

Brain-like Computing – Scalable Low-Power Chips for Learning & Optimal Control

Prof. Pinaki Mazumder
University of Michigan
Ann Arbor, MI 48105

**Acknowledgement: National Science
Foundation (CCF & ECCS)**

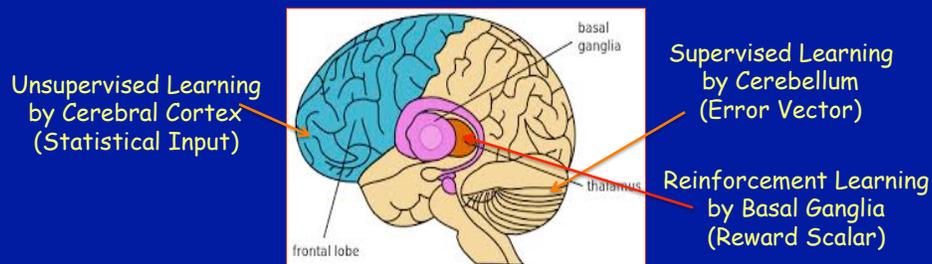
1

Brain-like Computing

Facets of Brain-like Computing:

1. Self-Healing
2. Associative Memory
3. Cognition
4. Learning & Plasticity

Transcribing Basal Ganglia on Silicon
for Optimal Control Algorithms

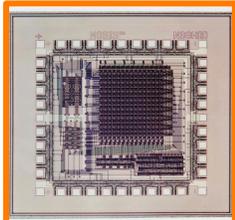


Adaptive Hardware Platform for Nonlinear Optimal Control, Swarm Intelligence, Robot Control, and Markov Decision Process (MDP)

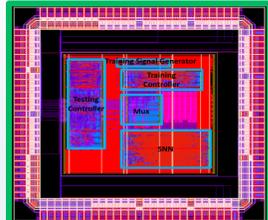
Mazumder Group's Neuromorphic Research

<h3 style="text-align: center; color: yellow;">Self-Healing VLSI Design (1989-1996)</h3> <p>Hopfield Neural Net as Algorithmic Hardware for Spare Allocation by Node Cover over Bipartite Graph</p> <ul style="list-style-type: none"> IEEE Trans. on CAS, 1993 IEEE Trans. on CAS, 1993 IEEE Trans. on Computer, 1996 	<h3 style="text-align: center; color: yellow;">Learning based VLSI Chips (2010- Now)</h3> <p>STDP Learning for Position Detector STDP Learning for Virtual Bug Navigation STDP Learning for XOR/Edge Detection Deep Learning for Pattern recognition</p> <p>Q-Learning for Maze Search Algorithm on Memristor Array</p> <ul style="list-style-type: none"> Proceedings of the IEEE, 2012 Nano Letters Journal, 2010 IEEE Nanotechnology, 2011 IEEE Nanotechnology, 2014 IEEE Cellular Neural Networks, 2012
<h3 style="text-align: center; color: yellow;">Cognition Chips using Cellular Neural Networks (2008-2013)</h3> <ul style="list-style-type: none"> Color Image Processing Velocity Tuned Filter Memristor/RRAM based CNN RTD+HEMT based CNN <ul style="list-style-type: none"> IEEE Trans. on VLSI, 2009 IEEE Trans. on Nanotechnology, 2008 IEEE Trans. on Neural Nets, 2014 IEEE Trans. on Nanotechnology, 2013 ACM Journal on Emerging Technologies in Computing Systems, 2013 	<h3 style="text-align: center; color: yellow;">Reinforcement Learning/Actor-Critic NN (2016 – Now)</h3> <ul style="list-style-type: none"> IEEE Trans. on Computer, 2016 IEEE Trans. on Neural Nets, 2018 IEEE Trans. on Circuits & Systems, 2018

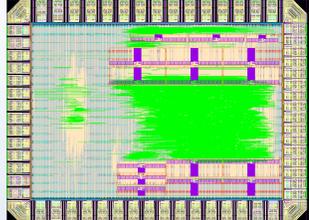
Neural-Inspired CMOS Chips Designed by Mazumder Group



