# IMPALA: Scalable Distributed Deep-RL with **Imp**ortance Weighted **A**ctor-**L**earner **A**rchitectures
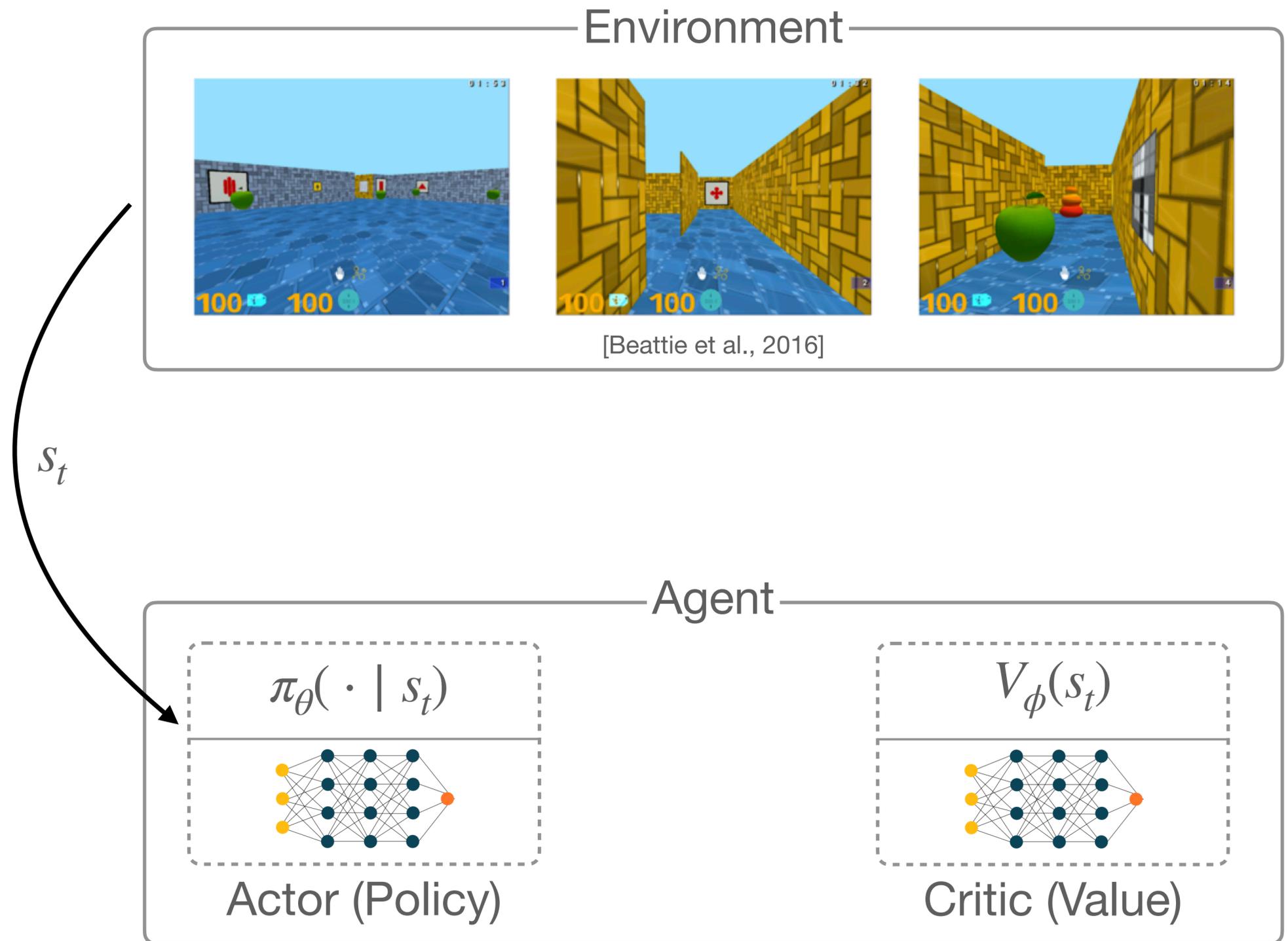
## ICML 2018

Lasse Espeholt *, Hubert Soyer *, Remi Munos *, Karen Simonyan, Volodymyr Mnih , Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, Koray Kavukcuoglu

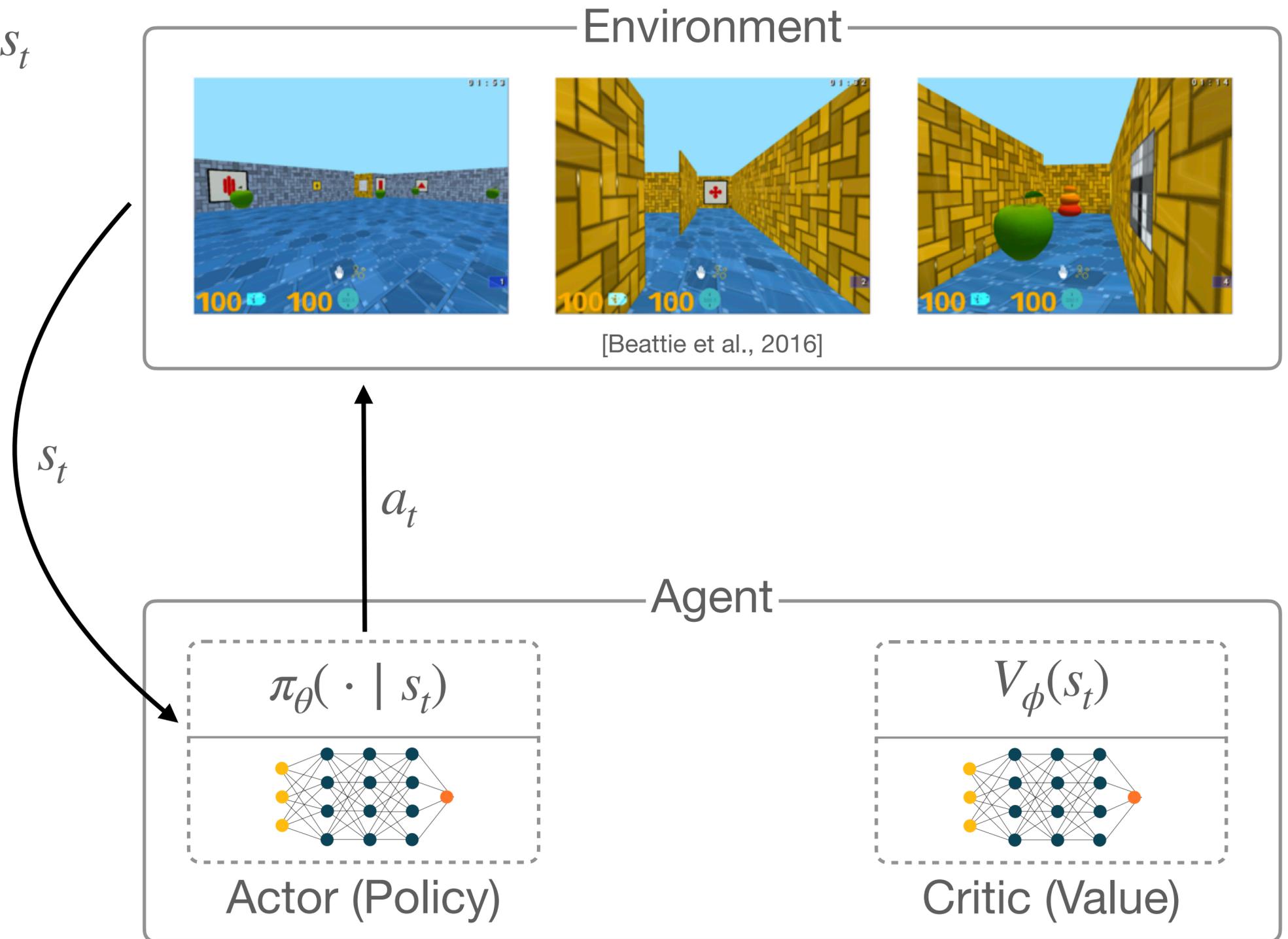Google DeepMind

Presenter - Aditya Thimmaiah

* Equal Contribution

# Motivation | Multi-Task RL



Environment

[Beattie et al., 2016]

$s_t$

Agent

$\pi_\theta( \cdot \mid s_t)$

Actor (Policy)

$V_\phi(s_t)$

Critic (Value)

# Motivation | Multi-Task RL

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

Agent

$\pi_\theta( \cdot \mid s_t)$

Actor (Policy)
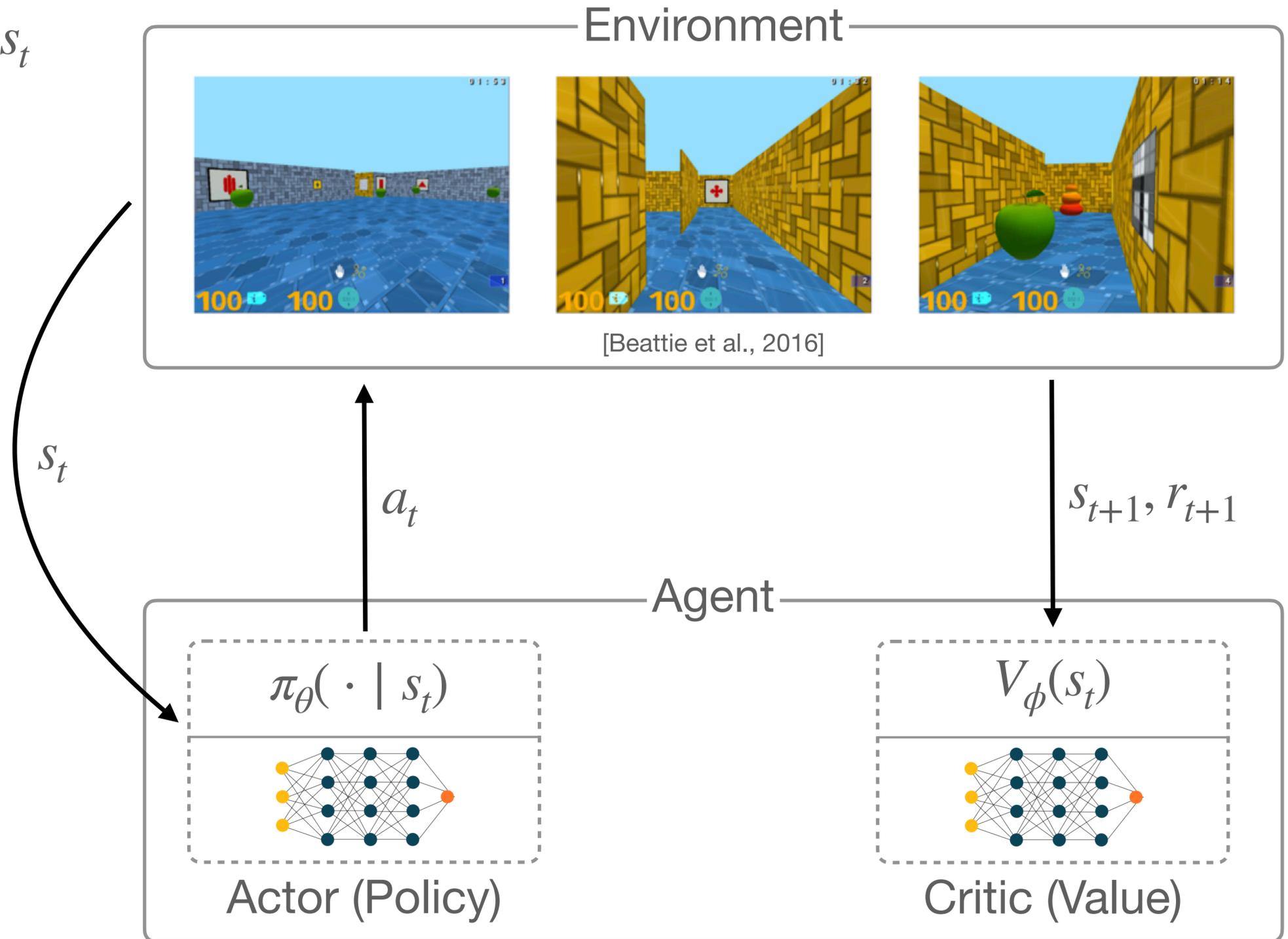
$V_\phi(s_t)$

Critic (Value)

# Motivation | Multi-Task RL

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \,|\, s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

$s_{t+1}, r_{t+1}$

Agent

$\pi_\theta( \cdot \,|\, s_t)$

Actor (Policy)

$V_\phi(s_t)$

Critic (Value)

# Motivation | Multi-Task RL

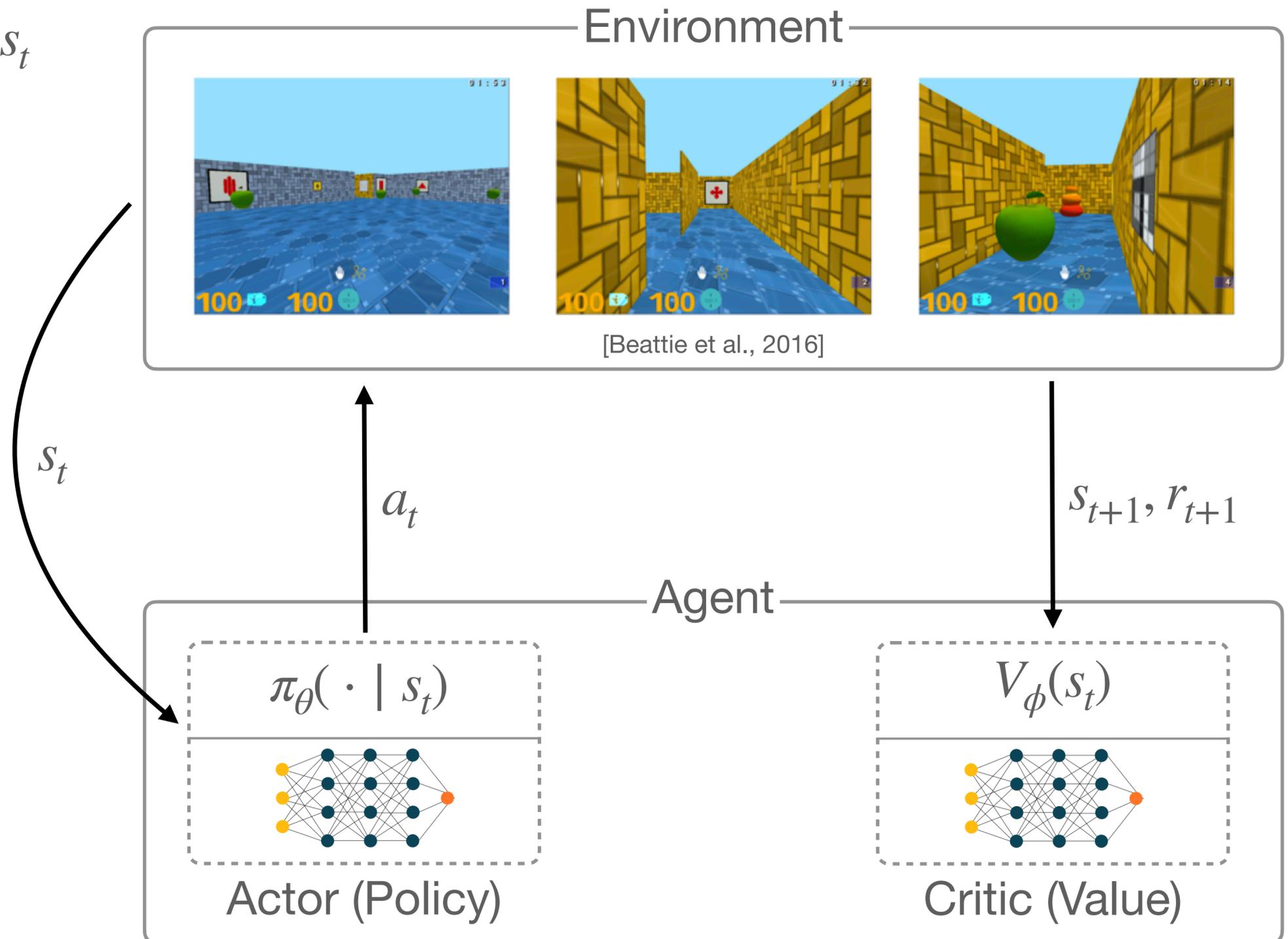**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \,|\, s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \,.\, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

