# Understanding the Training and Inference of Reinforcement Learning
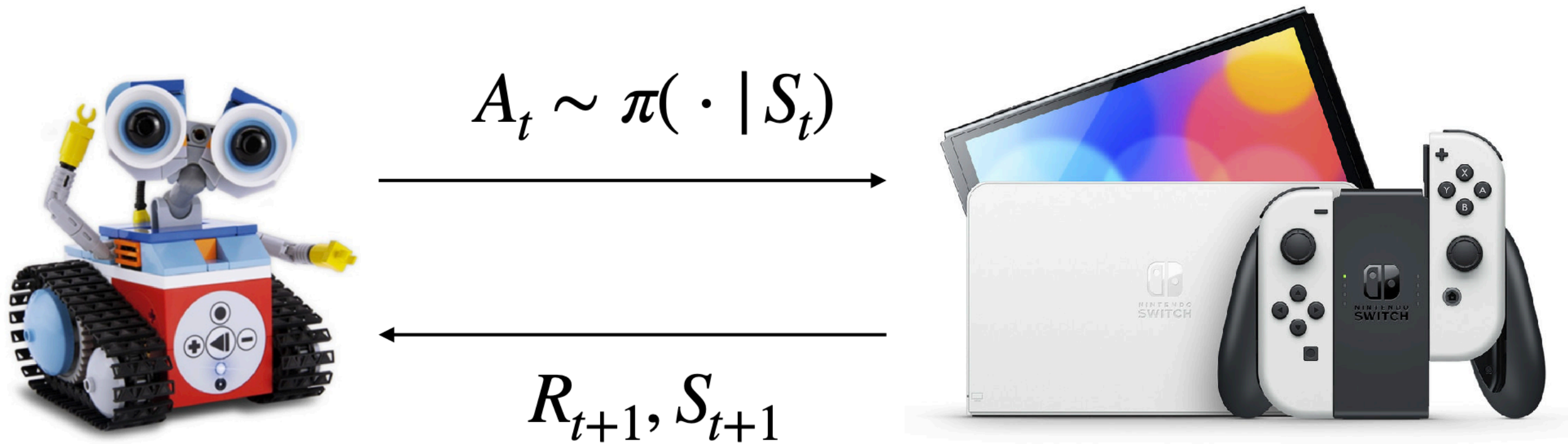
Shangtong Zhang
University of Virginia

# What is RL?

# What is RL?

- RL is PPO!

# What is RL?

$$A_t \sim \pi(\,\cdot\,|\,S_t)$$

$$R_{t+1}, S_{t+1}$$

$$v_\pi(s) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots \,|\, S_t = s]$$

# RL = SA + DP

- RL is to use Stochastic Approximation (SA) method to solve Dynamic Programming (DP) problem

- Bellman operator $Tv = r_\pi + \gamma P_\pi v$

# RL = SA + DP

- $v_{k+1} = Tv_k = r_\pi + \gamma P_\pi v_k$

- Challenge: unknown $P_\pi$

- Solution: use a sample

$$(P_\pi v_k)(s) = \sum_{s'} p(s' \mid s) v_k(s')$$
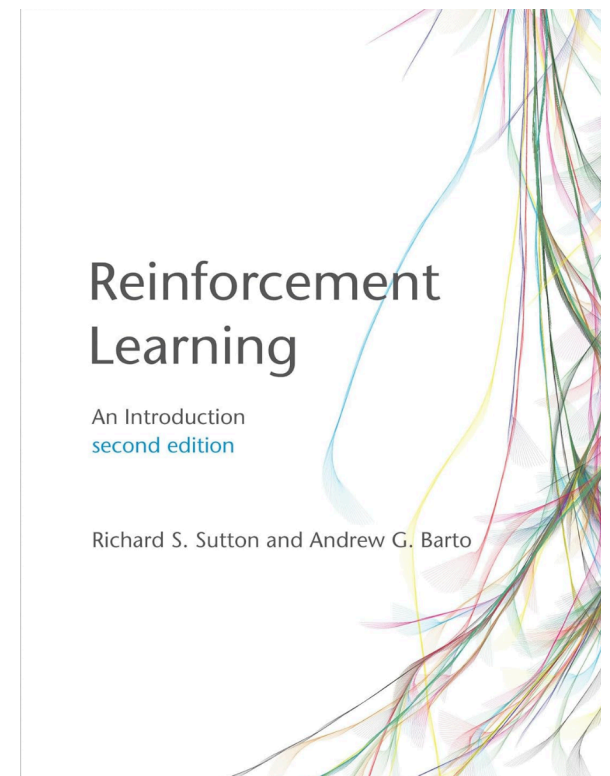
$$(P_\pi v_k)(s) \approx v_k(s')$$

