



# Energy Efficient and Scalable Neuromemristive Computing Substrates

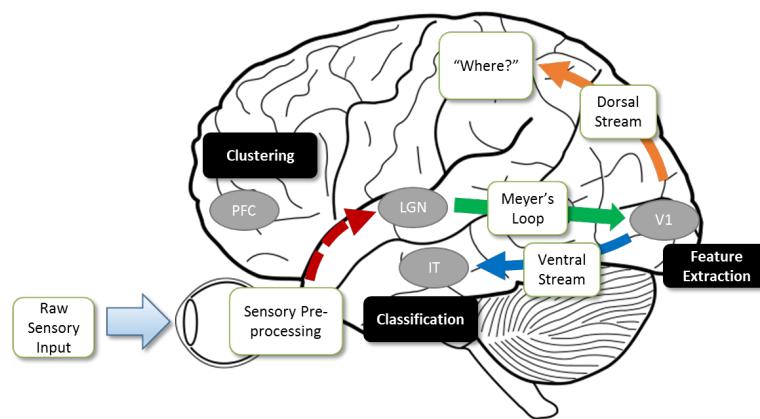
Dhireesha Kudithipudi<sup>†</sup>, Cory Merkel<sup>‡</sup>, James Mnatzaganian<sup>†</sup>, Nicholas Soures<sup>†</sup>, Qutaiba Saleh<sup>†</sup>

<sup>†</sup>NanoComputing Research Lab  
Rochester Institute of Technology

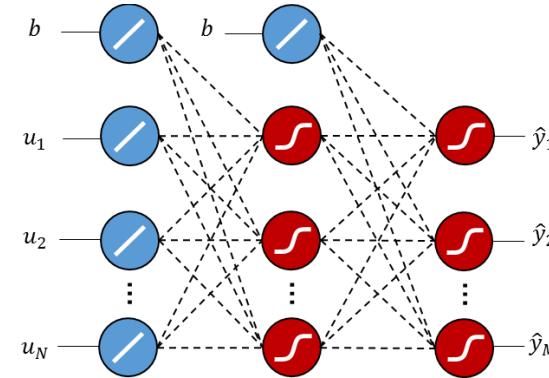
<sup>‡</sup>Information Directorate  
Air Force Research Laboratory

# Neuromemristive Systems

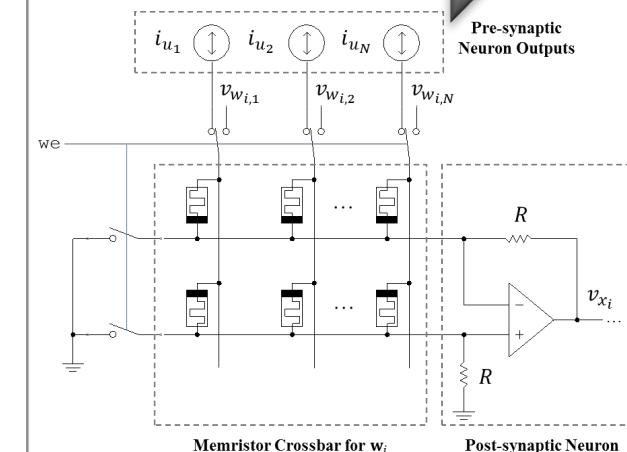
Study



Model



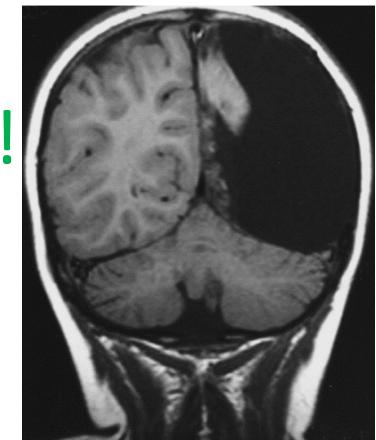
Implement



- **Brain-inspired adaptive computing platforms based on nanoscale resistive memory (memristors)**
- Memristor characteristics facilitate efficient computation and learning
- Improve the efficiency (over conventional computers) of *natural* processing tasks

# Why the Brain?

- Incredibly versatile
  - Can learn anything!
- Energy efficient
  - $\sim 10^{16}$  ops/sec @ a few watts!
- Robust/Resilient
  - Functions with noise!
  - Unreliable and damaged components!



[1] *Scientific American*

[2] [www.transhumanist.com](http://www.transhumanist.com)

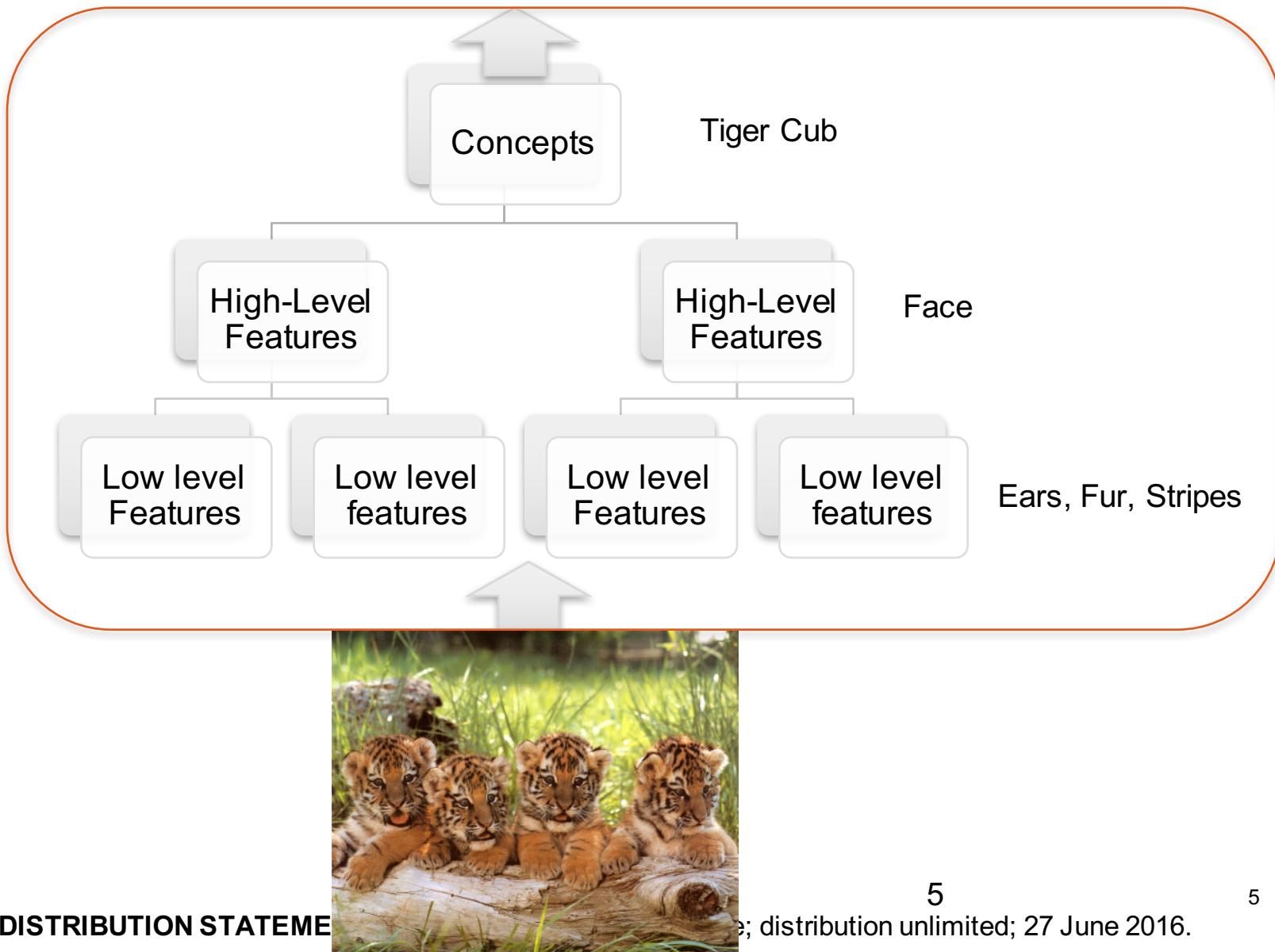
[3] <http://sites.psu.edu/cigerber02141993/2014/04/14/are-two-halves-better-the-one-whole/>