$s_{t+1}, r_{t+1}$

Agent

$\pi_\theta( \cdot \,|\, s_t)$

Actor (Policy)

$V_\phi(s_t)$

Critic (Value)

# Motivation | Multi-Task RL

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \,|\, s_t)$$
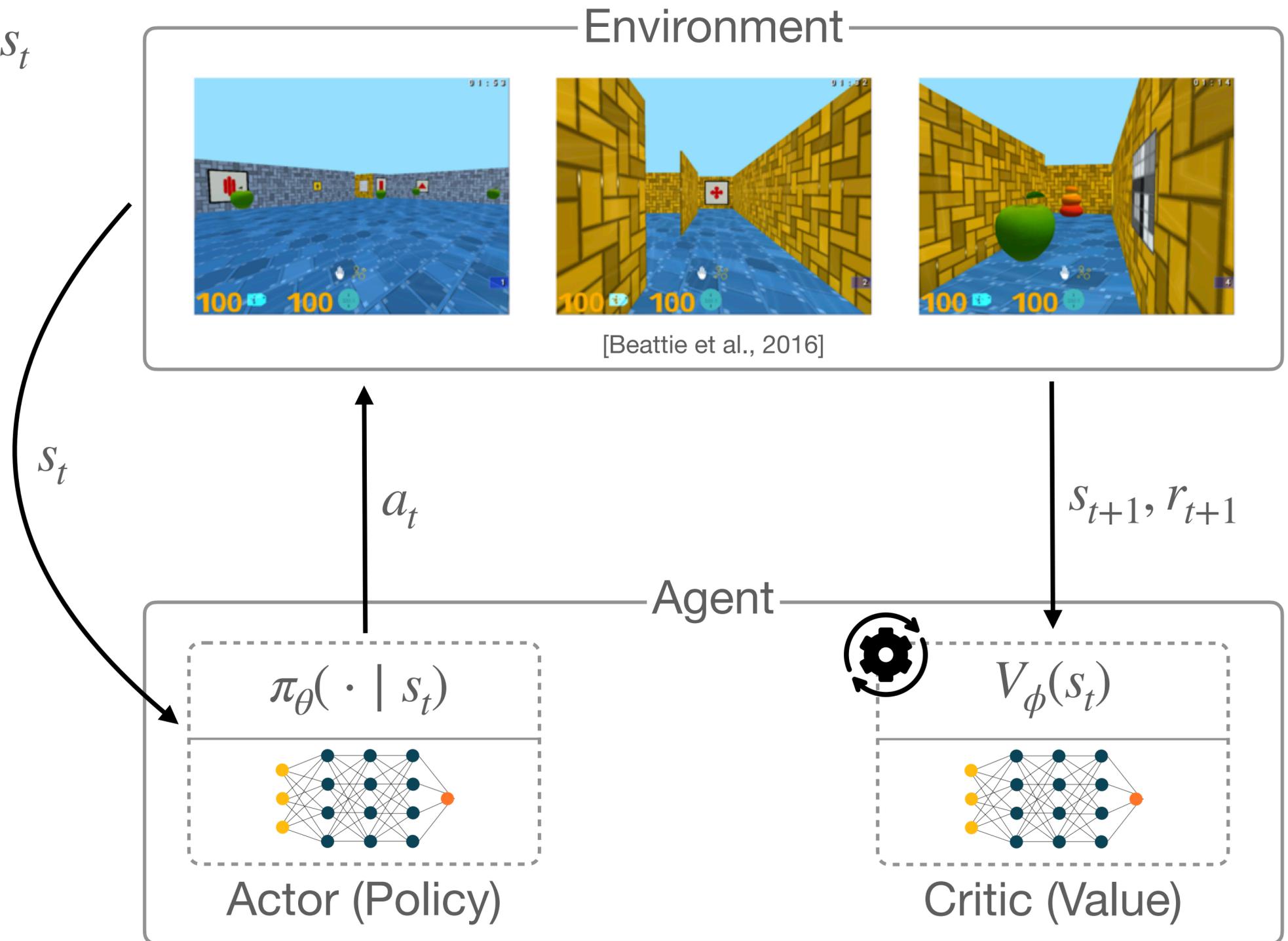
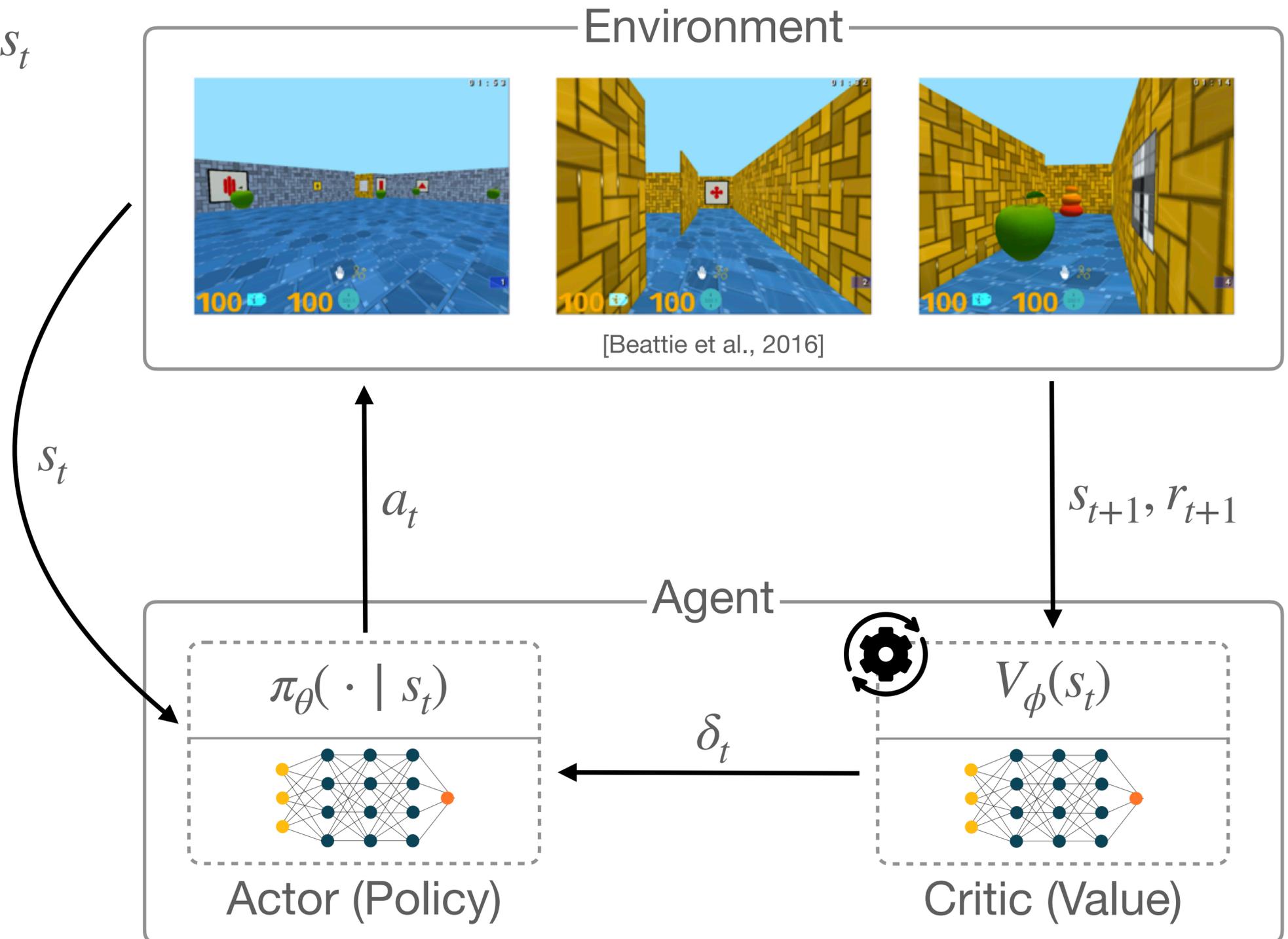**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \,.\, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

$s_{t+1}, r_{t+1}$

Agent

$\pi_\theta( \cdot \,|\, s_t)$

Actor (Policy)

$V_\phi(s_t)$

Critic (Value)

# Motivation | Multi-Task RL

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$
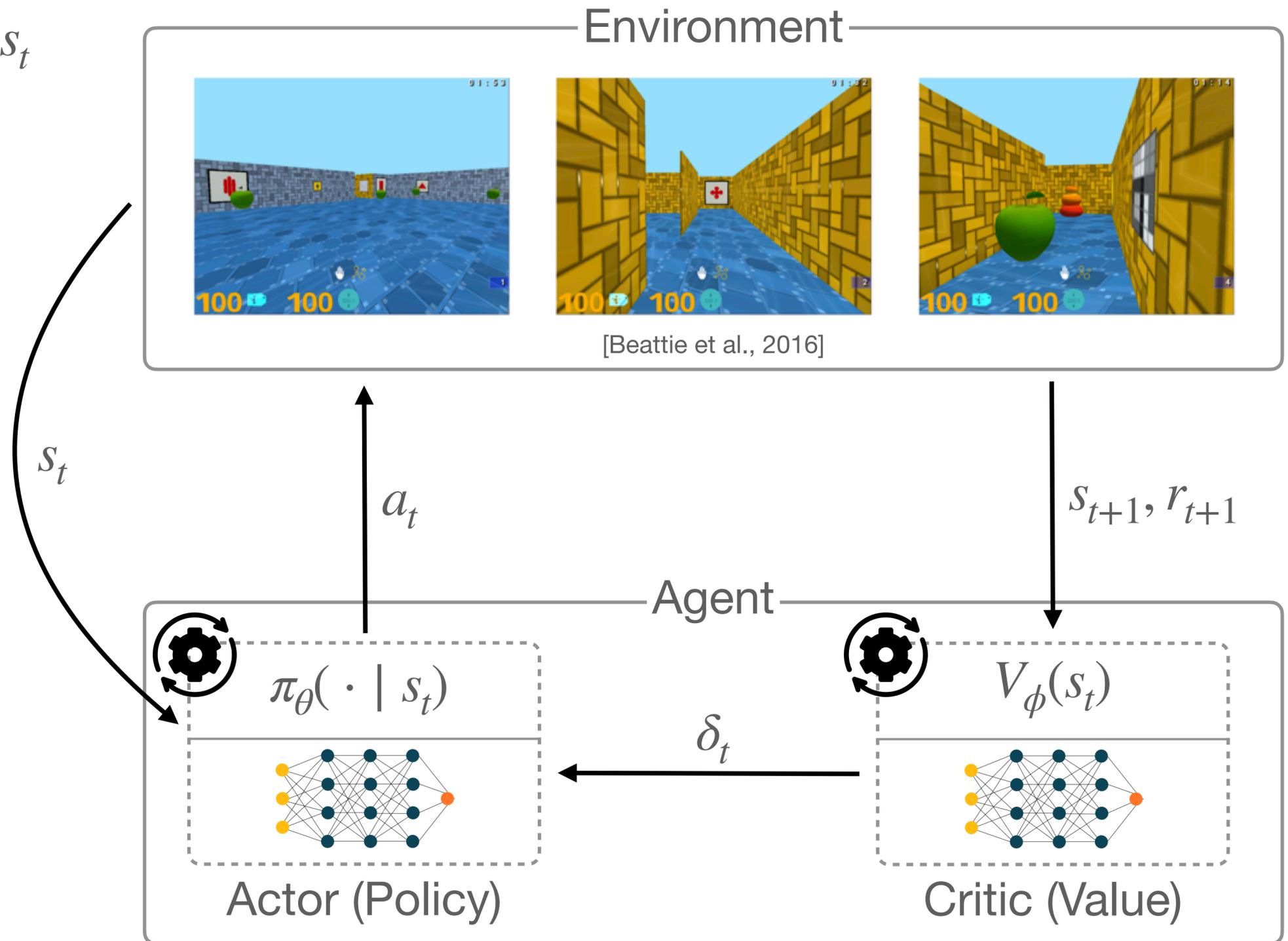
**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

$s_{t+1}, r_{t+1}$

Agent

$\pi_\theta( \cdot \mid s_t)$

$\delta_t$

$V_\phi(s_t)$

Actor (Policy)

Critic (Value)

# Motivation | Multi-Task RL

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**5** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Environment

[Beattie et al., 2016]

$s_t$

$a_t$

$s_{t+1}, r_{t+1}$

Agent

$\pi_\theta( \cdot \mid s_t)$
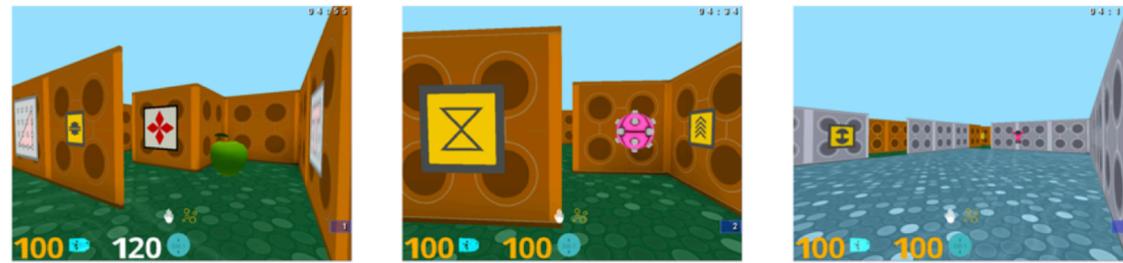
$\delta_t$