# RL = SA + DP

- $v_{k+1} = Tv_k = r_\pi + \gamma P_\pi v_k$

- Challenge: full update is too aggressive
$v_{k+1}(s) = r_\pi(s) + \gamma v_k(s')$

- Solution: incremental update
$v_{k+1}(s) = v_k(s) + \alpha_k(r_\pi(s) + \gamma v_k(s') - v_k(s))$

# RL = SA + DP

- $v_{k+1} = Tv_k = r_\pi + \gamma P_\pi v_k$

- Challenge: where to get $s'$?
  $\ldots, S_k, S_{k+1}, \ldots,$

- Solution: asynchronous update
  $v_{k+1}(S_k) = v_k(S_k) + \alpha_k(r_\pi(S_k) + \gamma v_k(S_{k+1}) - v_k(S_k))$

# RL = SA + DP

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$

- Different $h$ realizes different RL algorithms, e.g.,
  TD, Q-learning, linear TD, Gradient TD, Emphatic TD,
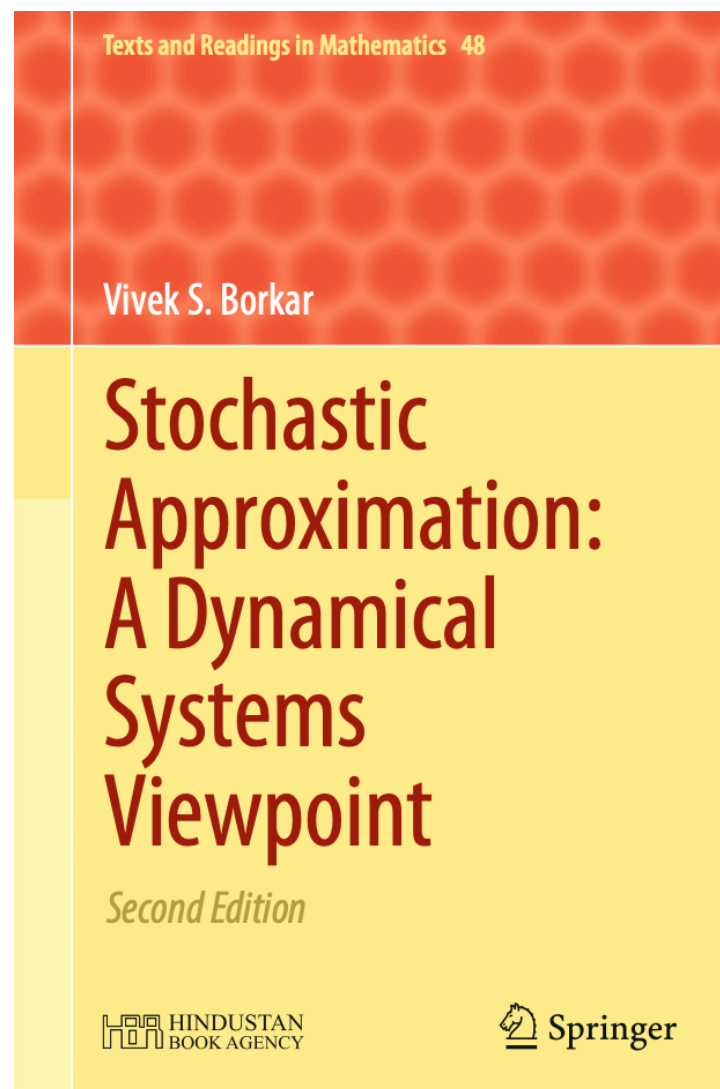  average reward TD, differential TD, differential Q-learning

Reinforcement
Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

# Does SA converge?

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$

- $v_k \to v_*$ almost surely?