# Brain vs. Conventional Computing

	Brain	Computer
<b>Signals</b>	Mixed signal	Digital
<b>Precision</b>	Low	High
<b>Parallelism</b>	Very high	Low
<b>Information Density</b>	High	Low
<b>Processing-Memory</b>	Closely-coupled	Separated
<b>Mutability</b>	Plastic	Constant

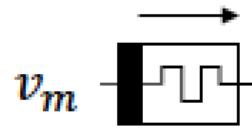
Brain-like computing is better for massively parallel applications with noisy data and relaxed precision requirements

# Neuromemristive Systems



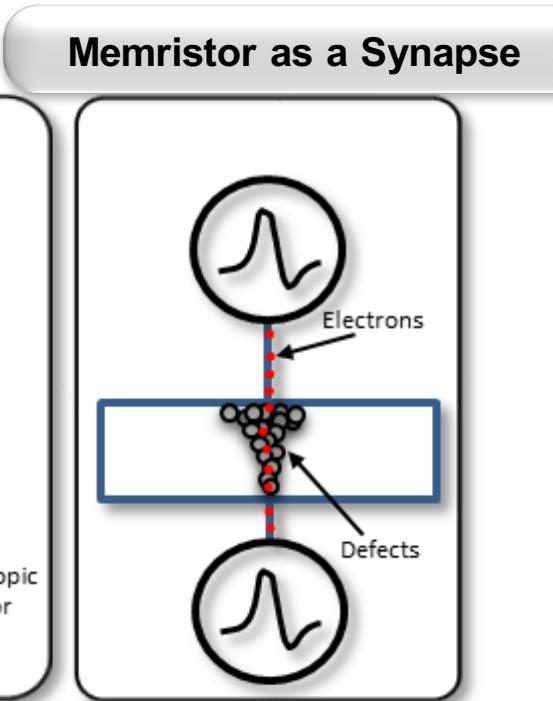
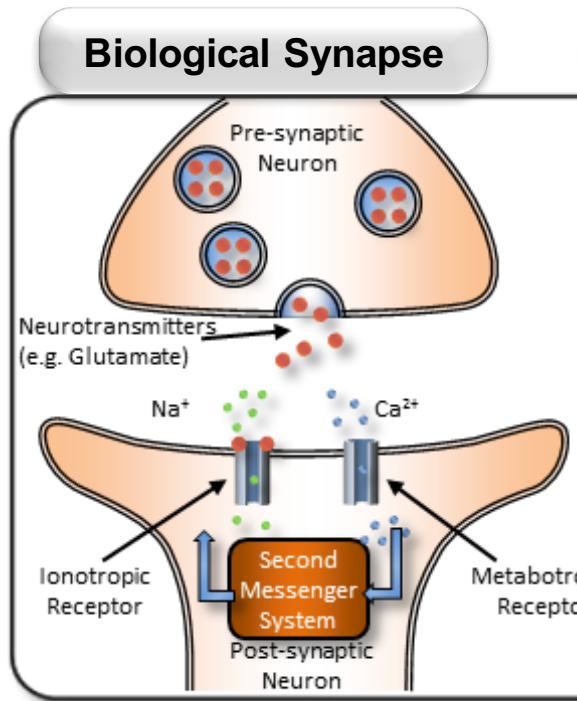
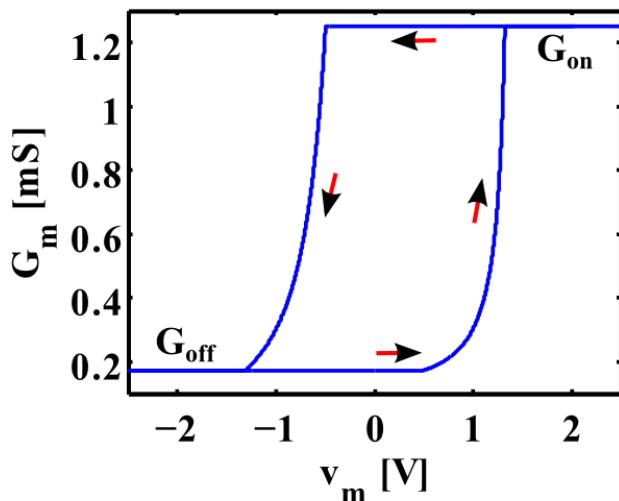
# Memristors for Plasticity

$$i_m = G_m(\gamma)v_m$$



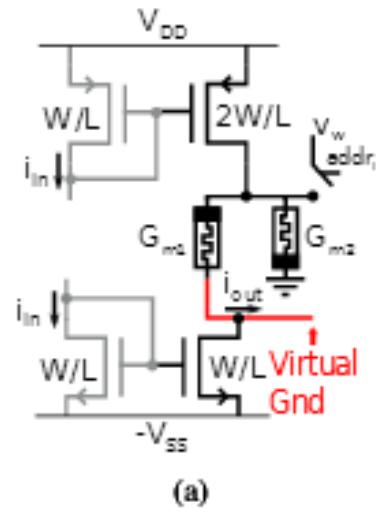
$$G_m(\gamma)$$

2-terminal device with state-dependent Ohm's Law

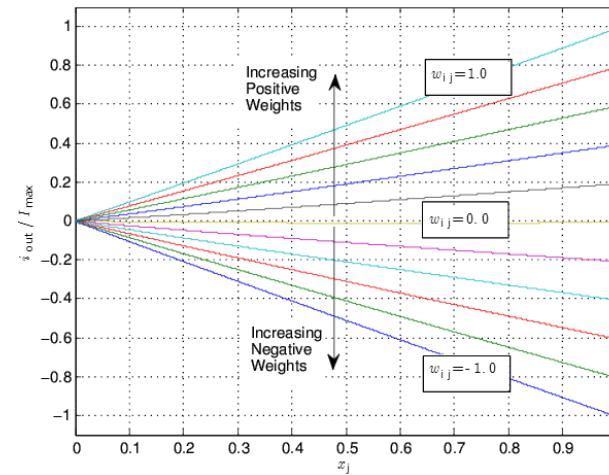


- Compatibility with CMOS
- Memristor characteristics facilitate efficient computation and learning

# Reconfigurable Synapses



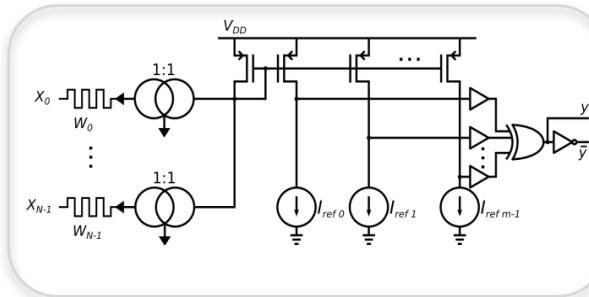
(a)



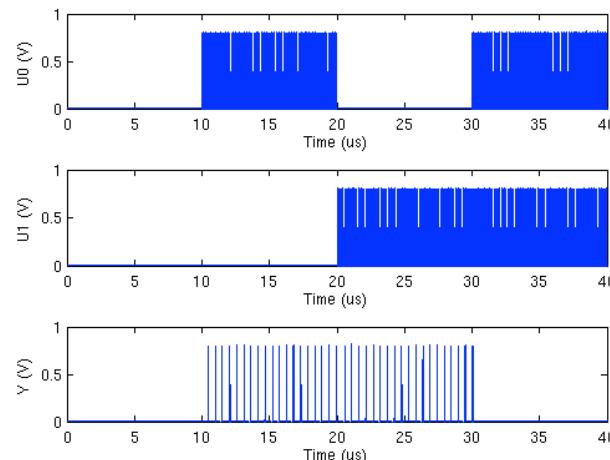
$$s_{ij} = \left( 2 \frac{G_{m1j}}{G_{m1j} + G_{m2j}} - 1 \right) x_j = w_{ij} x_j$$