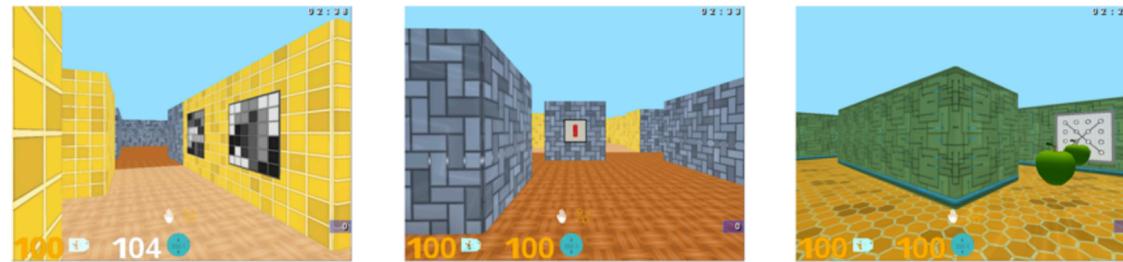$V_\phi(s_t)$
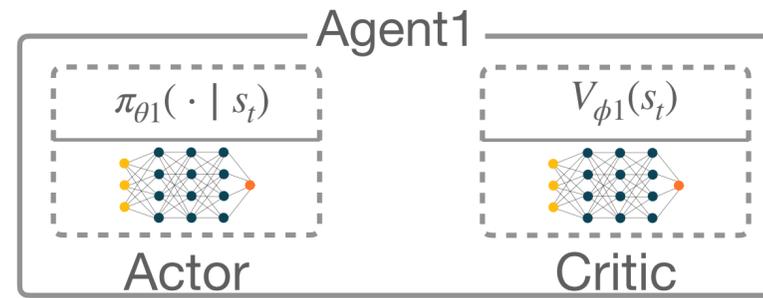
Actor (Policy)

Critic (Value)

# Motivation | Multi-Task RL
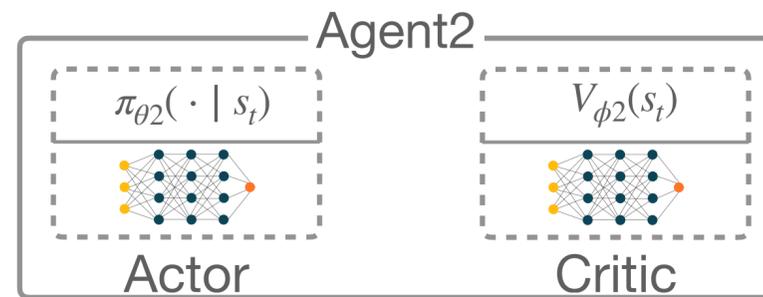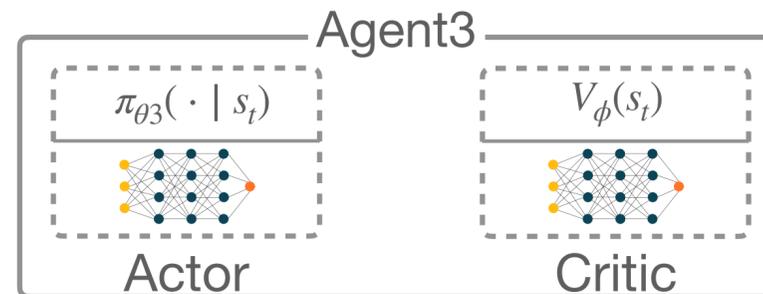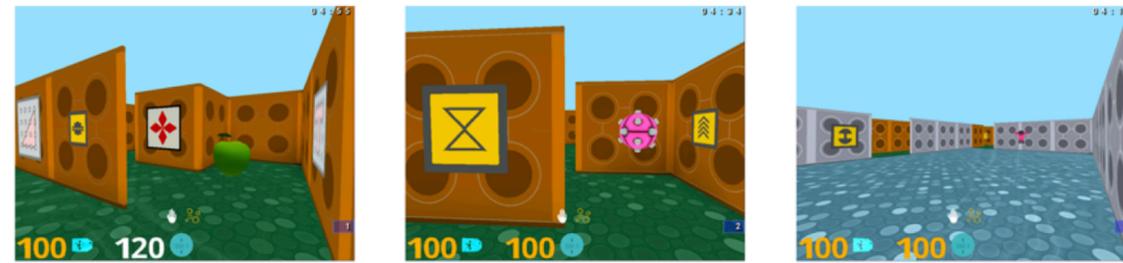
nav_maze*02



nav_maze*03



random_maze



[Beattie et al., 2016]

# Motivation | Multi-Task RL



[Beattie et al., 2016]

# Motivation | Multi-Task RL



nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

# Motivation | Multi-Task RL



Agent1
$\pi_{\theta 1}(\cdot \mid s_t)$
Actor
$V_{\phi 1}(s_t)$
Critic

Agent2
$\pi_{\theta 2}(\cdot \mid s_t)$
Actor
$V_{\phi 2}(s_t)$
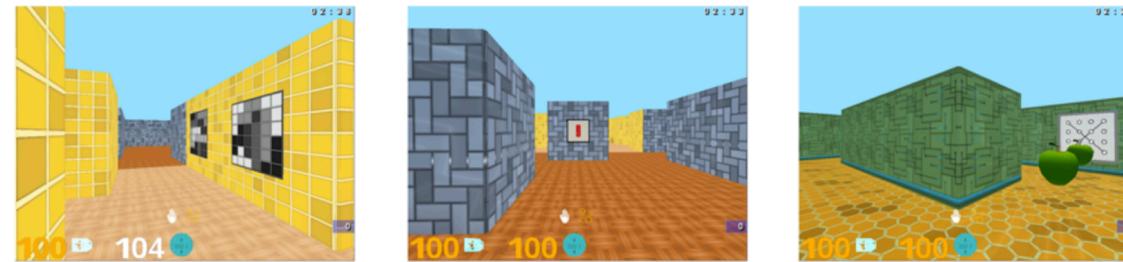Critic

Agent3
$\pi_{\theta 3}(\cdot \mid s_t)$
Actor
$V_{\phi}(s_t)$
Critic

nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

Agent
Actor1
Actor2
Actor3
$\pi_{\theta}(\cdot \mid s_t)$
$V_{\phi}(s_t)$
Shared Policy
Shared Critic

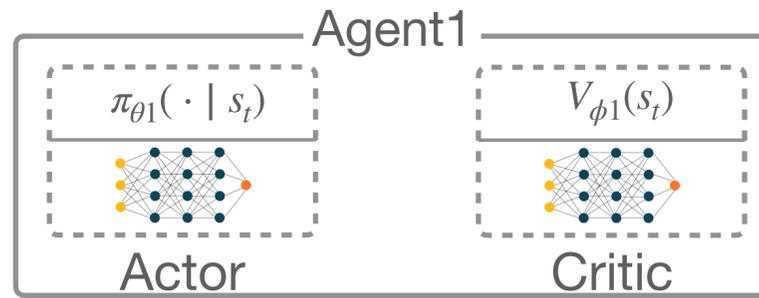✓ Positive transfer of learned features across tasks

# Motivation | Multi-Task RL



nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

Agent1
$\pi_{\theta 1}(\cdot \mid s_t)$ — Actor
$V_{\phi 1}(s_t)$ — Critic

Agent2
$\pi_{\theta 2}(\cdot \mid s_t)$ — Actor
$V_{\phi 2}(s_t)$ — Critic

Agent3
$\pi_{\theta 3}(\cdot \mid s_t)$ — Actor
$V_{\phi}(s_t)$ — Critic

Agent

Actor1
Actor2
Actor3
$\pi_{\theta}(\cdot \mid s_t)$
Shared Policy

$V_{\phi}(s_t)$
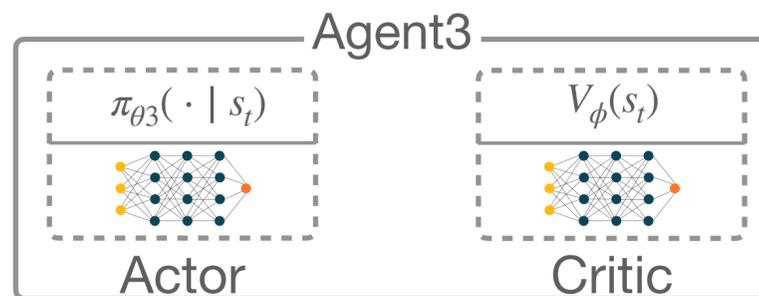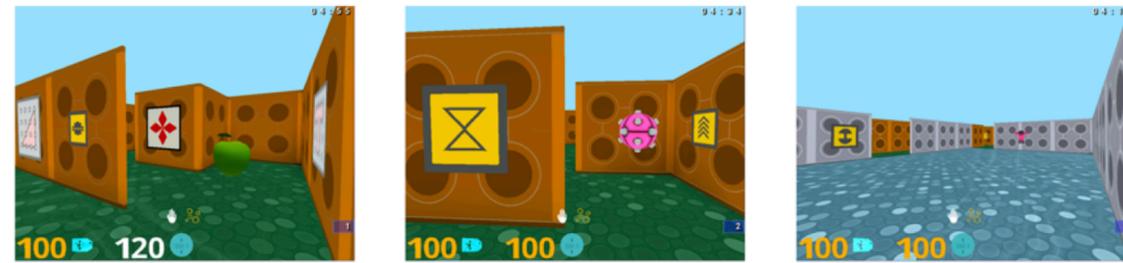Shared Critic

✅ Positive transfer of learned features across tasks

✅ Better compute efficiency

# Motivation | Multi-Task RL



Agent1
$\pi_{\theta 1}(\cdot \mid s_t)$ — Actor
$V_{\phi 1}(s_t)$ — Critic

Agent2
$\pi_{\theta 2}(\cdot \mid s_t)$ — Actor
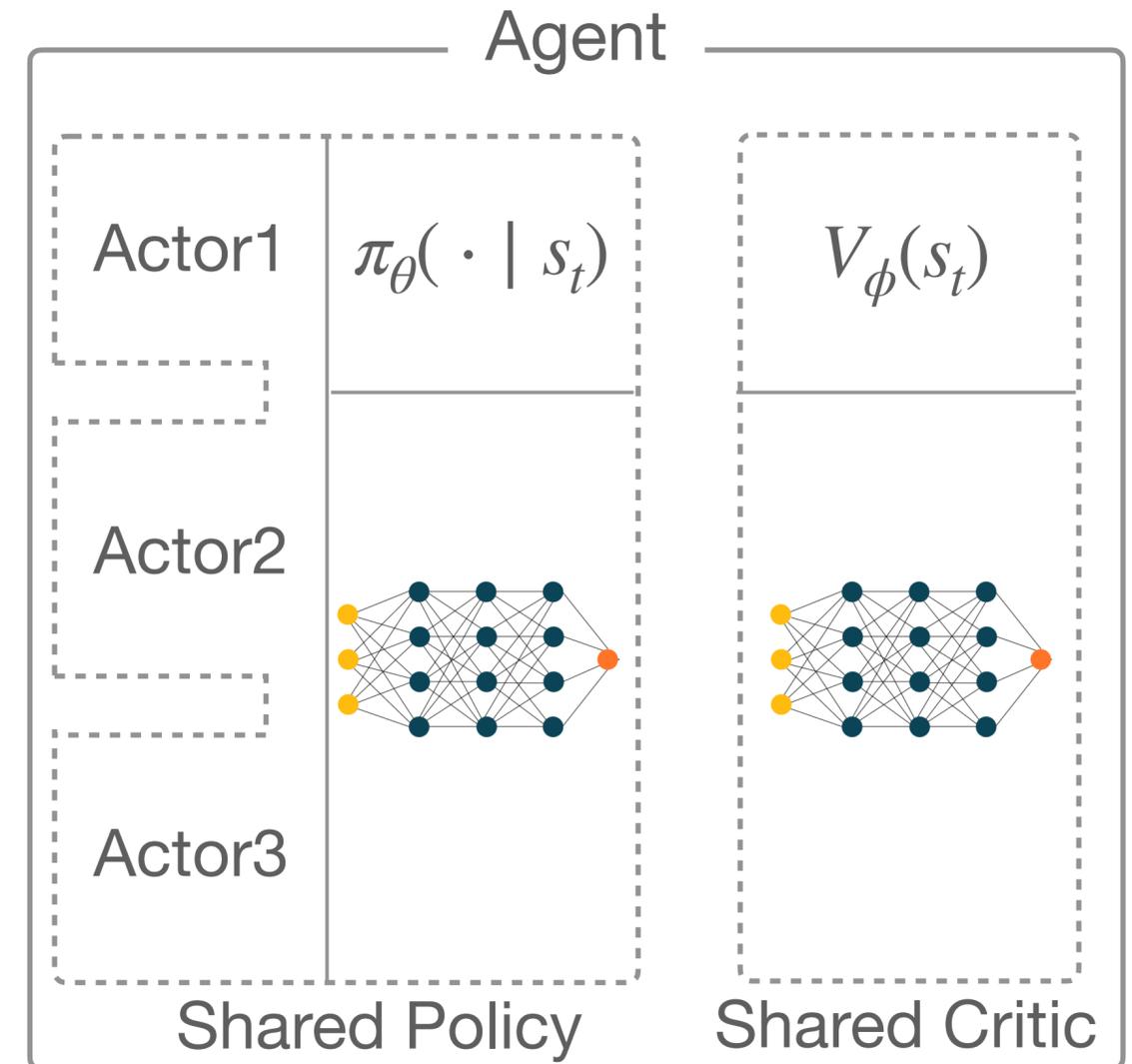$V_{\phi 2}(s_t)$ — Critic

Agent3
$\pi_{\theta 3}(\cdot \mid s_t)$ — Actor
$V_{\phi}(s_t)$ — Critic

nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

Agent
Actor1
Actor2
Actor3
$\pi_{\theta}(\cdot \mid s_t)$
$V_{\phi}(s_t)$
Shared Policy
Shared Critic

✅ Positive transfer of learned features across tasks
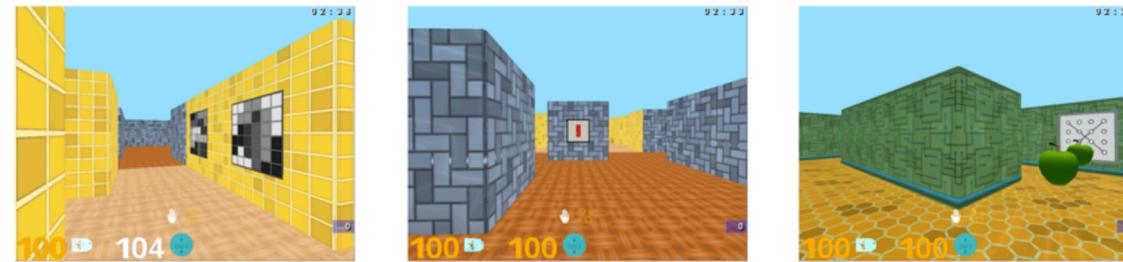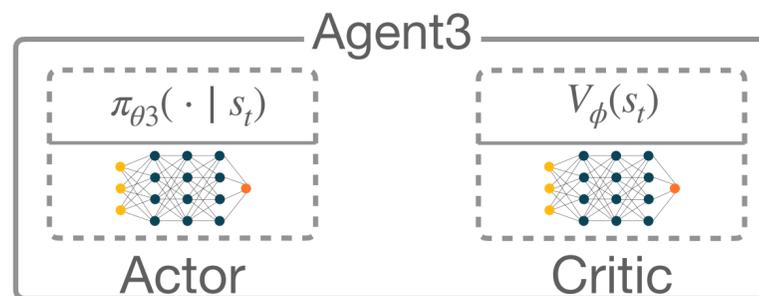
✅ Better compute efficiency