- Early RL pioneers borrow results from SA community

- Now RL theorists shift to fancier problems, e.g., offline RL, RLHF, etc.

# Tools in SA community do not apply to RL well

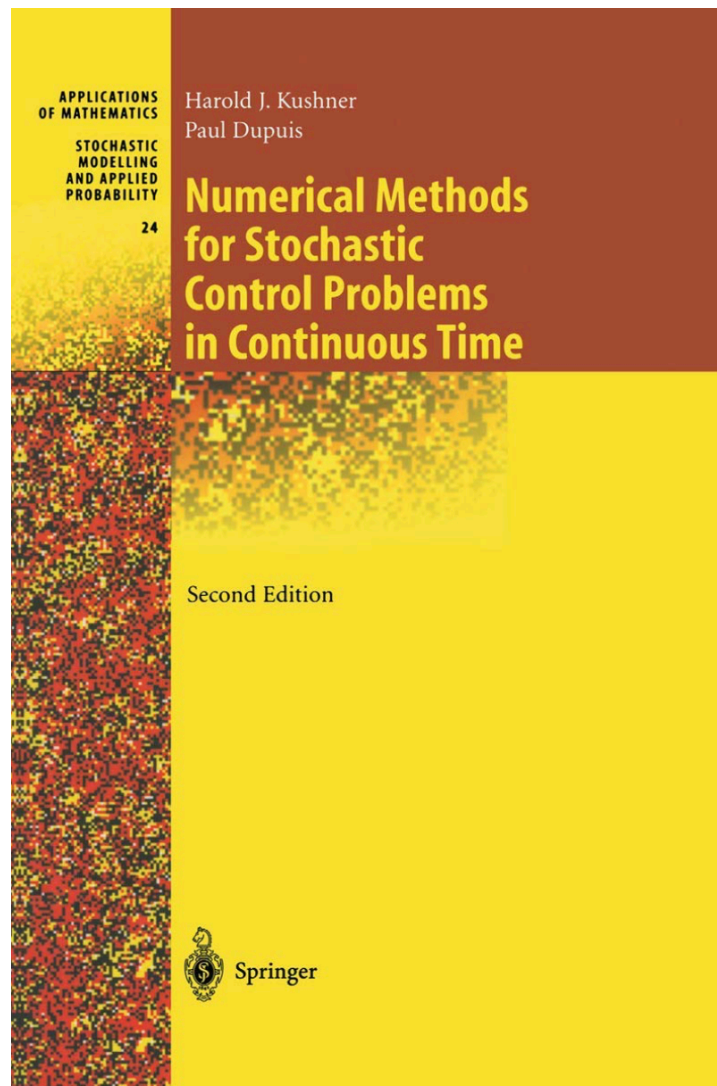- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$

