# Integrated Cognitive Architectures (Stanford CS379C)
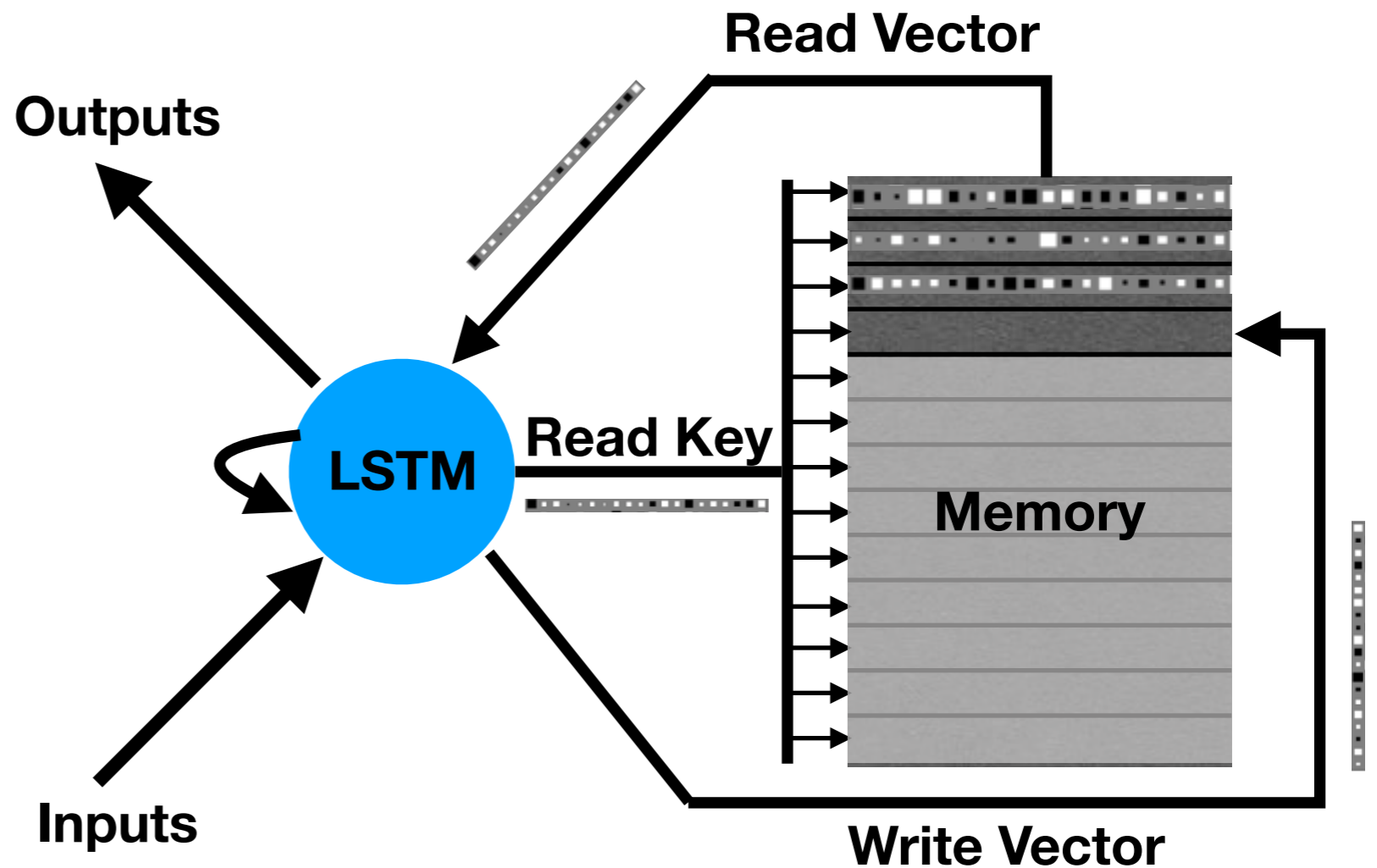
Greg Wayne
DeepMind Technologies Limited

# Talk Outline

- Neural Networks with External Memory

- Reinforcement Learning with Memory Systems

  - Some Limitations and Weaknesses Therein

- The MEmory, INference, and Reinforcement Learning Agent (MERLIN)

- Its Behavior on Interesting Tasks Characterized by Partial Observability

- Re The Programmer's Apprentice: Imitation-Learning for Complex Skills (Motor Skills)

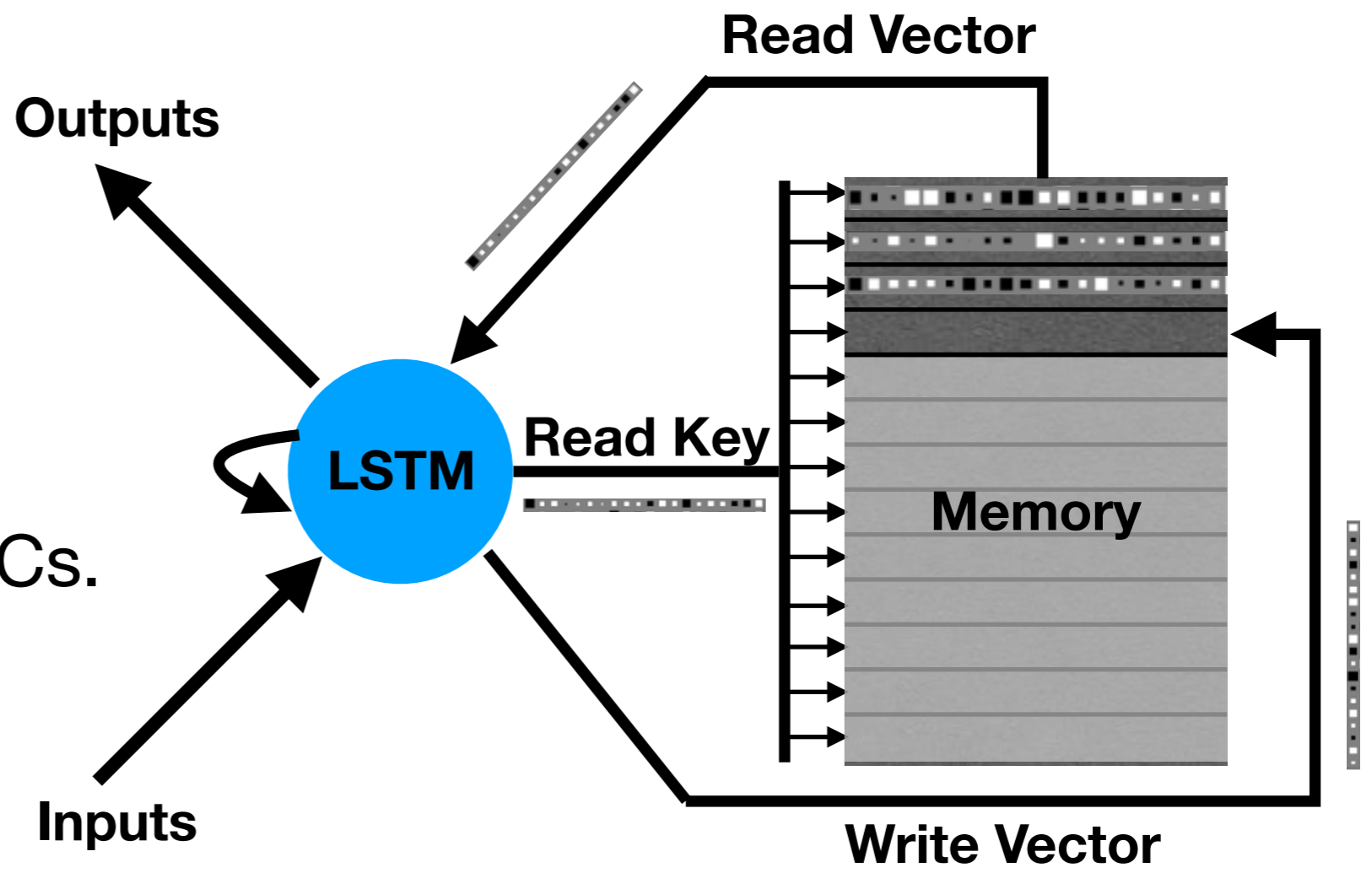# Neural Turing Machines (NTMs) and Differentiable Neural Computers (DNCs)

- Trainable neural networks that can read and write to external memory.

- Can instantiate simple algorithms operating over simple data structures like lists and graphs.

- Have higher capacity memory than LSTMs (Hochreiter and Schmidhuber, 1997) alone.

**Outputs**

**Read Vector**

**LSTM**

**Read Key**

**Memory**

**Inputs**

**Write Vector**

**(Graves, Wayne, Danihelka, arXiv 2014)**　　　**(Graves, Wayne et al., Nature 2016)**

# Neural Turing Machines (NTMs) and Differentiable Neural Computers (DNCs)



NTMs are a little more primitive / kludgy than DNCs.

**(Graves, Wayne, Danihelka, arXiv 2014)**     **(Graves, Wayne et al., Nature 2016)**
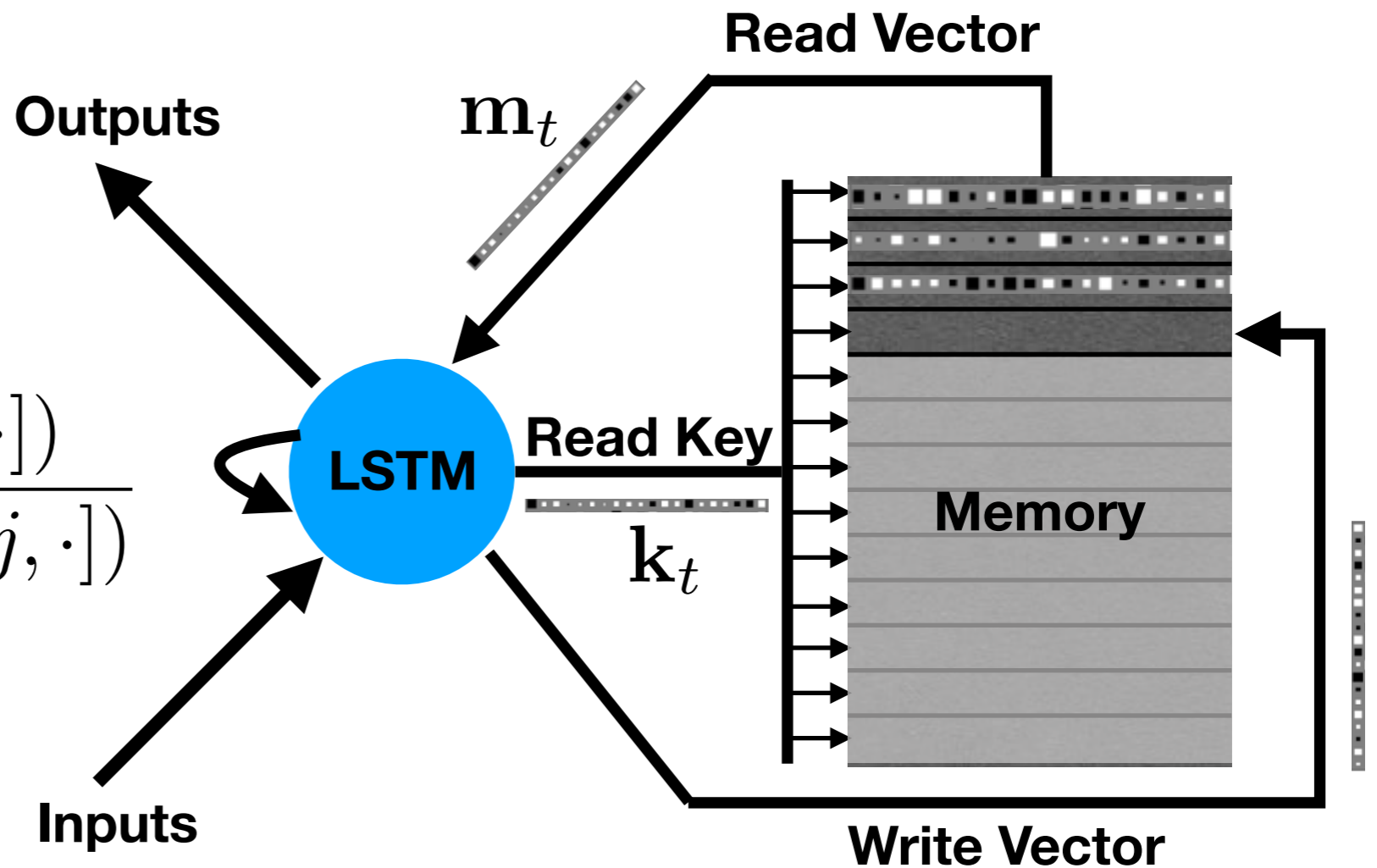
# A Simple DNC