❌ Compatible action and state space

# Motivation | Multi-Task RL



Agent1
$\pi_{\theta 1}(\,\cdot\mid s_t)$
$V_{\phi 1}(s_t)$
Actor
Critic

Agent2
$\pi_{\theta 2}(\,\cdot\mid s_t)$
$V_{\phi 2}(s_t)$
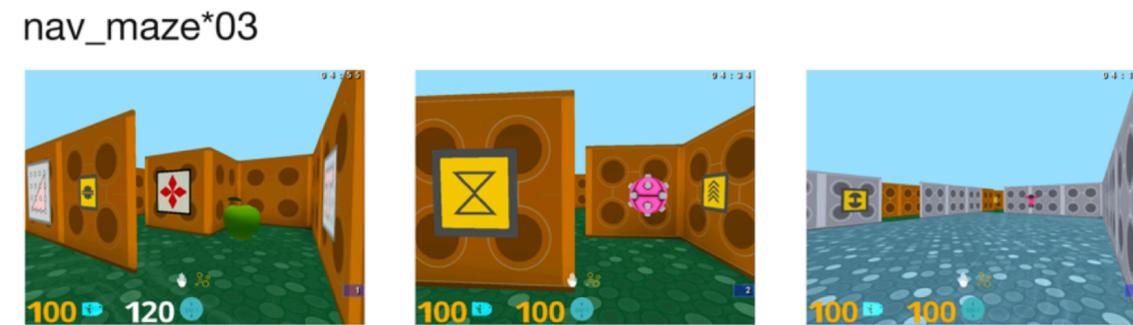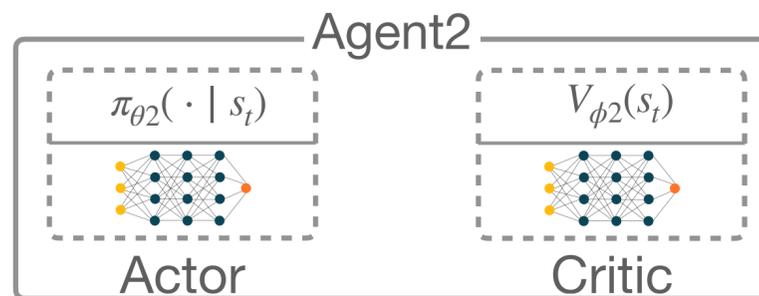Actor
Critic

Agent3
$\pi_{\theta 3}(\,\cdot\mid s_t)$
$V_{\phi}(s_t)$
Actor
Critic

nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

Agent

Actor1
Actor2
Actor3

$\pi_{\theta}(\,\cdot\mid s_t)$

$V_{\phi}(s_t)$

Shared Policy

Shared Critic

✅ Positive transfer of learned features across tasks

✅ Better compute efficiency

❌ Compatible action and state space

# Motivation | Multi-Task RL



nav_maze*02

nav_maze*03

random_maze

[Beattie et al., 2016]

Agent

Actor1

Actor2

Actor3

$\pi_\theta( \cdot \mid s_t)$

$V_\phi(s_t)$

Shared Policy

Shared Critic

Agent1

$\pi_{\theta 1}( \cdot \mid s_t)$

$V_{\phi 1}(s_t)$

Actor

Critic

Agent2

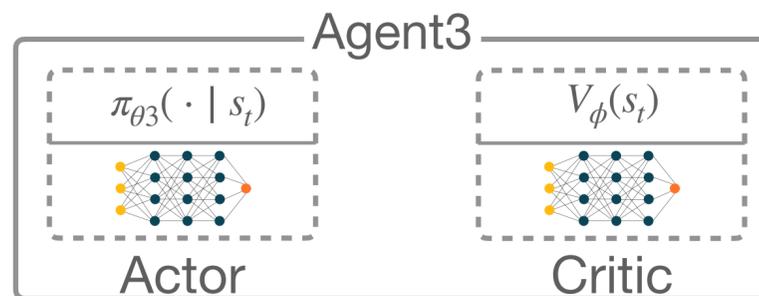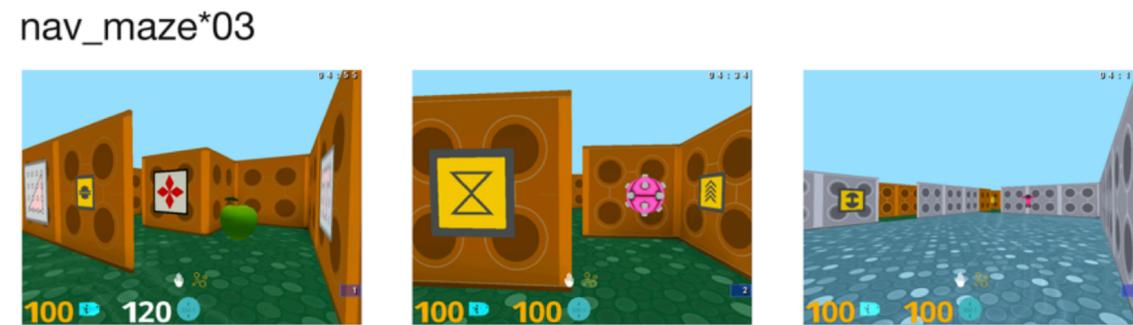$\pi_{\theta 2}( \cdot \mid s_t)$

$V_{\phi 2}(s_t)$

Actor

Critic

Agent3

$\pi_{\theta 3}( \cdot \mid s_t)$

$V_{\phi}(s_t)$

Actor

Critic

✅ Positive transfer of learned features across tasks

✅ Better compute efficiency

❌ Compatible action and state space

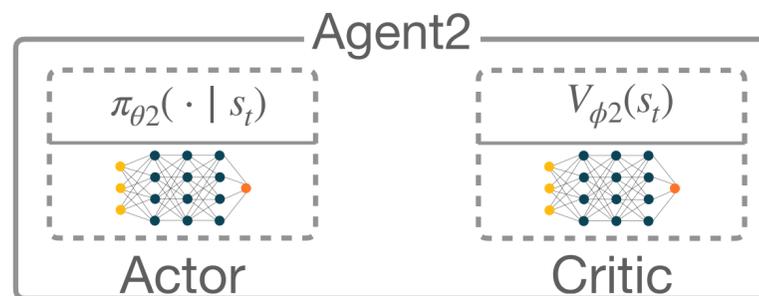One agent learns to perform optimally across several different tasks sharing similar characteristics

# Motivation | Multi-Task RL | Division of Labor

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# Motivation | Multi-Task RL | Division of Labor

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

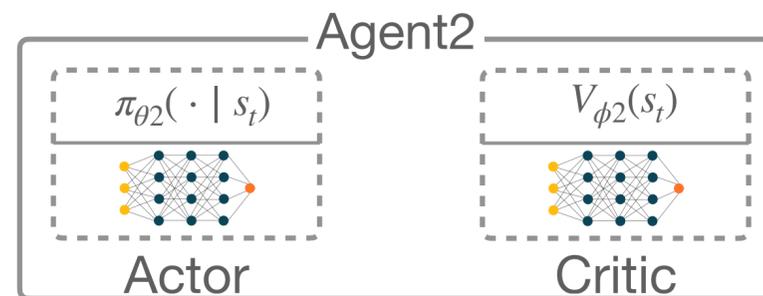$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# Motivation | Multi-Task RL | Division of Labor

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

A1 $\theta_1 \; \phi_1$

A2 $\theta_2 \; \phi_2$

Learner $\theta^\star \; \phi^\star$

A4 $\theta_4 \; \phi_4$

A3 $\theta_3 \; \phi_3$

Initial

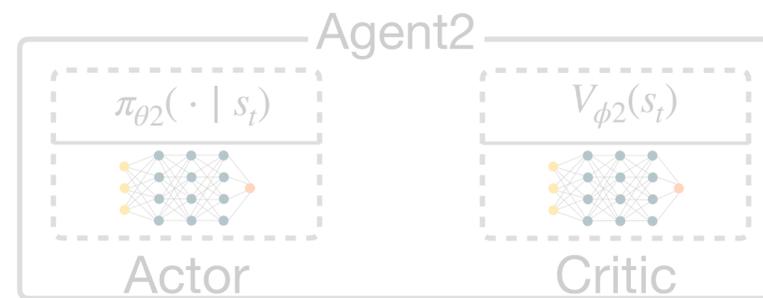$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Division of Labor

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\,\cdot\,|\,s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\,.\,step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$$

Increasing rollout time

A1 $\theta_1\ \phi_1$

A2 $\theta_2\ \phi_2$

Learner $\theta^\star\ \phi^\star$

A4 $\theta_4\ \phi_4$

A3 $\theta_3\ \phi_3$

Initial
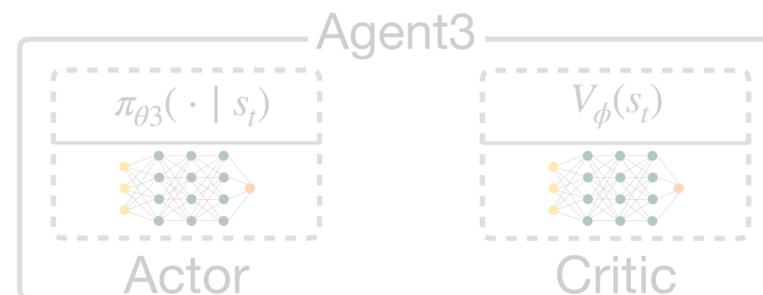$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star$$
$$\forall i \in [1,4]$$

# Motivation | Multi-Task RL | Division of Labor

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$
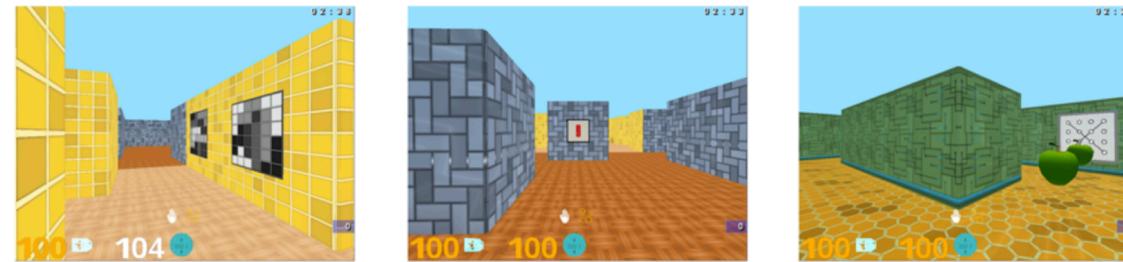
**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Increasing rollout time

A1 $\theta_1 \ \phi_1$

A2 $\theta_2 \ \phi_2$

Learner $\theta^\star \ \phi^\star$

A4 $\theta_4 \ \phi_4$

A3 $\theta_3 \ \phi_3$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

Update global parameters for individual actor rollout completions

OR

Update global parameters for combined actor rollout completions

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

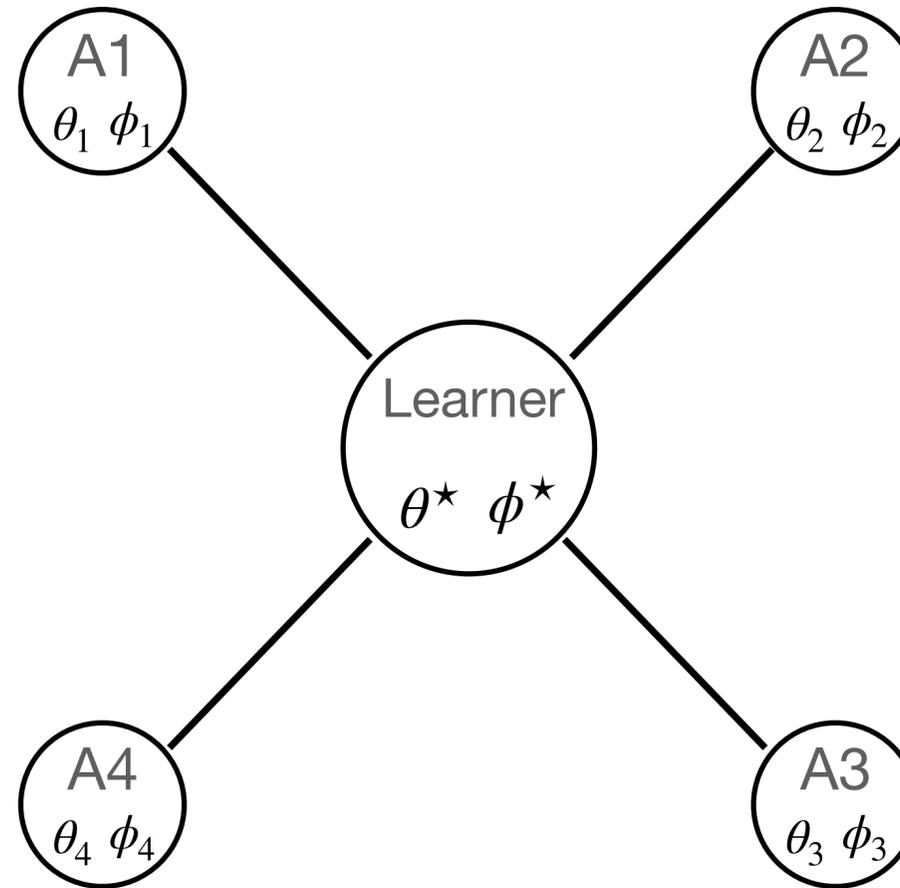**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Increasing rollout time

