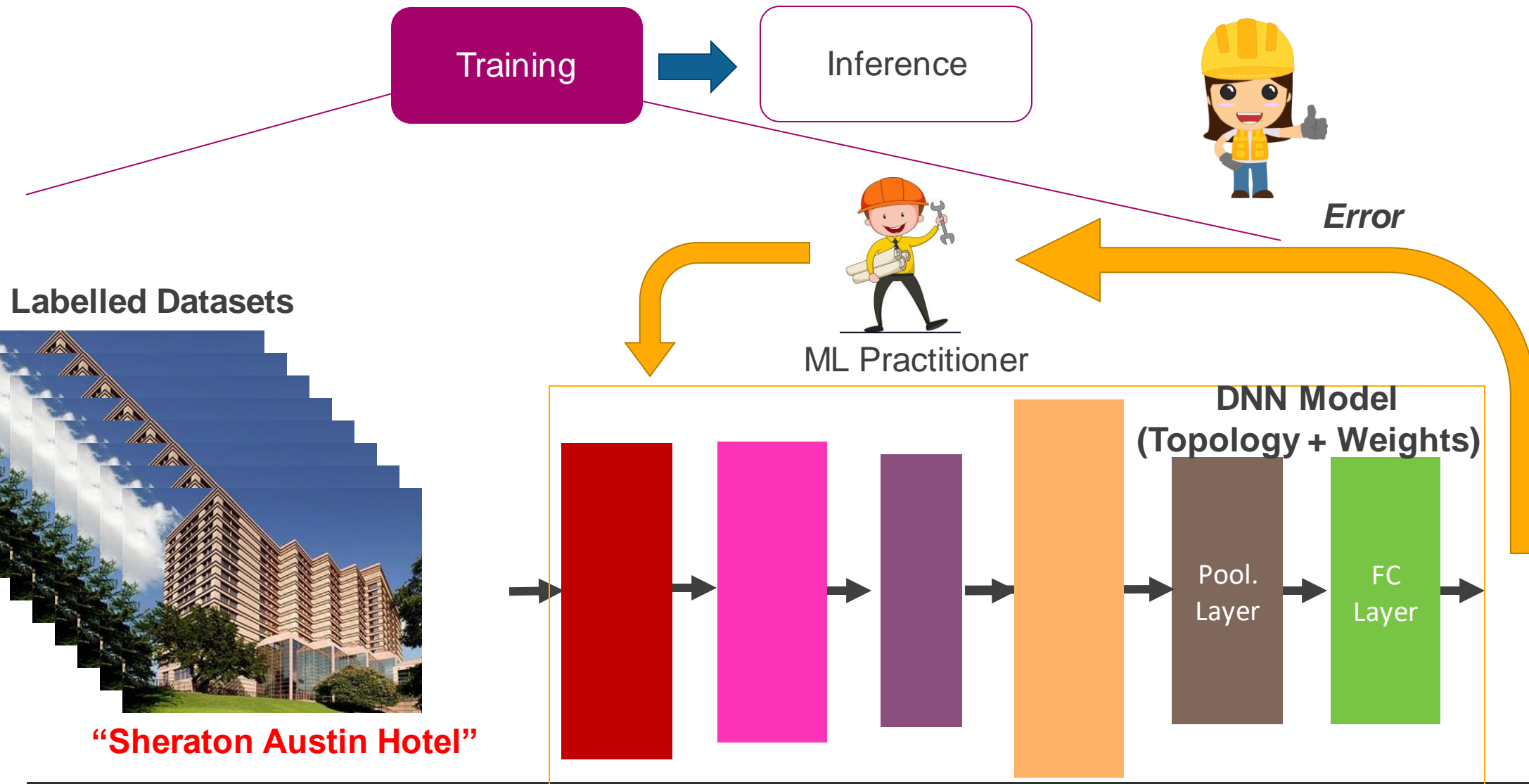
Deep Learning Landscape



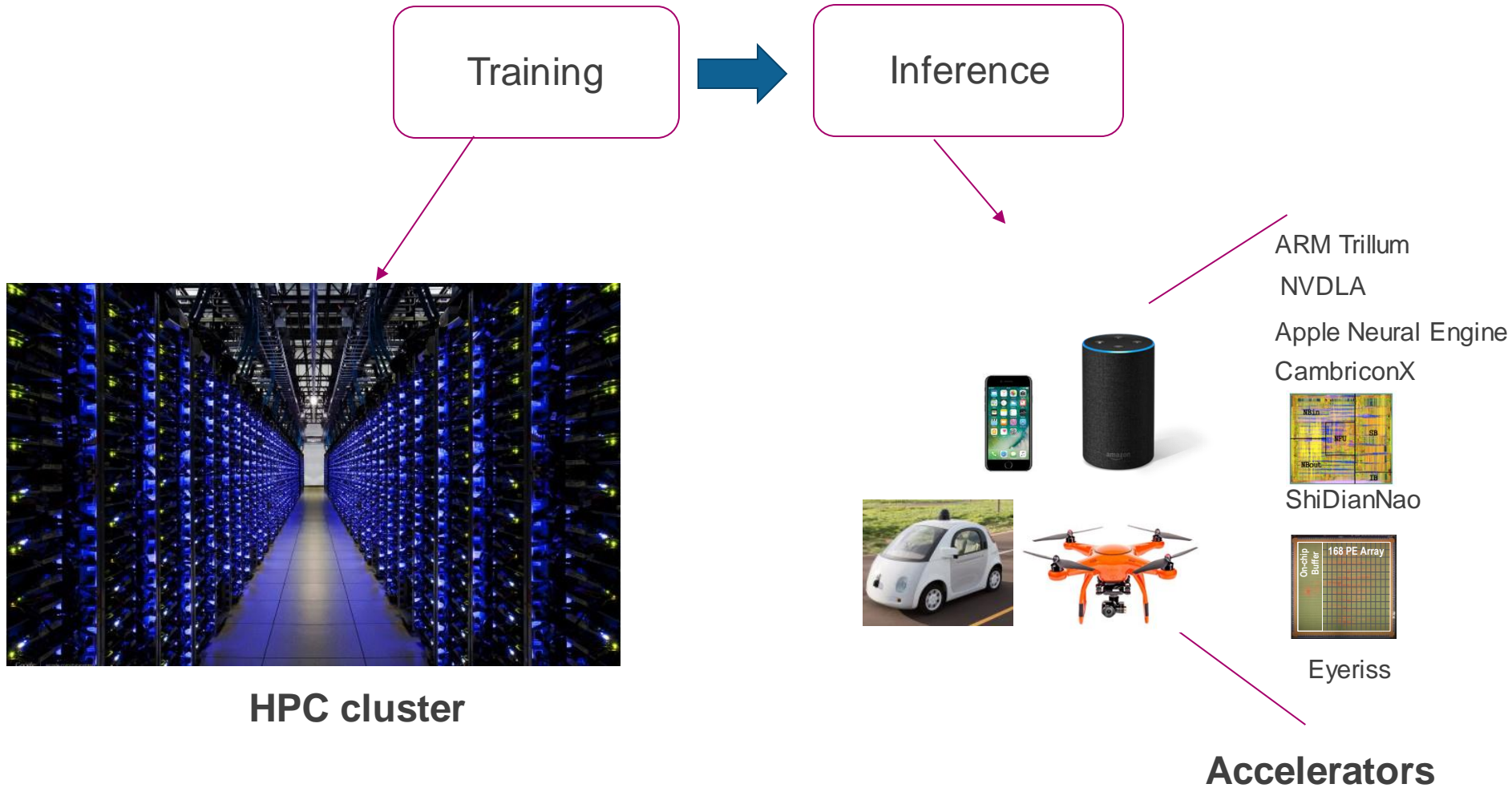
Deep Learning Landscape



Deep Learning Landscape



Computation Platforms



Challenges for Continuous Learning

- **Algorithm Challenge**

- No labeled datasets
- No human-in-the-loop
- Continuously Changing DNN Model

- **System Challenge**

- On-Device Learning => Low-energy Hardware Platform

Candidates for Continuous Learning

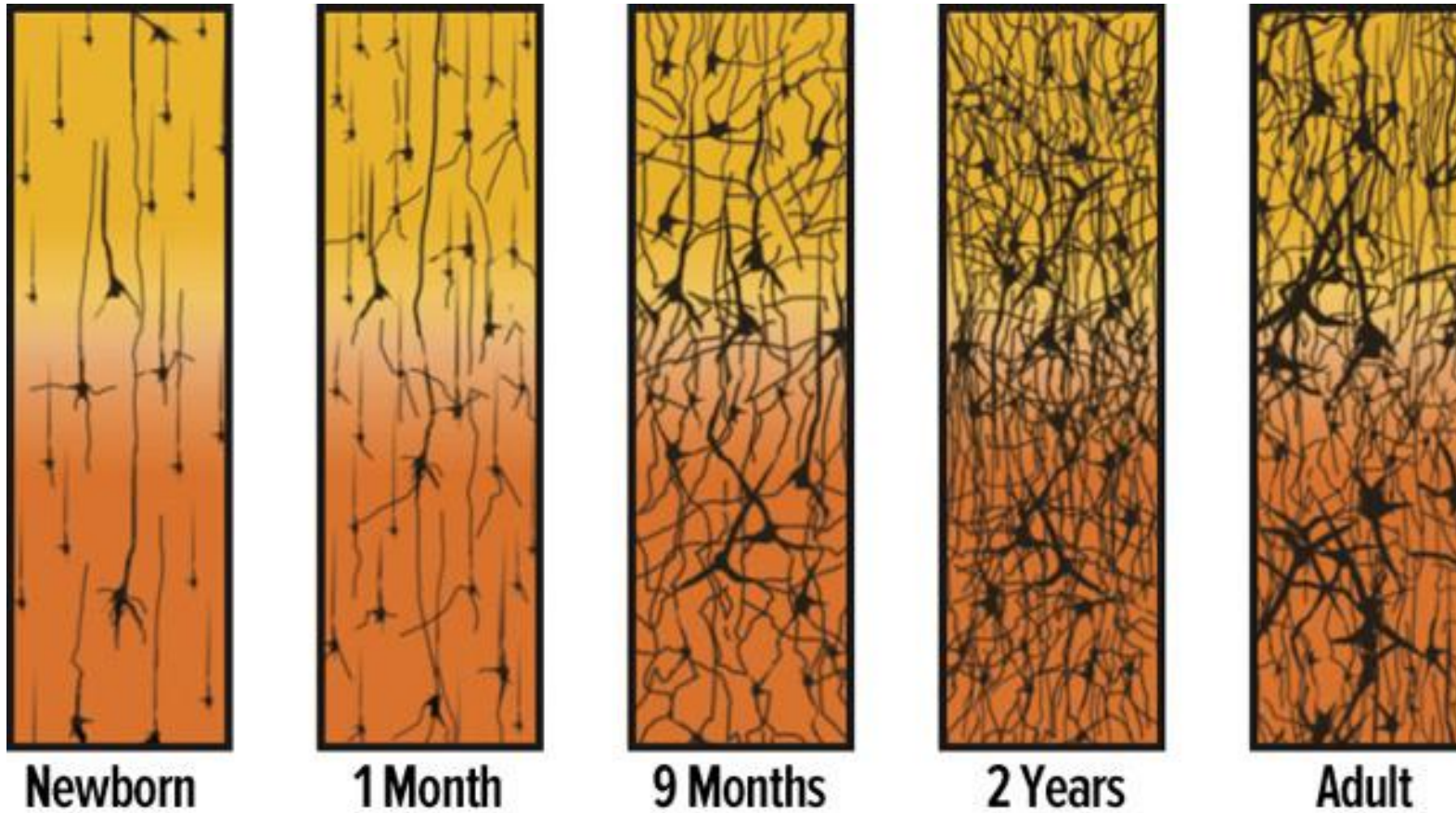
	Data	Hyper-param Tuning	DNN Plasticity	Compute	Memory
Supervised DL	Labeled X	Manual X	Designed for one problem X	Backprop (gradients) X	Backprop (gradients) X