**Reading from Memory**

**Read "Attention Weighting"**

$$w_t[i] = \frac{\exp(\mathbf{k}_t \cdot M_t[i, \cdot])}{\sum_j \exp(\mathbf{k}_t \cdot M_t[j, \cdot])}$$
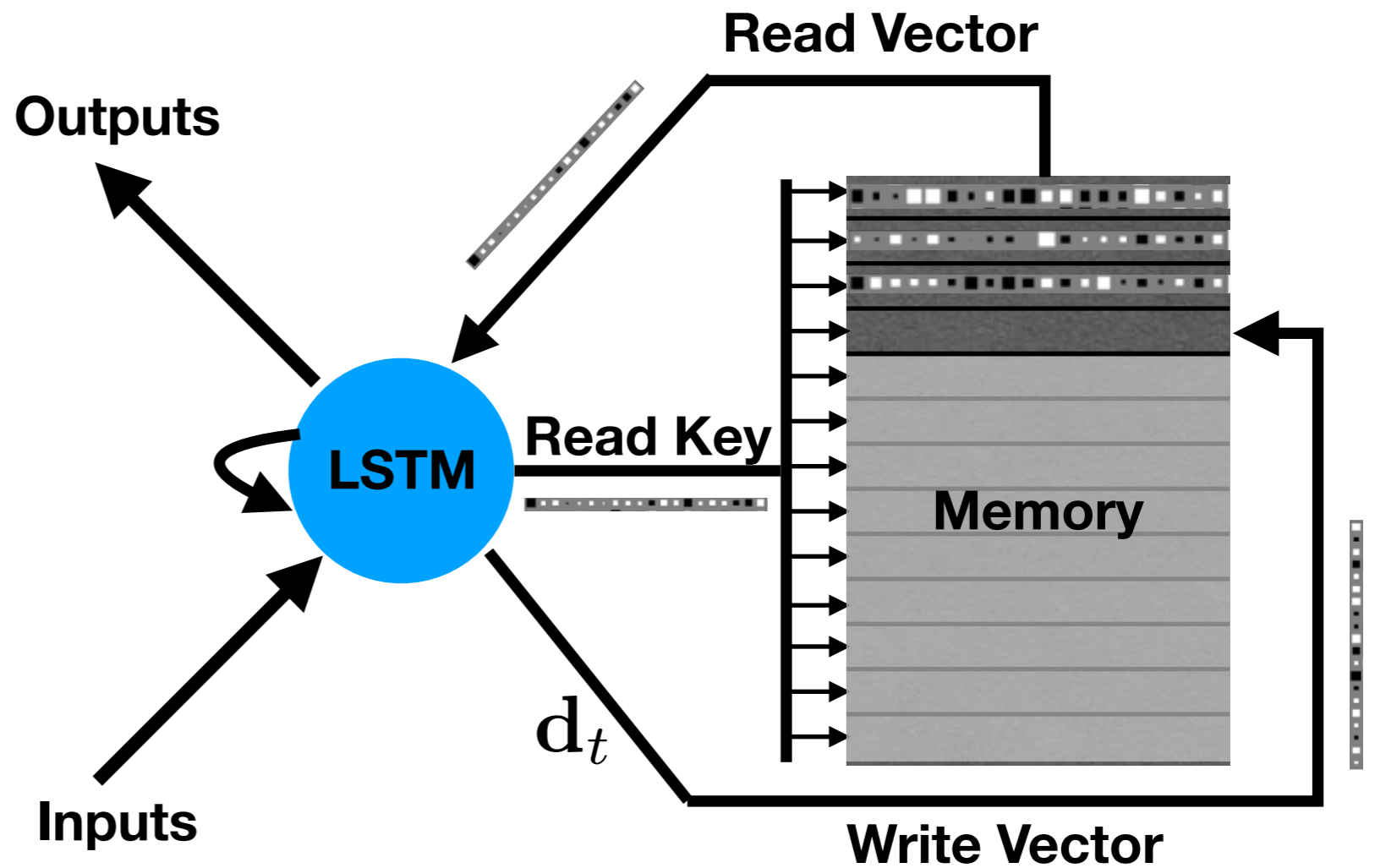
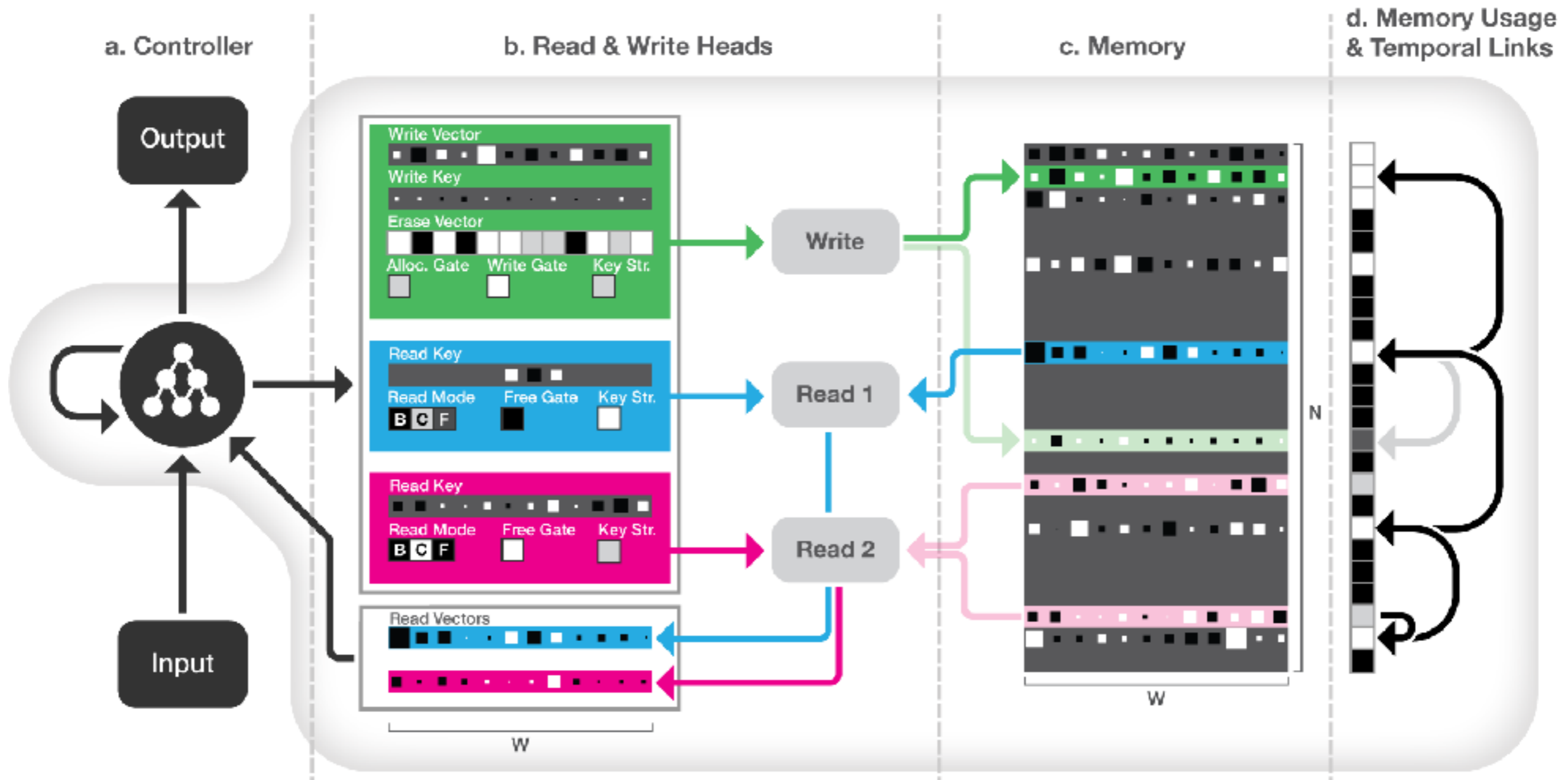**Read Vector**