**Texts and Readings in Mathematics  48**

**Vivek S. Borkar**

**Stochastic Approximation: A Dynamical Systems Viewpoint**

*Second Edition*

HINDUSTAN BOOK AGENCY

Springer

**Assumption:** $\{Y_k\}$ **are i.i.d.**

**Reality in RL:** $\{Y_k\}$ **are Markov chain**

# Tools in SA community do not apply to RL well

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$



APPLICATIONS OF MATHEMATICS

STOCHASTIC MODELLING AND APPLIED PROBABILITY

24

Harold J. Kushner
Paul Dupuis

**Numerical Methods for Stochastic Control Problems in Continuous Time**

Second Edition

Springer

**Assumption:** $\sup_k ||v_k|| < \infty$ **a.s.**
**Reality in RL: very hard to verify**
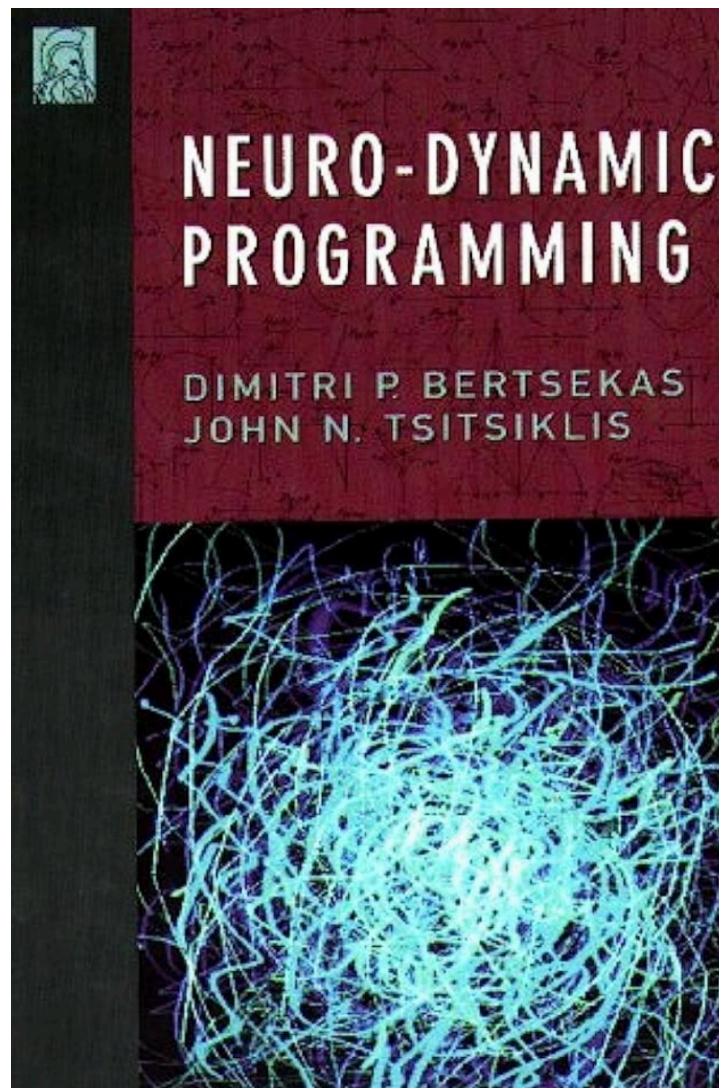
# Tools in SA community do not apply to RL well

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$



22

Albert Benveniste  Michel Métivier
Pierre Priouret

Adaptive Algorithms
and Stochastic
Approximations

Applications of Mathematics

Springer-Verlag

**Assumption:  Poisson's equation**
**Lyapunov function**
**Reality in RL: hard to prove existence**

# Tools in SA community do not apply to RL well

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$



**Assumption:  linear $h$**
**with a negative definite matrix**
**Reality in RL: not true in many RL algorithms**

# Why can't the RL community have its own SA theory?

- $v_{k+1} = v_k + \alpha_k h(v_k, Y_{k+1})$

**Assumption: Law of large numbers (LLN)**
**Reality in RL: if LLN does not hold,**
**good luck with your RL alg.**

# LLN on Markov chains

- $$\frac{X_1 + X_2 + \ldots X_n}{n} \to E[X_1]$$

- $$\frac{1}{n} \sum_{i=1}^{n} h(Y_i) \to \sum_{y} d(y)h(y)$$

  time average $\to$ space average (ergodicity)

- automatically hold in finite chains
  especially powerful in off-policy RL algorithms with
  eligibility traces

# Training is all before deep RL

- Tabular TD
$$v_{k+1}(S_k) = v_k(S_k) + \alpha_k(r_\pi(S_k) + \gamma v_k(S_{k+1}) - v_k(S_k))$$
$$s \to v_*(s)$$

- Linear TD
$$w_{k+1} = w_k + \alpha_k(r_\pi(S_k) + \gamma x(S_{k+1})^\top w_k - x(S_k)^\top w_k)x(S_k)$$
$$s \to x(s)^\top w$$

# Inference matters after deep RL

- Deep TD
$$\theta_{t+1} = \theta_t + \alpha_t(r_\pi(S_t) + \gamma v(S_{t+1}; \theta_t) - v(S_t; \theta_t)) \nabla v(S_t; \theta_t)$$

- Deep-Q-Networks
$$\theta_{t+1} = \theta_t + \alpha_t(R_{t+1} + \gamma \max_a q(S_{t+1}, a; \bar{\theta}_t) - q(S_t, A_t; \theta_t)) \nabla q(S_t, A_t; \theta_t)$$

- Training: $\theta_t \rightarrow \theta_{t+1} \rightarrow \ldots \rightarrow \theta_*$