Outline of Talk

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

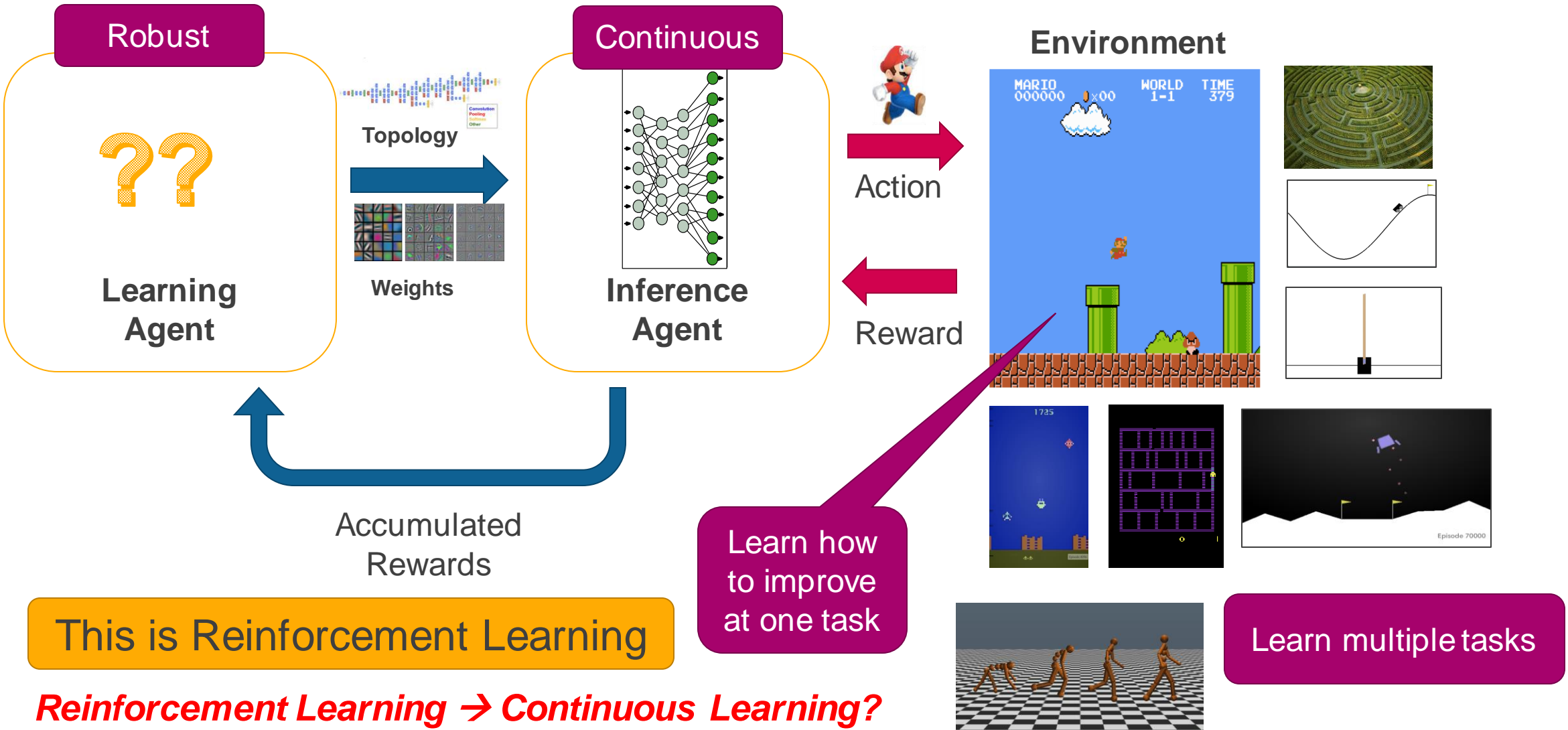
*Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna,
GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware,
In Proc of 51st Annual IEEE/ACM International Symposium on Microarchitecture(MICRO), 2018*

Continuous Learning in Brains



Constant synapse formation and pruning

Template for Continuous Learning



Candidates for Continuous Learning

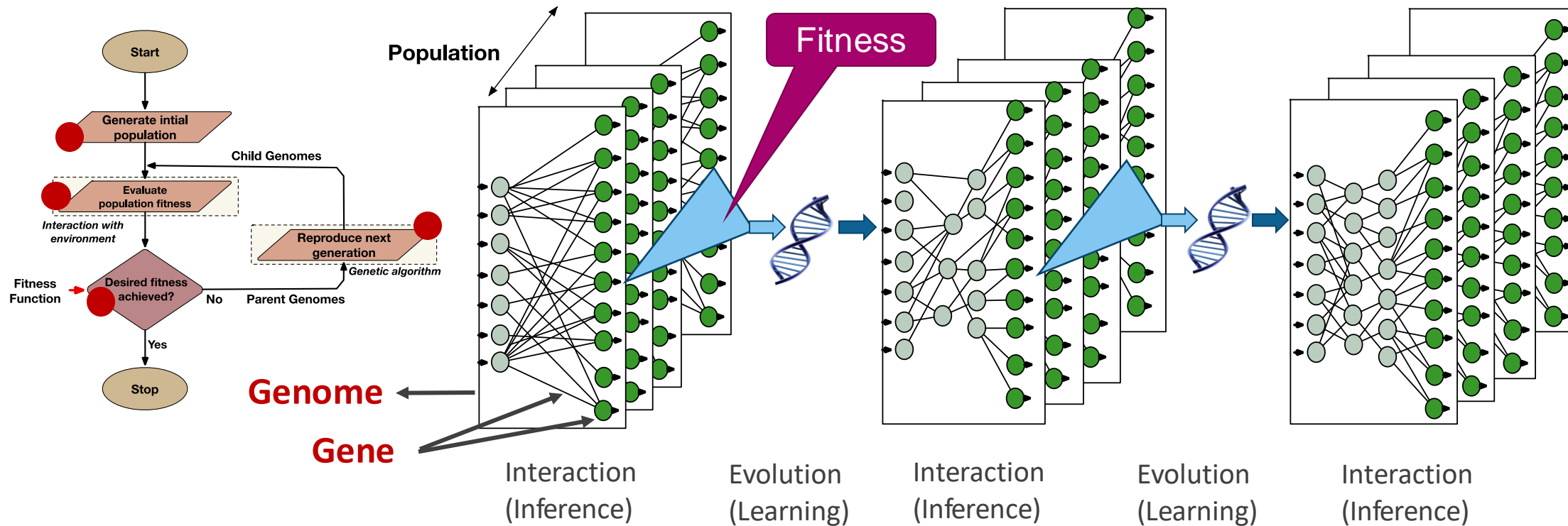
	Data	Hyper-param Tuning	DNN Plasticity	Compute	Memory
Supervised DL	Labeled X	Manual X	Designed for one problem X	Backprop (gradients) X	Backprop (gradients) X
Reinforcement Learning	Unlabeled ✓	Manual X	Reward Function ✓	Backprop (gradients) X	Backprop (gradients) X
	RL is not viable for continuous learning on the edge				

Outline of Talk

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

*Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna,
GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware,
In Proc of 51st Annual IEEE/ACM International Symposium on Microarchitecture(MICRO), 2018*