Inhibitory and Excitatory Synapses

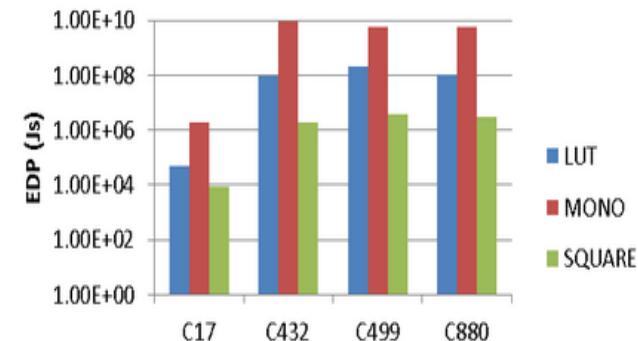
# Reconfigurable Neurons



Non-Monotonic Neuron



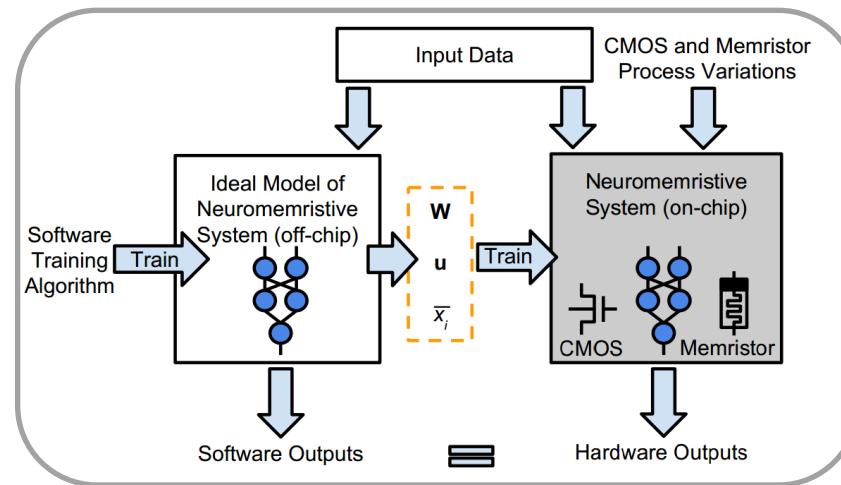
Edge Detection



Energy-Delay Product

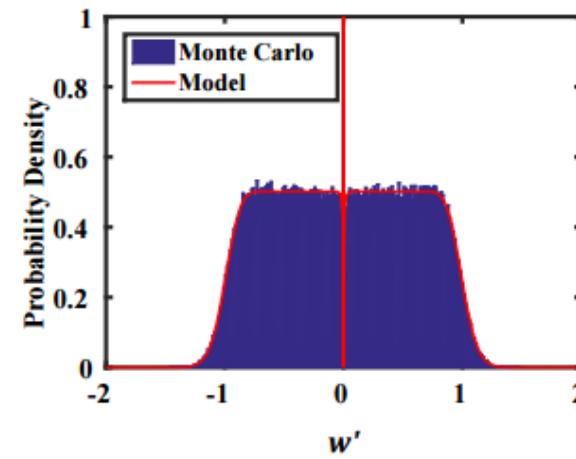
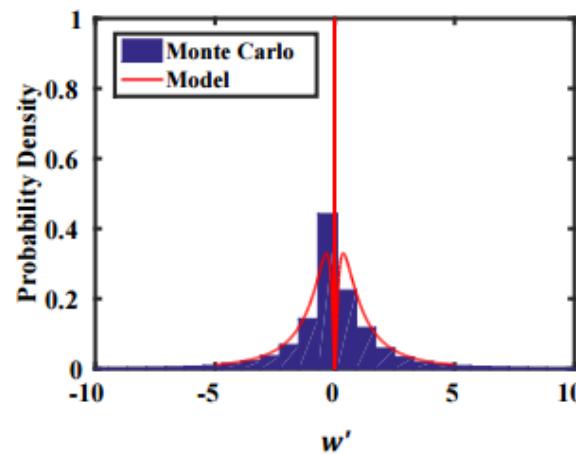
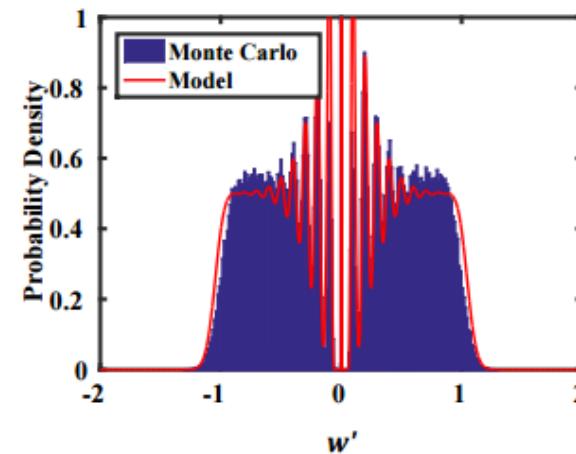
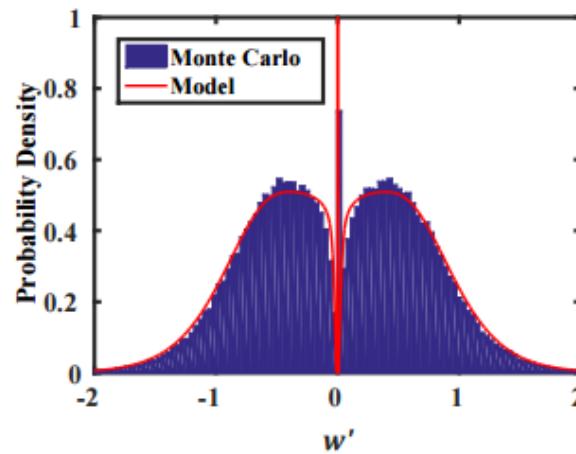
# On-Chip Training

- Variation is exploited in the training process



# Random Weight Synapses

- Exploit random mismatch in current mirrors
- Control distribution with sizing



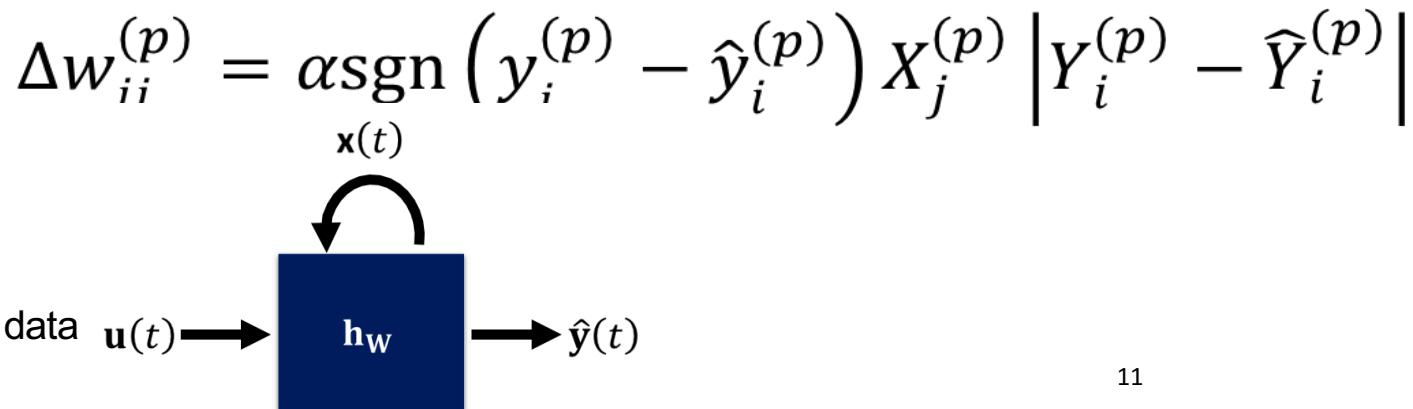
# On-Chip Training

- Least-mean squares (LMS) training algorithm:

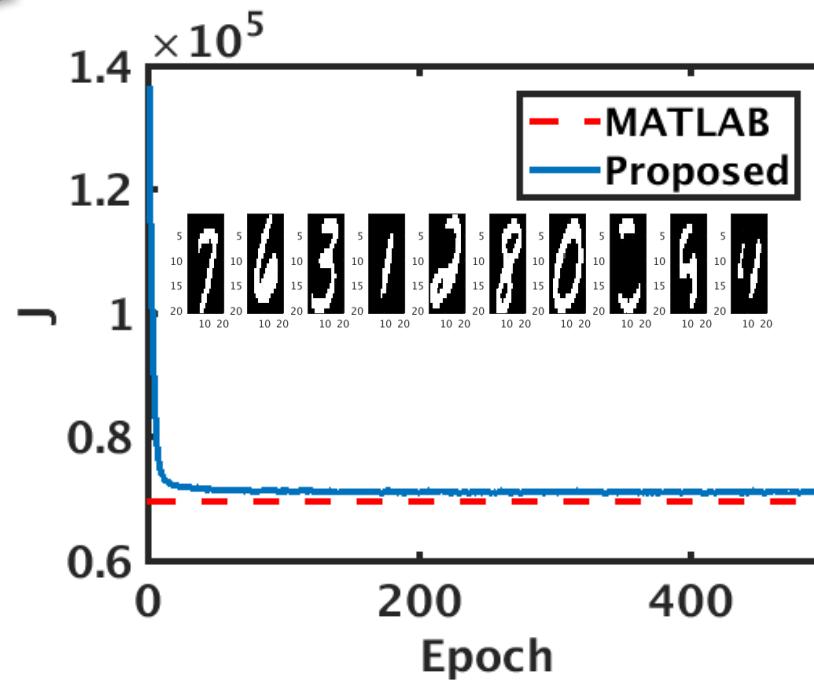
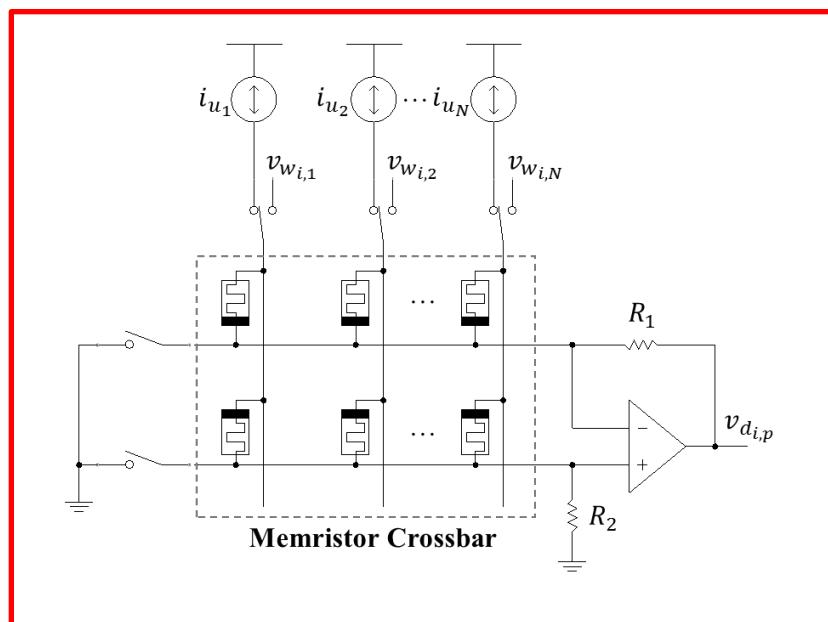
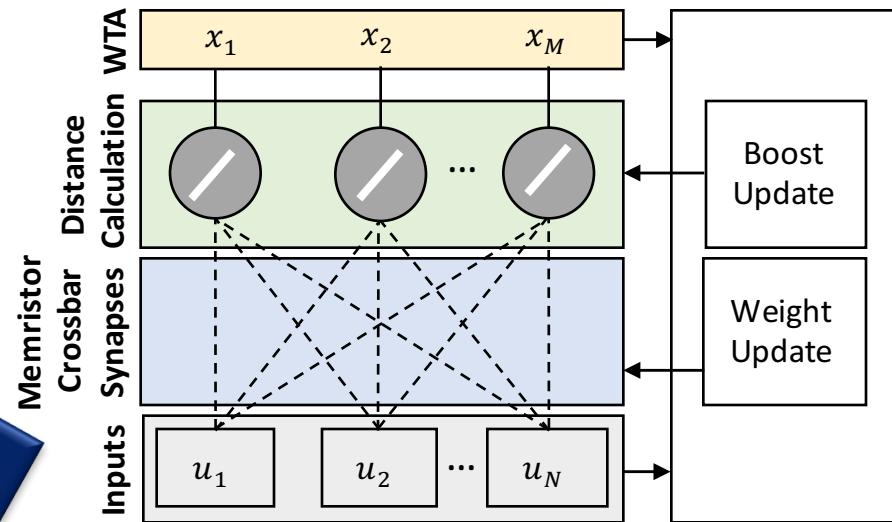
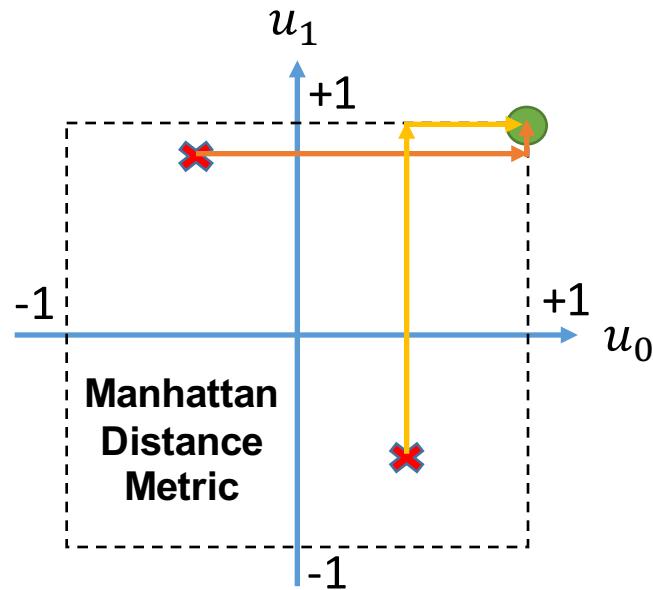
$$\Delta w_{ij}^{(p)} = \alpha u_j^{(p)} \left( y_i^{(p)} - \hat{y}_i^{(p)} \right)$$

Expected output

- Converted to stochastic LMS (SLMS) using proposed method:



# Unsupervised Clustering

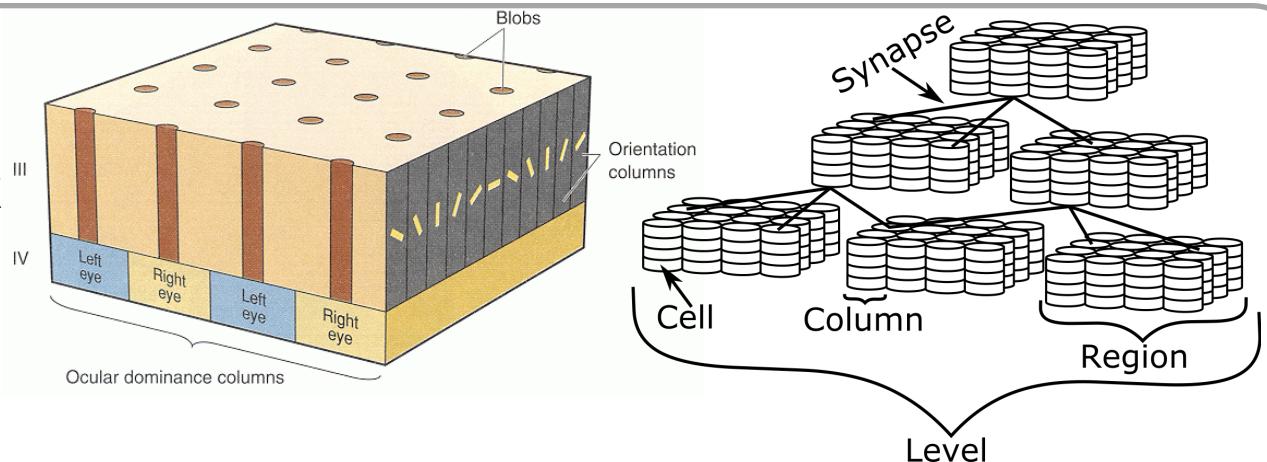


**DISTRIBUTION STATEMENT A.** Approved for public release; distribution unlimited; 27 June 2016.

# Hierarchical Temporal Memory

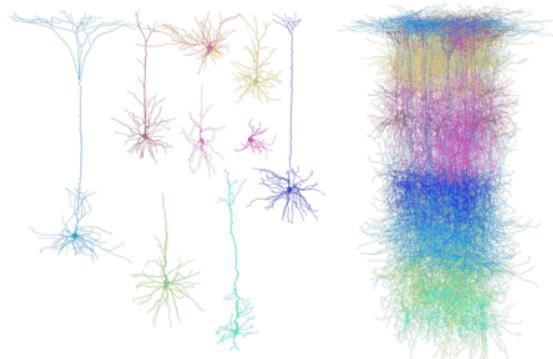
## Inspiration/Motivation

- Inspired by the neocortex
- Highly parallelizable
- Suitable for hardware design