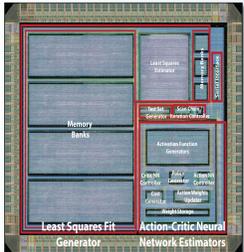
Self-Healing Chip, 1991



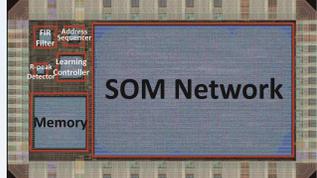
Synaptic Plasticity Chip, 2013



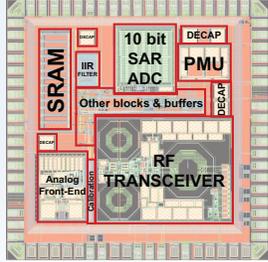
Deep Learning Chip, 2016



Actor-Critic Reinforcement Learning Chip, 2016



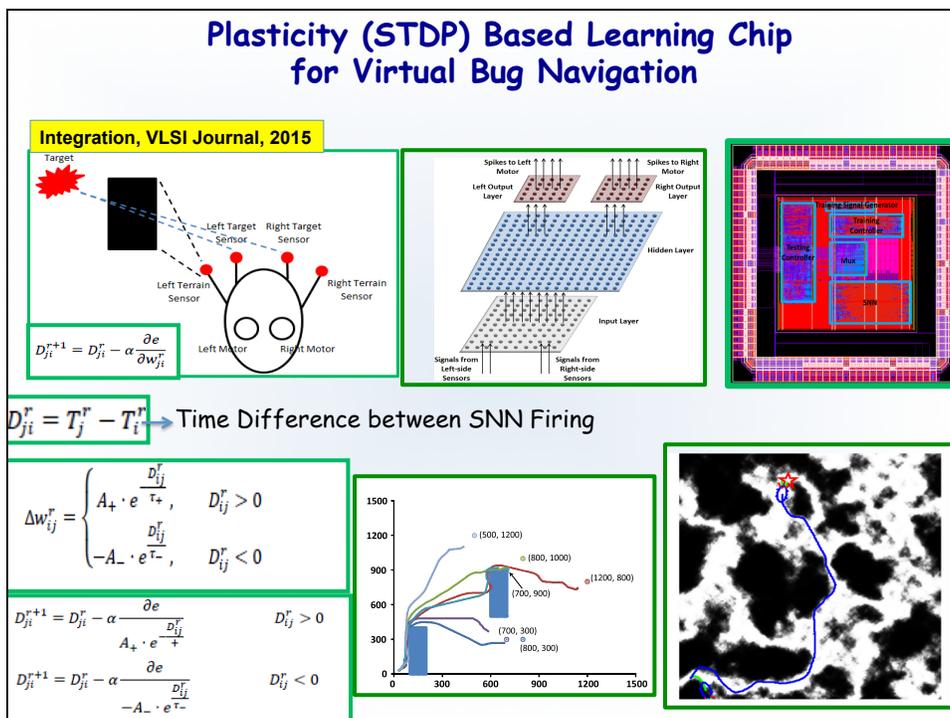
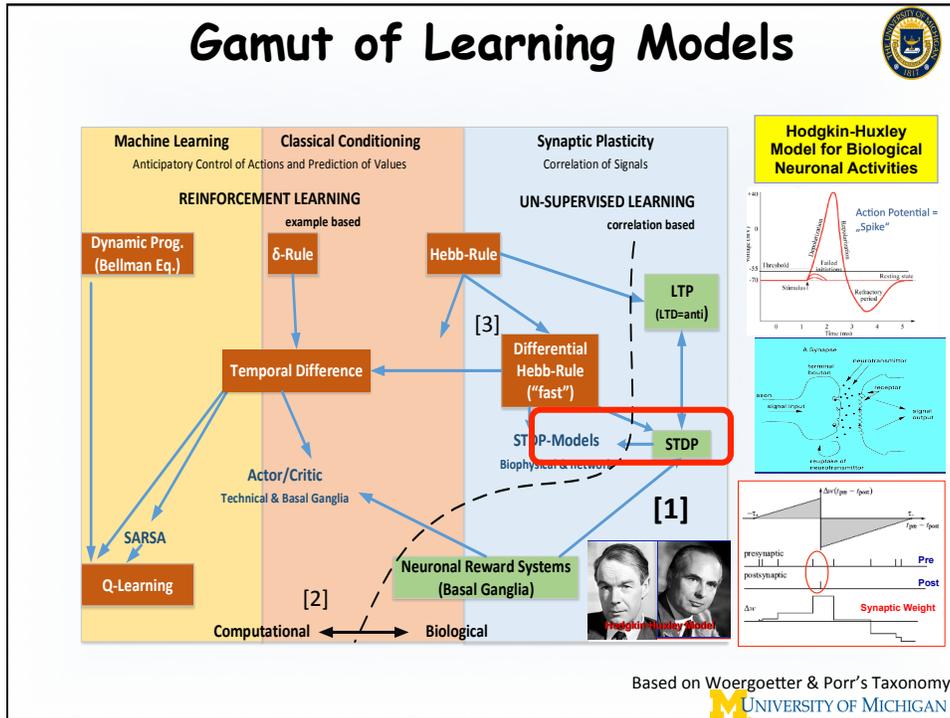
Self-Organizing Map Chip for ECG Clustering 2016



