# 7. Motor Control and Reinforcement Learning

## Outline

A. Action Selection and Reinforcement
B. Temporal Difference Reinforcement Learning
C. PVLV Model
D. Cerebellum and Error-driven Learning

3/5/17                    COSC 494/594 CCN                    2

## Sensory-Motor Loop

- Why animals have nervous systems but plants do not: *animals move*
  — a nervous system is needed to coordinate the movement of an animal's body
  — movement is fundamental to understanding cognition
- Perception conditions action
- Action conditions perception
  — profound effect of action on structuring perception is often neglected

3/5/17                    COSC 494/594 CCN                    3

## Overview

- Subcortical areas:
  o basal ganglia
    ➤ reinforcement learning (reward/punishment)
    ➤ connections to "what" pathway
  o cerebellum
    ➤ error-driven learning
    ➤ connections to "how" pathway
  o disinhibitory output dynamic
- Cortical areas:
  o frontal cortex
    ➤ connections to basal ganglia & cerebellum
  o parietal cortex
    ➤ maps sensory information to motor outputs
    ➤ connections to cerebellum

3/5/17                    COSC 494/594 CCN                    4

## Learning Rules Across the Brain

| Area | Learning Signal | | | | Dynamics | |
| | Reward | Error | Self Org | Separator | Integrator | Attractor |
|---|---|---|---|---|---|---|
| *Primitive* Basal Ganglia | +++ | - - - | - - - | ++ | - | - - - |
| Cerebellum | - - - | +++ | - - - | +++ | - - - | - - - |
| *Advanced* Hippocampus | + | + | +++ | +++ | - - - | +++ |
| Neocortex | ++ | +++ | ++ | - - - | +++ | +++ |

+ = has to some extent   …  +++ = defining characteristic – definitely has
- = not likely to have   …  - - - = definitely does not have

3/5/17                    COSC 494/594 CCN                    5
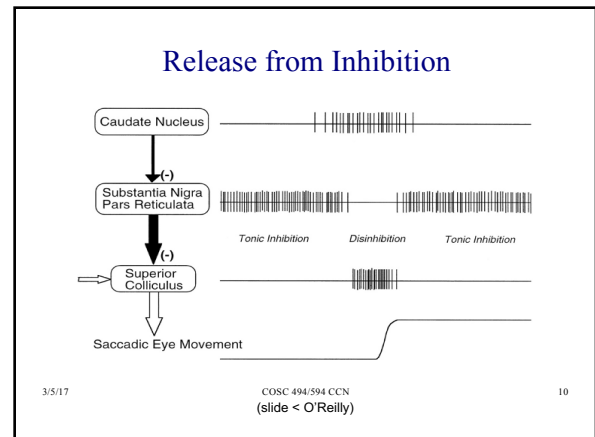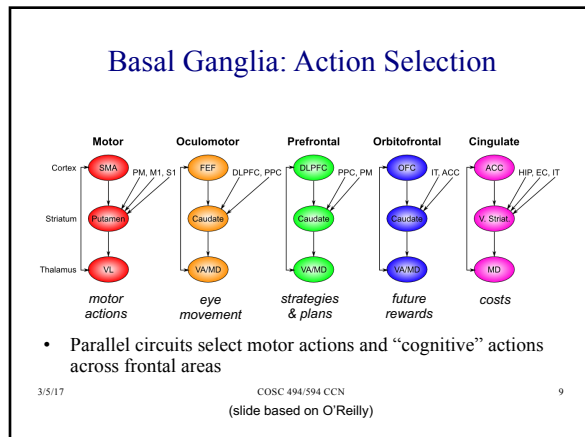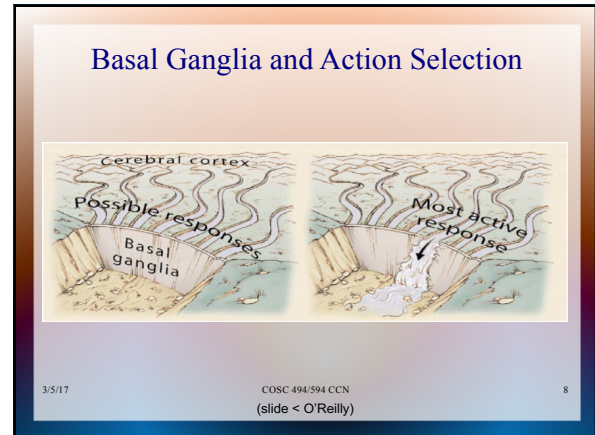(slide < O'Reilly)