Neuro-Evolutionary (NE) Algorithm



Neural Network (NN) expressed as a graph

Gene: Vertex or Edge in the graph

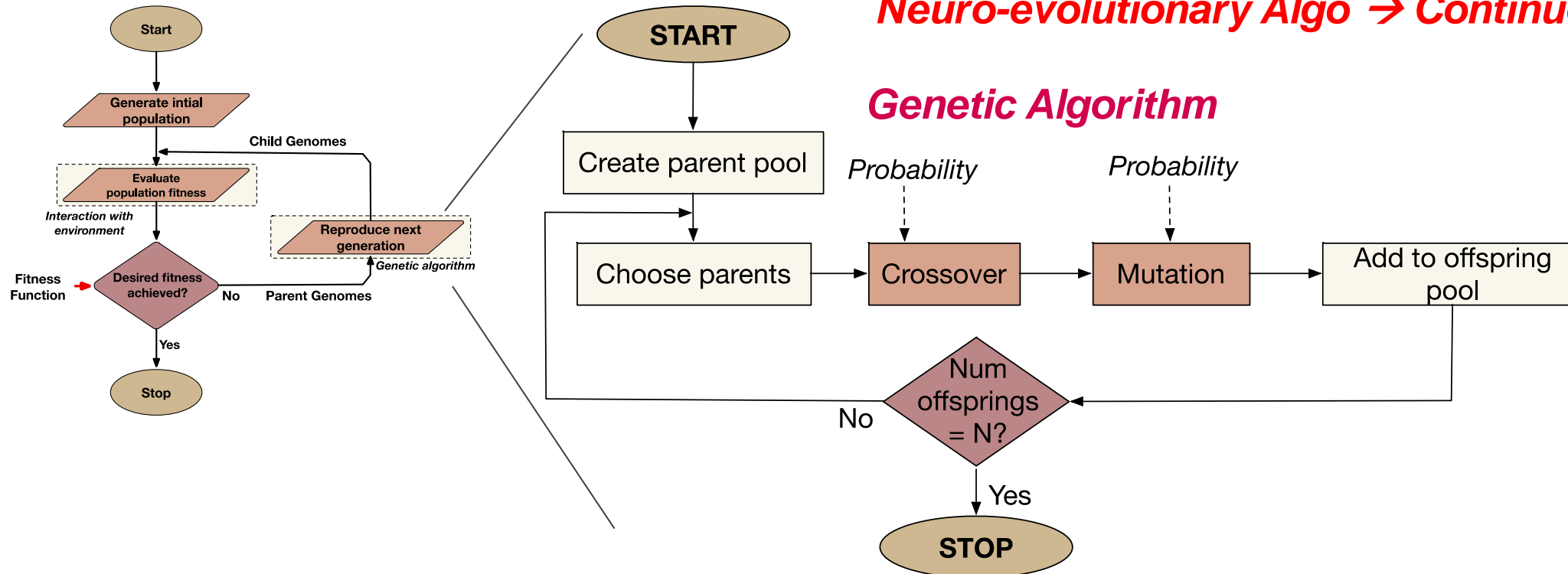
Genome: Collection of all genes (i.e., a NN)

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

Neuro-Evolutionary (NE) Algorithm

Neuro-evolutionary Algo → Continuous Learning?

Genetic 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.

Candidates for Continuous Learning

	Data	Hyper-param Tuning	DNN Plasticity	Compute	Memory
Supervised DL	Labeled X	Manual X	Designed for one problem X	Backprop (gradients) X	Backprop (gradients) X
Reinforcement Learning	Unlabeled ✓	Manual X	Reward Function ✓	Backprop (gradients) X	Backprop (gradients) X
Evolutionary	Unlabeled ✓	Automated ✓	Reward Function ✓	??	??

Outline of Talk


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

Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna,
GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware,
In Proc of 51st Annual IEEE/ACM International Symposium on Microarchitecture(MICRO),
2018

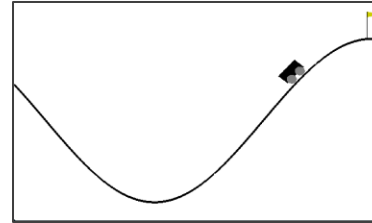
Characterization of NEAT

NEAT - Python

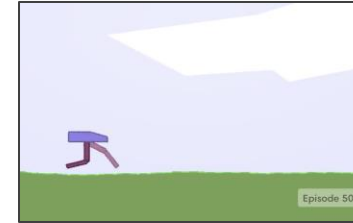
Codebase

 OpenAI Gym

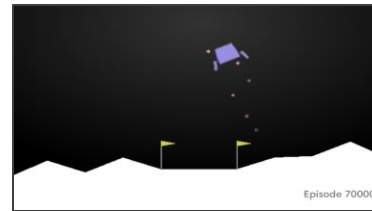
Environments



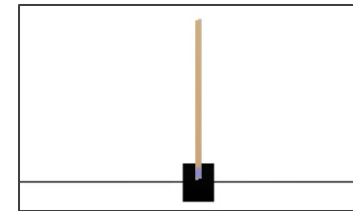
Mountain car



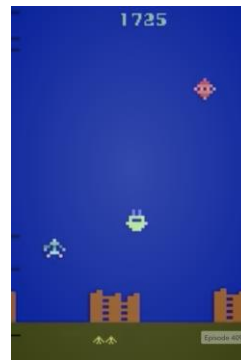
Bipedal



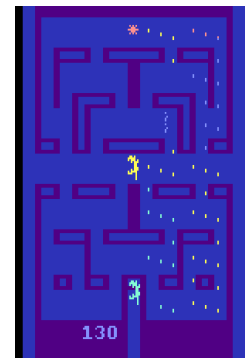
Lunar Lander



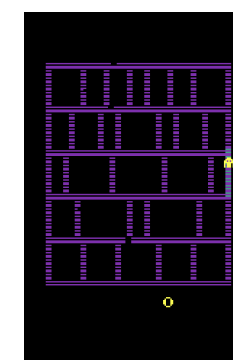
Cart pole



Airraid-RAM



Alien-RAM



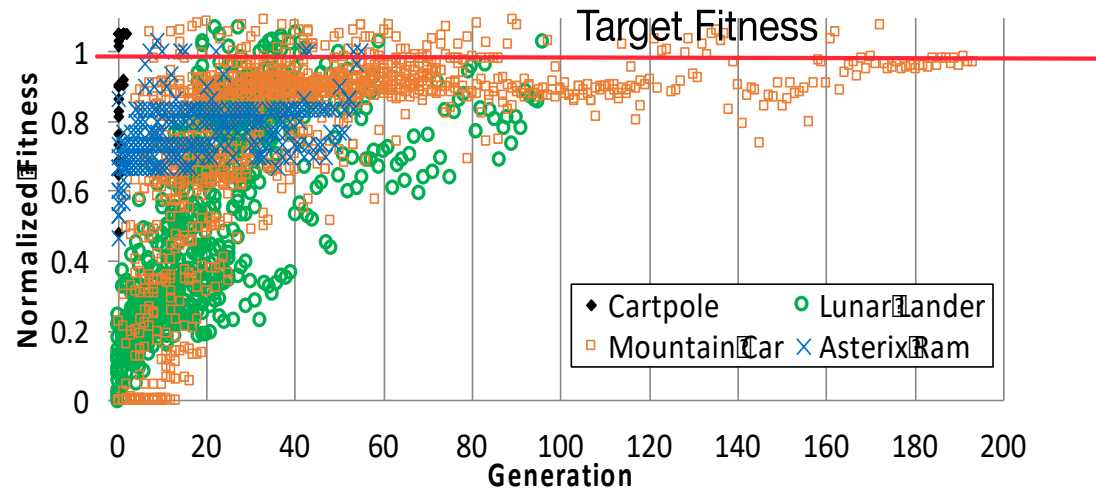
Amidar-RAM

Ran each environment
till convergence,
multiple times

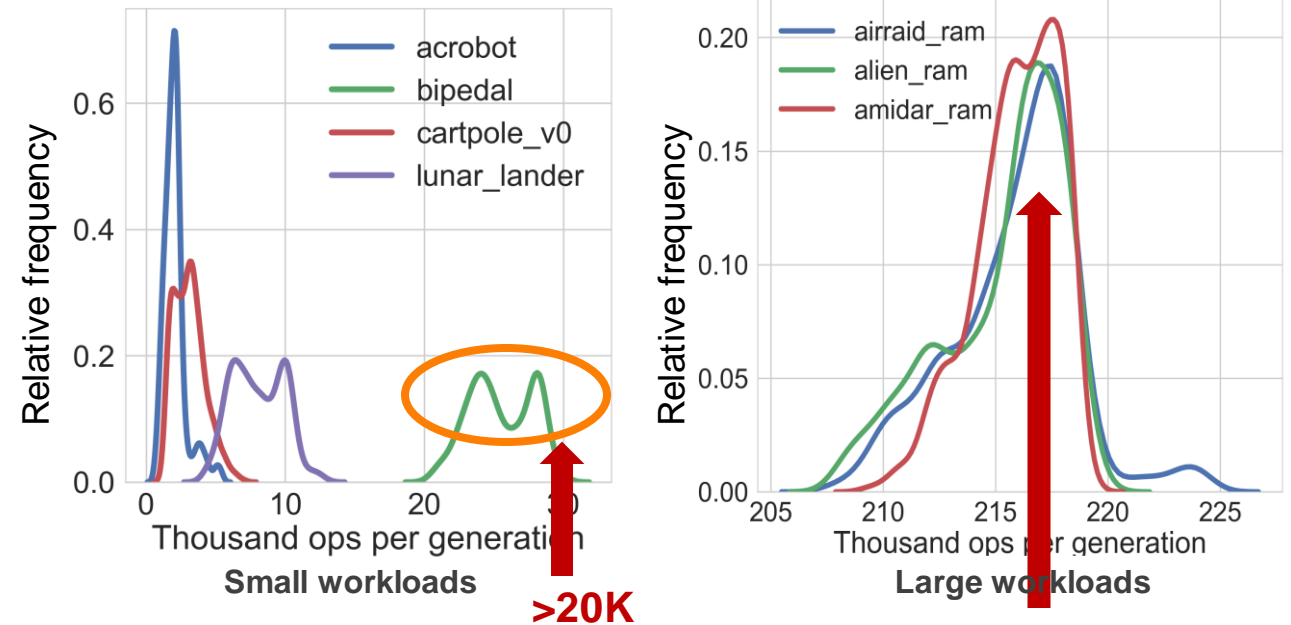
Only changed fitness
function between
workloads

Characterization of NEAT

Computations



Distribution of Operations/Generation



Inference:
Population level parallelism (PLP)

Evolution:
Gene level parallelism (GLP)