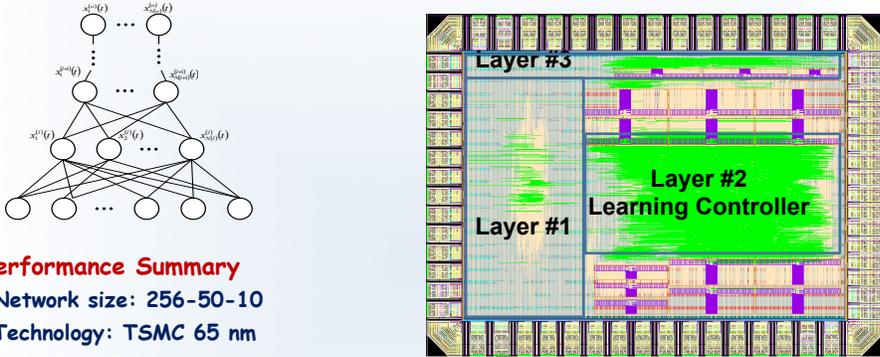
Body-Sensing Network with Wireless Transceiver

Figure 4.4: The proposed ECG clustering SOM chip layout.

© 2016 Prof. P. Mazumder of University of Michigan ⁴

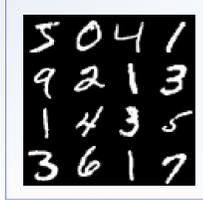


Discretized STDP Based Deep Neural Net



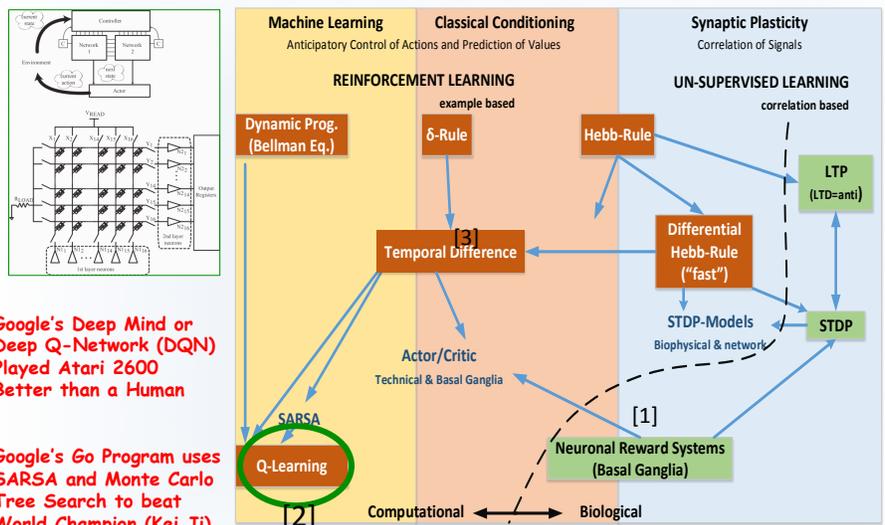
Performance Summary

- Network size: 256-50-10
- Technology: TSMC 65 nm
- Bit width of synapse: 24 bits
- Clock frequency: ~100 MHz
- Power: ~100 mW
- Time needed for one learning iteration: 30 μ s
- Speedup : ~700x v. digital computer
- Energy: ~500 pJ/spike @ VDD = 1.2



© 2016 Prof. P. Mazumder of University of Michigan

Gamut of Learning Models



Machine Learning
Anticipatory Control of Actions and Prediction of Values

REINFORCEMENT LEARNING
example based

Dynamic Prog. (Bellman Eq.)