A1 $\theta_1 \ \phi_1$

A2 $\theta_2 \ \phi_2$

Learner $\theta^\star \ \phi^\star$

A4 $\theta_4 \ \phi_4$

A3 $\theta_3 \ \phi_3$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \qquad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
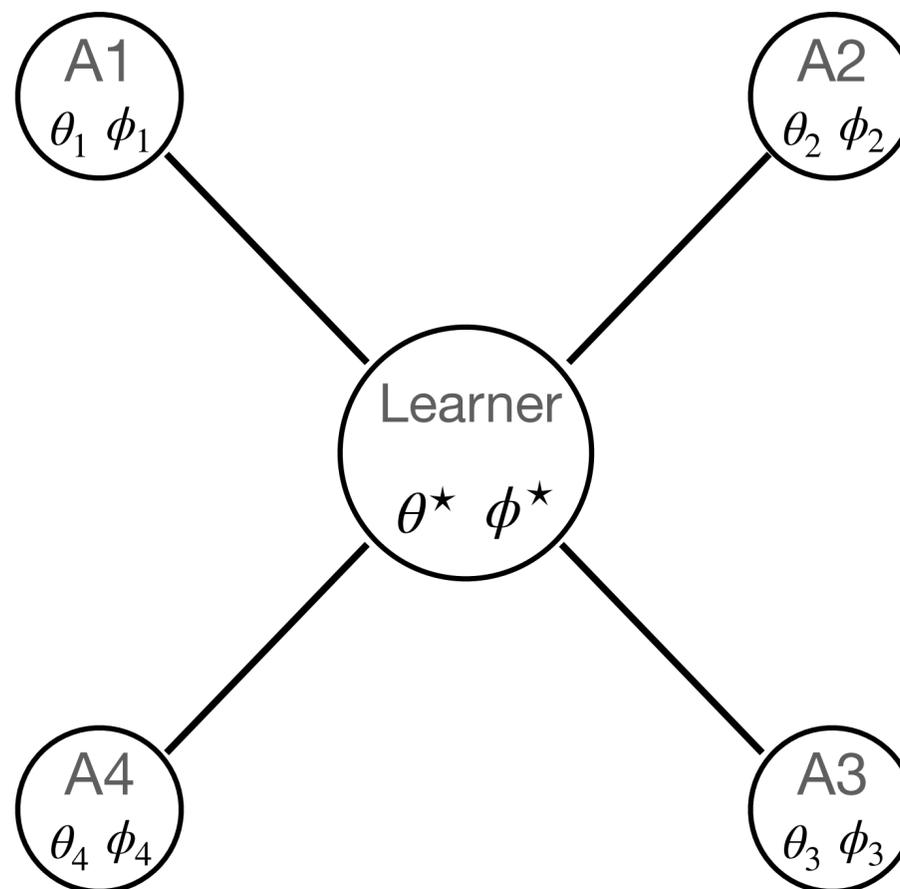
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Increasing rollout time

A1 $\theta_1 \ \phi_1$
A2 $\theta_2 \ \phi_2$
A3 $\theta_3 \ \phi_3$
A4 $\theta_4 \ \phi_4$

Learner $\theta^\star \ \phi^\star$

Initial
$\theta_i \leftarrow \theta^\star$
$\phi_i \leftarrow \phi^\star$ $\quad \forall i \in [1,4]$

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
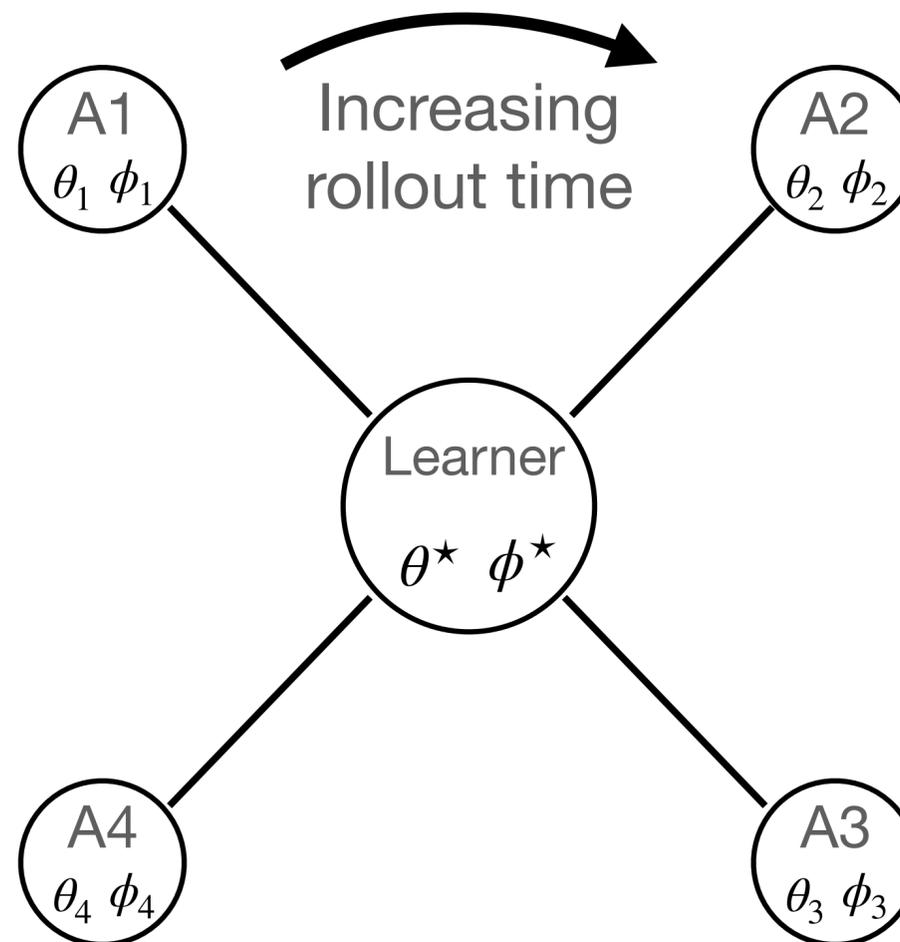
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Increasing rollout time

A1 $\theta_1 \ \phi_1$

A2 $\theta_2 \ \phi_2$

A4 $\theta_4 \ \phi_4$

A3 $\theta_3 \ \phi_3$

Learner $\theta^\star \ \phi^\star$

Initial

$\theta_i \leftarrow \theta^\star$

$\phi_i \leftarrow \phi^\star$ $\quad \forall i \in [1,4]$

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \cdot step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
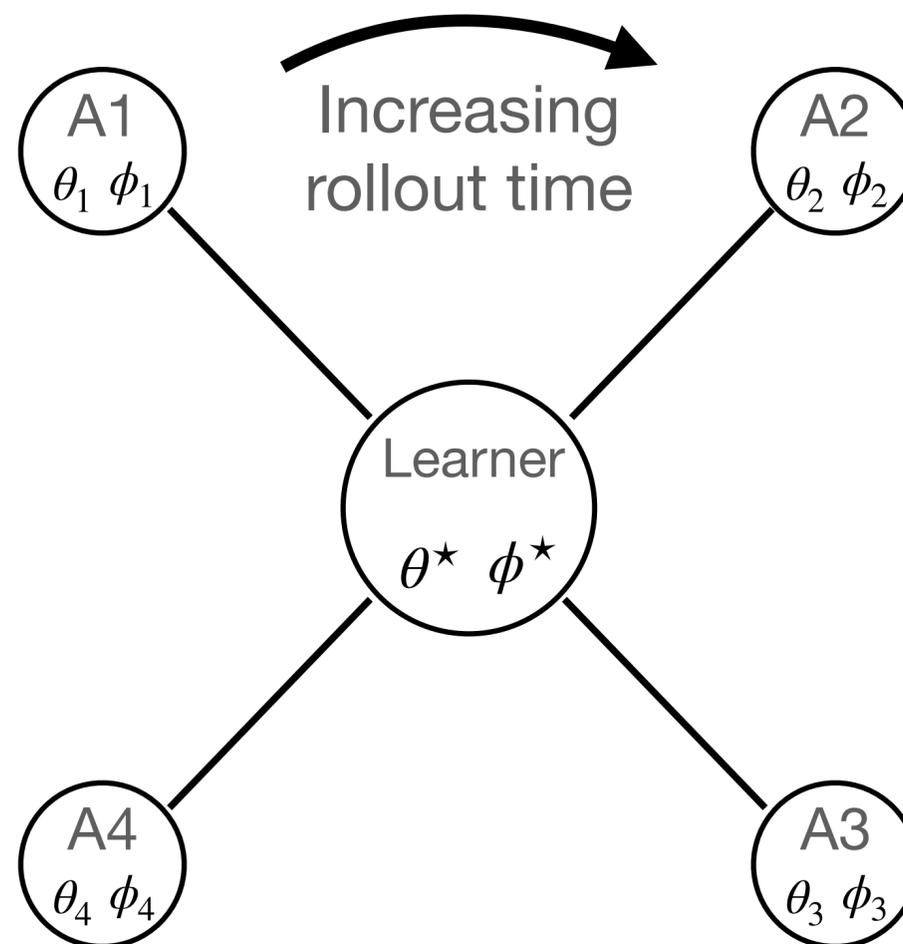
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \boxed{\delta_t \nabla_\phi V_\phi(s_t)}$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \boxed{\delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)}$$

Increasing rollout time

A1 $\theta_1 \phi_1$  A2 $\theta_2 \phi_2$  A3 $\theta_3 \phi_3$  A4 $\theta_4 \phi_4$

Learner $\theta^\star \phi^\star$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
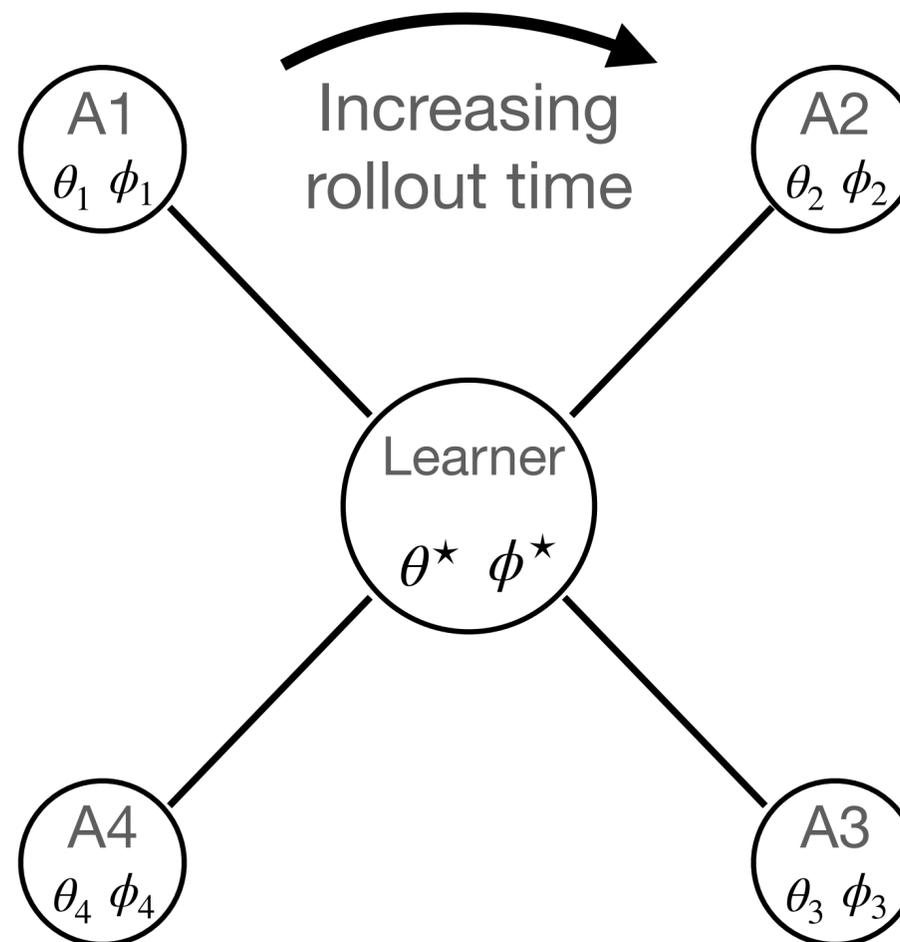
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

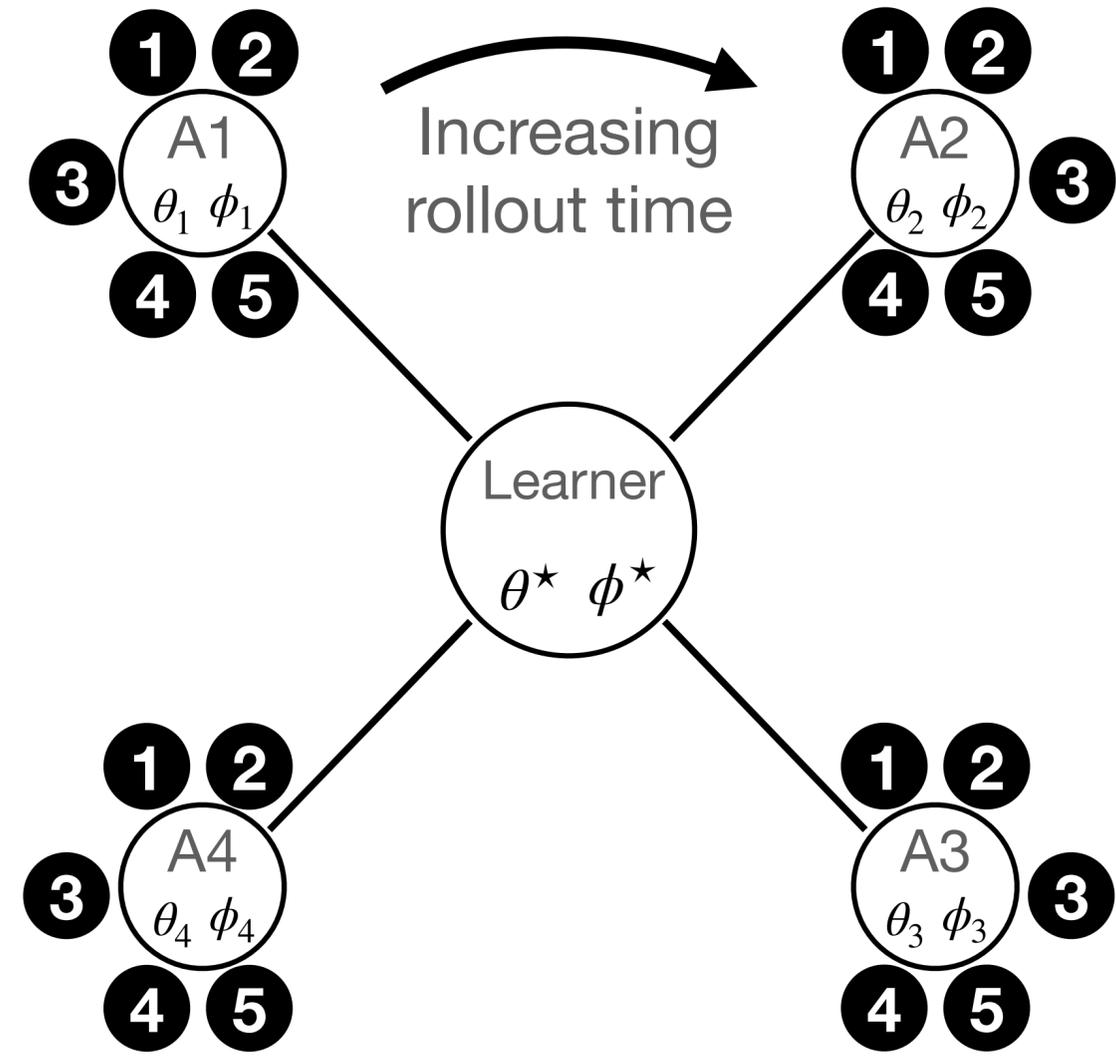**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \boxed{\delta_t \nabla_\phi V_\phi(s_t)}$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \boxed{\delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)}$$



Increasing rollout time

Learner $\theta^\star \ \phi^\star$

A1 $\theta_1 \ \phi_1$

A2 $\theta_2 \ \phi_2$

A3 $\theta_3 \ \phi_3$

A4 $\theta_4 \ \phi_4$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)\Big|_{\phi_1, \theta_1}$$