$$\mathbf{m}_t = \sum_i w_t[i] M_t[i, \cdot]$$

**Outputs**

**Read Vector**

$\mathbf{m}_t$

**LSTM**

**Read Key**

$\mathbf{k}_t$

**Memory**

**Inputs**

**Write Vector**

# A Simple DNC

**Read Vector**

**Outputs**

**Read Key**

**Memory**

**Inputs**

**Write Vector**

*Writing to Memory*

$$M_{t+1}[t, \cdot] = \mathbf{d}_t$$

**LSTM**

$\mathbf{d}_t$

# A More Complicated DNC



**Supports reading items based on the order they were written as well.**
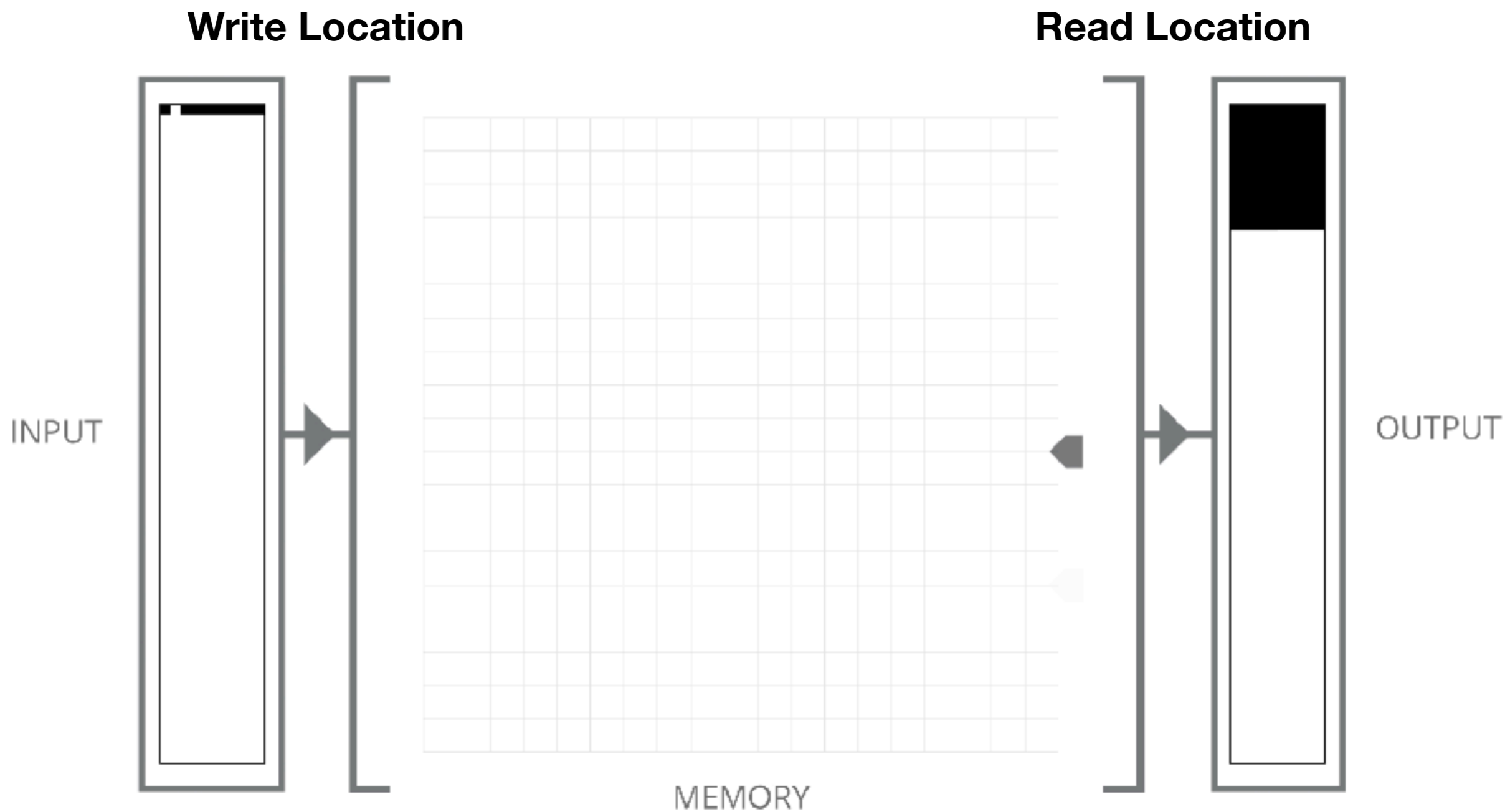
# Supervised Training of RNNs / DNCs
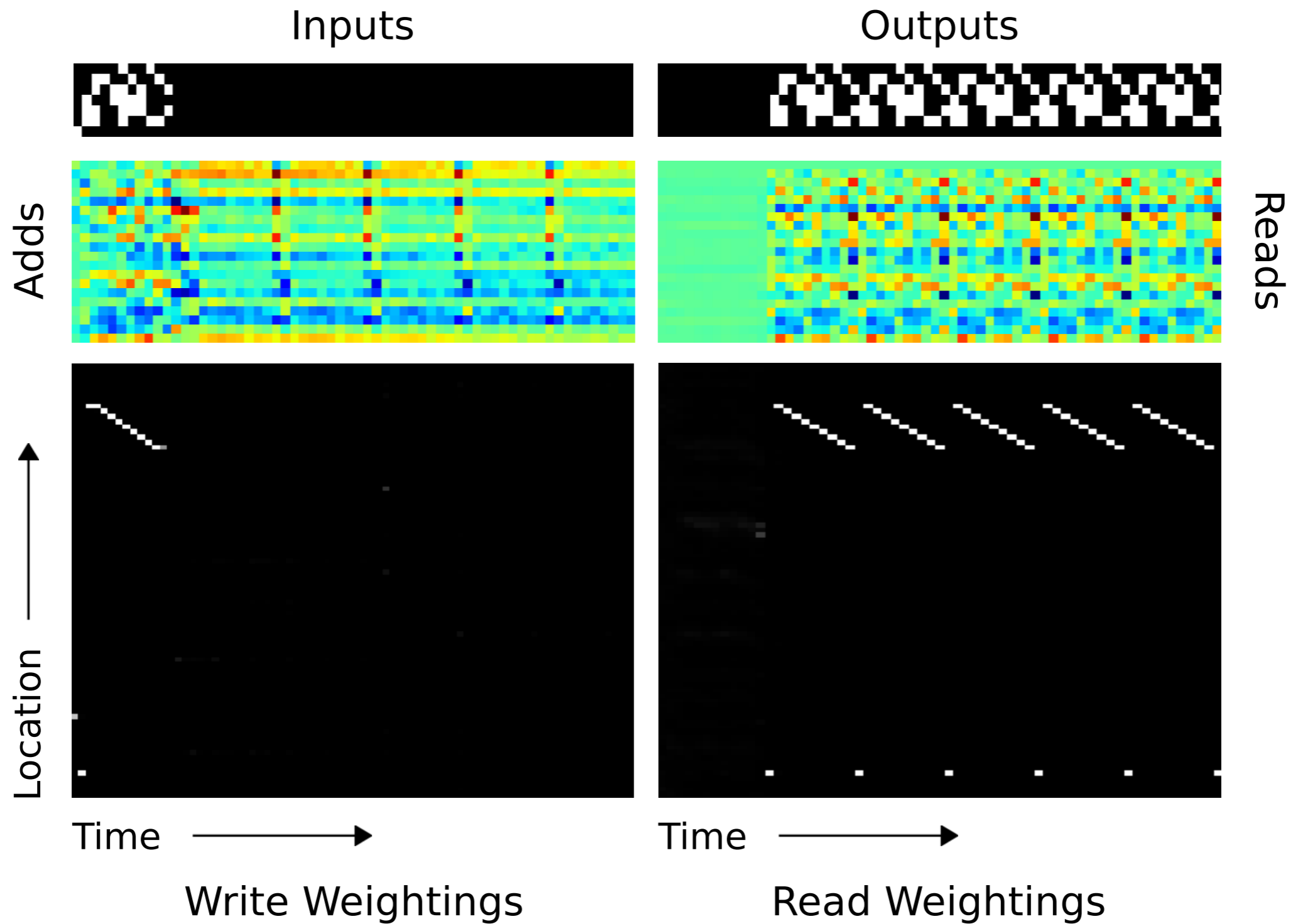
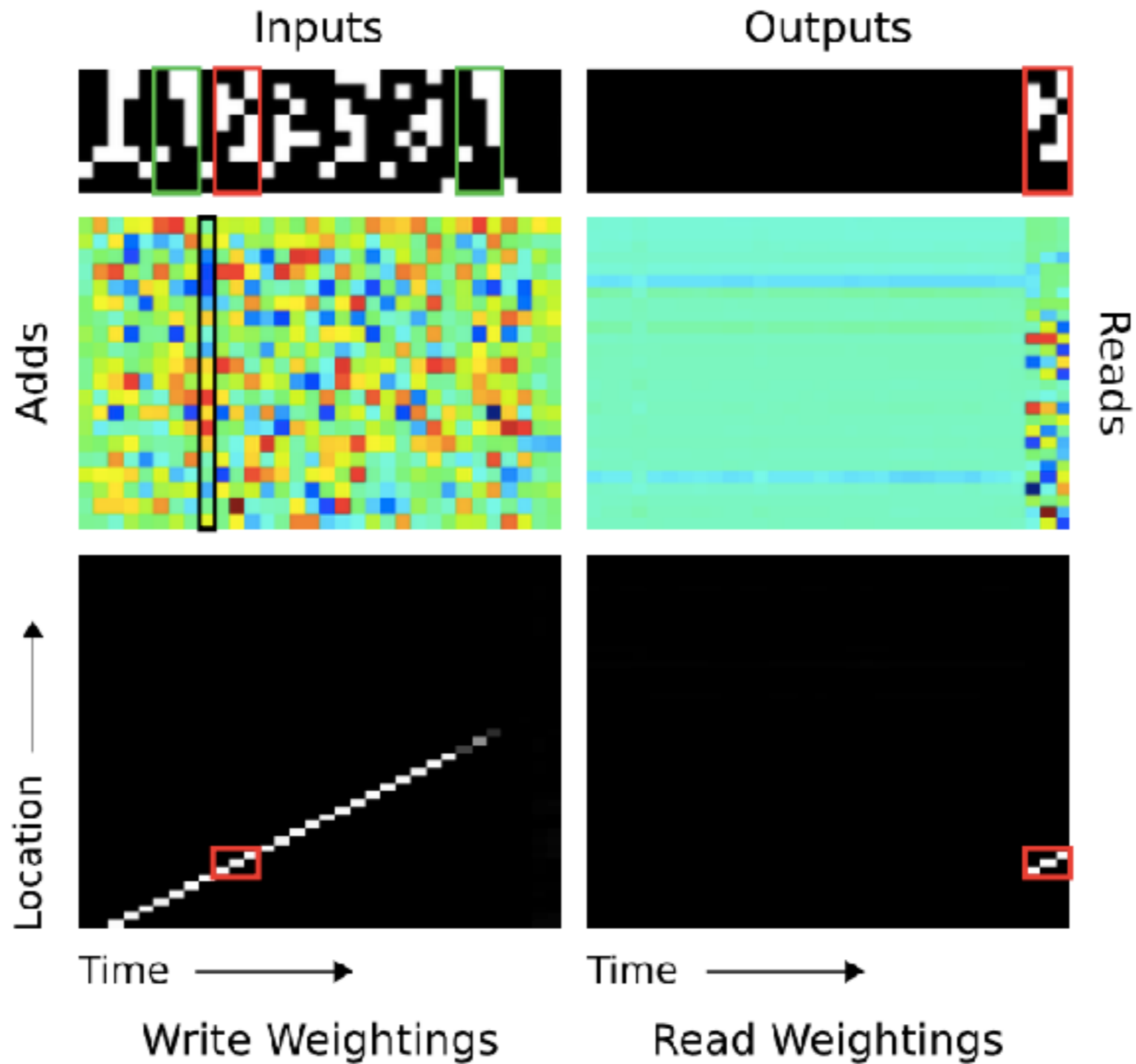# Supervised Training of RNNs / DNCs



**Backpropagation through Time**

**Write Location**　　　　　　　　　　　**Read Location**

INPUT　　　　　　　　　　　　　MEMORY　　　　　　　　　　　OUTPUT

**Example of an NTM (circa 2014) on a looped copy problem.**

**(LSTMs struggle on this.)**

# Looped Copy

# Associative List



Inputs      Outputs

Adds    Reads

Location

Time      Time

Write Weightings      Read Weightings

# Learning Complex Data Structures (Random Graphs)

Time

1. Source Node Label - Edge Label - Destination Node Label

2. Source Node Label - Edge Label - Destination Node Label

3. Source Node Label - Edge Label - Destination Node Label

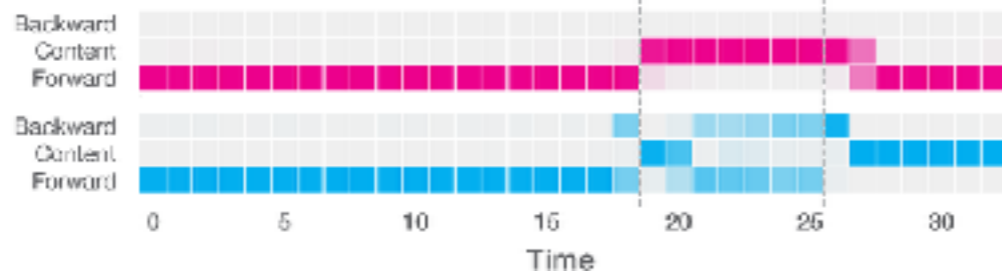Labels:
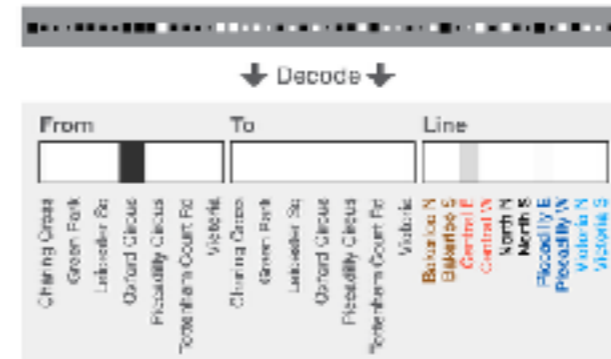Digital representations of numbers between 0 and 999.

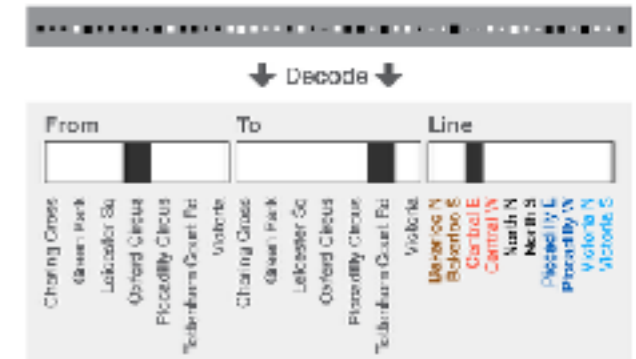# Tube Traversal Problem



a. Read and Write Weightings
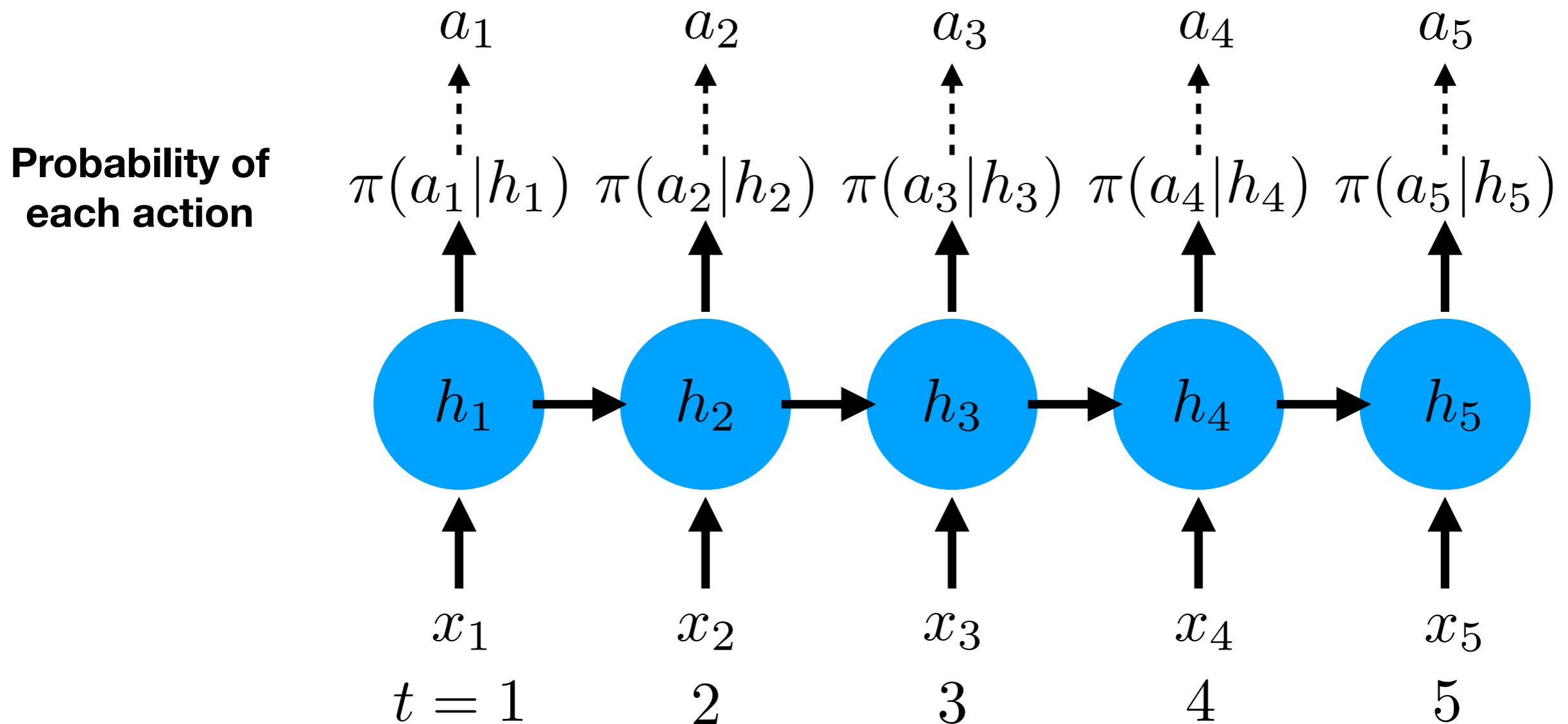
b. Read Mode