δ -Rule

Temporal Difference [3]

Actor/Critic
Technical & Basal Ganglia

SARSA

Q-Learning [2]

Classical Conditioning
Correlation of Signals

UN-SUPERVISED LEARNING
correlation based

Hebb-Rule

Differential Hebb-Rule ("fast")

STDP-Models
Biophysical & network

STDP

Neuronal Reward Systems (Basal Ganglia) [1]

LTP (LTD=anti)

Synaptic Plasticity
Correlation of Signals

Computational \longleftrightarrow Biological

Based on Woergoetter & Porr's Taxonomy

Google's Deep Mind or Deep Q-Network (DQN) Played Atari 2600 Better than a Human

Google's Go Program uses SARSA and Monte Carlo Tree Search to beat World Champion (Kei Ji) in a Go contest.

Q-Learning Hardware - Reinforcement Learning

SARSA: State-Action-Reward-State-Action

$NewEstimate \leftarrow OldEstimate + StepSize * (Target - OldEstimate)$

$$Q(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha_t}_{\text{learning rate}} \cdot \left(\underbrace{r_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)$$

$$\tilde{Q}(s_t, a_t) \leftarrow \tilde{Q}(s_t, a_t) \times (1 - \alpha_t(s_t, a_t)) + \alpha_t(s_t, a_t) \times r_t + \alpha_t(s_t, a_t) \times \max_{a_{t+1}} [\tilde{Q}(s_{t+1}, a_{t+1})]$$

$$\max_{a_{t+1}} [\tilde{Q}(s_{t+1}, a_{t+1})] = \tilde{Q}(s_t, a_t) + \delta_t$$

$$\tilde{Q}(s_t, a_t) \leftarrow \tilde{Q}(s_t, a_t) + \alpha_t(s_t, a_t) \times (r_t + \delta_t)$$

Evaluative Feedback (Rewards)

© 2016 Prof. P. Mazumder of University of Michigan

Performance of Memristor Q-Learning Hardware

(a) (b)

Memristor Model Used in Matlab Simulation

steps vs iteration

(15, 28)

steps vs iteration

(26, 45)

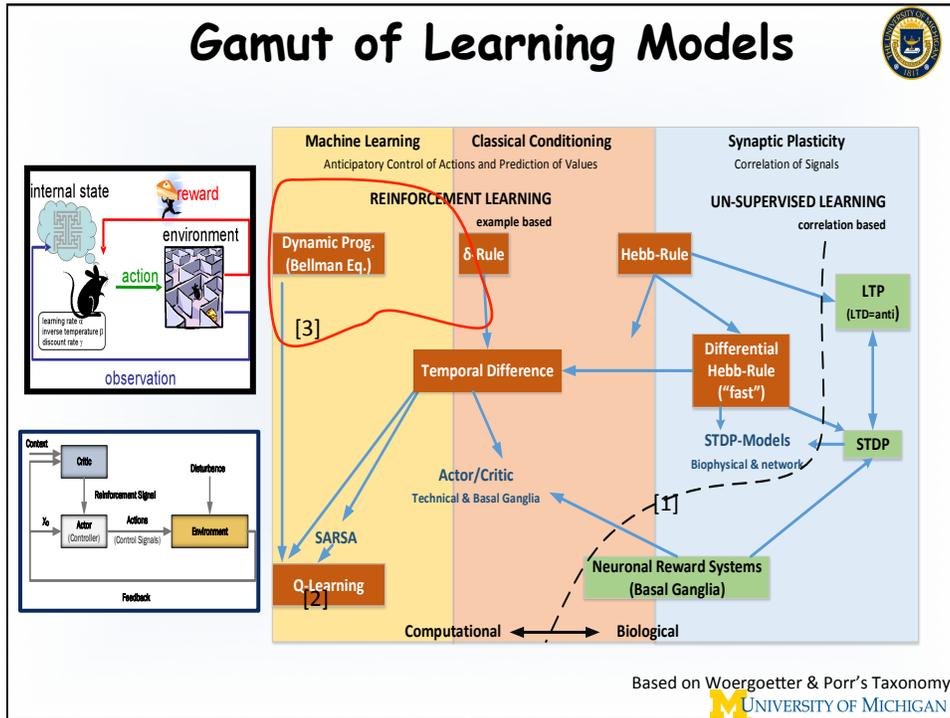
Synapse States after 1st and 2nd Iterations

Ebong & Mazumder, IEEE Nano 2014

Superior Application of Q Learning by Google, UK

Human-level control through deep reinforcement learning
Nature, v. 518, pp. 529-533, Feb 2015.