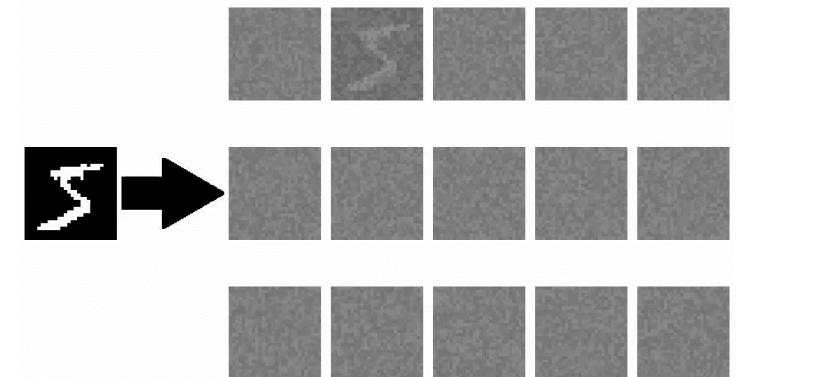


## Critical Aspects

- Spatiotemporal data
- Online, unsupervised learning
- Classification & prediction
- Distinct learning components
- Customizable architecture



## Applications/Results

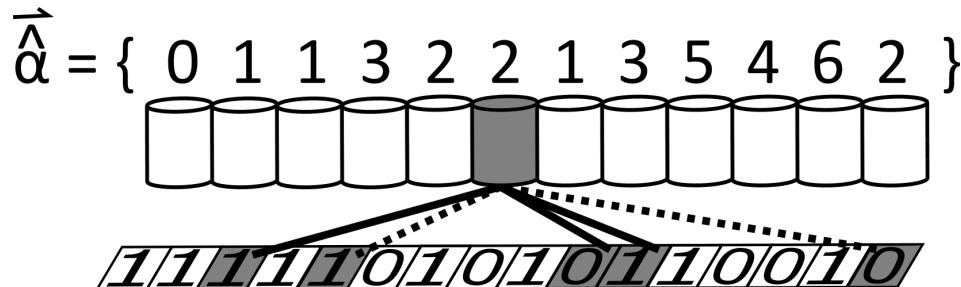


**Given** : (A) 11111-.....-.....  
**Predicted** : (?) .....-.....-.....

# Mathematical Formalization of the Spatial Pooler


 Example column  
 Example column's input  
 — Connected synapse  
 ... Unconnected synapse

<http://arxiv.org/abs/1601.06116>



<https://github.com/tehtechguy/mHTM>

$$\vec{\alpha} \equiv \begin{cases} \vec{\hat{\alpha}}_i \vec{b}_i & \vec{\hat{\alpha}}_i \geq \rho_d, \\ 0 & \text{otherwise} \end{cases} \quad \forall i \quad \vec{\hat{\alpha}}_i \equiv \mathbf{X}_i \bullet \mathbf{Y}_i$$

Overlap

$$\hat{\mathbf{c}} \equiv I(\vec{\alpha}_i \geq \vec{\gamma}_i) \quad \forall i \quad \text{inhibition}$$

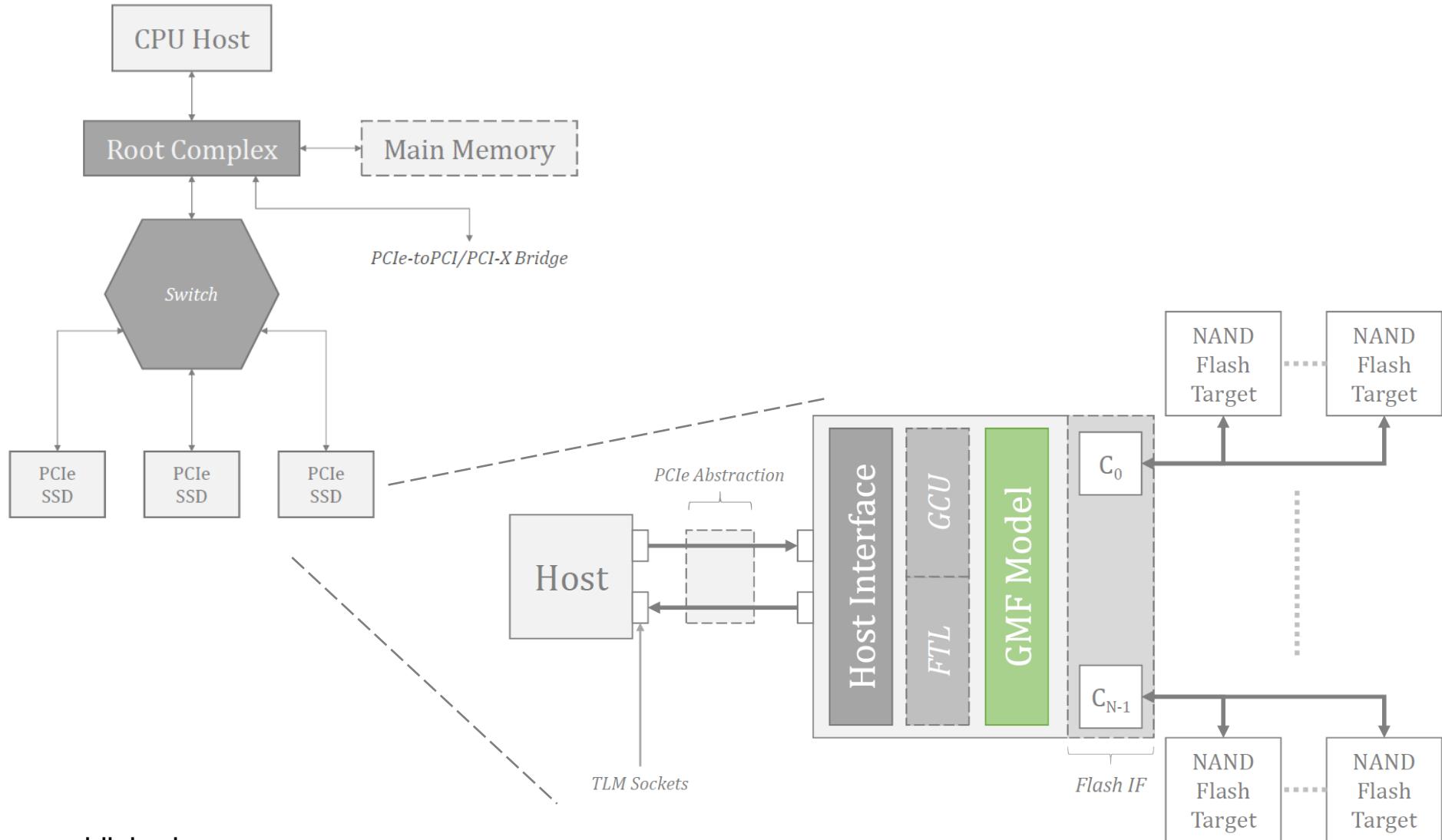
$$\vec{\gamma} \equiv \max(\text{kmax}(\mathbf{H}_i \odot \vec{\alpha}, \rho_c), 1) \quad \forall i$$

$$\delta\Phi \equiv \vec{\hat{\mathbf{c}}}^T \odot (\phi_+ \mathbf{X} - (\phi_- \neg \mathbf{X})) \quad \text{Learning}$$