c. London Underground Map

d. Read Key

e. Location Content

Content-based lookup

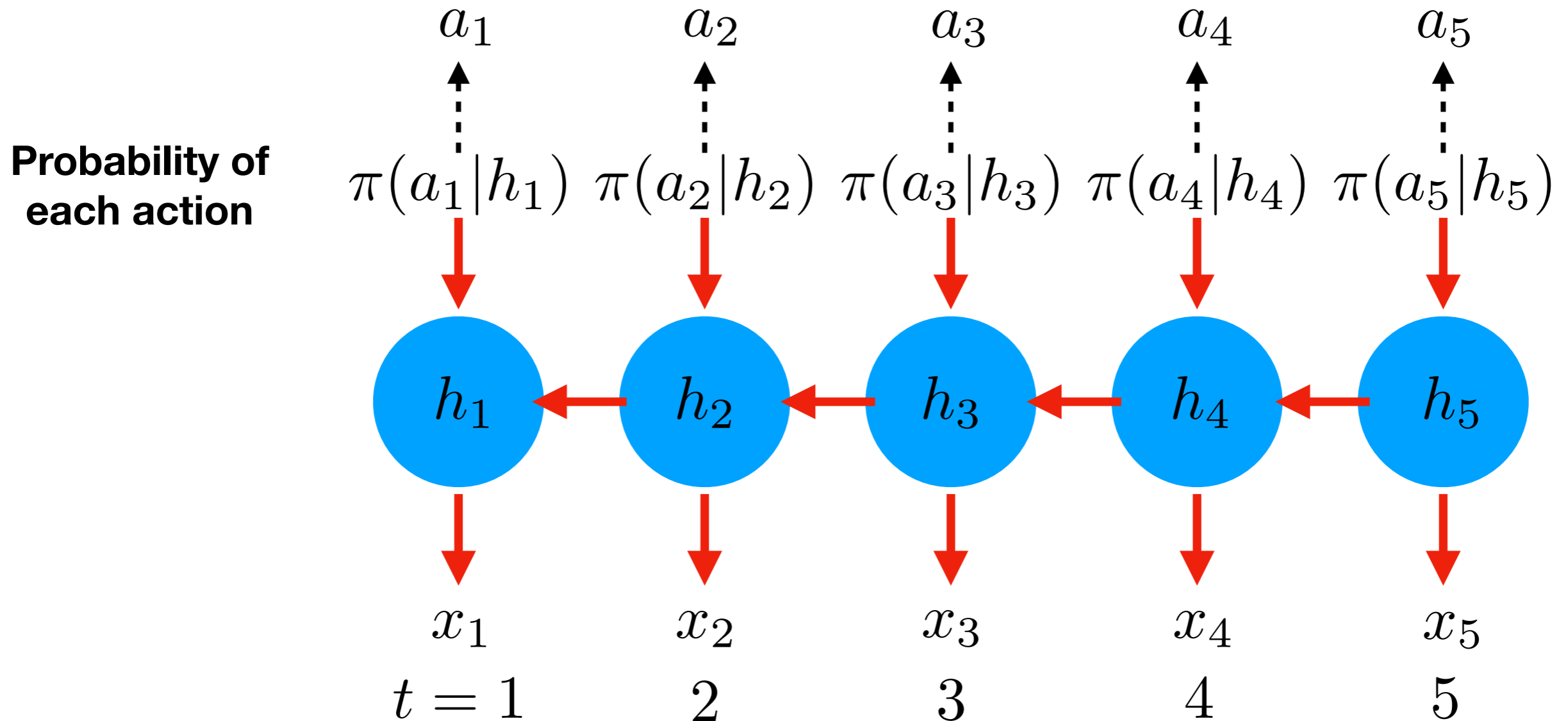**Queries of form: From Victoria Station go North on Victoria Line, North on Victoria Line, East on Central Line, ….**

# Reinforcement Learning

**Probability of each action**

$$a_1 \qquad a_2 \qquad a_3 \qquad a_4 \qquad a_5$$

$$\pi(a_1|h_1) \quad \pi(a_2|h_2) \quad \pi(a_3|h_3) \quad \pi(a_4|h_4) \quad \pi(a_5|h_5)$$

$$h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow h_4 \rightarrow h_5$$

$$x_1 \qquad x_2 \qquad x_3 \qquad x_4 \qquad x_5$$

$$t = 1 \qquad 2 \qquad 3 \qquad 4 \qquad 5$$

**Maximize expected "return"** $\quad R_t = r_t + r_{t+1} + r_{t+2} + \cdots + r_T$

# Reinforcement Learning

$a_1$ $a_2$ $a_3$ $a_4$ $a_5$

**Probability of each action**

$\pi(a_1|h_1)$ $\pi(a_2|h_2)$ $\pi(a_3|h_3)$ $\pi(a_4|h_4)$ $\pi(a_5|h_5)$

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

$t = 1$ $\quad$ 2 $\quad$ 3 $\quad$ 4 $\quad$ 5

**Policy Gradient**

$$\Delta\theta \propto \sum_{t=0}^{T} R_t \nabla_\theta \log \pi(a_t|h_t)$$

# A Simple Memory-Based RL Challenge (Memory Game)


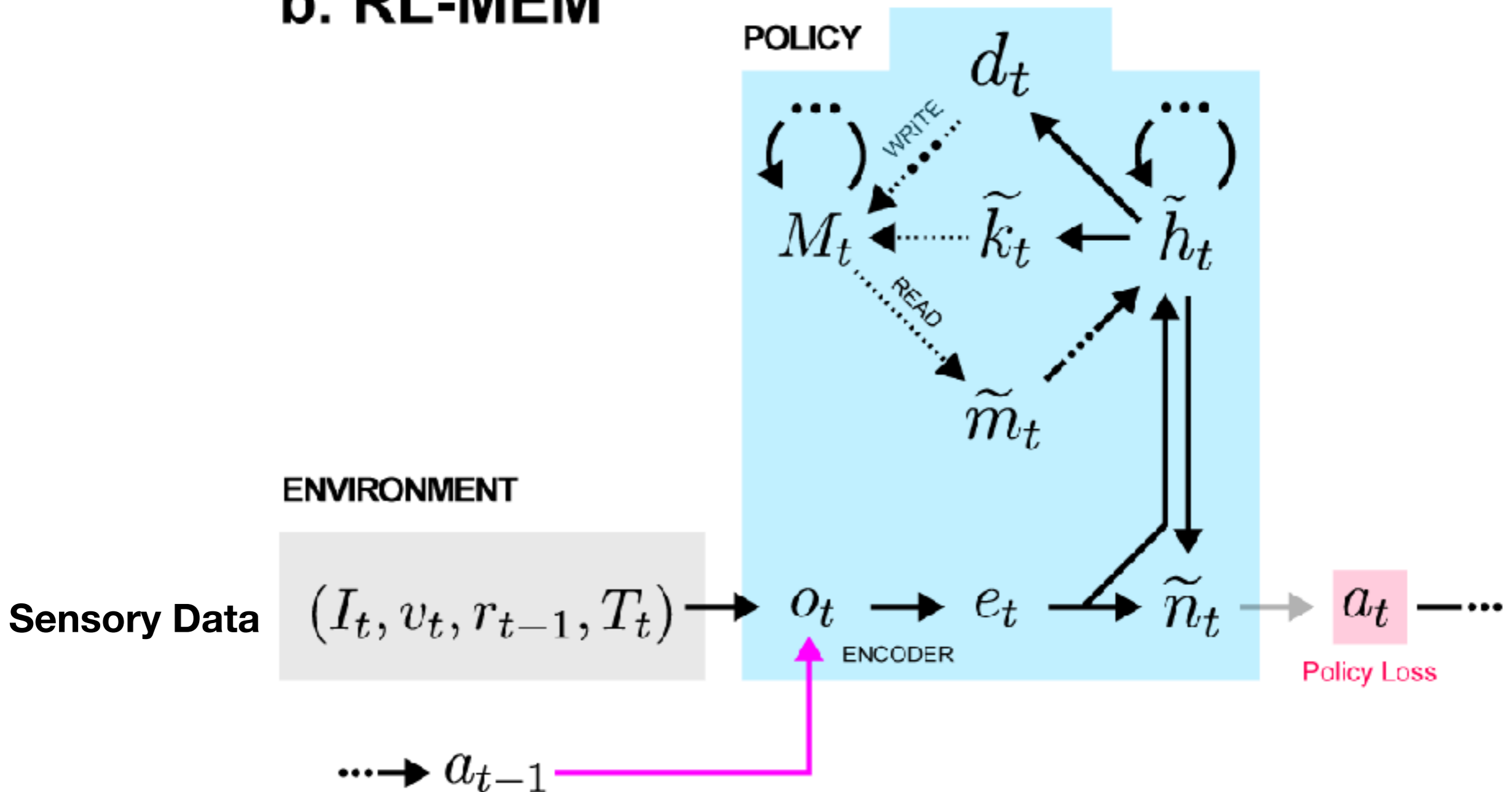
**Partially observed,
need to know "what is where"**

**(a kind of state estimation)**

# Control Architecture "a"

# Control Architecture "b"

# Remarkably

- Both these models fail to solve this very simple task. Why?

- RL-LSTM: lacks **capacity for storing** and **search mechanism for recalling** positions and images.

- RL-MEM: has capacity and search mechanism, but (we argue) policy gradient RL does not provide the right **objective function** to determine what is stored in memory.

# Policy Gradient

- Policy gradient: relatively high variance (noisy) estimator of performance improvement direction.

  - trajectory dependent

- Intuitively, feels wrong: memory for, say, the memory game should not only develop through gradients from **success or failure** due to **trial and error** action choices ("motor twitches").

# Prediction, Compression, and State Estimation

Classical Control: Separate State Estimation from Control Design

Our strategy: instead build perceptual representations and memories by unsupervised predictive modeling / state estimation

# Unsupervised Predictive Memory in a Goal-Directed Agent

# Unsupervised Predictive Memory in a Goal-Directed Agent



Alden Hung     Tim Lillicrap

# Unsupervised Learning

Modern framework for unsupervised learning:

Variational Autoencoders (Kingma and Welling, 2013)
(Rezende and Mohamed, 2014)



$$q(z|x) = \mathcal{N}(\mu(x), \Sigma(x))$$

$$\log p(x) \geq \mathbb{E}_{q(z|x)}\left[\log p(x|z) - \mathrm{KL}[q(z|x)\|p(z)]\right]$$

# Unsupervised Predictive Modeling

Sequential Variational Autoencoders or Variational RNNs (Gregor, 2015), (Chung, 2015).