## Primitive, Basic Learning…

| Area | Learning Signal | | | | Dynamics | |
| | Reward | Error | Self Org | Separator | Integrator | Attractor |
|---|---|---|---|---|---|---|
| *Primitive* Basal Ganglia | +++ | - - - | - - - | ++ | - | - - - |
| Cerebellum | - - - | +++ | - - - | +++ | - - - | - - - |

- Reward & Error = most basic learning signals (self organized learning is a luxury…)

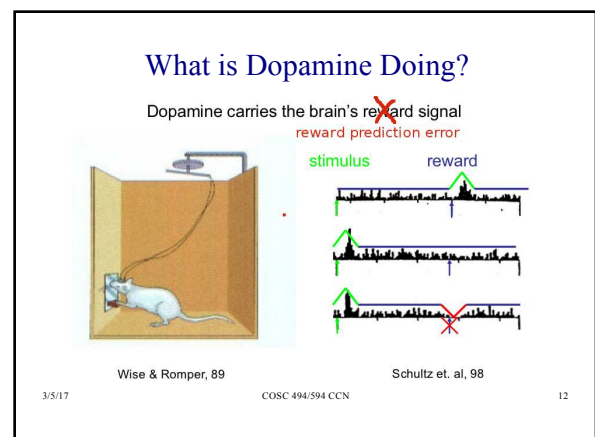- Simplest general solution to any learning problem is a *lookup table* = separator dynamics

3/5/17                    COSC 494/594 CCN                    6
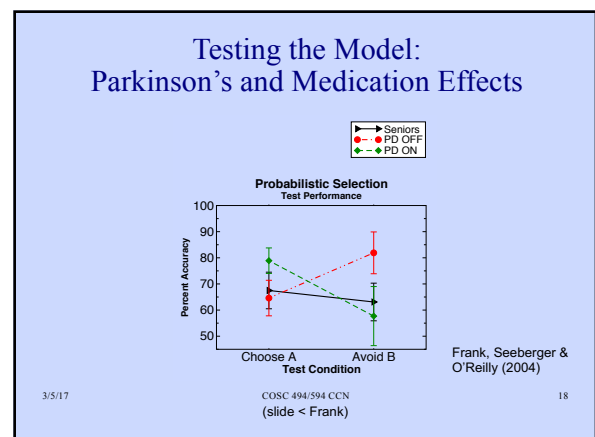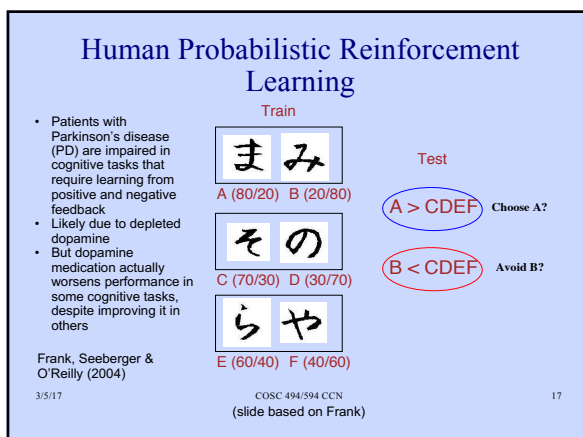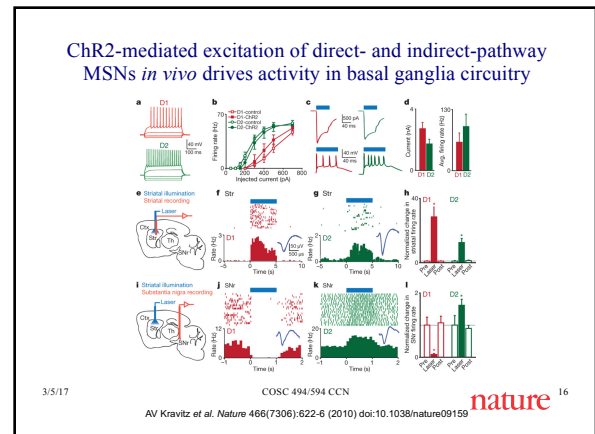(slide < O'Reilly)