Global update for A1

# Motivation | Multi-Task RL | Asynchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

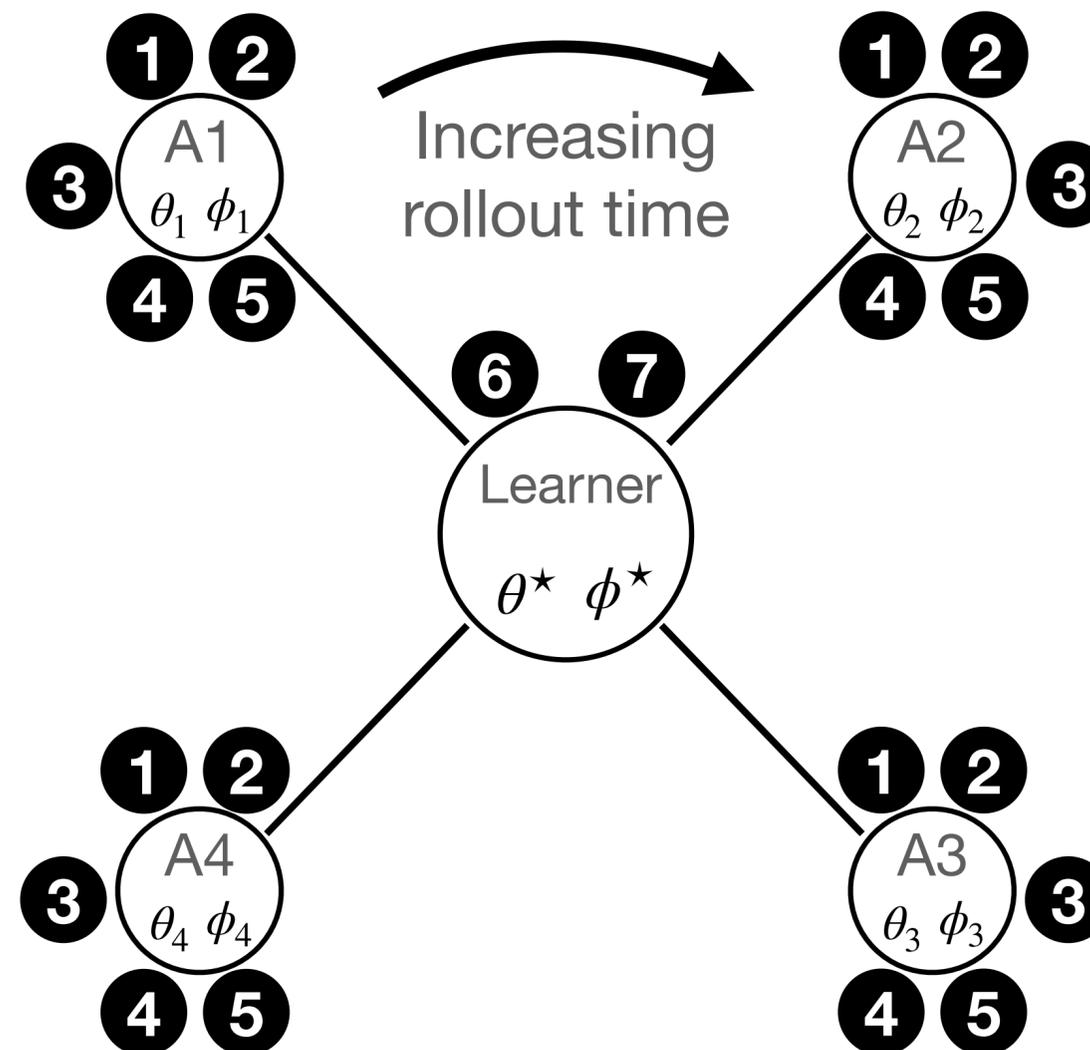**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \boxed{\delta_t \nabla_\phi V_\phi(s_t)}$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \boxed{\delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)}$$



Increasing rollout time

Learner $\theta^\star \ \phi^\star$

A1 $\theta_1 \ \phi_1$  A2 $\theta_2 \ \phi_2$  A3 $\theta_3 \ \phi_3$  A4 $\theta_4 \ \phi_4$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi_1, \theta_1}$$

$$\phi_2^\star \leftarrow \phi_1^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_2^\star \leftarrow \theta_1^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi_2, \theta_2}$$

$\cdots$

Global update for A1          Global update for A2

# Motivation | Multi-Task RL | Asynchronous Update

**❶** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\ \cdot\ |\ s_t)$$

**❷** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\ .\ step(a_t, s_t)$$

**❸** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
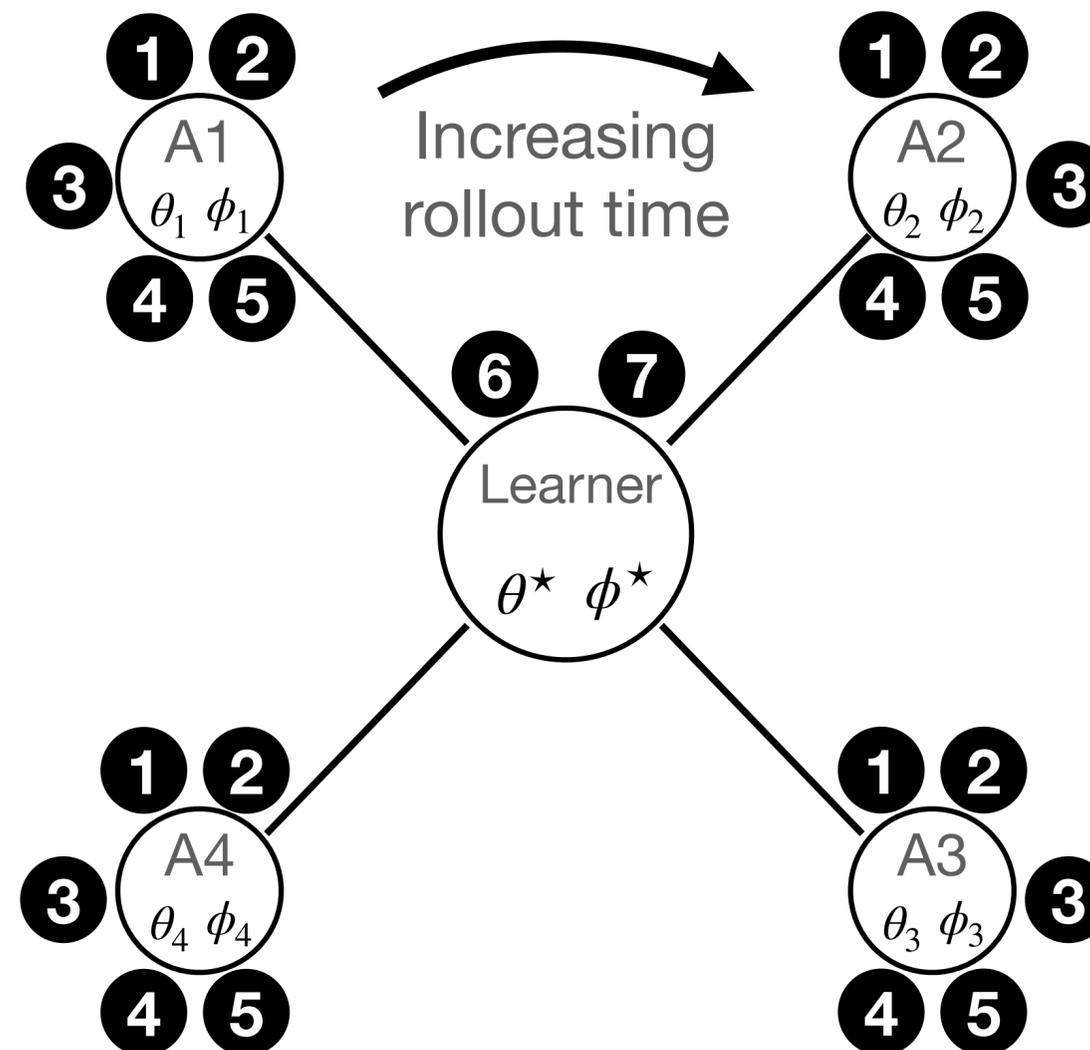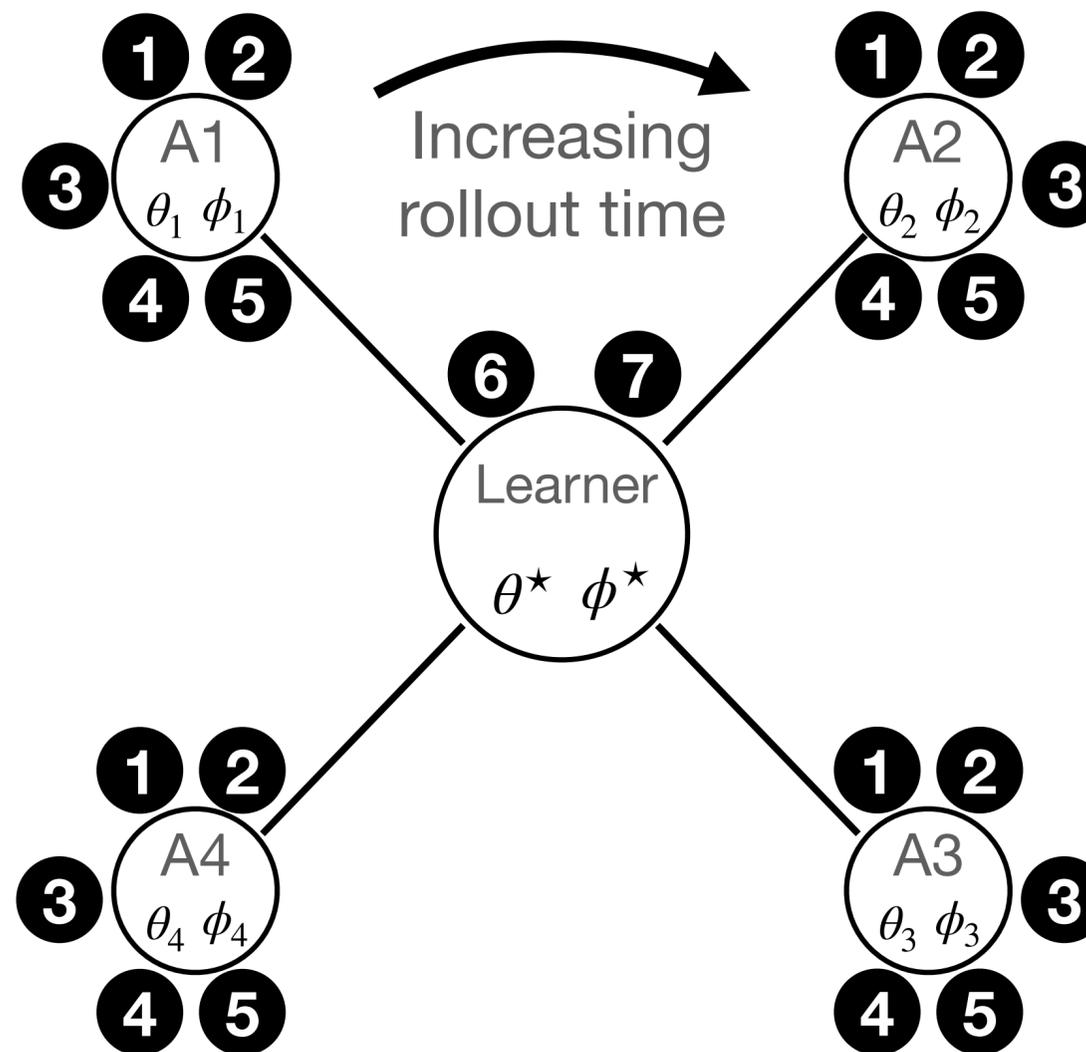
**❹** Value Gradient $\nabla_\phi V_\phi(s_t)$

**❺** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t\ |\ s_t)$

**❻** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \boxed{\delta_t \nabla_\phi V_\phi(s_t)}$$

**❼** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \boxed{\delta_t \nabla_\theta \log \pi_\theta(a_t\ |\ s_t)}$$



**❶ ❷**
A1
$\theta_1\ \phi_1$
**❸**
**❹ ❺**

Increasing rollout time

**❶ ❷**
A2
$\theta_2\ \phi_2$ **❸**
**❹ ❺**

**❻ ❼**

Learner
$\theta^\star\ \phi^\star$

**❶ ❷**
A4
$\theta_4\ \phi_4$
**❸**
**❹ ❺**

**❶ ❷**
A3
$\theta_3\ \phi_3$ **❸**
**❹ ❺**

A3C
(Mnih et al., 2016)

Initial

$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t\ |\ s_t) \Big|_{\phi_1, \theta_1}$$

$$\phi_2^\star \leftarrow \phi_1^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_2^\star \leftarrow \theta_1^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t\ |\ s_t) \Big|_{\phi_2, \theta_2}$$

$\cdots$

Global update for A1

Global update for A2

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \, \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \,.\, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

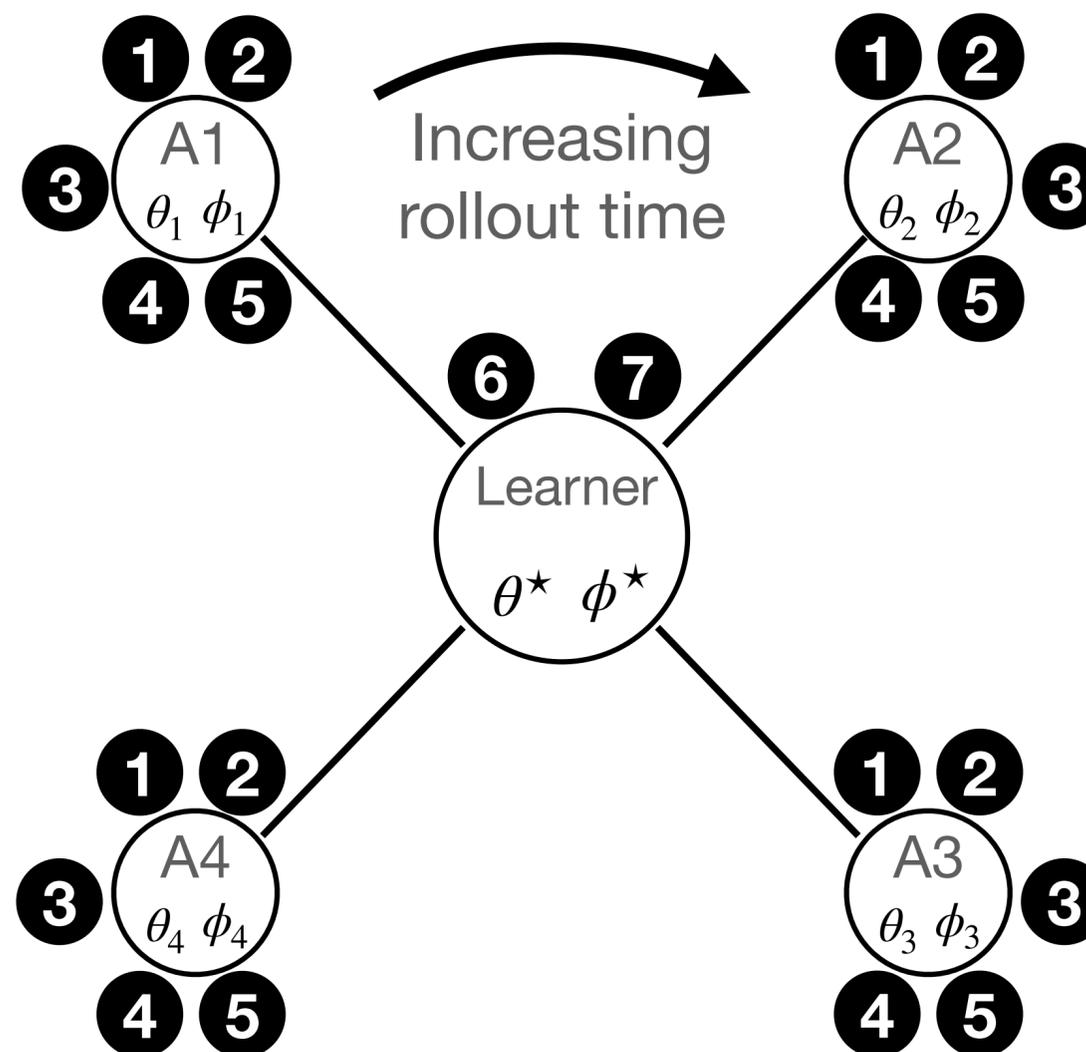**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