- Inference: $(s, a) \rightarrow q(s, a; \theta_*)$

# In-context learning is perhaps the most trending inference problem

```
6 -> number
a -> letter
7 ->
```

forward pass of LLM / Transformer →

number

# In-context learning is perhaps the most trending inference problem

$$
\begin{array}{l}
S_1 \rightarrow r_\pi(S_1) \\
S_2 \rightarrow r_\pi(S_2) \\
\ldots \\
S_n \rightarrow r_\pi(S_n) \\
s \rightarrow
\end{array}
$$

**forward pass of LLM / Transformer**

$\longrightarrow$

# In-context learning is perhaps the most trending inference problem

$$\boxed{\begin{aligned} S_1 &\to r_\pi(S_1) \\ S_2 &\to r_\pi(S_2) \\ &\dots \\ S_n &\to r_\pi(S_n) \\ s &\to \end{aligned}}$$

forward pass of LLM / Transformer

$\longrightarrow$

$\boxed{r_\pi(s)}$

# What about predicting value?

$$
\begin{array}{l}
S_1 \rightarrow r_\pi(S_1) \\
S_2 \rightarrow r_\pi(S_2) \\
\ldots \\
S_n \rightarrow r_\pi(S_n) \\
s \rightarrow
\end{array}
$$

**forward pass of LLM / Transformer**

$\longrightarrow$

$v_\pi(s)$

# Humans predict value via TD

$$S_1 \to r_\pi(S_1)$$
$$S_2 \to r_\pi(S_2)$$
$$\dots$$
$$S_n \to r_\pi(S_n)$$
$$s \to$$

$$w_{k+1} = w_k + \alpha_k(r_\pi(S_k) + \gamma x(S_{k+1})^\top w_k - x(S_k)^\top w_k)x(S_k)$$

$$w_0 \to w_1 \to w_2 \to w_3 \to \dots$$

$$w_0^\top x(s) \to w_1^\top x(s) \to w_2^\top x(s) \to w_3^\top x(s) \to \dots$$

# Transformers **CAN** mimic what humans do!

$$S_1 \rightarrow r_\pi(S_1)$$
$$S_2 \rightarrow r_\pi(S_2)$$
$$\ldots$$
$$S_n \rightarrow r_\pi(S_n)$$
$$s \rightarrow 0$$

**Attention Layer** $\longrightarrow$

$$S_1 \rightarrow r_\pi(S_1)$$
$$S_2 \rightarrow r_\pi(S_2)$$
$$\ldots$$
$$S_n \rightarrow r_\pi(S_n)$$
$$s \rightarrow v_1$$

**Attention Layer** $\longrightarrow$

$$S_1 \rightarrow r_\pi(S_1)$$
$$S_2 \rightarrow r_\pi(S_2)$$
$$\ldots$$
$$S_n \rightarrow r_\pi(S_n)$$
$$s \rightarrow v_2$$