**State variable (compressed representation)** $z_t$

**Prior** $p(z_t | z_1, a_1, z_2, a_2, \ldots, z_{t-1}, a_{t-1})$

**Posterior** $q(z_t | z_1, a_1, z_2, a_2, \ldots, z_{t-1}, a_{t-1}, \textcolor{red}{x_t})$

**Variational Lower Bound (VLB)** $\log p(x_t | z_t) - \mathrm{KL}[q || p]$

**Prior is trained to predict posterior.**

**Posterior is trained to make minimal deviation from prior whilst propagating information about observation.**

# What the Mechanism Looks Like

**Sensory Data**

**ENVIRONMENT**

$(I_t, v_t, r_{t-1}, T_t)$



$\cdots \rightarrow a_{t-1}$

PRIOR

$e_t \rightarrow n_t \leftarrow p \leftarrow \cdots - h_t$

$o_t$

ENCODER

KL Loss

$q \rightarrow z_t$

POSTERIOR

DECODER

$(\hat{I}_t, \hat{R}_t, \hat{v}_t, \hat{a}_{t-1}, \hat{r}_{t-1}, \hat{T}_t)$

Reconstruction Loss

**Two Modes:**
1)      prediction
2)   inference / state estimation

# Memory-Based Predictor



**MEMORY-BASED PREDICTOR**

**Sensory Data**

**ENVIRONMENT**

$(I_t, v_t, r_{t-1}, T_t)$

$o_t$

ENCODER

$e_t \rightarrow n_t \leftarrow p \leftarrow h_t \leftarrow m_t$

PRIOR

$k_t \cdots\rightarrow M_t$

READ

WRITE

KL Loss

$q \rightarrow z_t$

POSTERIOR

$a_{t-1}$

DECODER

$(\hat{I}_t, \hat{R}_t, \hat{v}_t, \hat{a}_{t-1}, \hat{r}_{t-1}, \hat{T}_t)$

Reconstruction Loss

**Two-fold advantage: state variables are good representations; also, memory helps prediction.**

# MERLIN

# MERLIN on the Memory Game

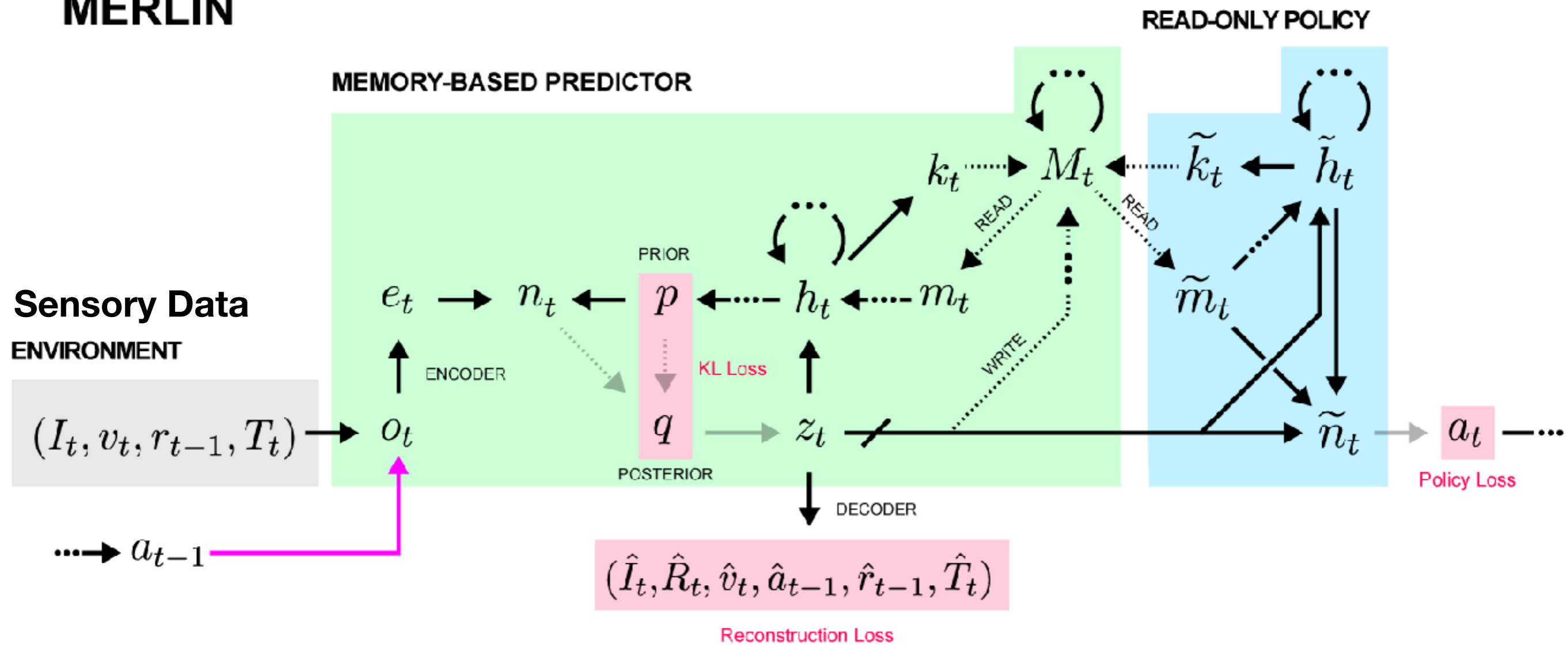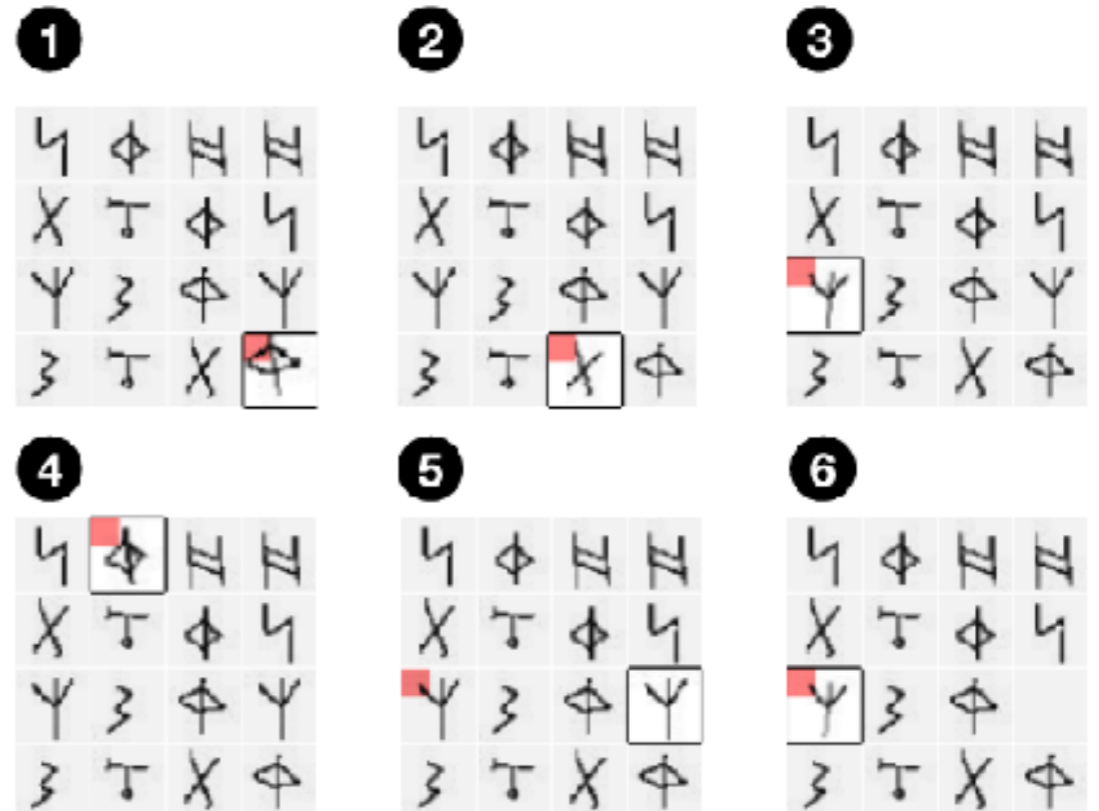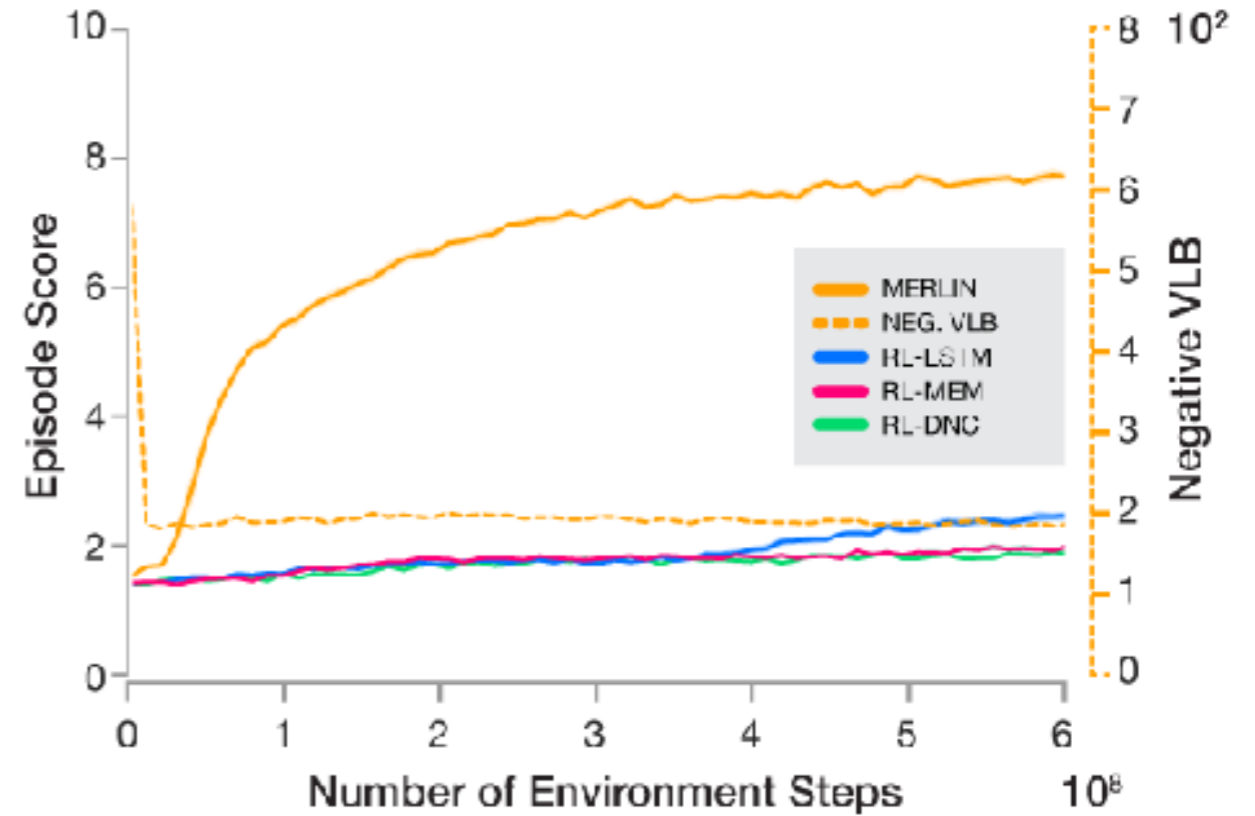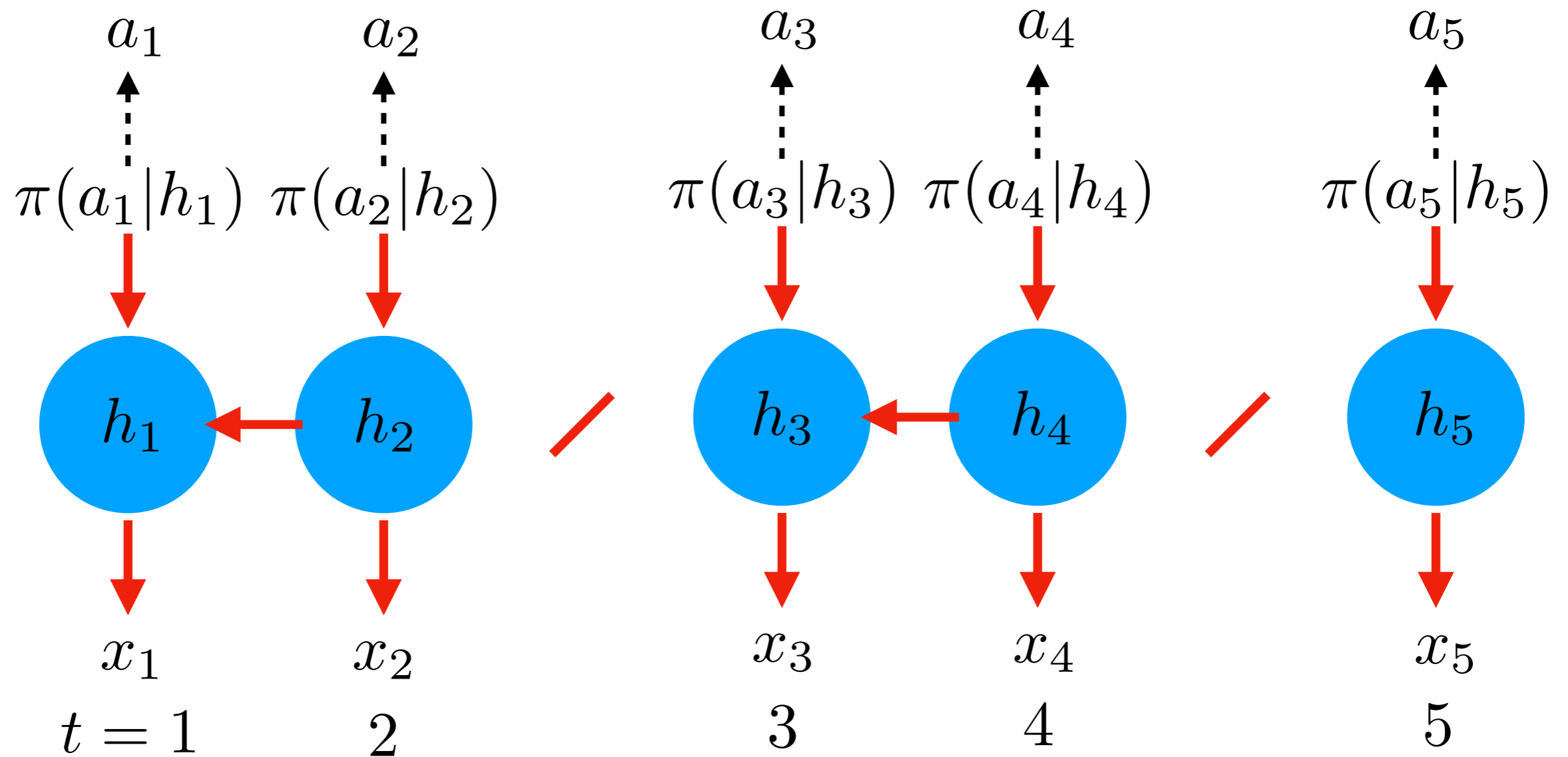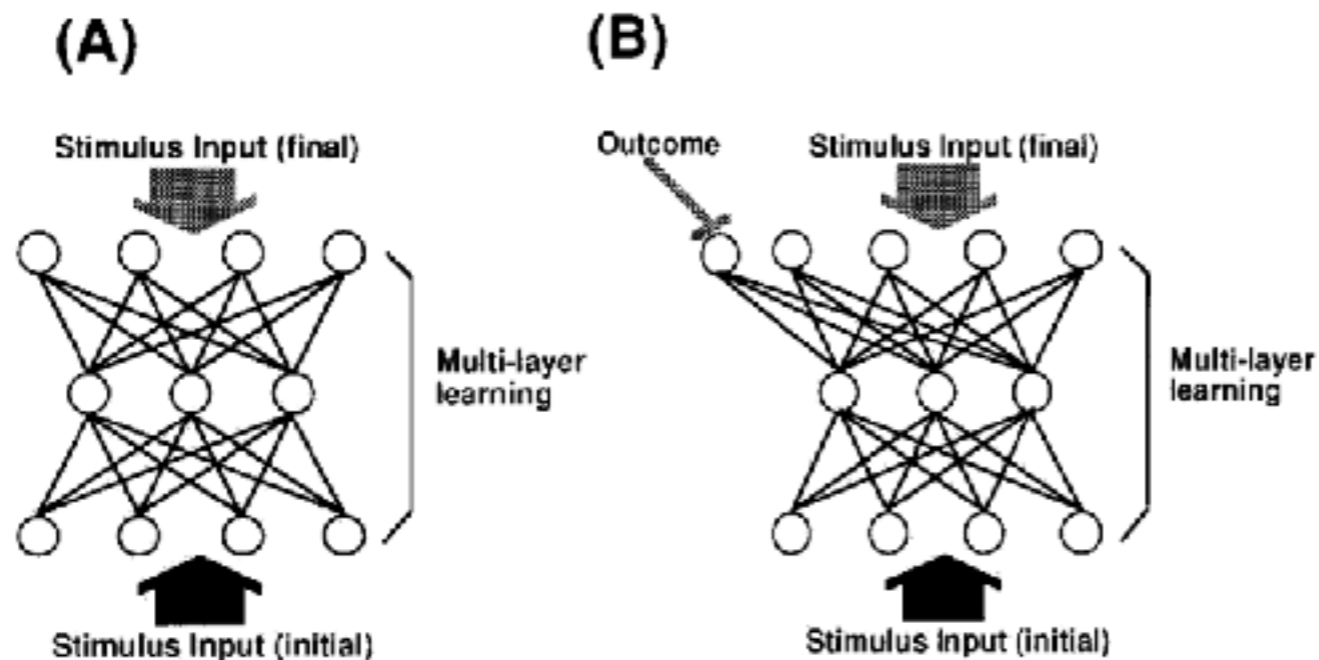# Truncation for Very Long Tasks