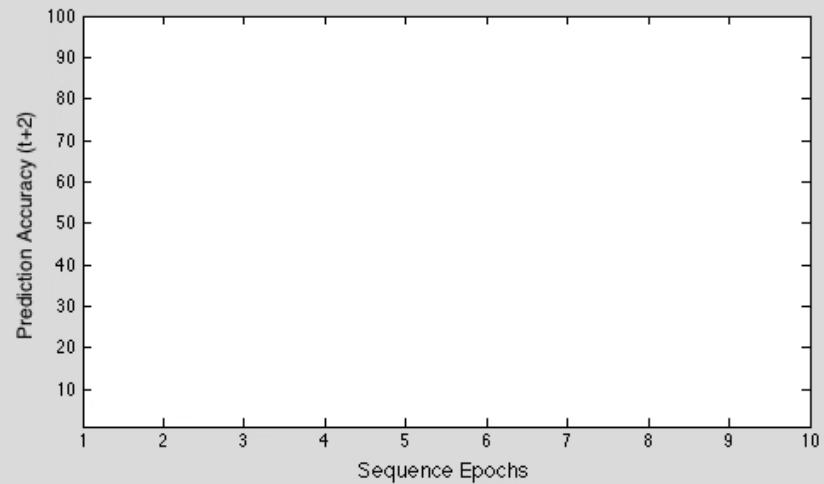
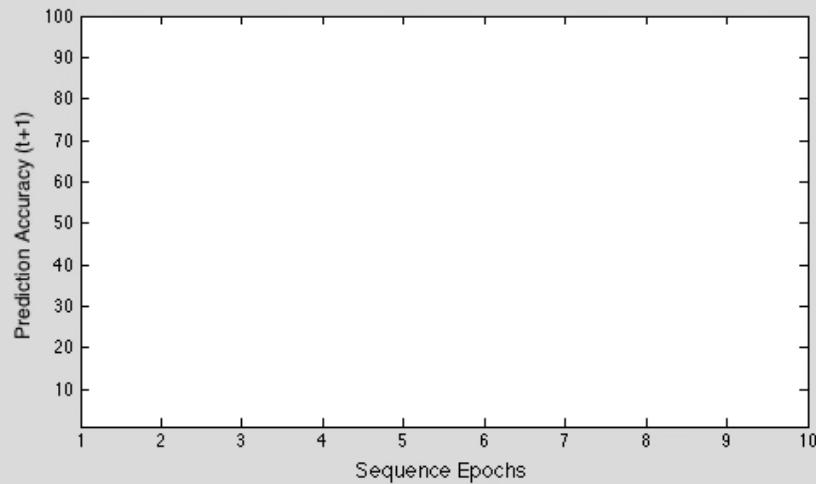
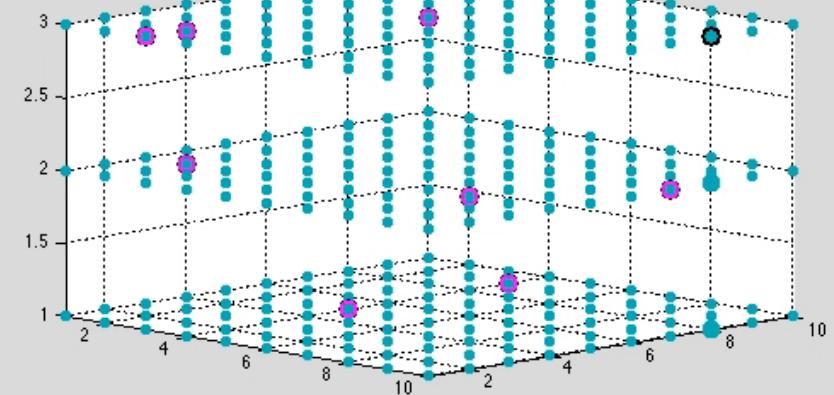
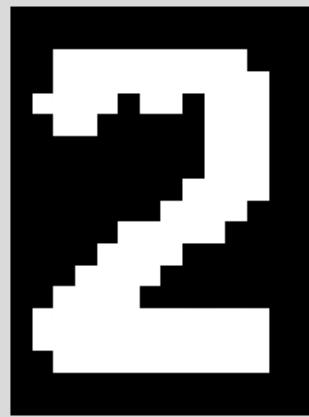
$$\Phi \equiv \text{clip}(\Phi \oplus \delta\Phi, 0, 1)$$

# Reconfigurable HTM Architecture

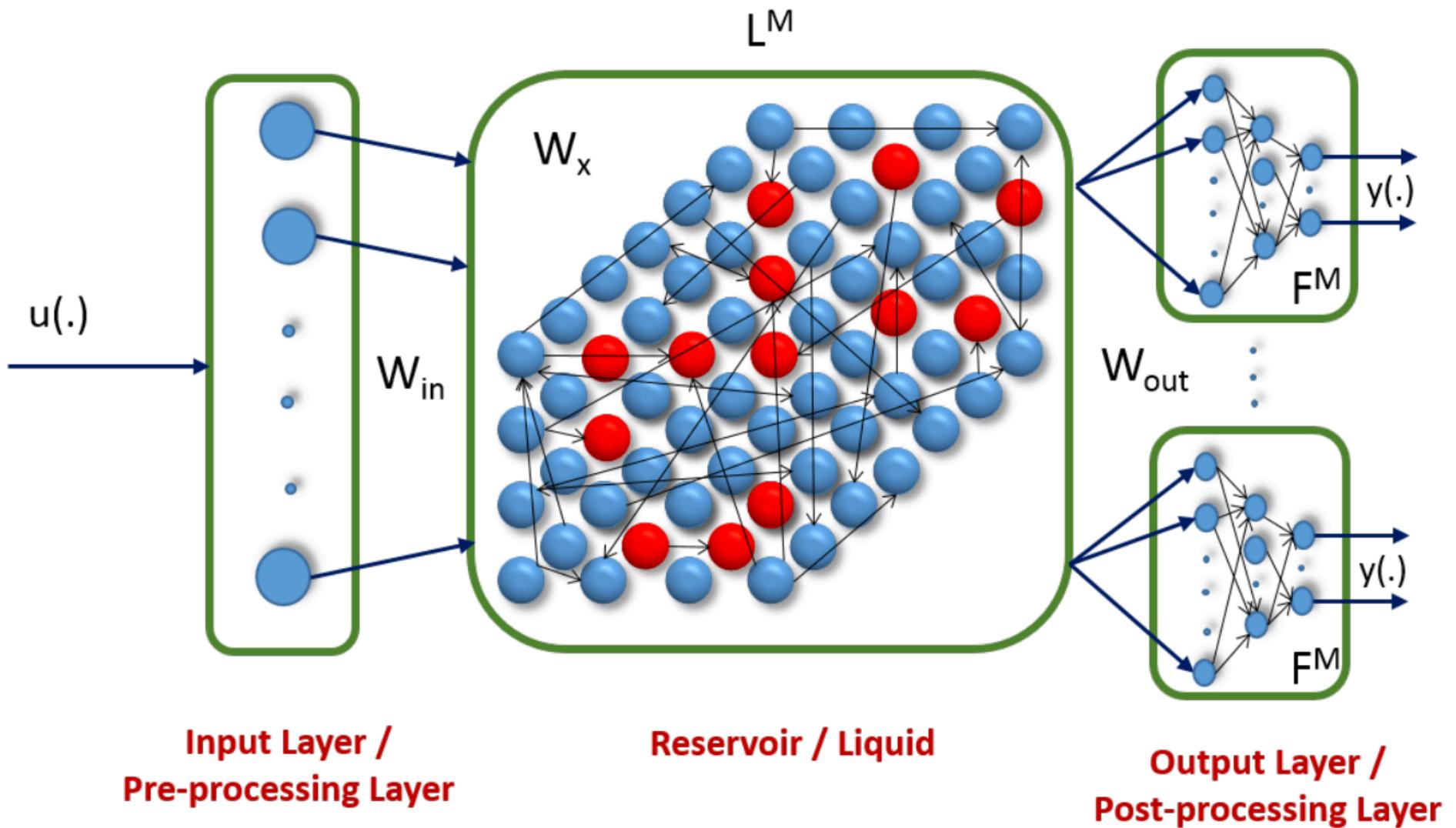
Storage processor units may leverage PCIe SSD technology