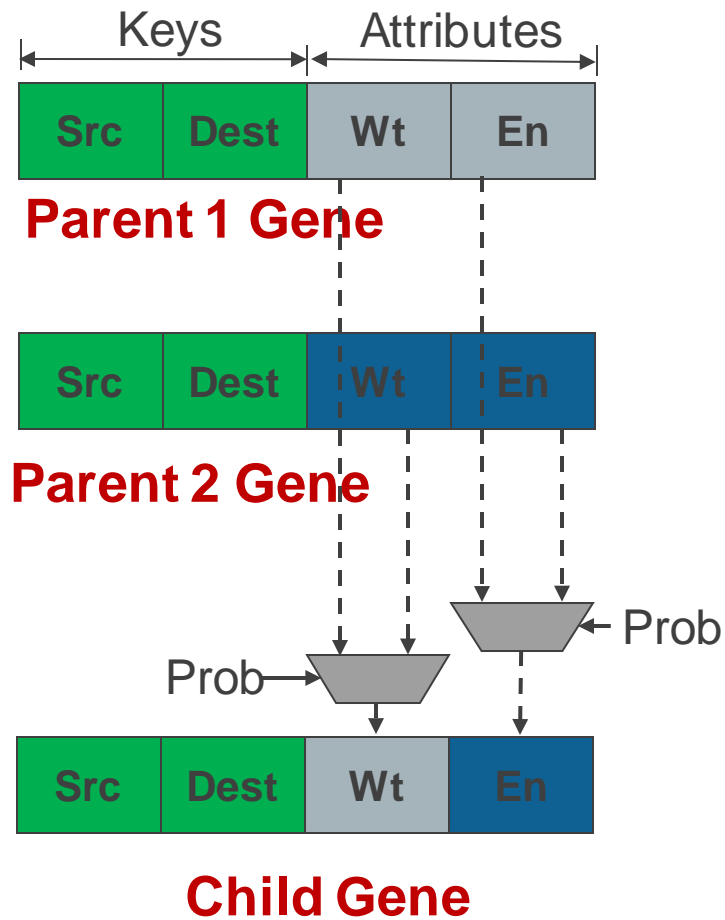
All operations are independent

Large operation level Parallelism

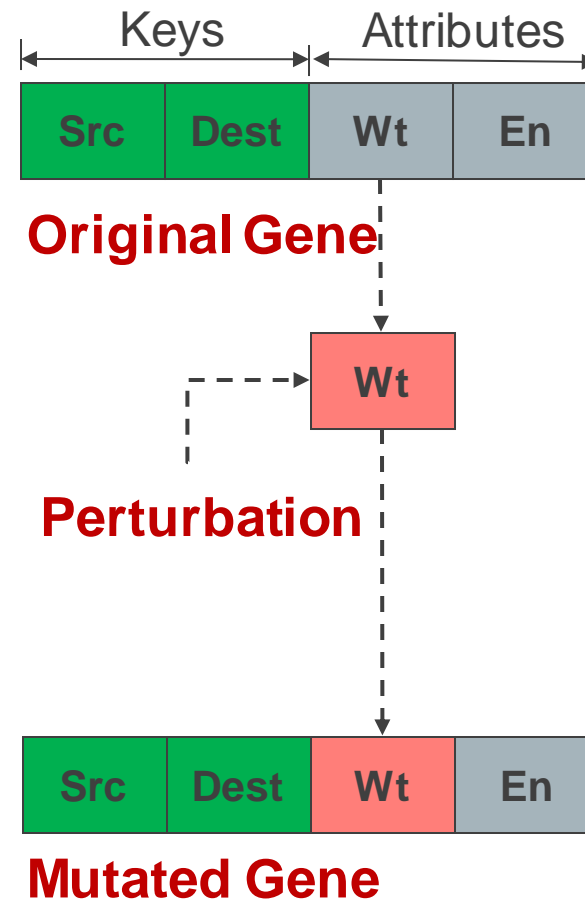
Operations in NEAT

Evolution

Crossover



Mutation



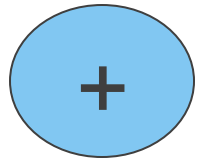
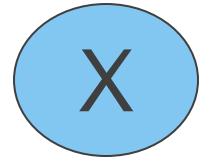
Addition mutation