**Here, unless otherwise specified, we truncate every 20 steps (15 steps per second).**

# The Bullet Problem



**Saliency and "The Bullet Problem"**

**Gluck and Myers, 1993**

496    HIPPOCAMPUS  VOL. 3, NO. 4, OCTOBER 1993

(A)                                    (B)

Stimulus Input (final)        Outcome        Stimulus Input (final)

Multi-layer learning                           Multi-layer learning

Stimulus Input (initial)                       Stimulus Input (initial)
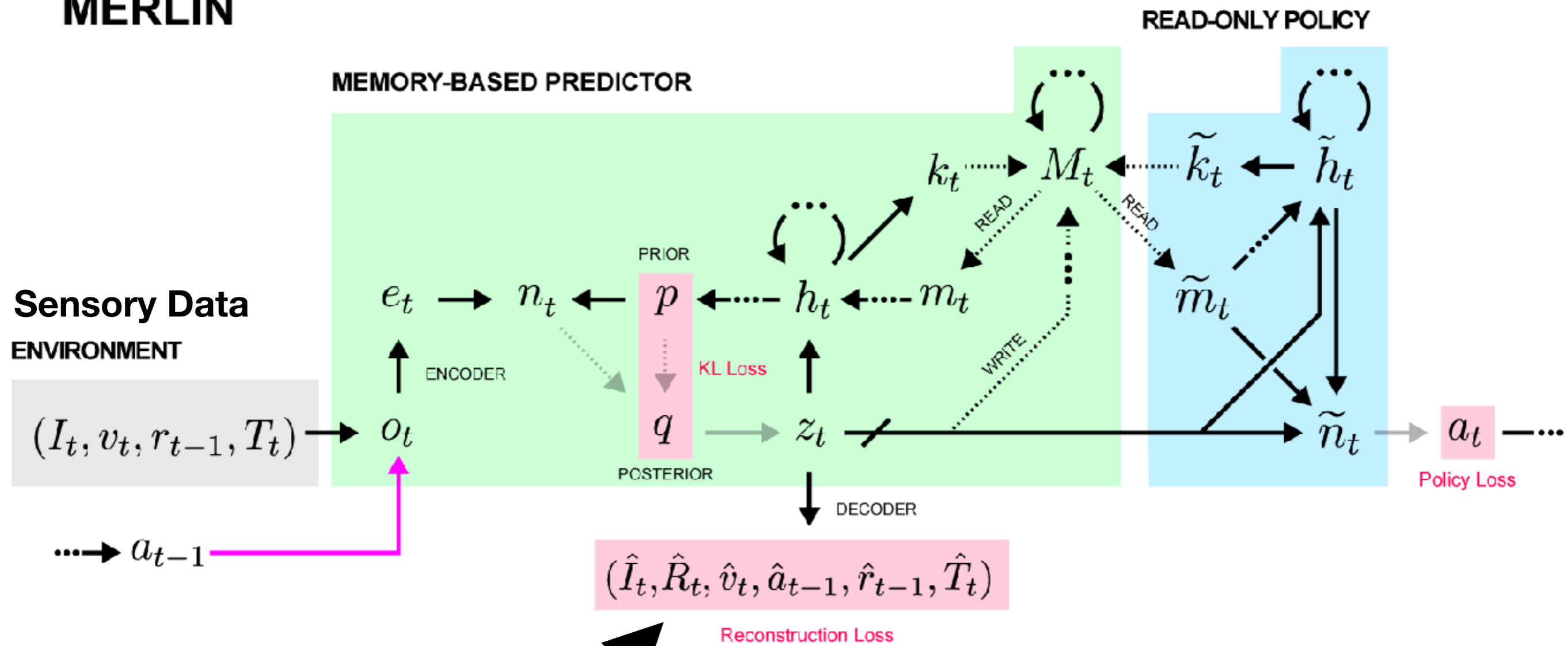
Fig. 4. (A) An autoencoder (Hinton, 1989). This network learns to reproduce its inputs at the output layer, using a multilayer learning algorithm. It contains a narrow internal layer, and so is forced to generate an internal representation that compresses redundancies in the input pattern. (B) A predictive autoencoder. This network is an autoencoder augmented with the constraint that it must also output a classification (or prediction of outcome) for the output pattern. The representation formed in the narrow internal layer must still compress redundancies; now, however, it must also differentiate representations of input features that are especially useful in predicting the outcome.
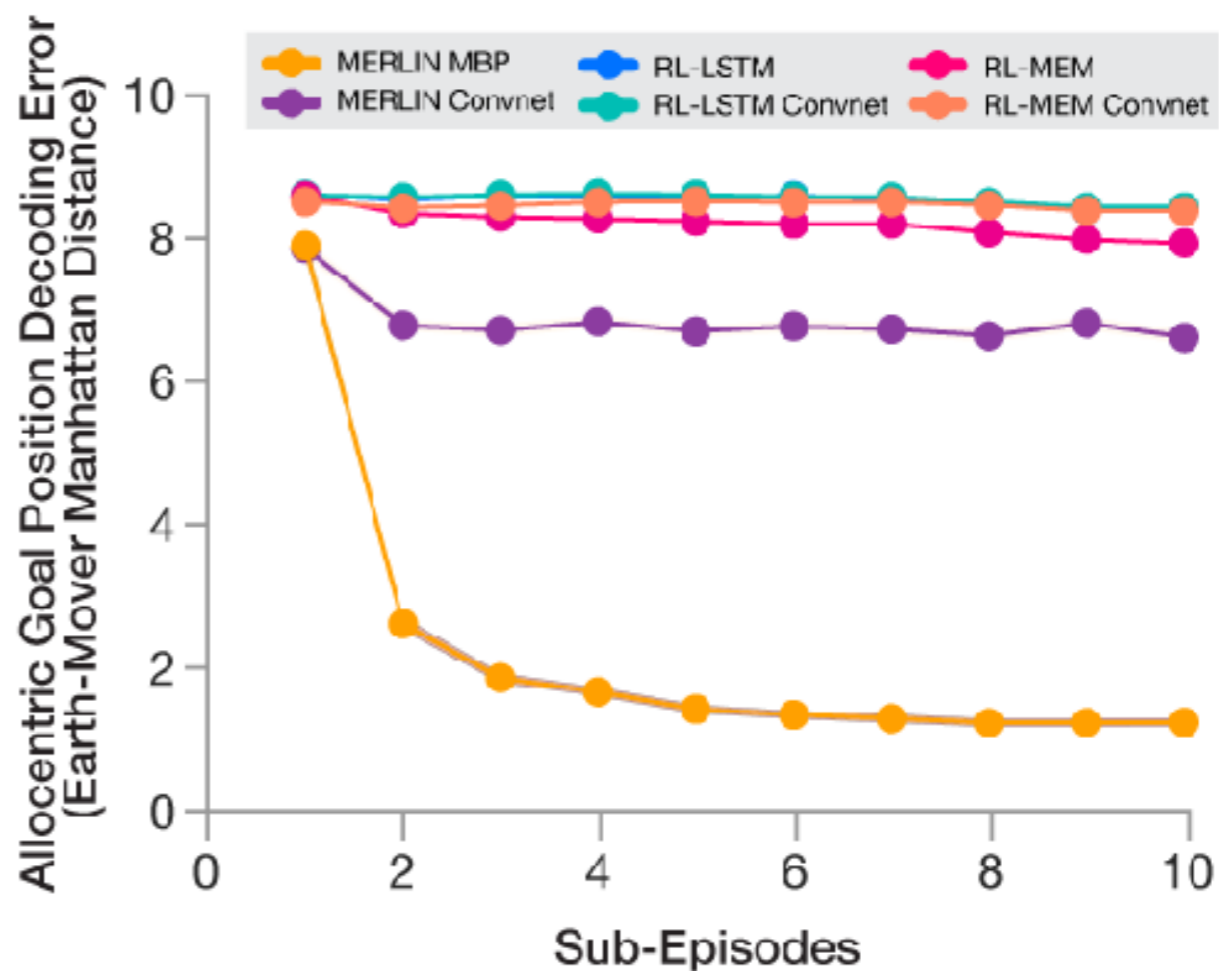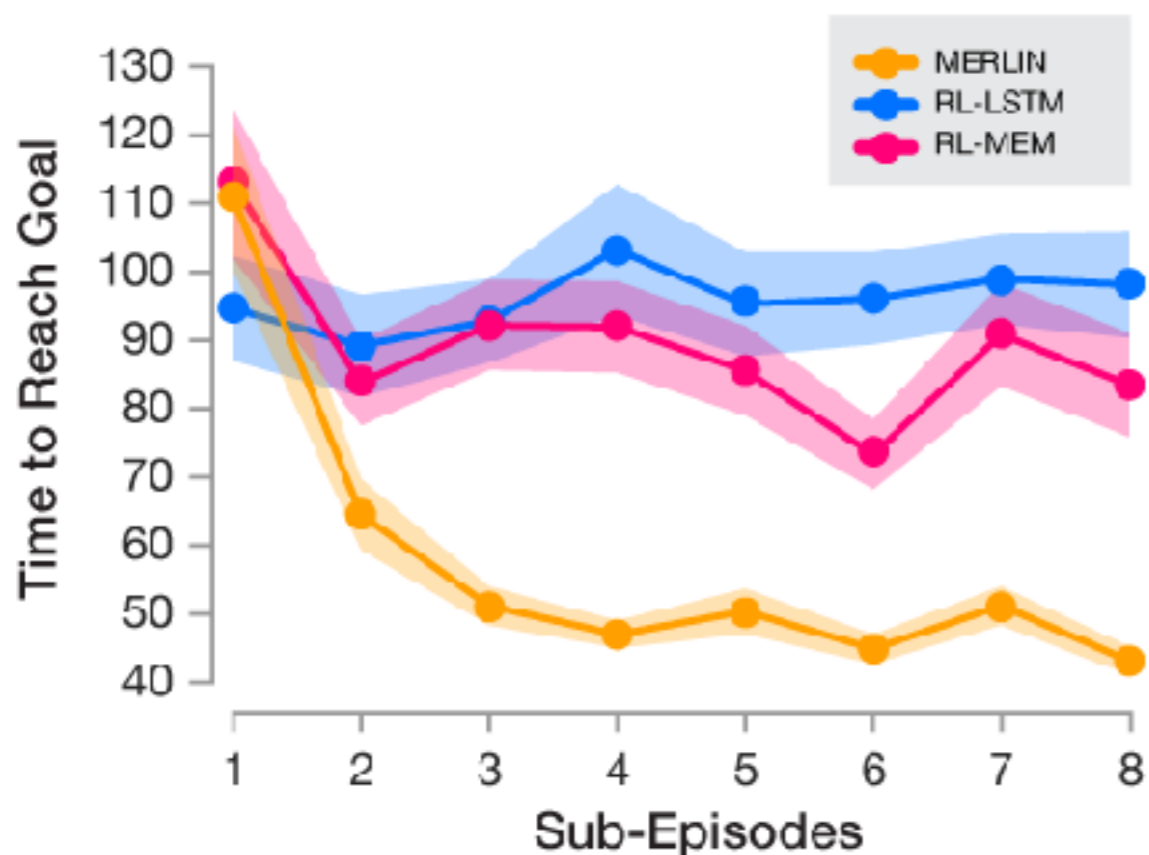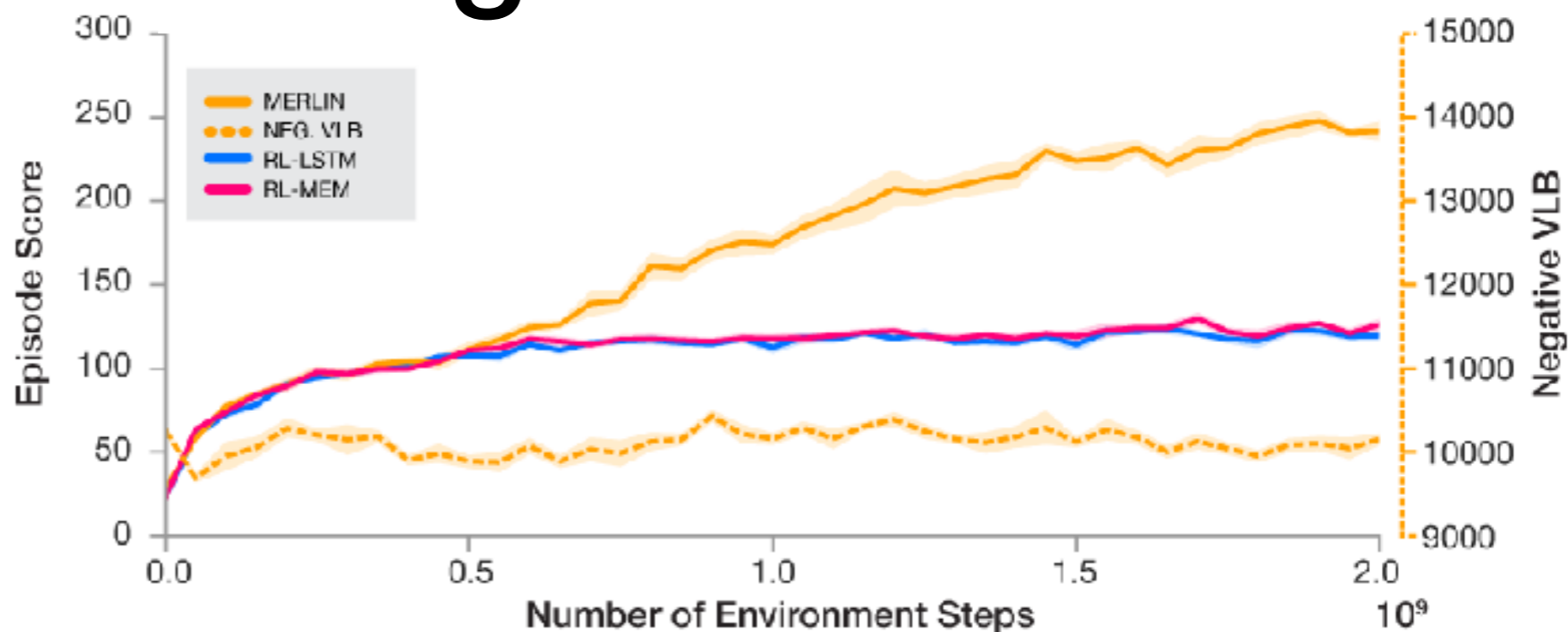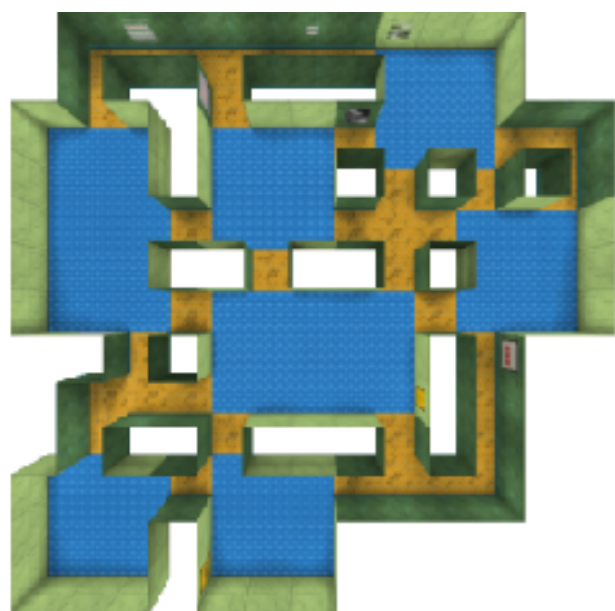
# MERLIN