## A. Action Selection and Reinforcement

3/5/17    COSC 494/594 CCN    7

## Basal Ganglia and Action Selection



3/5/17    COSC 494/594 CCN    8
(slide < O'Reilly)

## Basal Ganglia: Action Selection



| | Motor | Oculomotor | Prefrontal | Orbitofrontal | Cingulate |
|---|---|---|---|---|---|
| Cortex | SMA — PM, M1, S1 | FEF — DLPFC, PPC | DLPFC — PPC, PM | OFC — IT, ACC | ACC — HIP, EC, IT |
| Striatum | Putamen | Caudate | Caudate | Caudate | V. Striat. |
| Thalamus | VL | VA/MD | VA/MD | VA/MD | MD |
| | motor actions | eye movement | strategies & plans | future rewards | costs |

- Parallel circuits select motor actions and "cognitive" actions across frontal areas

3/5/17    COSC 494/594 CCN    9
(slide based on O'Reilly)

## Release from Inhibition



Caudate Nucleus

(-)

Substantia Nigra Pars Reticulata

Tonic Inhibition    Disinhibition    Tonic Inhibition

(-)

Superior Colliculus

Saccadic Eye Movement

3/5/17    COSC 494/594 CCN    10
(slide < O'Reilly)

## Basal Ganglia System

- Striatum
  - matrix clusters (inhib.)
    - direct (Go) pathway → GPi
    - indirect (NoGo) path → GPe
  - patch clusters
    - to dopaminergic system
- Globus pallidus, int. segment (GPi)[*]
  - tonically active
  - inhibit thalamic cells
- Globus pallidus, ext. segment (GPe)
  - tonically active
  - inhibits corresponding GPi neurons

  *and substantia nigra pars reticulata (SNr)

- Thalamus[*]
  - cells fire when both:
    - excited (cortex)
    - disinhibited (GPi)
  - disinhibits FC deep layers
- Substantia nigra pars compacta (SNc)
  - releases dopamine (DA) into striatum
  - excites D1 receptors (Go)
  - inhibits D2 receptors (NoGo)
- Subthalamic nucleus (STN)
  - hyperdirect pathway
  - input from cortex
  - diffuse excitatory output to GPi
  - global NoGo delays decision

  *and superior colliculus (SC)

3/5/17    COSC 494/594 CCN    11

## What is Dopamine Doing?

Dopamine carries the brain's reward signal
reward prediction error



stimulus    reward

Wise & Romper, 89    Schultz et. al, 98

3/5/17    COSC 494/594 CCN    12

## Basal Ganglia Reward Learning
### (Frank, 2005…; O'Reilly & Frank 2006)



a) Dopamine Burst          b) Dopamine Dip

- Feedforward, modulatory (disinhibition) on cortex/motor (same as cerebellum)
- Co-opted for higher level cognitive control → PFC

3/5/17          COSC 494/594 CCN          13
(slide < O'Reilly)

## Basal Ganglia Architecture: Cortically-based Loops



3/5/17          COSC 494/594 CCN          14
(slide < Frank)

## Fronto-basal Ganglia Circuits in Motivation, Action, & Cognition



3/5/17          COSC 494/594 CCN          15
(slide < Frank)

## ChR2-mediated excitation of direct- and indirect-pathway MSNs *in vivo* drives activity in basal ganglia circuitry



3/5/17          COSC 494/594 CCN          16
AV Kravitz et al. Nature 466(7306):622-6 (2010) doi:10.1038/nature09159

## Human Probabilistic Reinforcement Learning

- Patients with Parkinson's disease (PD) are impaired in cognitive tasks that require learning from positive and negative feedback
- Likely due to depleted dopamine
- But dopamine medication actually worsens performance in some cognitive tasks, despite improving it in others

Frank, Seeberger & O'Reilly (2004)



Train

A (80/20)  B (20/80)

C (70/30)  D (30/70)

E (60/40)  F (40/60)

Test

A > CDEF    Choose A?