- Add new node
- Add new connection

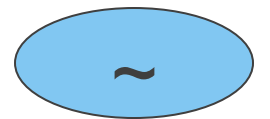
Deletion mutation

- Delete connection
- Delete node

Inference



MAC

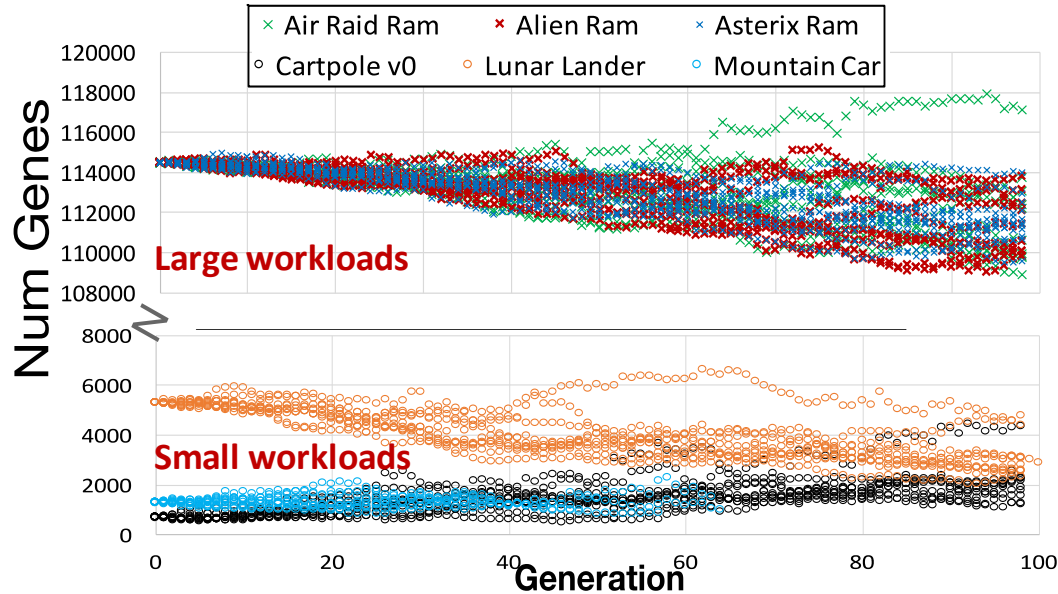


Activation

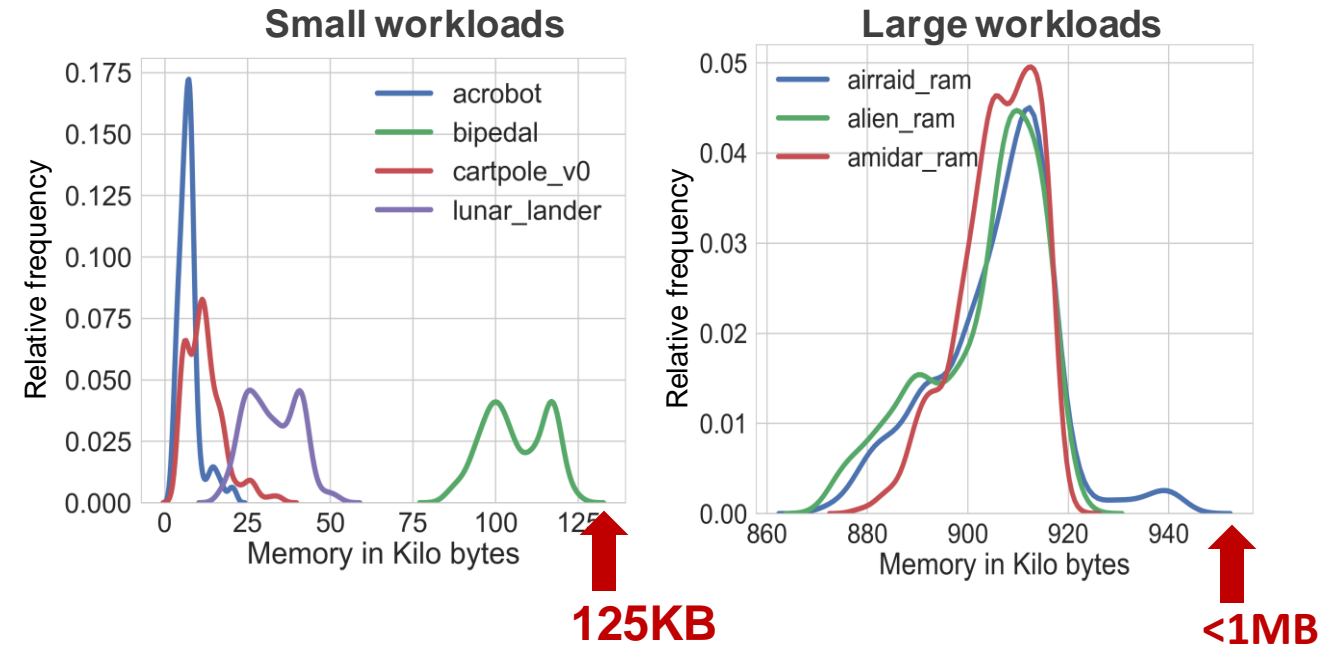
Simple operations

Characterization of NEAT

Memory



Distribution of Memory footprint/Generation



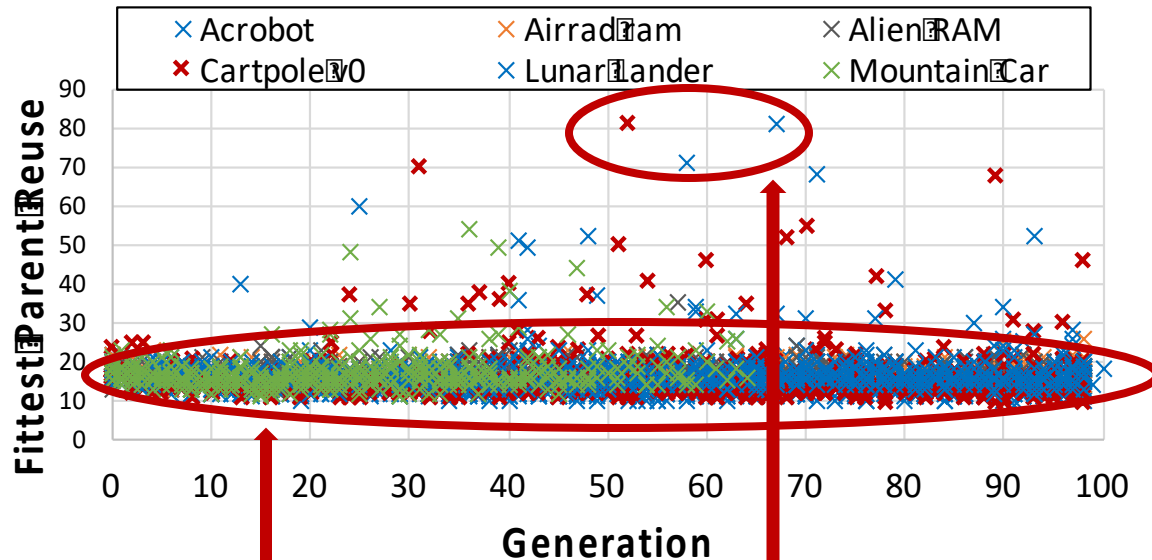
Entire population can fit on-chip

Only need to store the weights and node info

Characterization of NEAT

Memory

Opportunity for Reuse

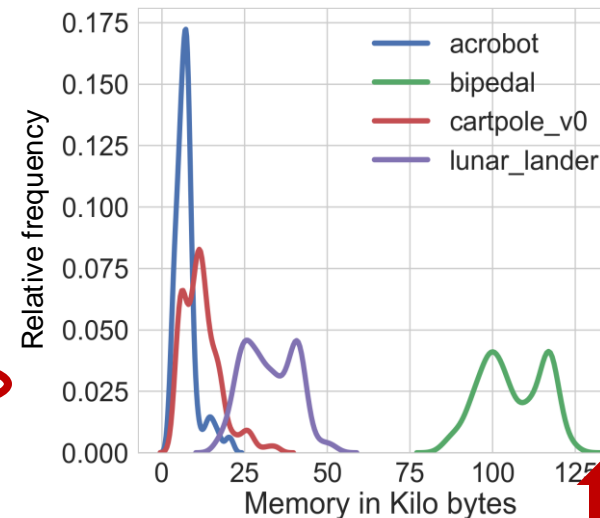


Fittest parent genome is used about ~10-20 times each generation

Even higher in certain cases

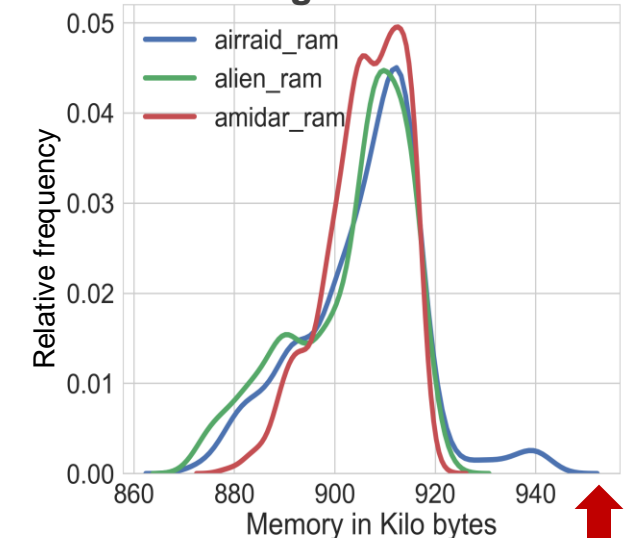
Distribution of Memory footprint/Generation

Small workloads



125KB

Large workloads



<1MB

Entire population can fit on-chip

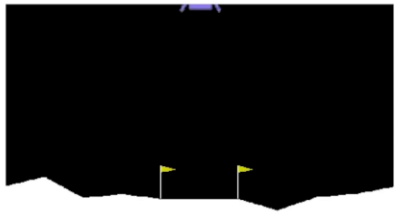
Only need to store the weights and node info

Properties of NE algorithms

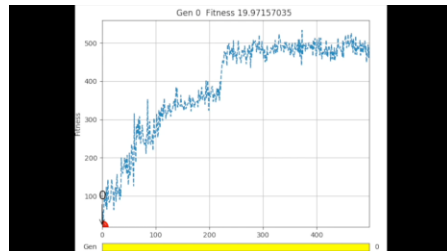
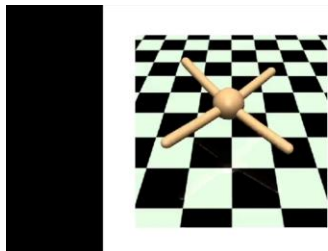
Algorithmic

No Explicit Training

No Manual Tuning



Change fitness function



Systems

Massive Parallelism

Genomes within Population

Genes within a Genome

No Backprop

No gradient calculations or storage

Low Memory Footprint

Only store genomes in current generation

Simple HW-friendly Ops

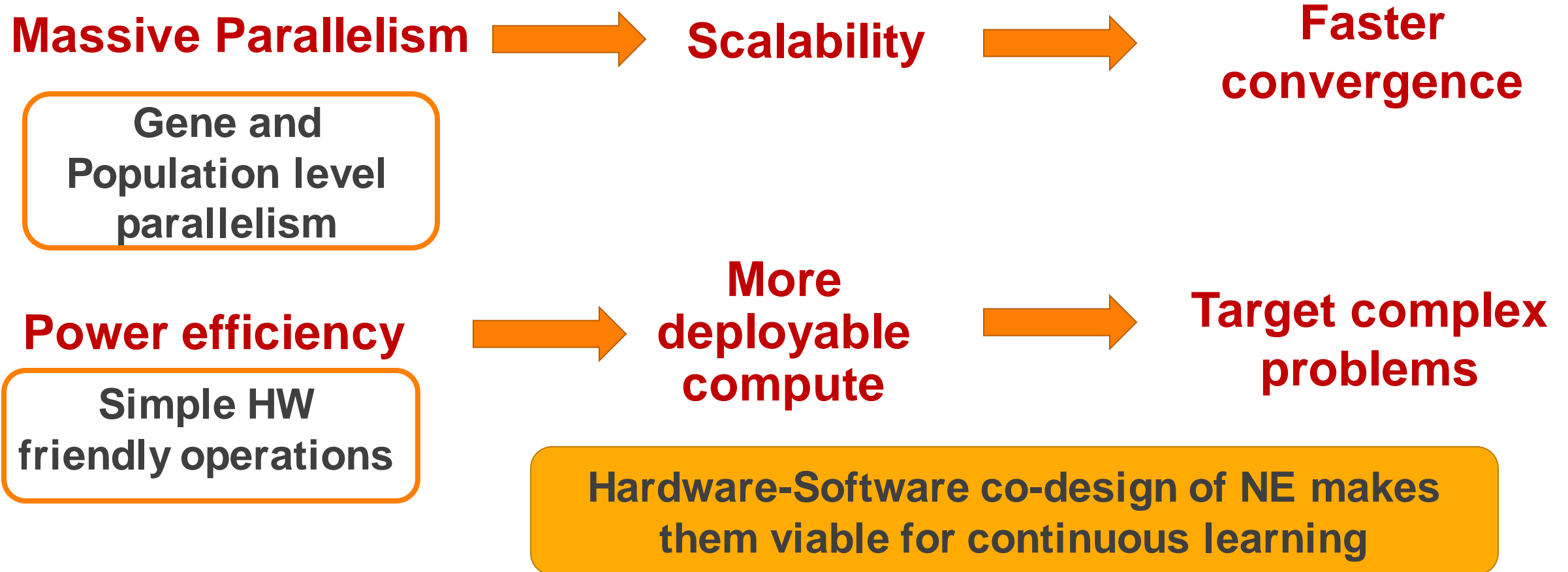
MACs in Inference
Crossover and Mutation in Evolution

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

Candidates for Continuous Learning

	Data	Hyper-param Tuning	DNN Plasticity	Compute	Memory
Supervised DL	Labeled X	Manual X	Designed for one problem X	Backprop (gradients) X	Backprop (gradients) X
Reinforcement Learning	Unlabeled ✓	Manual X	Reward Function ✓	Backprop (gradients) X	Backprop (gradients) X
Evolutionary	Unlabeled ✓	Automated ✓	Reward Function ✓	Massive Parallelism ✓	Only store model ✓