# Navigation

# Navigation

# Return Prediction

# Gluck and Myers



Episode Score (Large Environment) vs. Return Cost Coefficient

% Variance Explained vs. Return Cost Coefficient

Observation / MERLIN / No Return Prediction Gradient to z at t=-20, t=-10, t=-5

Agent Steps to Goal

# Predictive Model (Prior vs. Posterior) for Different Values of Return Cost Coefficient

# Hierarchical Behavior



Memory-Based Predictor Reading from Memory



Read Head 1

Read Head 2

Read Head 3

# Arbitrary Visuomotor Mapping (Wise and Murray, 2000)

# Arbitrary Visuomotor Mapping (Wise and Murray, 2000)

# Rapid Reward Valuation (One-Shot)

# Rapid Reward Valuation (One-Shot)

# An Episodic Task that Breaks Gradient Flow

**Block 1**   **Block 2**   **Block 3**   **Block 1'**



**Episodic Water Mazes**

Red Room   Yellow Room   Green Room

# An Episodic Task that Breaks Gradient Flow

Episodic Water Mazes

# Learning Curves with Longer Temporal Credit Assignment Windows

MERLIN — MERLIN (tau = 200) — RL-LSTM — RL-LSTM (tau = 200) — RL-MEM — RL-MEM (tau = 200)

## Episodic Water Mazes

# Learning Curves with Longer Temporal Credit Assignment Windows

Legend: MERLIN | MERLIN (tau = 200) | RL-LSTM | RL-LSTM (tau = 200) | RL-MEM | RL-MEM (tau = 200)

## a. Arbitrary Visuomotor Mapping



## b. Executing Transient Instructions



## c. Episodic Water Mazes

# Latent Learning



**Learning in the absence of direct task reward.**    **Tolman and Honzik, 1930**

# Latent Learning

# Learning from Demonstrations

Programmer's Apprentice: (Waters and Rich, 1987)

The objective for programming is hard to write down!

Imitation is a useful proxy when the objective is hard to articulate.

Maybe that's why it was called the programmer's *apprentice*.

# Learning Human Behaviors

Given human input-output measurements, how to build a good imitation learner?

Simple approach is supervised learning or behavioral cloning:

$$\tau = (s_1, a_1, s_2, a_2, s_3, a_3, \ldots) \qquad \max_\theta \sum_k \log p_\theta(a_k | s_k)$$

Problem with this approach: if you find yourself in a state unlike one the human demonstrated, the policy may do anything.

# Learning Human Behaviors

Another approach: make the achieved states / trajectories look like one's from demonstration data.

$$\max_{\theta} \sum_k \log[p(s_{k+1}|s_k, a_k)p_\theta(a_k|s_k)]$$

Advantage: if the agent gets away from the distribution of states, there is a motivation to get back.

Problem with this approach: It typically implies having a model of transitions that is good. This can be hard for complicated environments.

# Learning Human Behaviors

Another approach: make marginal distributions over states similar to the expert's.

$$\max_\theta \sum_k p(s_k)$$

Inverse Reinforcement Learning

Generative Adversarial Imitation Learning (Ho and Ermon, 2016) based on GANs (Goodfellow, 2014). Here, fool the discriminator:

$$\max_\theta p_{\text{discriminator}}(\text{demonstration}|s_k)$$

# Learning Human Behaviors

Generative Adversarial Imitation Learning (Ho and Ermon, 2016)

Learning Human Behaviors from Motion Capture by Adversarial Imitation (Merel, Tassa, Dhruva TB, Sriram Srinivasan, Jay Lemmon, Ziyu Wang, Greg Wayne, Nicolas Heess, 2017)

# Conditional GAIL

Note: no action information strictly needed; important, as explained next

# Algorithm

---

**Algorithm 1**

---

**Input:** Set of demonstration observations $\{z_t^d, c_t^d\}_{t=1...T^d}$
Randomly initialize policy ($\pi_\theta$) and discriminator ($D_\phi$)
// Perform N training iterations of policy & discriminator updating
**for** $i$ in $1...N$ **do**
    Execute policy rollouts to collect $T^g$ timestep observations, $\{z_t^g, c_t^g\}_{t=1...T^g}$
    Compute rewards $\{r_t = -\log(1 - D_\phi(z_t^g, c_t^g))\}_{t=1...T^g}$
    Update $\theta$ (e.g. by TRPO)
    // Perform M discriminator updates steps
    **for** $j$ in $1...M$ **do**
        $\ell(\phi) = \sum_{t=1...T^g} \log(1 - D_\phi(z_t^g, c_t^g)) - \sum_{t=1...T^d} \log(D_\phi(z_t^d, c_t^d))$
        Update $\phi$ by a gradient method w.r.t. $\ell(\phi)$
**Return:** $\pi$