A1 $\theta_1 \; \phi_1$

Increasing rollout time

A2 $\theta_2 \; \phi_2$

Learner $\theta^\star \; \phi^\star$

A4 $\theta_4 \; \phi_4$

A3 $\theta_3 \; \phi_3$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

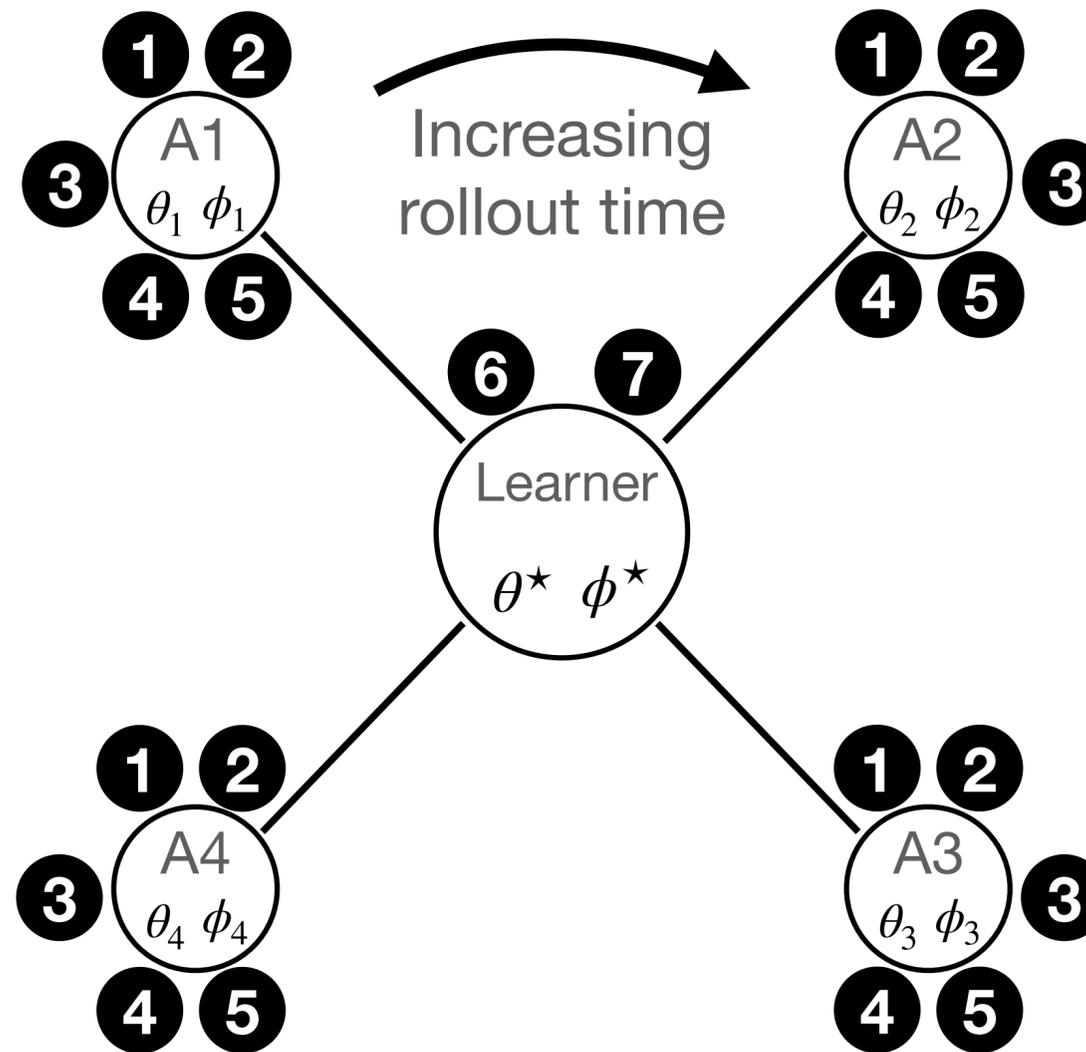**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Increasing rollout time

A1 $\theta_1 \ \phi_1$   A2 $\theta_2 \ \phi_2$   A4 $\theta_4 \ \phi_4$   A3 $\theta_3 \ \phi_3$   Learner $\theta^\star \ \phi^\star$

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\,\cdot\,|\,s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
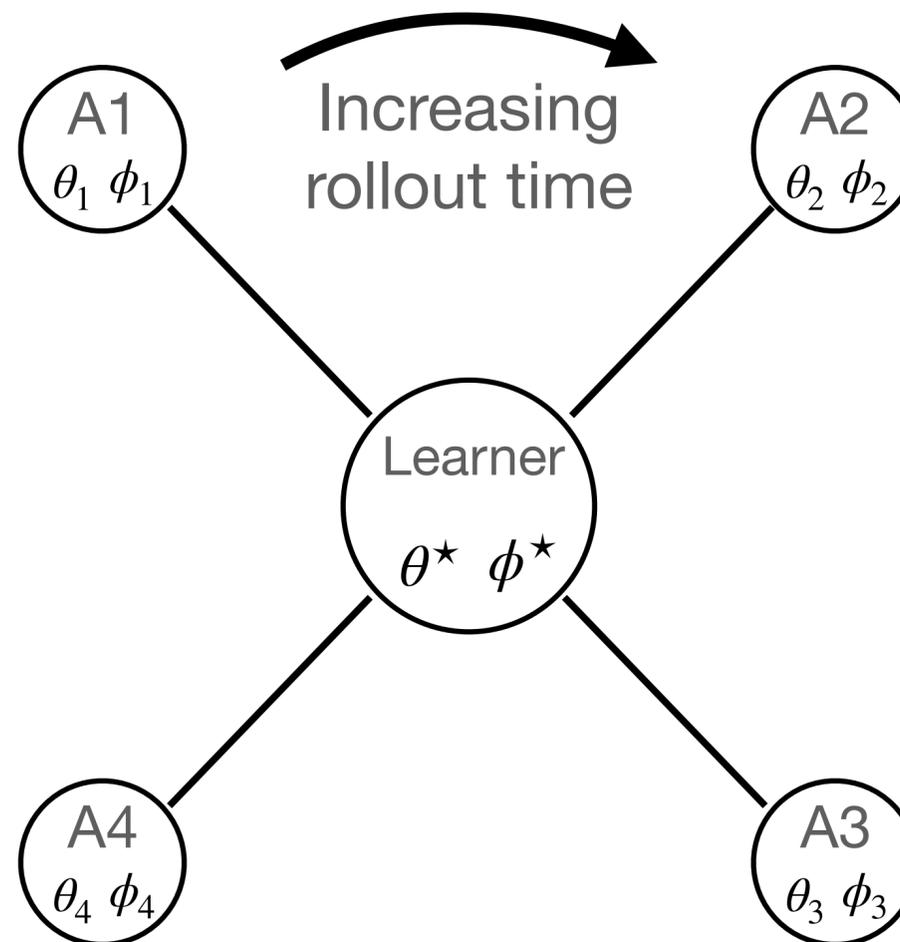
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$$



Increasing rollout time

**1** **2** A1 $\theta_1$ $\phi_1$  **1** **2** A2 $\theta_2$ $\phi_2$

**3** **4** **5** Learner $\theta^\star$ $\phi^\star$

**6** **7**

**1** **2** A4 $\theta_4$ $\phi_4$  **1** **2** A3 $\theta_3$ $\phi_3$

Initial

$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

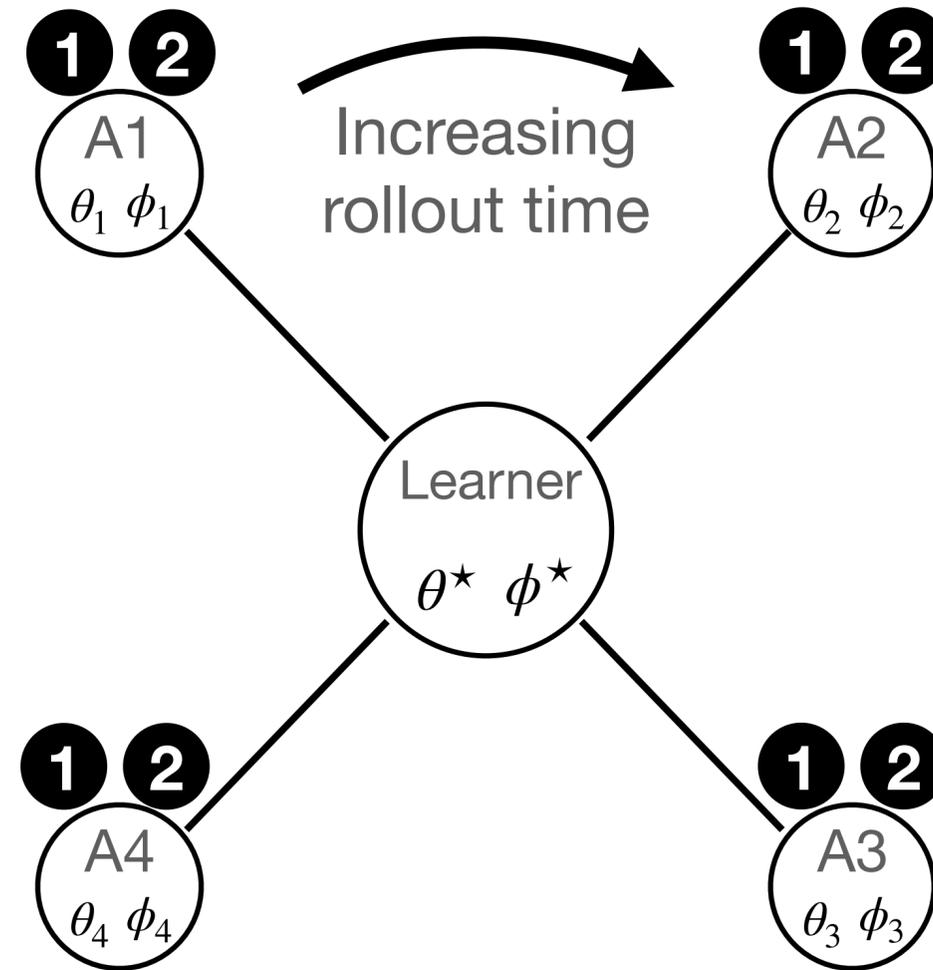**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

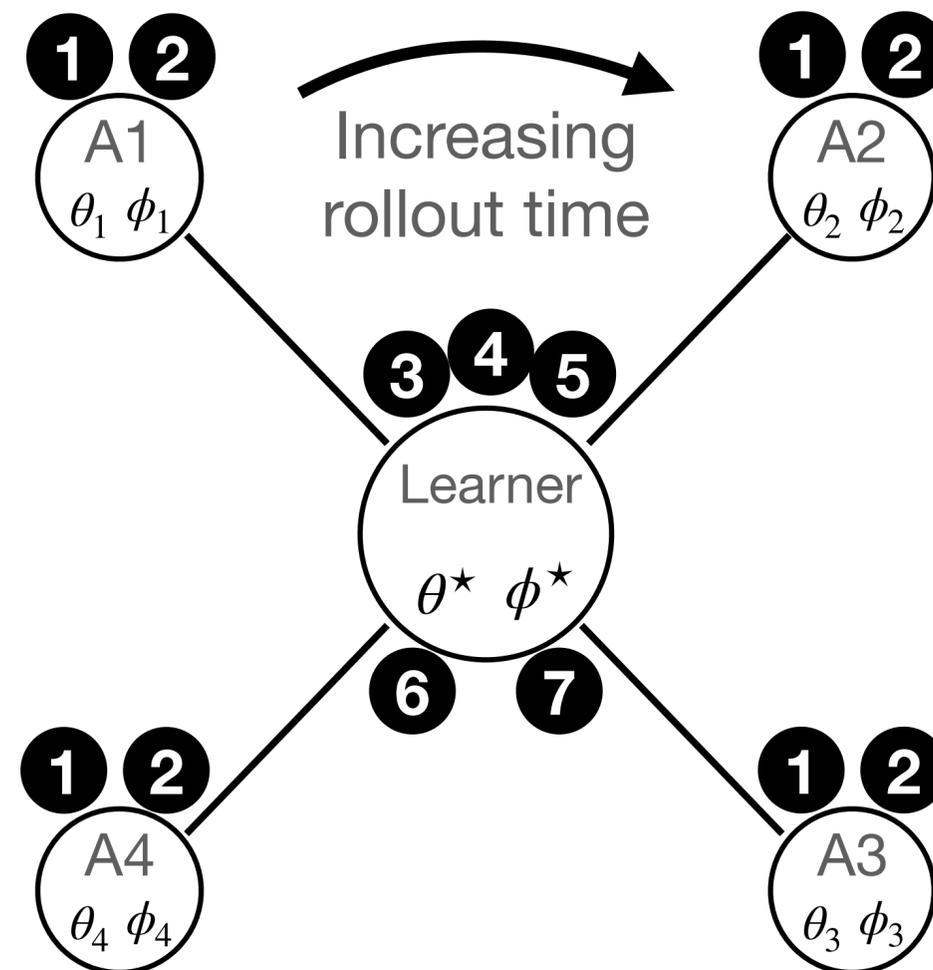**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Increasing rollout time

**1** **2** A1 $\theta_1$ $\phi_1$

**1** **2** A2 $\theta_2$ $\phi_2$

**3** **4** **5**

Learner $\theta^\star$ $\phi^\star$

**6** **7**

**1** **2** A4 $\theta_4$ $\phi_4$

**1** **2** A3 $\theta_3$ $\phi_3$

{ **Batched** **Trajectory** }

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star$$
$$\forall i \in [1,4]$$