Motivation for HW-SW Co-Design

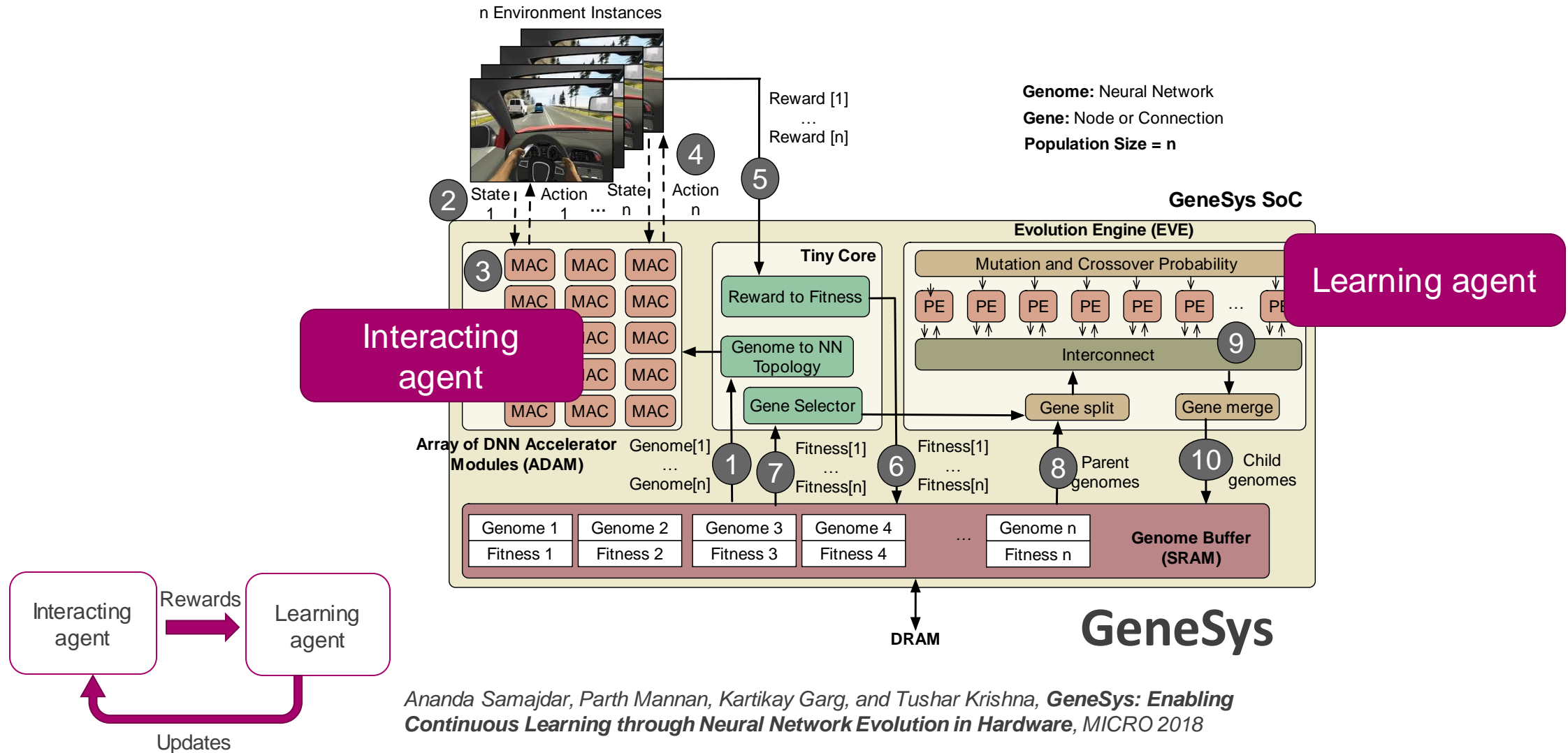


Outline of Talk

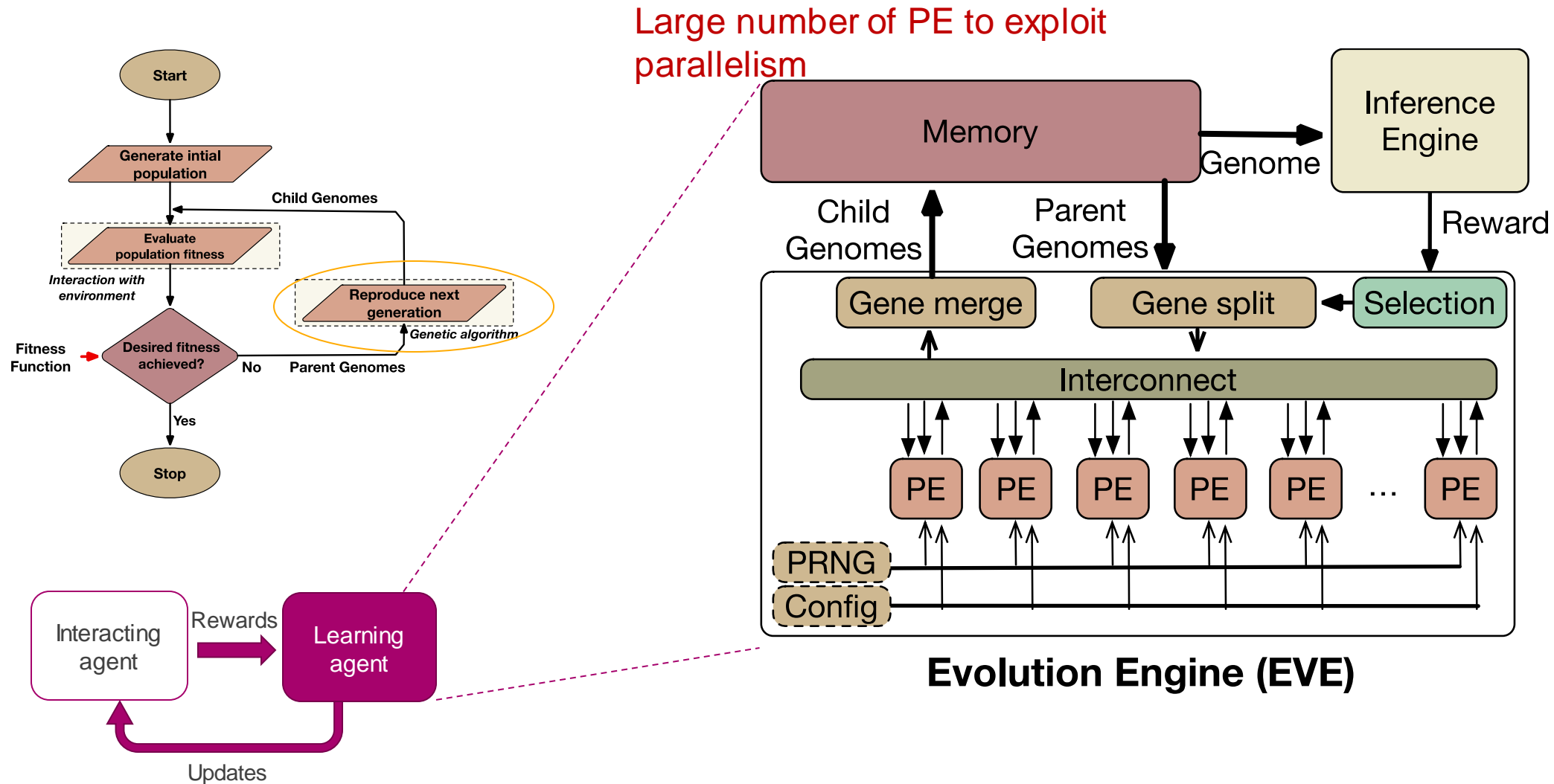
- Continuous Learning Template
- Neuro-Evolutionary Algorithms
 - Algorithm Description
 - Characterizing NEAT
- Microarchitecture
- Evaluations

*Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna,
GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware,
In Proc of 51st Annual IEEE/ACM International Symposium on Microarchitecture(MICRO), 2018*