Played Atari 2600
Better than a Human



Actor-Critic NN CMOS Chip

- **Post-Layout Metrics**
 - Technology: 65 nm
 - Area: 550 $\mu\text{m} \times 550 \mu\text{m}$
 - Arithmetic precision: 24-bit fixed-point
 - Supply voltage: 1.2 V
 - Clock frequency: 175 MHz
 - Power Consumption: 25 mW
 - 270x speed up over software ADP

Controlling balance, rotation, movement of microrobots that operate on batteries

Provide smart sensing strategies for energy-constrained sensors

550 μm

550 μm

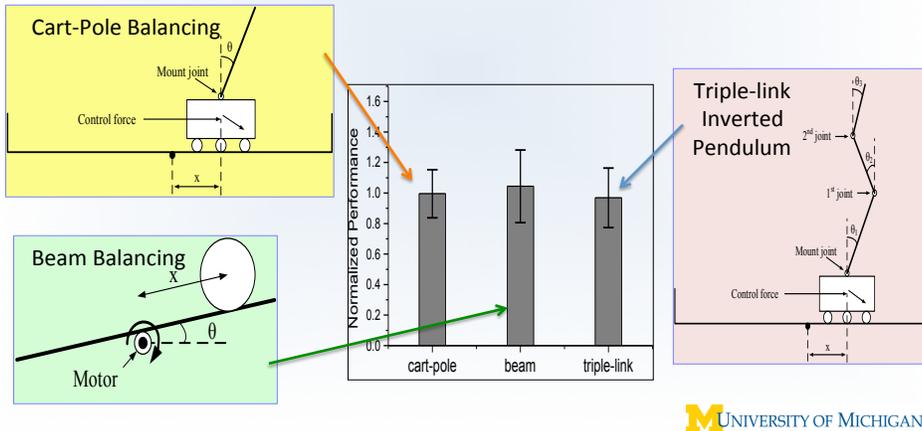
UNIVERSITY OF MICHIGAN

Benchmark Results



▪ **Benchmark Tests**

- Cart-Pole Balancing
- Beam Balancing
- Triple-link Inverted Pendulum
- Results obtained from the accelerator are normalized w.r.t software results
- Hardware solution consumes 25 mW and is faster than software by 225 times



Problem Statement



▪ **Optimal Control/Decision-Making Problem**

- Looking for the optimal policy/strategy
- Target is to maximize reward & minimize cost
- Solve using Bellman equation: The optimal equality.
- Continuous-space problem is handled by Hamilton-Jacobi-Bellman equation



System Model: $x(t+1) = f[x(t), a(t)]$

Reward to be maximized: $J[x(t)] = \sum_{k=1}^{\infty} \gamma^{k-1} r[x(t+k)]$

Bellman equation: $J^*[x(t)] = \max_{a(t)} (r[x(t)] + \gamma J^*[x(t+1)])$

- $x(t)$ -state of the system at time t
- $a(t)$ -action performed at time t
- $J[x(t)]$ -state value function at state time $x(t)$
- $J^*[x(t)]$ -optimal state value function
- $r[x(t+1)]$ -reward received at time t
- γ -discount factor used to encourage near-term reward



Traditional Applications of ADP in Software

Autopilot

- Used to control flight of a plane
- Cope with unexpected conditions

Power Systems

- Used as static compensators to regulate voltage in power systems

Communication Systems

- Self-learning call admission control scheme that can adapt for new communication environment

Autonomous Robot

- Navigation in unknown environment
- Control of a moving robotic arm

UNIVERSITY OF MICHIGAN

Action-Dependent Heuristic Dynamic Programming

- **Essentials**
- Model-free learning
- Learning by minimizing temporal difference error
- Critic: estimate the reward-to-go
- Actor: pick the right action to maximize future reward

$$\delta(t) = \hat{J}[x(t-1)] - \gamma \hat{J}[x(t)] - r[x(t)]$$

$$\hat{a}[x(t)] = \max_{\hat{a}[x(t)]} (r[x(t+1)] + \gamma \hat{J}[x(t+1)])$$

UNIVERSITY OF MICHIGAN