# Reconfigurable HTM Architecture

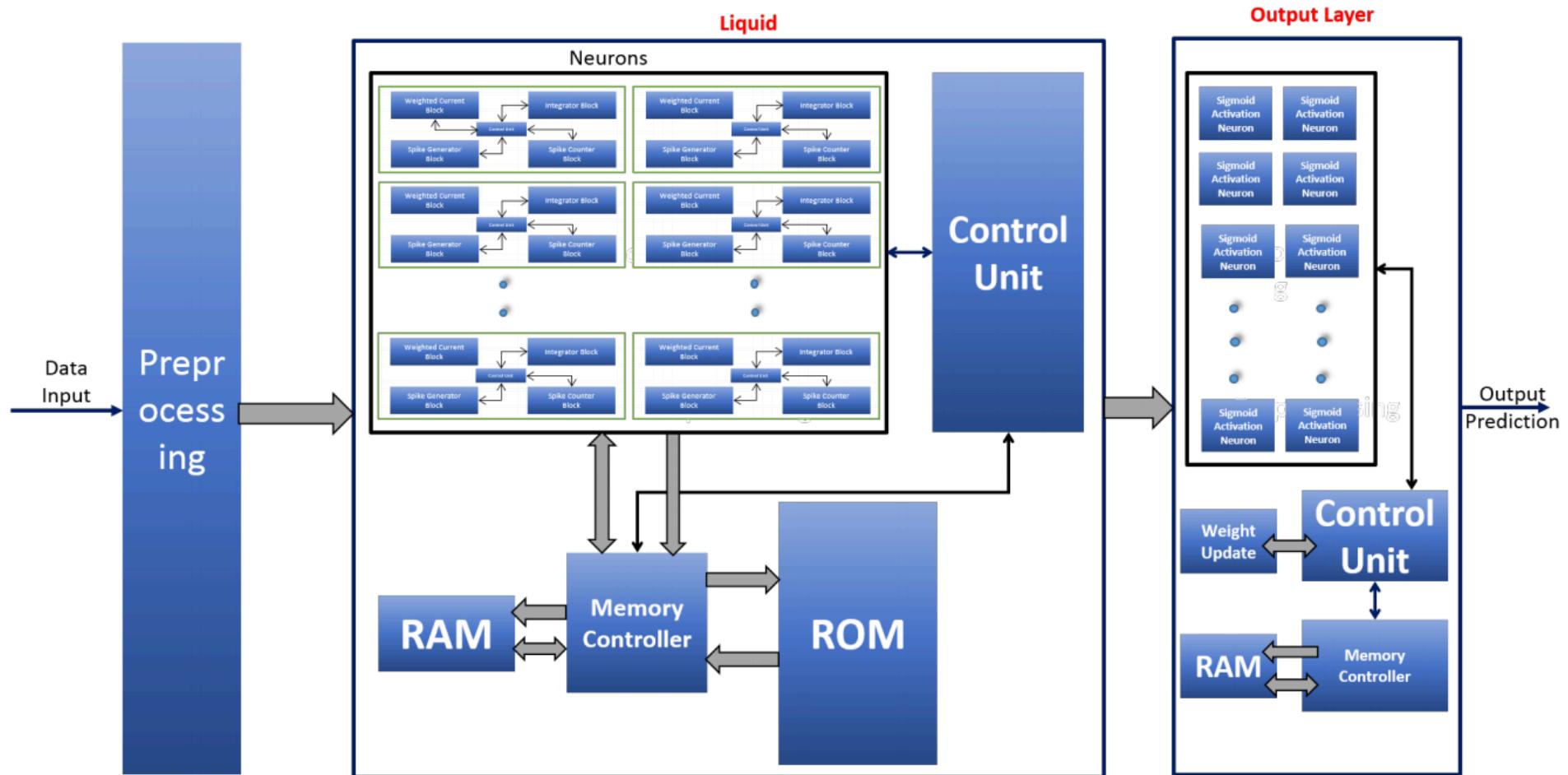


# Generalizable Intelligence Engine

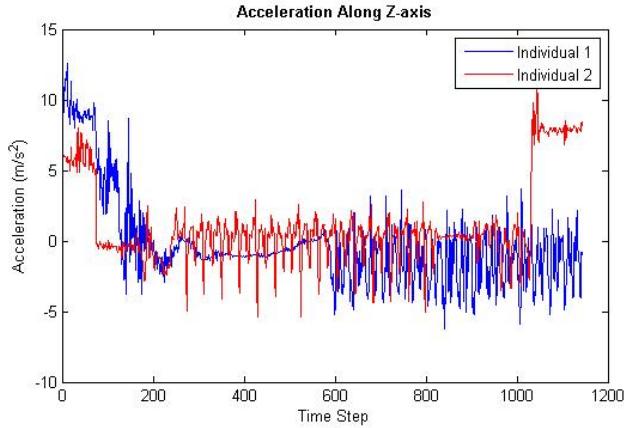
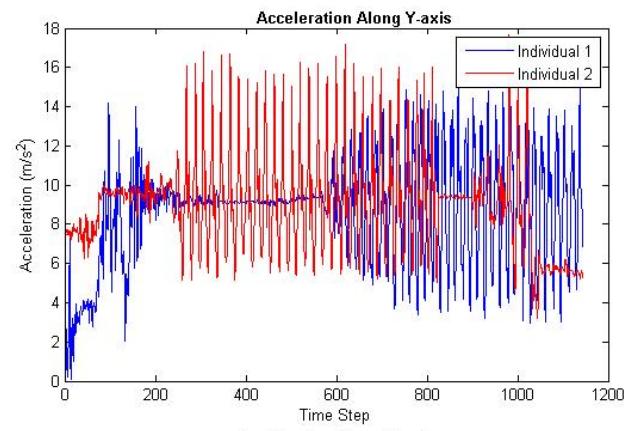
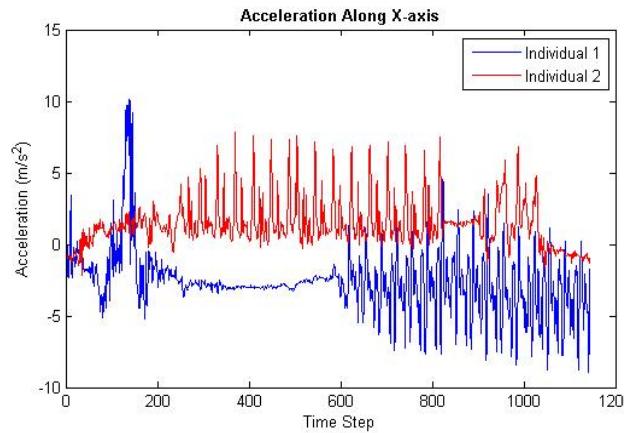


DISTRIBUTION STATEMENT A. Approved for public release; distribution unlimited; 27 June 2016.