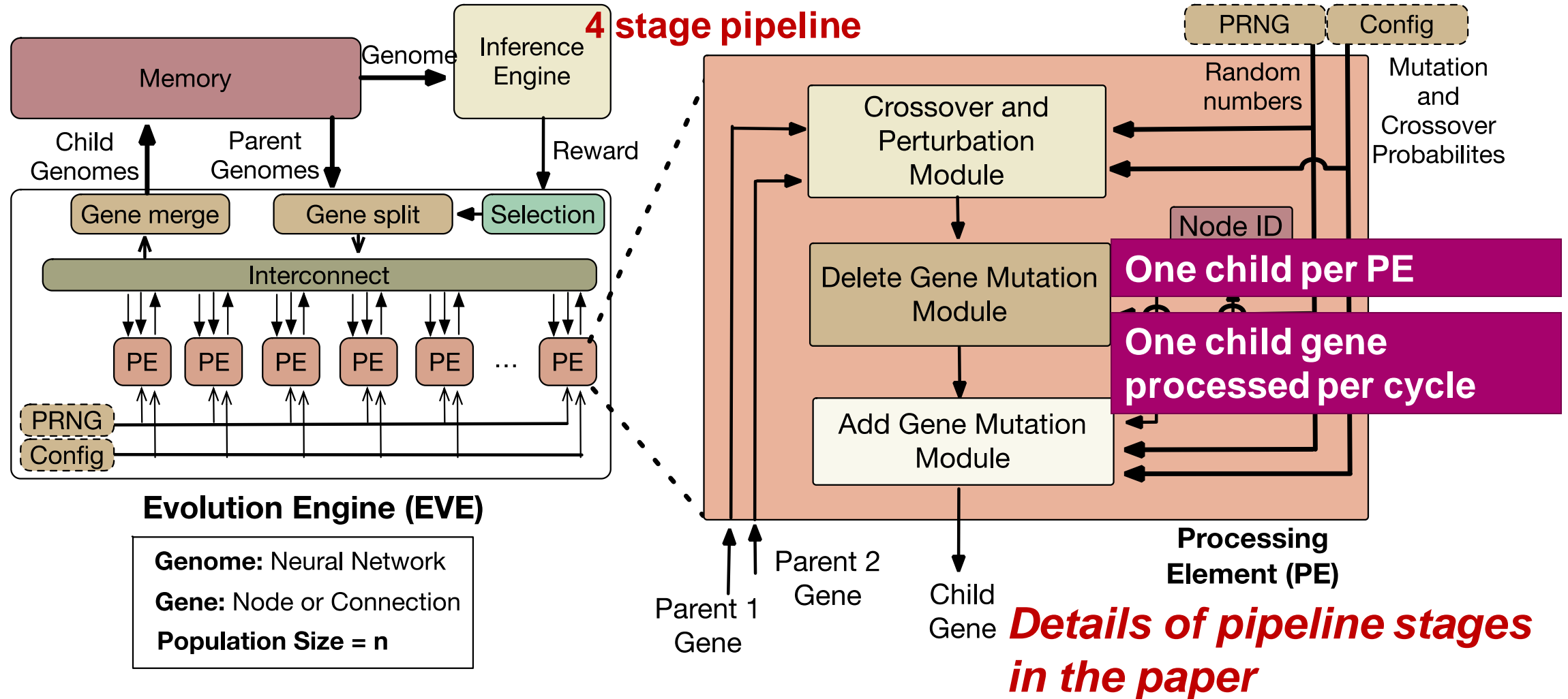
GeneSys SoC



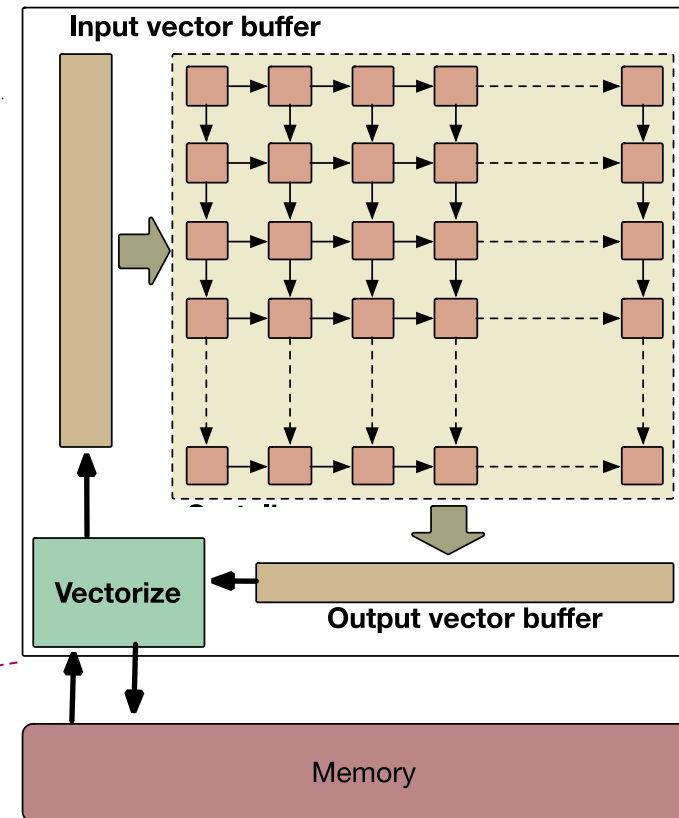
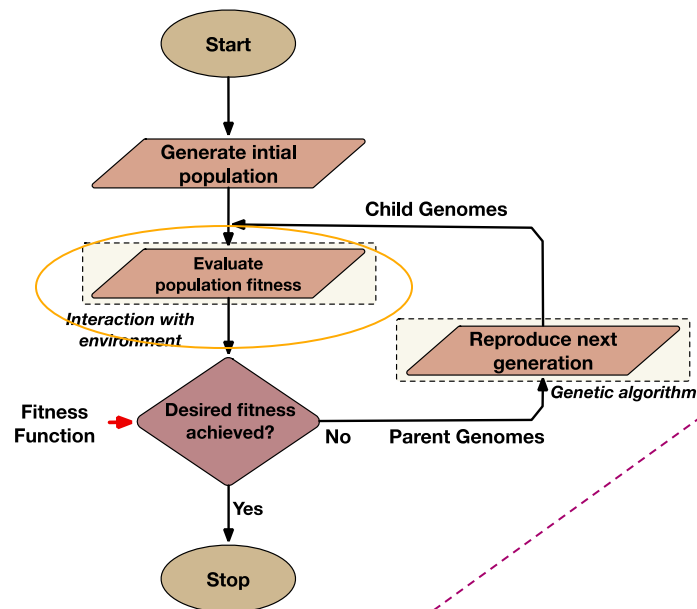
Evolution Engine: EvE Microarchitecture



PE Microarchitecture



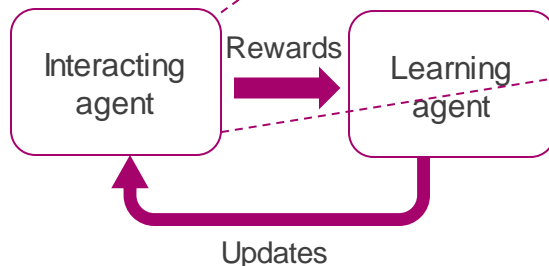
Inference Engine: ADAM Microarchitecture



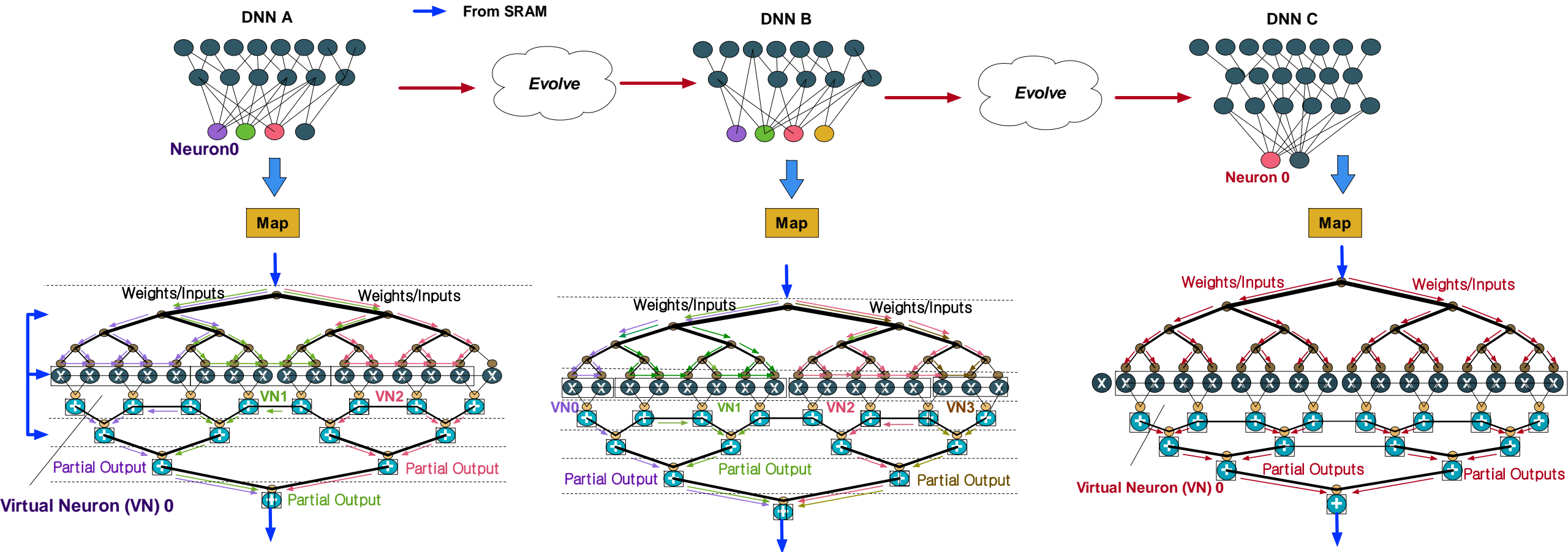
Conventional DNN
Inference Accelerator

Exploit Population
Level Parallelism

Networks generated
by NEAT are irregular
(thus sparse)



Mapping Evolving DNNs over ADAM



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

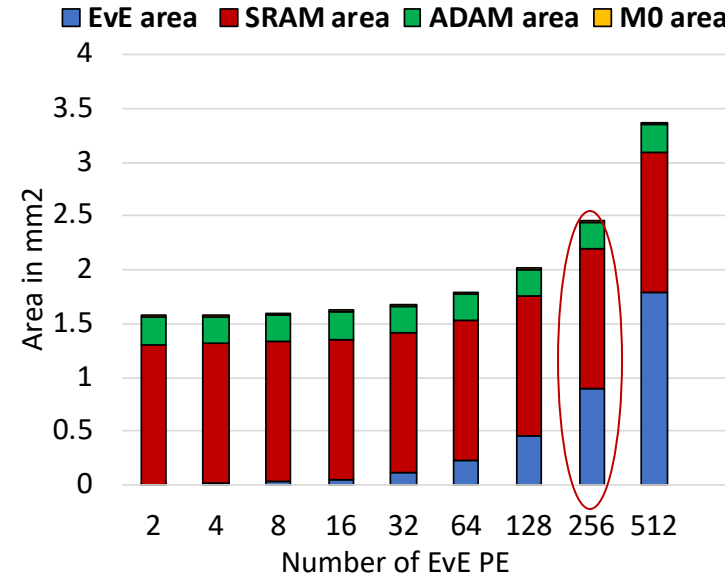
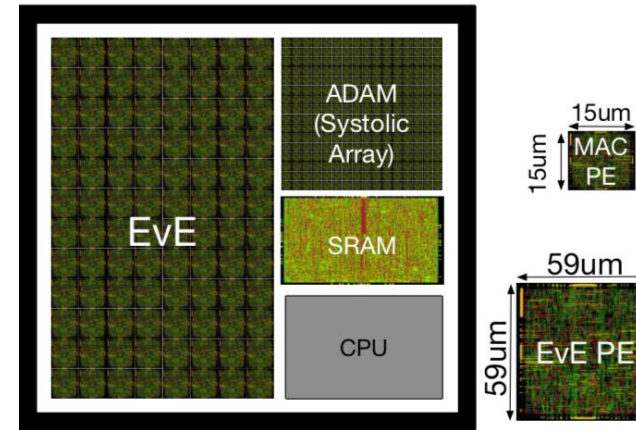
- Continuous Learning Template
- Neuro-Evolutionary Algorithms
 - Algorithm Description
 - Characterizing NEAT
- Microarchitecture
- Evaluations

*Ananda Samajdar, Parth Mannan, Kartikay Garg, and Tushar Krishna,
GeneSys: Enabling Continuous Learning through Neural Network Evolution in Hardware,
In Proc of 51st Annual IEEE/ACM International Symposium on Microarchitecture(MICRO), 2018*