B < CDEF    Avoid B?

3/5/17          COSC 494/594 CCN          17
(slide based on Frank)

## Testing the Model: Parkinson's and Medication Effects



Frank, Seeberger & O'Reilly (2004)

3/5/17          COSC 494/594 CCN          18
(slide < Frank)

## Slide: BG Model: DA Modulates Learning from Positive/Negative Reinforcement



**BG Model: DA Modulates Learning from Positive/Negative Reinforcement**

(A) The corticostriato-thalamo-cortical loops, including the direct (Go) and indirect (NoGo) pathways of the basal ganglia.
(B) The Frank (in press) neural network model of this circuit.
(C) Predictions from the model for the probabilistic selection task

Michael J. Frank et al. Science 2004;306:1940-1943

Published by AAAS

## Slide: emergent Demonstration: BG



**emergent Demonstration: BG**

A simplified model compared to Frank, Seeberger, & O'Reilly (2004)

3/5/17                COSC 494/594 CCN                20

## Slide: Anatomy of BG Gating Including Subthalamic Nucleus (STN)



**Anatomy of BG Gating Including Subthalamic Nucleus (STN)**

PFC-STN provides an override mechanism

3/5/17                COSC 494/594 CCN                21
(slide < Frank)

## Slide: Subthalamic Nucleus: Dynamic Modulation of Decision Threshold



**Subthalamic Nucleus: Dynamic Modulation of Decision Threshold**

Conflict (entropy) in choice prob ⇒ delay decision!

3/5/17                COSC 494/594 CCN                22
(slide < Frank)

## Slide: B. Temporal Difference Reinforcement Learning

**B. Temporal Difference Reinforcement Learning**

3/5/17                COSC 494/594 CCN                23

## Slide: Reinforcement Learning: Dopamine



**Reinforcement Learning: Dopamine**

No prediction
Reward occurs

Reward predicted
Reward occurs

Reward predicted
No reward occurs

Rescorla-Wagner / Delta Rule:

- $\delta = r - \hat{r}$
- $\delta = r - \sum xw$

But no CS-onset firing – need to anticipate the future!

- $\delta = (r + f) - \hat{r}$

CS-onset = future reward = $f$

3/5/17                COSC 494/594 CCN                24
(slide < O'Reilly)

## Temporal Differences Learning

- $V(t) = r(t) + \gamma^1 r(t+1) + \gamma^2 r(t+2)...$

- $\hat{V}(t) = r(t) + \gamma \hat{V}(t+1)$

- $0 = \left(r(t) + \hat{V}(t+1)\right) - \hat{V}(t)$

- $\delta = \left(r(t) + \hat{V}(t+1)\right) - \hat{V}(t)$

- $f = \gamma \hat{V}(t+1)$  ← this is the future!

3/5/17      COSC 494/594 CCN      25
(slide < O'Reilly)

## Network Implementation



3/5/17      COSC 494/594 CCN      26
(slide < O'Reilly)

## The RL-cond Model

- ExtRew: external reward $r(t)$ (based on input)
- TDRewPred: learns to predict reward value
  - minus phase = prediction $V(t)$ from previous trial
  - plus phase = predicted $V(t+1)$ based on Input
- TDRewInteg: Integrates ExtRew and TDRewPred
  - minus phase = $V(t)$ from previous trial
  - plus phase = $V(t+1) + r(t)$
- TD: computes temporal dif. delta value $\approx$ dopamine signal
  - compute plus – minus from TDRewInteg

3/5/17      COSC 494/594 CCN      27

## Classical Conditioning

- Forward conditioning
  - unconditioned stimulus (US): doesn't depend on experience
  - leads to unconditioned response (UR)
  - preceding conditioned stimulus (CS) becomes associated with US
  - leads to conditioned response (CR)
- Extinction
  - after CS established, CS is presented repeatedly without US
  - CR frequency falls to pre-conditioning levels
- Second-order conditioning
  - CS1 associated with US through conditioning
  - CS2 associated with CS1 through conditioning, leads to CR