---

# Pipeline



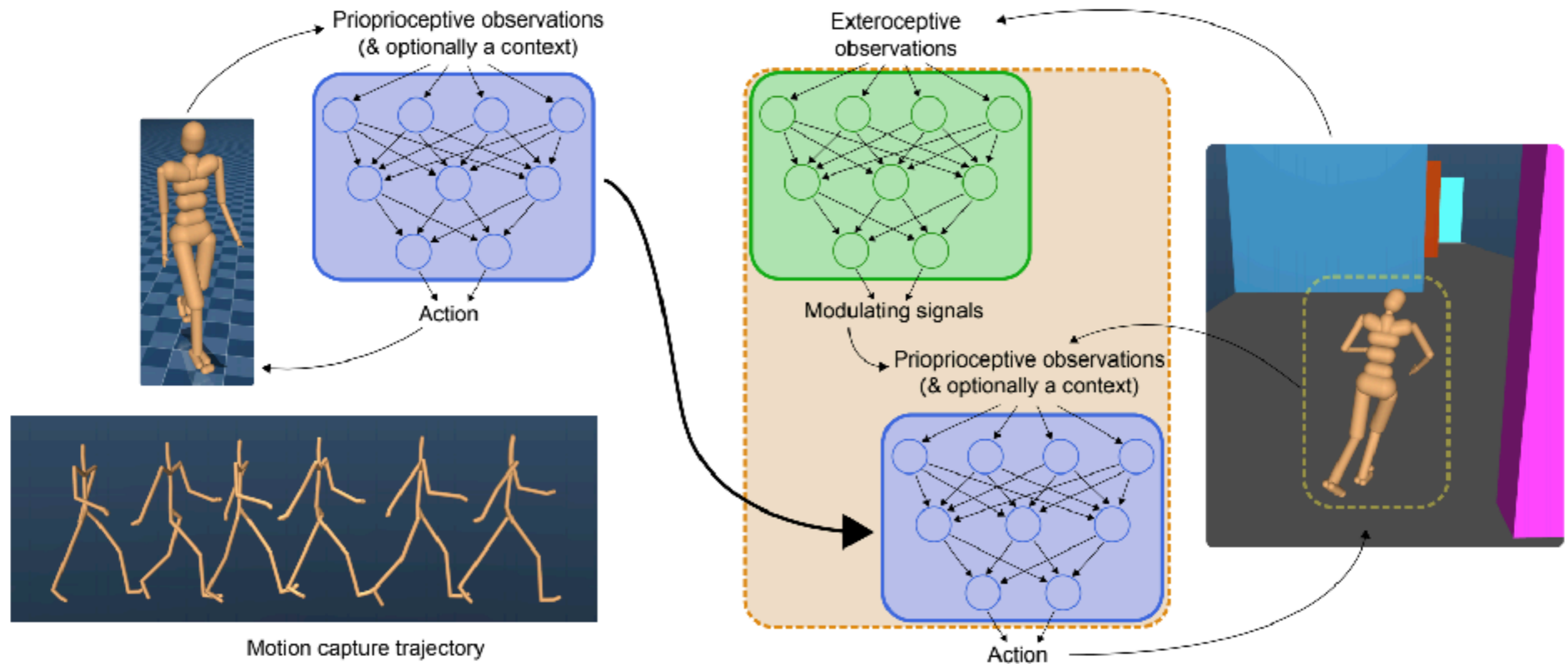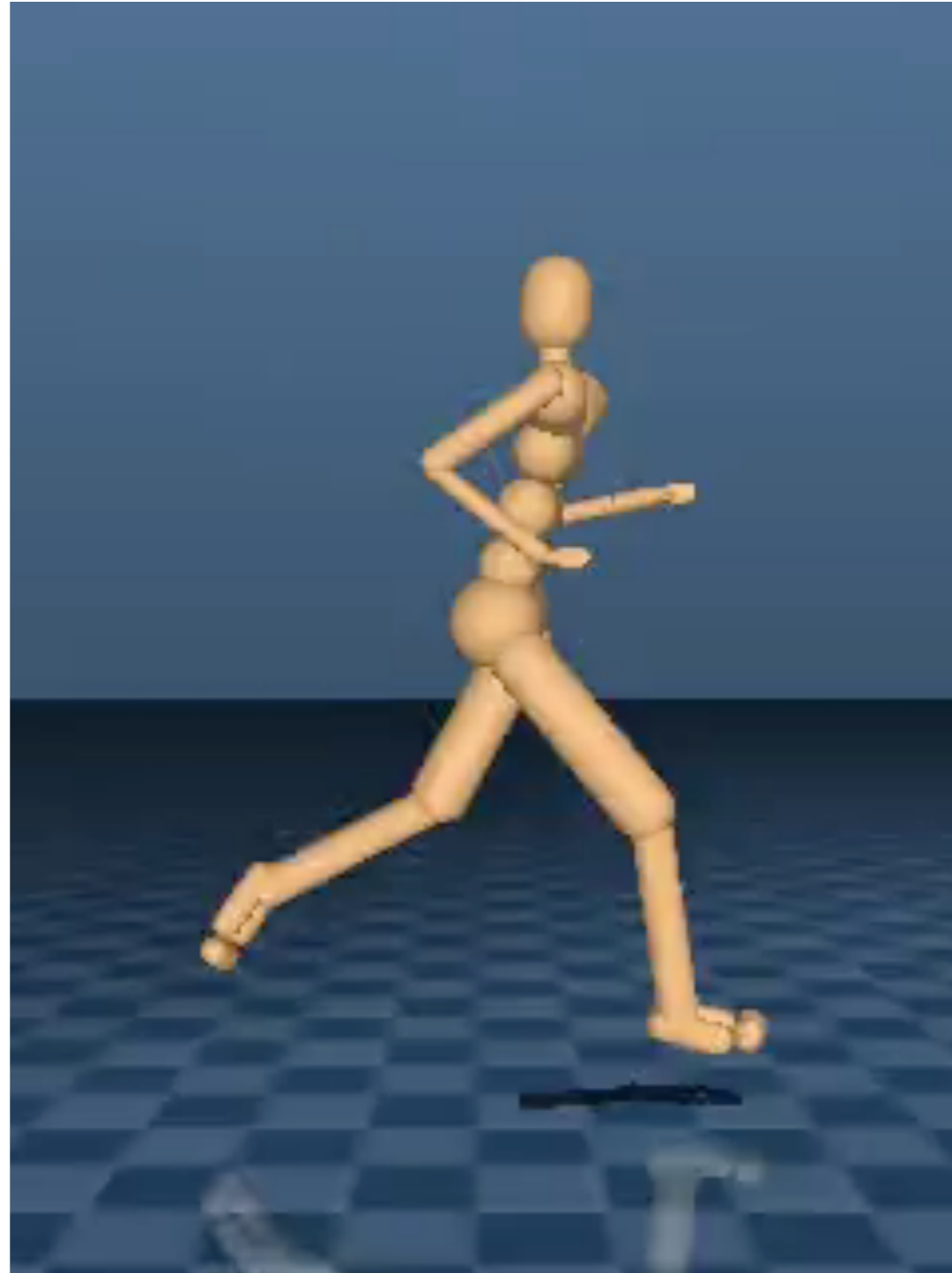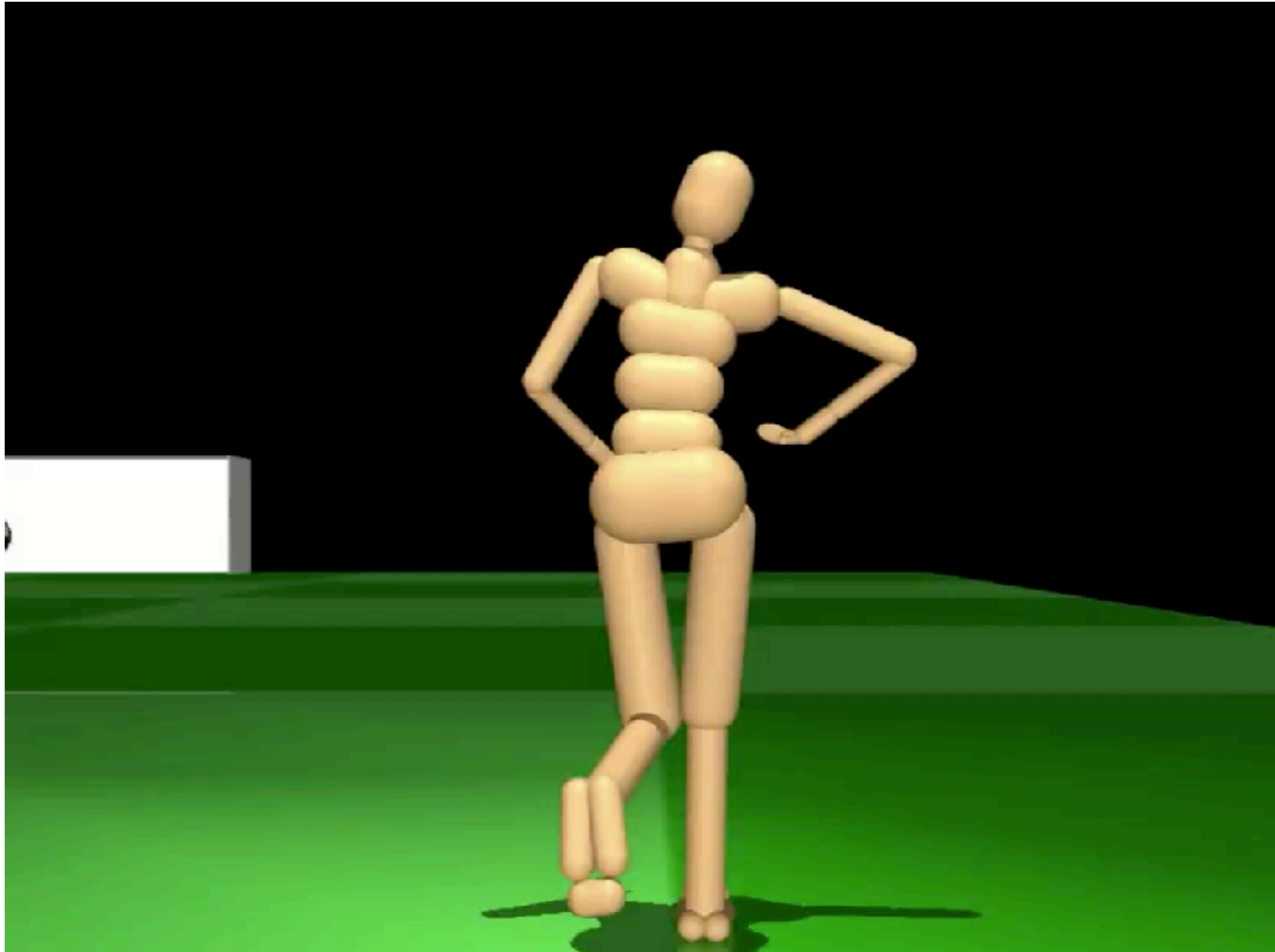Figure 1: Overview of our approach: (Left) First train specific skills into low-level controller (LLC) policies by imitation learning from motion capture data. (Right) Train a high-level controller (HLC) by RL to reuse pre-trained LLCs.

# Running and Turning

# Running and Turning

# Conclusions and Prospects

- MERLIN model can solve quite sophisticated tasks from raw sensory data in partially observed environments.

- Perception and memory formation are guided by a process of predictive modeling and compression, less by trial and error task success.

  - 200 dimensional z captures relevant features from approximately 10^4 sensory dimensions.

- Lessens the need for end-to-end gradient computation.

  - MERLIN uses a temporal credit assignment window of 1.3 seconds to solve memory tasks of 6 minutes.

- Can make use of information acquired without associated reward, exhibiting properties of latent learning.

- Provides a conceptual but **functional** model of the interaction of multiple neural systems in a complete, goal-directed cognitive architecture.

- GAIL can be scaled up to do imitation learning of non-obvious objectives.

# Noticed This

- Btw, sounds relevant: NEURAL SKETCH LEARNING FOR CONDITIONAL PROGRAM GENERATION (Murali et al., 2018)

# Backmatter

$(v_t, r_{t-1}, T_t) \rightarrow o_t \rightarrow e_t \rightarrow \tilde{n}_t \rightarrow a_t --\cdots$

ENCODER

Policy Loss

$(I_t, v_t, r_{t-1}, T_t) \rightarrow$

$\rightarrow a_{t-1}$

$\cdots \rightarrow a_{t-1}$

ERLIN

READ-ONLY POL

MEMORY-BASED PREDICTOR

$k_t \cdots\rightarrow M_t \leftarrow\cdots \tilde{k}_t \leftarrow$

READ

READ

PRIOR

$e_t \rightarrow n_t \leftarrow p \leftarrow\cdots h_t \leftarrow\cdots m_t$

$\tilde{m}_t$

ENCODER

KL Loss

WRITE

$(v_t, r_{t-1}, T_t) \rightarrow o_t$

$q \rightarrow z_t$

POSTERIOR

$\rightarrow a_{t-1}$

DECODER

$(\hat{I}_t, \hat{R}_t, \hat{v}_t, \hat{a}_{t-1}, \hat{r}_{t-1}, \hat{T}_t)$

Reconstruction Loss

unsupervised learning / compressing sensory data

state variables — prior and posterior; instantiation in memory model

variational autoencoder framework

problem with pure unsupervised learning: bullet problem (gluck and myers)

scaling to challenging rl problems — typically use truncated backpropagation through time

Knock on benefit: can cut backpropagation time scales

figures: back matter learning curves for navigation tasks, one step prediction, unroll length ext. fig 10