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



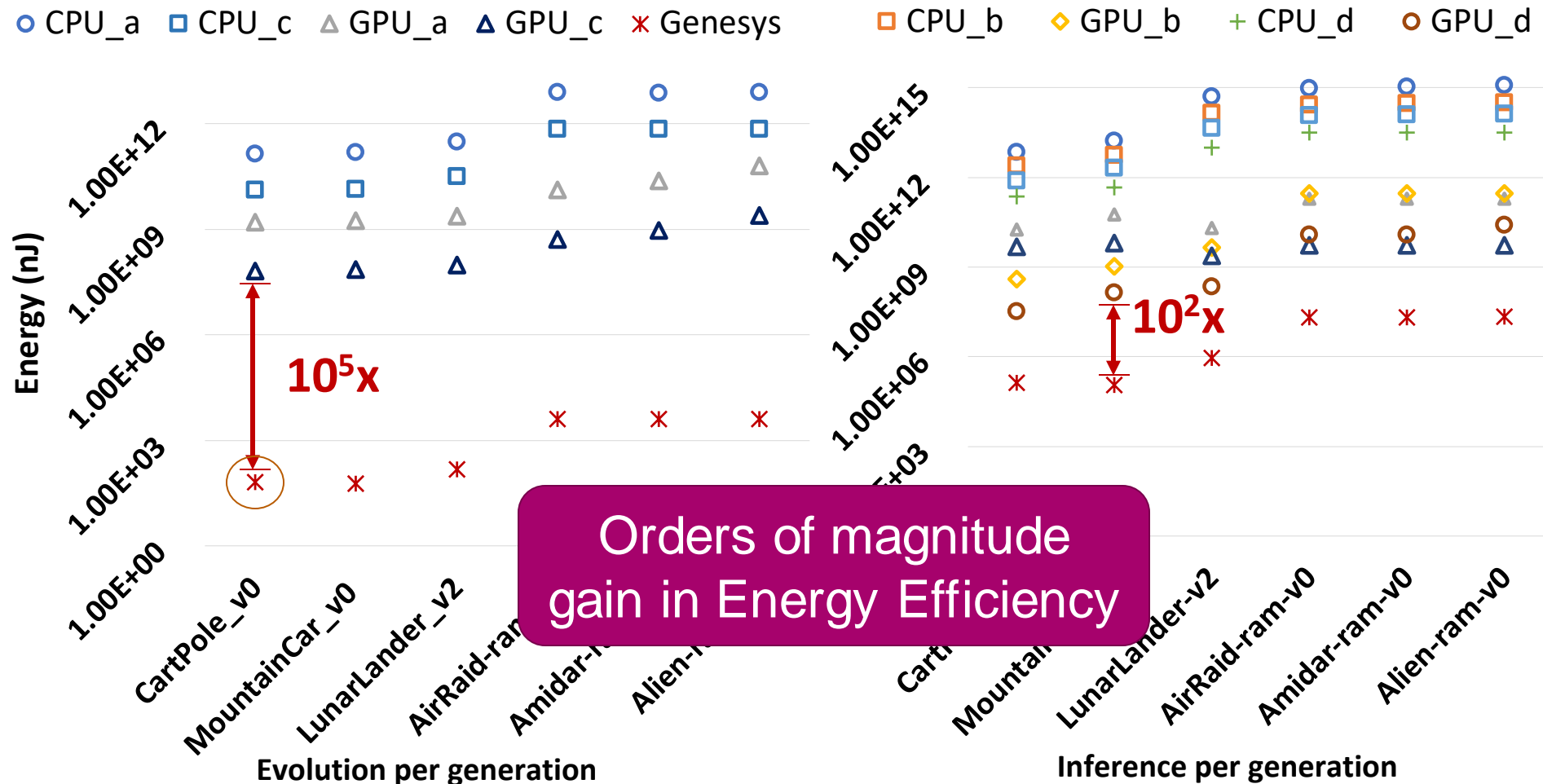
Evaluations

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

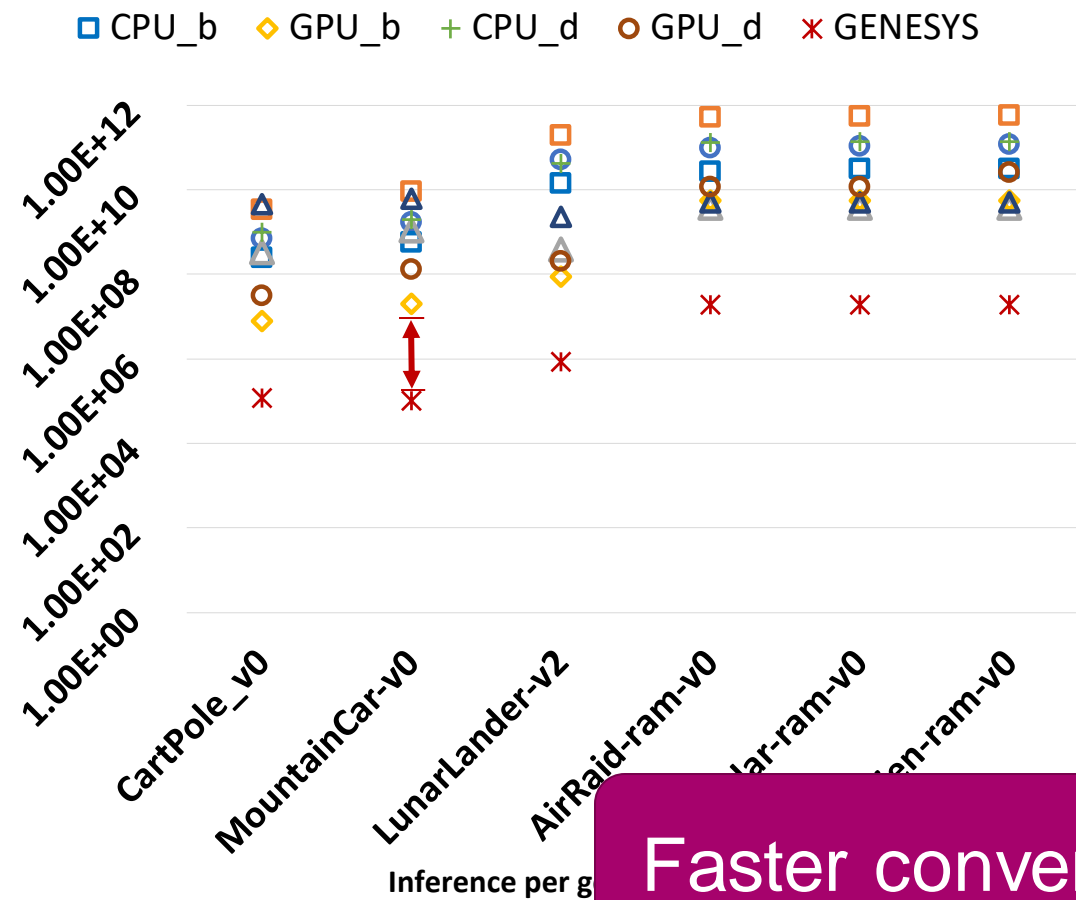
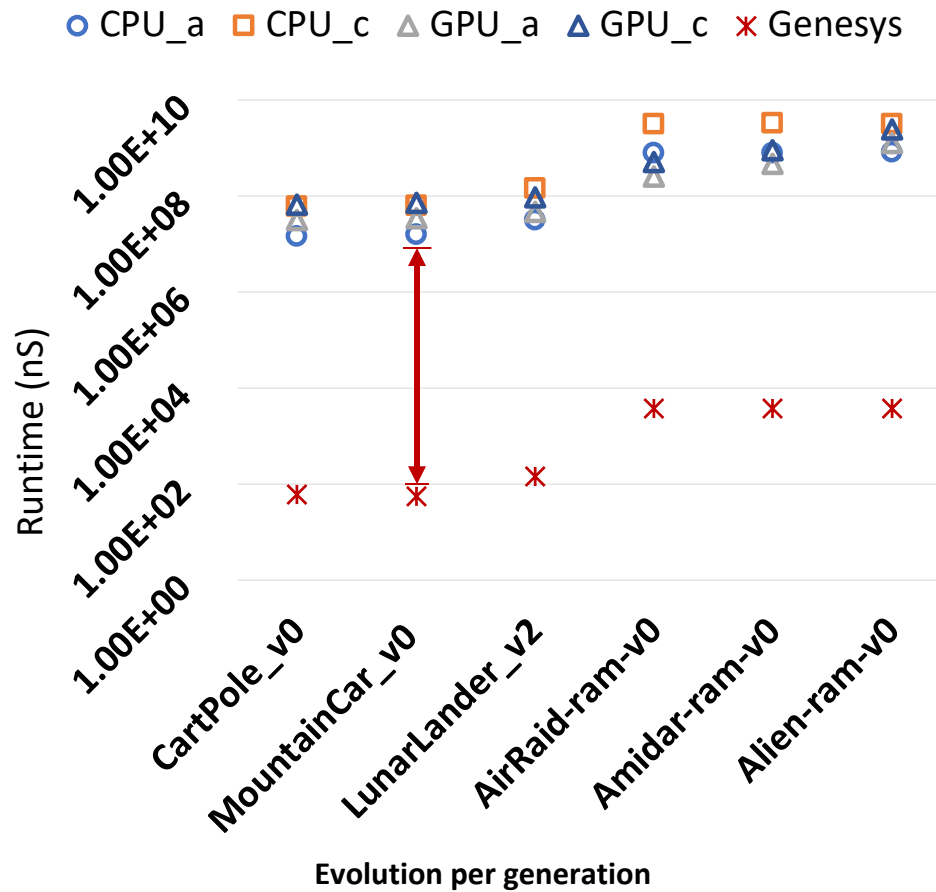
PLP (GLP) - Population (Gene) Level Parallelism

BSP - Bulk Synchronous Parallelism (GPU)

Evaluations: Energy



Evaluations: Runtime



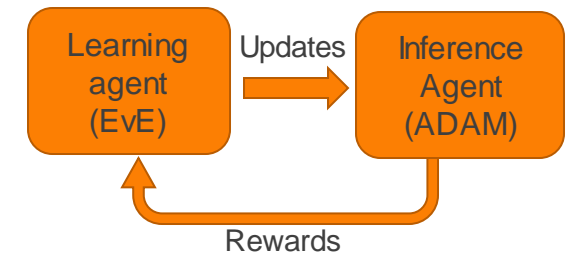
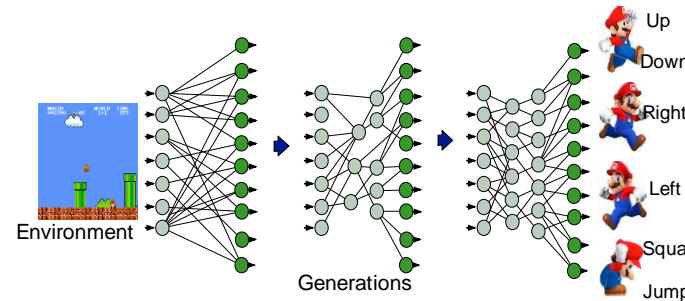
Faster convergence

Conclusion

- **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

Thank you!

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