**Attention Layer** $\longrightarrow$

$$v_1 = w_1^\top x(s) \qquad\qquad v_2 = w_2^\top x(s)$$

**(If the linear attention layer has special weights)**

# Transformers **DO** mimic what humans do!

- Deep TD
$$\theta_{t+1} = \theta_t + \alpha_t(r_\pi(S_t) + \gamma v(S_{t+1}; \theta_t) - v(S_t; \theta_t)) \nabla v(S_t; \theta_t)$$
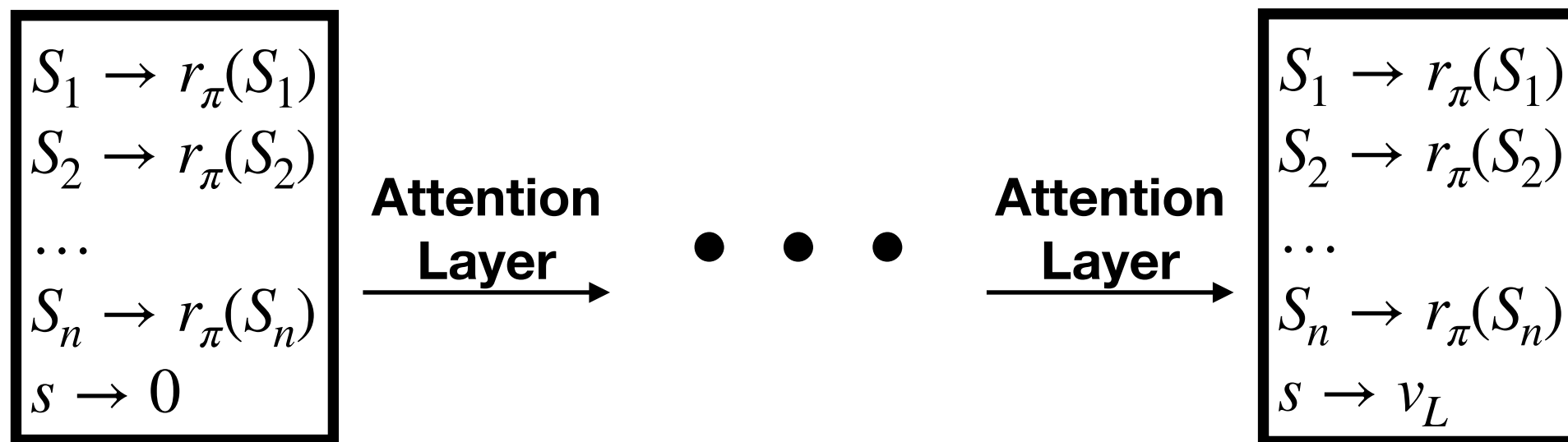
- Parameterize $v(\text{context}, s; \theta)$ as an $L$-layer transformer

$$\boxed{\begin{array}{l} S_1 \to r_\pi(S_1) \\ S_2 \to r_\pi(S_2) \\ \dots \\ S_n \to r_\pi(S_n) \\ s \to 0 \end{array}} \quad \xrightarrow{\textbf{Attention Layer}} \quad \bullet \ \bullet \ \bullet \quad \xrightarrow{\textbf{Attention Layer}} \quad \boxed{\begin{array}{l} S_1 \to r_\pi(S_1) \\ S_2 \to r_\pi(S_2) \\ \dots \\ S_n \to r_\pi(S_n) \\ s \to v_L \end{array}}$$

# Transformers **DO** mimic what humans do!

- Run deep TD to train the $L$-layer transformer $v(c, s; \theta)$

$$\theta_{t+1} = \theta_t + \alpha_t(r_\pi(S_t) + \gamma v(c, S_{t+1}; \theta_t) - v(c, S_t; \theta_t)) \nabla v(c, S_t; \theta_t)$$

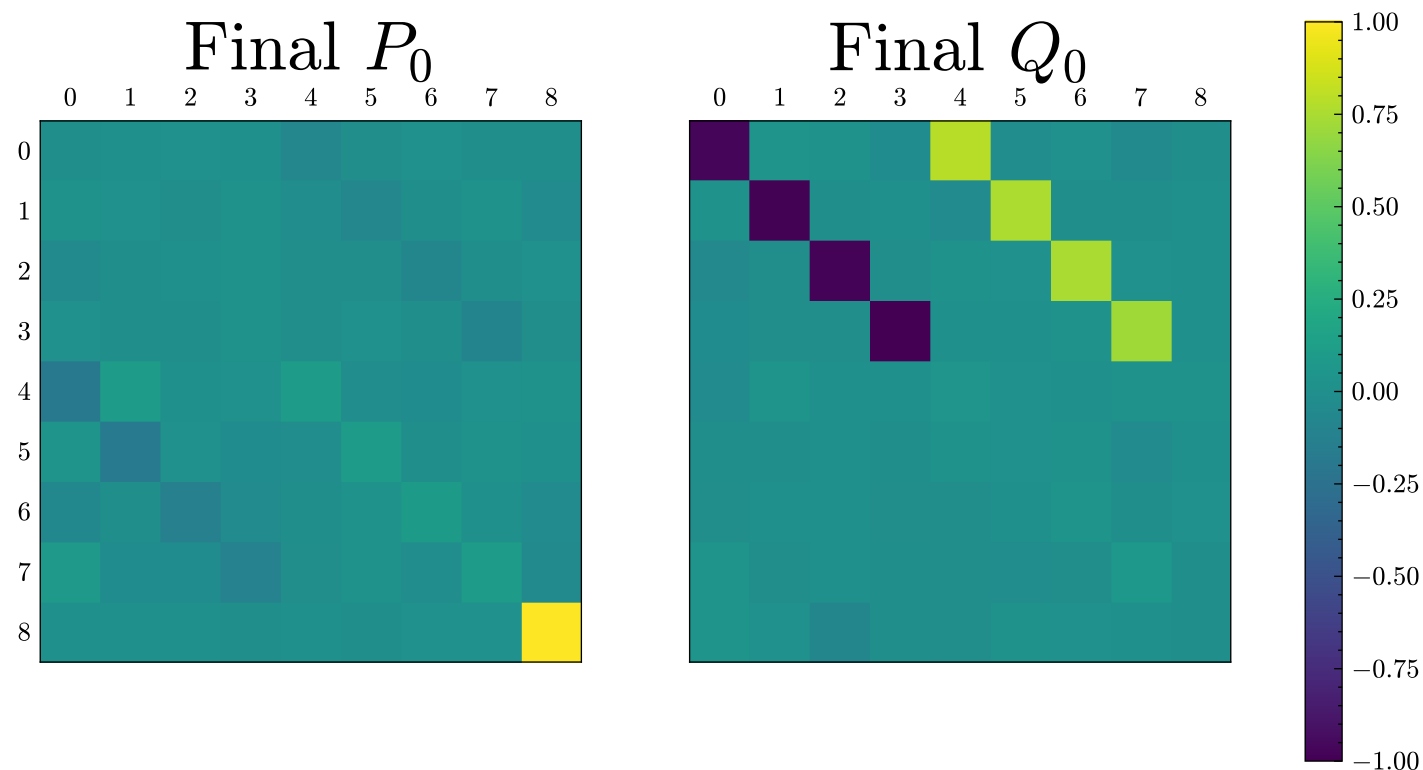- The weights that implement in-context TD emerge after training!