3/5/17      COSC 494/594 CCN      28

## CSC Experiment

- A <u>serial-compound stimulus</u> has a series of distinguishable components
- A <u>complete serial-compound (CSC) stimulus</u> has a component for every small segment of time before, during, and after the US
  - Richard S. Sutton & Andrew G. Barto, "Time-Derivative Models of Pavlovian Reinforcement," *Learning and Computational Neuroscience: Foundations of Adaptive Networks*, M. Gabriel and J. Moore, Eds., pp. 497–537. MIT Press, 1990
- RL-cond.proj implements this form of conditioning
  - somewhat unrealistic, since the stimulus or some trace of it must persist until the US
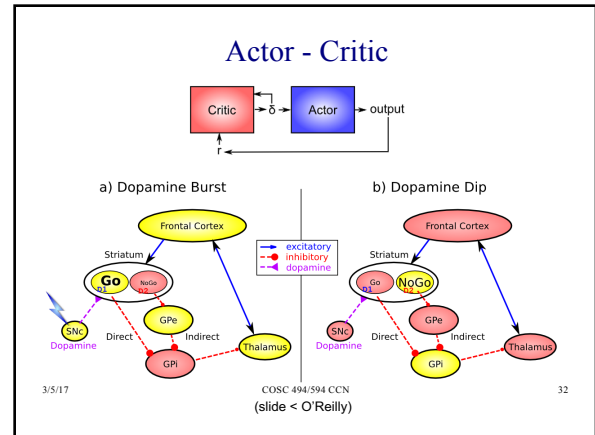
3/5/17      COSC 494/594 CCN      29

## RL-cond.proj



3/5/17      COSC 494/594 CCN      30

# emergent Demonstration: RL

A simplified model of temporal difference reinforcement learning

3/5/17                     COSC 494/594 CCN                     31

## Actor - Critic



a) Dopamine Burst                    b) Dopamine Dip

3/5/17                     COSC 494/594 CCN                     32
(slide < O'Reilly)

## Opponent-Actor Learning (OpAL)

- Actor has independent G and N weights
- Scaled by dopamine (DA) levels during choice
- Choice based on relative activation levels
- Low DA: costs amplified, benefits diminished $\Rightarrow$ choice 1
- High DA: benefits amplified, costs diminished $\Rightarrow$ choice 3
- Moderate DA $\Rightarrow$ choice 2
- Accounts for differing costs & benefits



3/5/17                     COSC 494/594 CCN                     33

# C. PVLV Model of DA Biology

A model of dopamine firing in the brain

3/5/17                     COSC 494/594 CCN                     34

## Brain Areas Involved in Reward Prediction

- Lateral hypothalamus (LHA): provides a primary reward signal for basic rewards like food, water etc.
- Patch-like neurons in ventral striatum (VS-patch)
  — have direct inhibitory connections onto dopamine neurons in VTA and SNc
  — likely role in canceling influence of primary reward signals when they're successfully predicted
- Central nucleus of amygdala (CNA)
  — important for driving dopamine firing at the onset of conditioned stimuli
  — receives input broadly from cortex
  — projects directly and indirectly to the VTA and SNc (DA neurons)
  — neurons in the CNA exhibit CS-related firing

3/5/17                     COSC 494/594 CCN                     35

## PVLV Model of Dopamine Firing

- Two distinct systems: Primary Value (PV) and Learned Value (LV)
- DA signal at time of external reward (US):
$$\delta_{pv} = PV_e - PV_i = r - \hat{r}$$
- DA signal for LV when PV not present/expected:
$$\delta_{lv} = LV_e - LV_i$$
- $LV_e$ is excitatory drive from CNA responding to CS (eventually canceled by $LV_i$)
- $LV_e$ and $LV_i$ values learned from $PV_e$ when rewards present/expected
- Hence, CS (or some trace) must still be present when US occurs
- CNA supports 1st order conditioning, but not 2nd order (that's in BLA)