# Reconfigurable Reservoir Architecture



# Reconfigurable Reservoir Architecture



		Training Confusion Matrix	
Output Class	Target Class	0	1
		0	1
0	0	42 46.7%	1 1.1%
	1	1 1.1%	46 51.1%
1	0	97.7% 2.3%	97.9% 2.1%
	1	97.7% 2.3%	97.8% 2.2%

		Validation Confusion Matrix	
Output Class	Target Class	0	1
		0	1
0	0	10 52.6%	0 0.0%
	1	0 0.0%	9 47.4%
1	0	100% 0.0%	100% 0.0%
	1	100% 0.0%	100% 0.0%

		Test Confusion Matrix	
Output Class	Target Class	0	1
		0	1
0	0	11 57.9%	0 0.0%
	1	0 0.0%	8 42.1%
1	0	100% 0.0%	100% 0.0%
	1	100% 0.0%	100% 0.0%

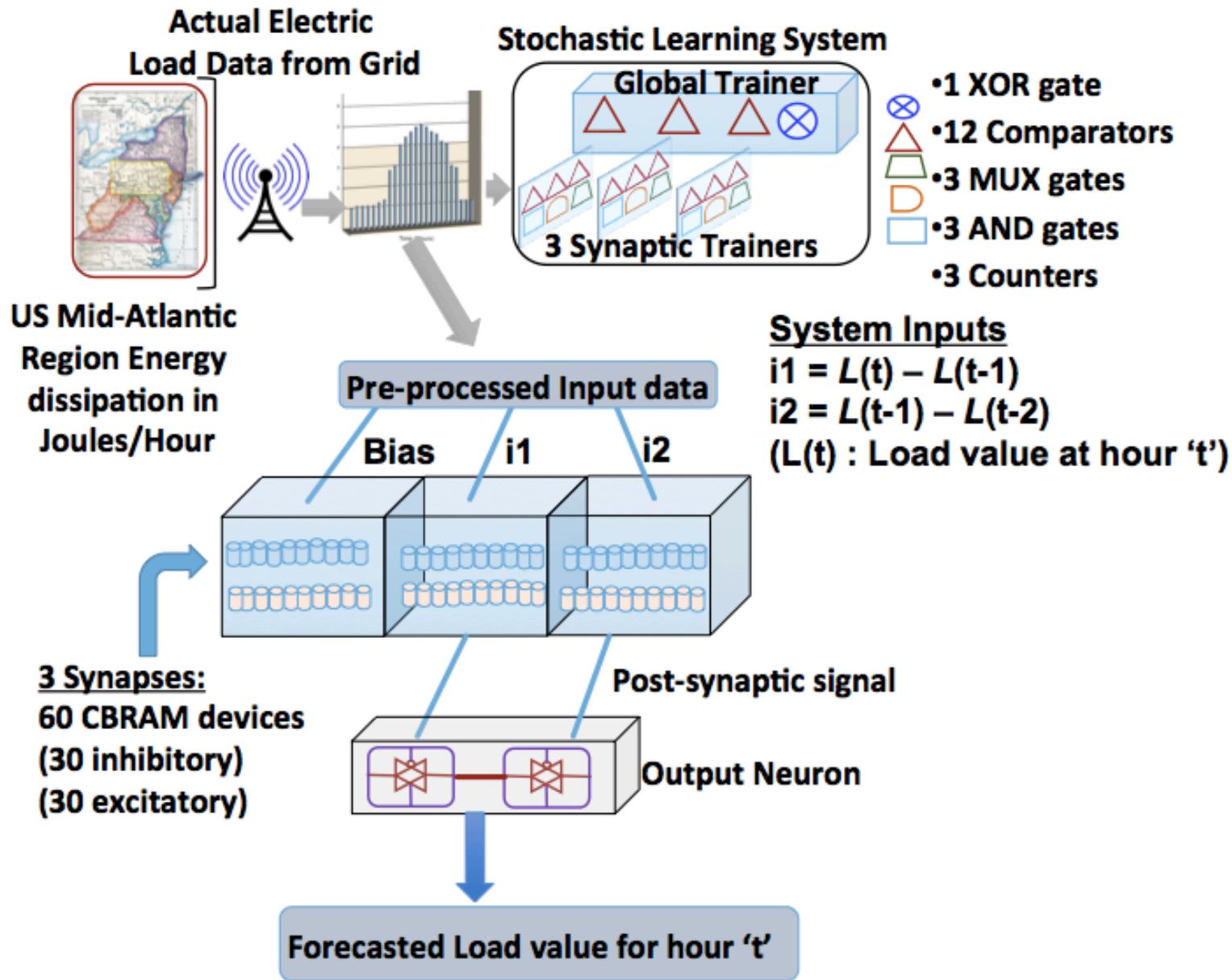
		All Confusion Matrix	
Output Class	Target Class	0	1
		0	1
0	0	63 49.2%	1 0.8%
	1	0 0.8%	63 49.2%
1	0	98.4% 1.6%	98.4% 1.6%
	1	98.4% 1.6%	98.4% 1.6%

User Authentication based on Gait Patterns

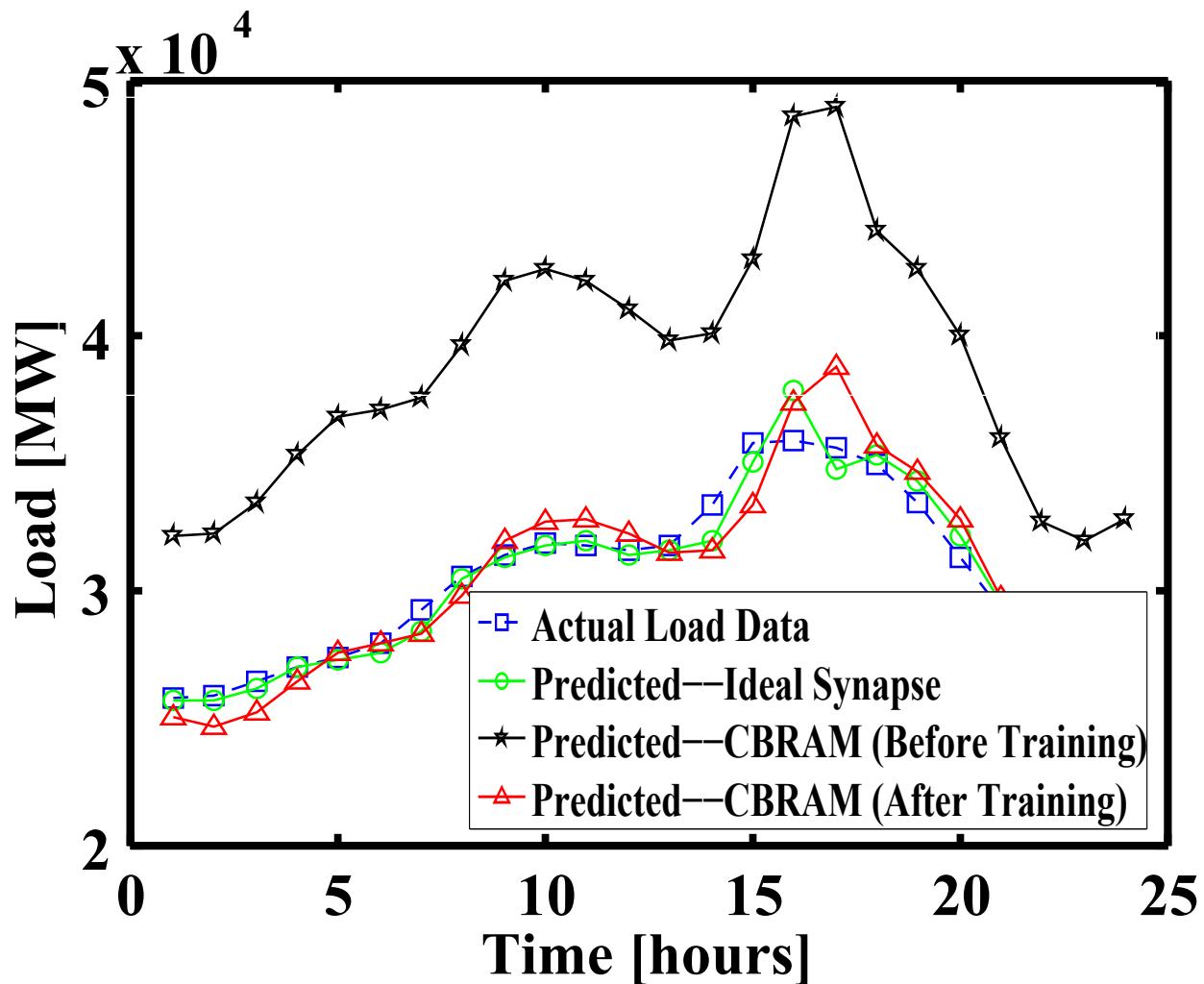
Rebooting'16

DISTRIBUTION STATEMENT A. Approved for public release; distribution unlimited; 27 June 2016.

# Smart Grid Load Forecasting



# Smart Grid Load Forecasting



# Summary

- Reconfigurability is integral to the nature of computation
  - Precomputation is occurring in communication channels
  - No standardized metrics/benchmarks to evaluate
    - Designing technology agnostic vs. technology aware systems
- Looking forward
  - one shot learning ....

# Team & Collaborators