$\theta_\infty = (P_0, Q_0)$

# WHY do transformers mimic what humans do?

- Run deep TD to train the $L$-layer transformer $v(c, s; \theta)$
  $$\theta_{t+1} = \theta_t + \alpha_t(r_\pi(S_t) + \gamma v(c, S_{t+1}; \theta_t) - v(c, S_t; \theta_t)) \nabla v(c, S_t; \theta_t)$$

- The weights that implement in-context TD form an invariant set of the deep TD update

# Transformers can implement more RL algorithms

- Residual gradient

- TD($\lambda$)

- Average-reward TD

# In-context regression as gradient descent (Ahn et.al. 2023)

- RL algorithm is **NOT** gradient descent
  The algorithm in inference is **NOT** gradient descent
  The training algorithm is **NOT** gradient descent

- RL prediction is inhomogeneous

# In-context regression as gradient descent (Ahn et.al.)

- To implement average-reward TD

  - multiple head linear attention

  - overparameterized prompt

$$
\begin{array}{|l|}
\hline
S_1 \rightarrow r_\pi(S_1) \\
S_2 \rightarrow r_\pi(S_2) \\
\ldots \\
S_n \rightarrow r_\pi(S_n) \\
s \rightarrow \\
\hline
\end{array}
\qquad
\begin{array}{|l|}
\hline
S_1 \rightarrow r_\pi(S_1)X \\
S_2 \rightarrow r_\pi(S_2)X \\
\ldots \\
S_n \rightarrow r_\pi(S_n)X \\
s \rightarrow \\
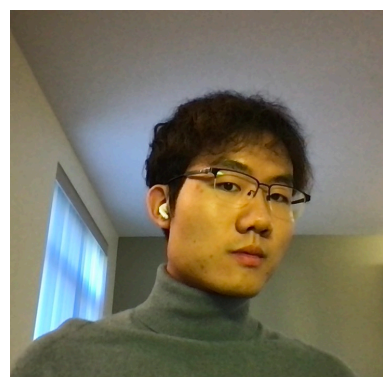\hline
\end{array}
$$

**Shuze Liu (UVA)**

**Shuhang Cheng (Scaled Foundation)**

**Ethan Blaser (UVA)**

**Jiuqi Wang (UVA)**

**Hadi Daneshmand (MIT)**

The ODE paper:
https://arxiv.org/abs/2401.07844

The in-context TD paper:
https://arxiv.org/abs/2405.13861

# Thanks!