3/5/17                     COSC 494/594 CCN                     36

### Slide 37

## Biology of Dopamine Firing



3/5/17 COSC 494/594 CCN 37

(slide < O'Reilly)

### Slide 38

## More Detailed Description of PVLV

- Major issue: Which of PV/LV systems should be in charge of overall dopamine system?
- PV and LV learning occur when PV present or expected (indicated by $PV_r > \Theta_{pv}$)
- PVr system learns: $\delta w_{pvr} = r_{present} - PV_r$ (improves prediction)
- Recall alternative DA signals:
$$\delta_{pv} = PV_e - PV_i, \qquad \delta_{lv} = LV_e - LV_i$$
- Novelty Value (NV) signal reflects stimulus novelty
- Overall dopamine signal:
$$\delta = \begin{cases} \delta_{pv}(t) - \delta_{pv}(t-1) & \text{if } PV_r > \Theta_{pv} \\ [\delta_{lv}(t) - \delta_{lv}(t-1)] + [NV(t) - NV(t-1)] & \text{otherwise} \end{cases}$$
- Note DA burst is phasic (ceases after CS onset)

3/5/17 COSC 494/594 CCN 38

### Slide 39

## More Detailed Description (ctu'd)

- Learning $PV_i$ weights:
$$\delta w_{pv} = \varepsilon(PV_e - PV_i)x$$

- Learning LV weights is conditional on PV filter:
$$\delta w_{lv} = \begin{cases} \varepsilon(PV_e - LV_e)x & \text{if } PV_r > \Theta_{pv} \\ 0 & \text{otherwise} \end{cases}$$

3/5/17 COSC 494/594 CCN 39

### Slide 40

## PVLV.proj Model

- PV in Ventral Striatum system
- LV in Amygdala system
- $VTA_l$ and VS adapt to US+
- Eventually $VTA_l$ bursts for CS onset
- LHB+RMTg and VS adapt to US–
- $VTA_m$ and VS adapt to US–
- Eventually DA dip for CS



simplified!

3/5/17 COSC 494/594 CCN 40

### Slide 41

## emergent Demonstration: PVLV

3/5/17 COSC 494/594 CCN 41

### Slide 42

## D. Cerebellum and Error-driven Learning

"The blessing of dimensionality"

3/5/17 COSC 494/594 CCN 42

## Functions of Cerebellum

- Maintenance of equilibrium and posture
- Timing of learned, skilled motor movement
  — any motor movement that improves with practice
  — timing, fluency, rhythm, coordination
  — involved in cognitive processes too
- Correction of errors during the execution of movements
  — error-driven learning
- Many inputs from cortical motor and sensory areas
- Influences cortical motor outputs to spinal chord

3/5/17 COSC 494/594 CCN 43

## Lookup Table & Pattern Separation



f(x)

Lookup Table -- store learned input/output pairs

x

Cerebellum
Cortex

3/5/17 COSC 494/594 CCN 44
(slide < O'Reilly)

## Cerebellum

- Inputs from parietal cortex and motor areas of frontal cortex
- Three layers, very many cortical maps
- Single basic circuit replicated throughout
- 200 million mossy fiber inputs (each to 500 granule cells)
  — projection of input into hyperdimensional space
  — separator learning and dynamics
- 40 billion granule cells (input from 4–5 mossy fibers)
- 15 million Purkinje cells (input from 200,000 granule cells)
  — matrix organization
  — enormous integration and cross connection
- Climbing fibers (one per Purkinje, from inferior olive)

3/5/17 COSC 494/594 CCN 45

## Cerebellar Error-driven Learning



Separation may be easier in higher dimensions

complex in low dimensions    simple in higher dimensions

Cerebellum = Support Vector Machine

- Granule cells = high-dimensional encoding (separation)
- Purkinje/Olive = delta-rule error-driven learning
- Classic ideas from Marr (1969) & Albus (1971)

3/5/17 COSC 494/594 CCN 46
(slide < O'Reilly)

## Cerebellum is Feed Forward

Feedforward circuit:

Input (PN) $\rightarrow$ granules $\rightarrow$ Purkinje $\rightarrow$ Output (DCN)

Inhibitory interactions – no attractor dynamics

Key idea: does delta-rule learning bridging small temporal gap:

$S(t-100) \rightarrow R(t)$

$\uparrow$ Error$(t+100)$



3/5/17 COSC 494/594 CCN 47
(slide < O'Reilly)

## Properties of Hyperdimensional Spaces

- Hyperdimensional spaces = spaces of very high dimension
- Consider vectors of 10,000 bits
  — measure distance by Hamming distance (HD)
  — or normalized Hamming distance (NHD)
- Mean HD = 5000, SD = 50 (binomial distribution)
- $< 10^{-9}$ of space closer than NHD = 0.47 or farther than 0.53 ($\pm 300 = \pm 6$ SD)
- Therefore random vectors almost surely have NHD = 0.5$\pm$0.03
- Vectors with < 3000 changed bits still accurately recognized
- Ref: Pentti Kanerva (2009), Hyperdimensional Computing: An Introduction to Computing in Distributed Representation with High-Dimensional Random Vectors, *Cognitive Computation*, 1(2)

3/5/17 COSC 494/594 CCN 48

## Orthogonality of Random Hyperdimensional Bipolar Vectors

- 99.99% probability of being within $4\sigma$ of mean
- It is 99.99% probable that random $n$-dimensional vectors will be within $\varepsilon = 4/\sqrt{n}$ orthogonal
- $\varepsilon = 4\%$ for $n = 10,000$
- Probability of being less orthogonal than $\varepsilon$ decreases exponentially with $n$
- The brain gets approximate orthogonality by assigning random high-dimensional vectors

$$|\mathbf{u}\cdot\mathbf{v}| < 4\sigma$$
$$\text{iff } \|\mathbf{u}\|\,\|\mathbf{v}\|\,|\cos\theta| < 4\sqrt{n}$$
$$\text{iff } n|\cos\theta| < 4\sqrt{n}$$
$$\text{iff } |\cos\theta| < 4/\sqrt{n} = \varepsilon$$

$$\Pr\left\{|\cos\theta| > \varepsilon\right\} = \operatorname{erfc}\left(\frac{\varepsilon\sqrt{n}}{\sqrt{2}}\right)$$
$$\approx \frac{1}{6}\exp\left(-\varepsilon^2 n/2\right) + \frac{1}{2}\exp\left(-2\varepsilon^2 n/3\right)$$

3/5/17                                                                                             49

## Hyperdimensional Pattern Associator

- Suppose $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_P$ are a set of random hyperdimensional bipolar vectors (inputs)
- Let $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_P$ be arbitrary bipolar vectors (outputs)
- Define Hebbian linear associator matrix

$$M = \sum_{k=1}^{P} \mathbf{q}_k \mathbf{p}_k^T$$

- Then $M\mathbf{p}_k \approx \mathbf{q}_k$ (table lookup)
- To encode a sequence of random vectors $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_P$:

$$M = \sum_{k=1}^{P-1} \mathbf{p}_{k+1} \mathbf{p}_k^T$$

- Then $M\mathbf{p}_k = \mathbf{p}_{k+1}$ (sequence readout)

3/5/17                              COSC 494/594 CCN                              50

## BG + Cerebellum Capacities

- Learn what satisfies basic needs, and what to avoid (BG reward learning)
  — And what information to maintain in working memory (PFC) to support successful behavior
- Learn basic Sensory → Motor mappings accurately (Cerebellum error-driven learning)
  — Sensory → Sensory mappings? (what is going to happen next)

3/5/17                              COSC 494/594 CCN                              51
(slide < O'Reilly)

## BG + Cerebellum Incapacities

- Generalize knowledge to novel situations
  — Lookup tables don't generalize well…
- Learn abstract semantics
  — Statistical regularities, higher-order categories, etc
- Encode episodic memories (specific events)
  — Useful for instance-based reasoning
- Plan, anticipate, simulate, etc…
  — Requires robust working memory

3/5/17                              COSC 494/594 CCN                              52
(slide < O'Reilly)

## emergent Demonstration: Cereb

3/5/17                              COSC 494/594 CCN                              53