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi^\star, \theta^\star}$$

Single Global Update

# Motivation | Multi-Task RL | Synchronous Update

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \cdot step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
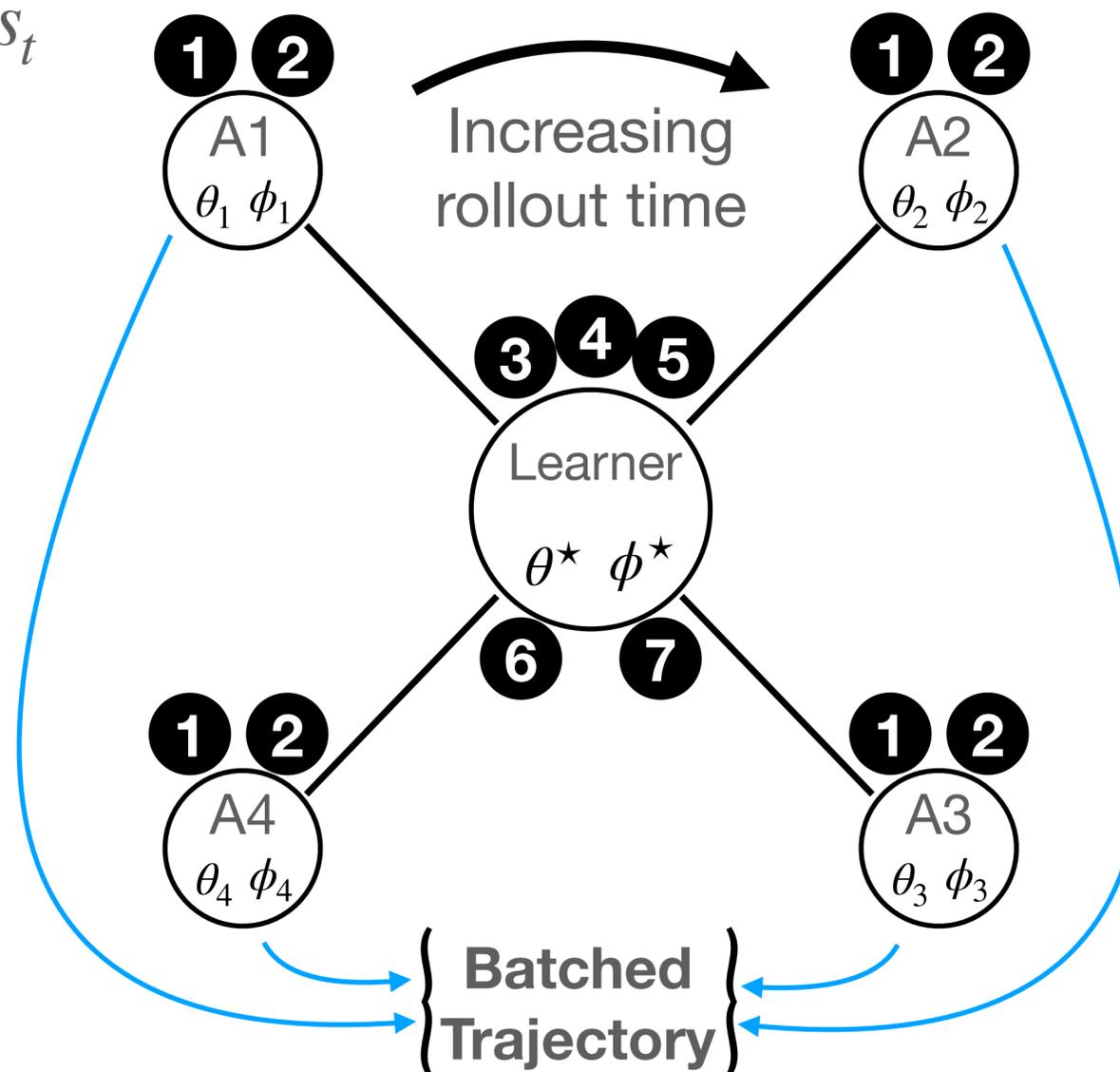
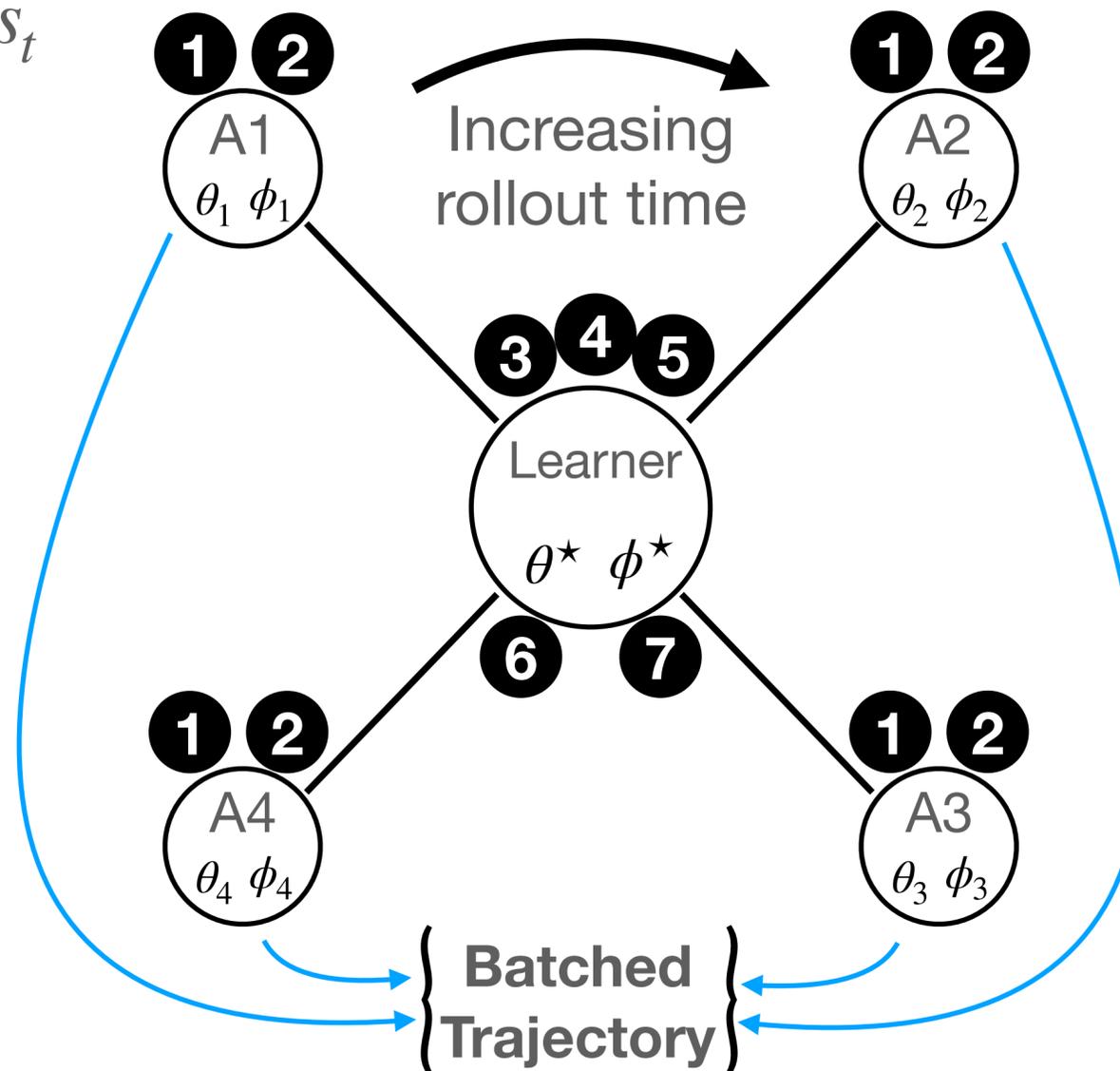**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



A2C
(Clemente et al., 2017)

Initial
$$\theta_i \leftarrow \theta^\star$$
$$\phi_i \leftarrow \phi^\star \quad \forall i \in [1,4]$$

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi^\star, \theta^\star}$$

Single Global Update

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



## A2C (Synchronous)

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



## A2C (Synchronous)



✅ Global update per actor rollout

❌ Global update bottlenecked by slowest actor

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



## A2C (Synchronous)



✅ Global update per actor rollout

❌ Local gradient computation - GPU underutilized

❌ Global update bottlenecked by slowest actor

✅ Batched gradient computation - GPU scalable

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



✅ Global update per actor rollout

❌ Local gradient computation - GPU underutilized

❌ Bandwidth Limited - Transfer massive gradients

## A2C (Synchronous)



❌ Global update bottlenecked by slowest actor

✅ Batched gradient computation - GPU scalable

✅ Better bandwidth utilization - Transfer trajectories

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



## A2C (Synchronous)



| | A3C (Asynchronous) | | A2C (Synchronous) |
|---|---|---|---|
| ✅ | Global update per actor rollout | ❌ | Global update bottlenecked by slowest actor |
| ❌ | Local gradient computation - GPU underutilized | ✅ | Batched gradient computation - GPU scalable |
| ❌ | Bandwidth Limited - Transfer massive gradients | ✅ | Better bandwidth utilization - Transfer trajectories |
| ❌ | Unstable - Policy lag due to asynchronous updates | ✅ | Stable - No Policy Lag |

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



## A2C (Synchronous)



| | A3C | | A2C |
|---|---|---|---|
| ✅ | Global update per actor rollout | ❌ | Global update bottlenecked by slowest actor |
| ❌ | Local gradient computation - GPU underutilized | ✅ | Batched gradient computation - GPU scalable |
| ❌ | Bandwidth Limited - Transfer massive gradients | ✅ | Better bandwidth utilization - Transfer trajectories |
| ❌ | Unstable - Policy lag due to asynchronous updates | ✅ | Stable - No Policy Lag |

$$\phi_1^\star \leftarrow \phi^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_1^\star \leftarrow \theta^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi_1, \theta_1}$$

$$\phi_2^\star \leftarrow \phi_1^\star + \delta_t \nabla_\phi V_\phi(s_t) \Big|$$
$$\theta_2^\star \leftarrow \theta_1^\star + \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \Big|_{\phi_2, \theta_2}$$

$\bullet \bullet \bullet$

Global update for A1          Global update for A2

# Motivation | Multi-Task RL | Trade-Offs

## A3C (Asynchronous)



✅ Global update per actor rollout

❌ Local gradient computation - GPU underutilized

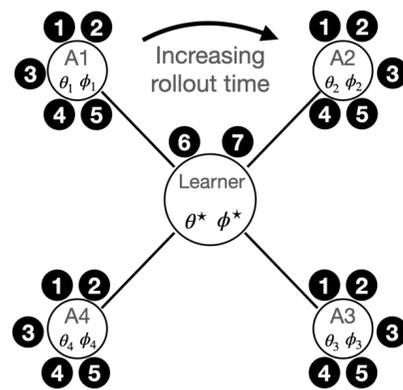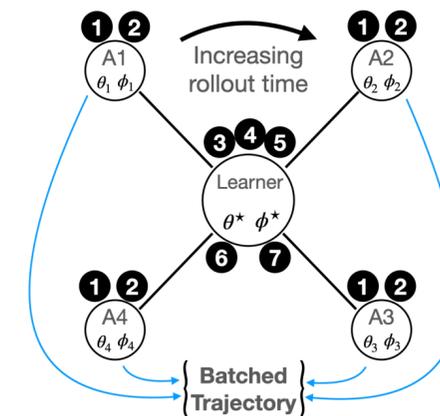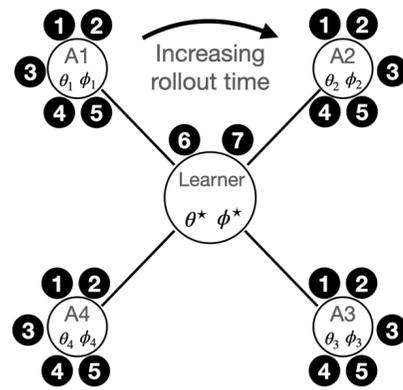❌ Bandwidth Limited - Transfer massive gradients

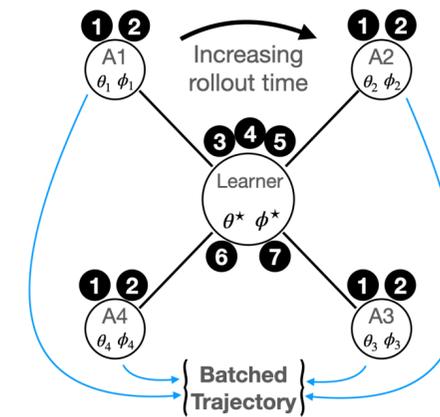❌ Unstable - Policy lag due to asynchronous updates

## A2C (Synchronous)



❌ Global update bottlenecked by slowest actor

✅ Batched gradient computation - GPU scalable

✅ Better bandwidth utilization - Transfer trajectories

✅ Stable - No Policy Lag

$$\phi_1^{\star} \leftarrow \phi^{\star} + \delta_t \nabla_{\phi} V_{\phi}(s_t) \Big|$$
$$\theta_1^{\star} \leftarrow \theta^{\star} + \delta_t \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \Big|_{\phi_1, \theta_1}$$

Global update for A1

$$\phi_2^{\star} \leftarrow \phi_1^{\star} + \delta_t \nabla_{\phi} V_{\phi}(s_t) \Big|$$
$$\theta_2^{\star} \leftarrow \theta_1^{\star} + \delta_t \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \Big|_{\phi_2, \theta_2}$$

Global update for A2

$\bullet\ \bullet\ \bullet$

Gradients computed for $\phi^{\star}, \theta^{\star}$

But update applied to $\phi_1^{\star}, \theta_1^{\star}$

# IMPALA | V-Trace

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# IMPALA | V-Trace

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\,\cdot\,|\,s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\,.\,step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \,|\, s_t)$$

# IMPALA | V-Trace

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \cdot step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

**1** **2**  A1 $\theta_1 \, \phi_1$

Increasing rollout time

**1** **2**  A2 $\theta_2 \, \phi_2$

**3** **4** **5**

Learner $\theta^\star \, \phi^\star$

**6** **7**

**1** **2**  A4 $\theta_4 \, \phi_4$

**1** **2**  A3 $\theta_3 \, \phi_3$

Looks similar to A2C but…..

$$\theta_i \neq \theta^\star \qquad \phi_i \neq \phi^\star$$

# IMPALA | V-Trace

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \cdot step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Increasing rollout time

**A1** $\theta_1\ \phi_1$

**A2** $\theta_2\ \phi_2$

**Learner** $\theta^\star\ \phi^\star$

**A4** $\theta_4\ \phi_4$

**A3** $\theta_3\ \phi_3$

Looks similar to A2C but…..

$$\theta_i \neq \theta^\star \qquad \phi_i \neq \phi^\star$$

We have trajectory from $\theta_i$ but we need to update $\theta^\star$

# IMPALA | V-Trace

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\,\cdot\,|\,s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\,.\,step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

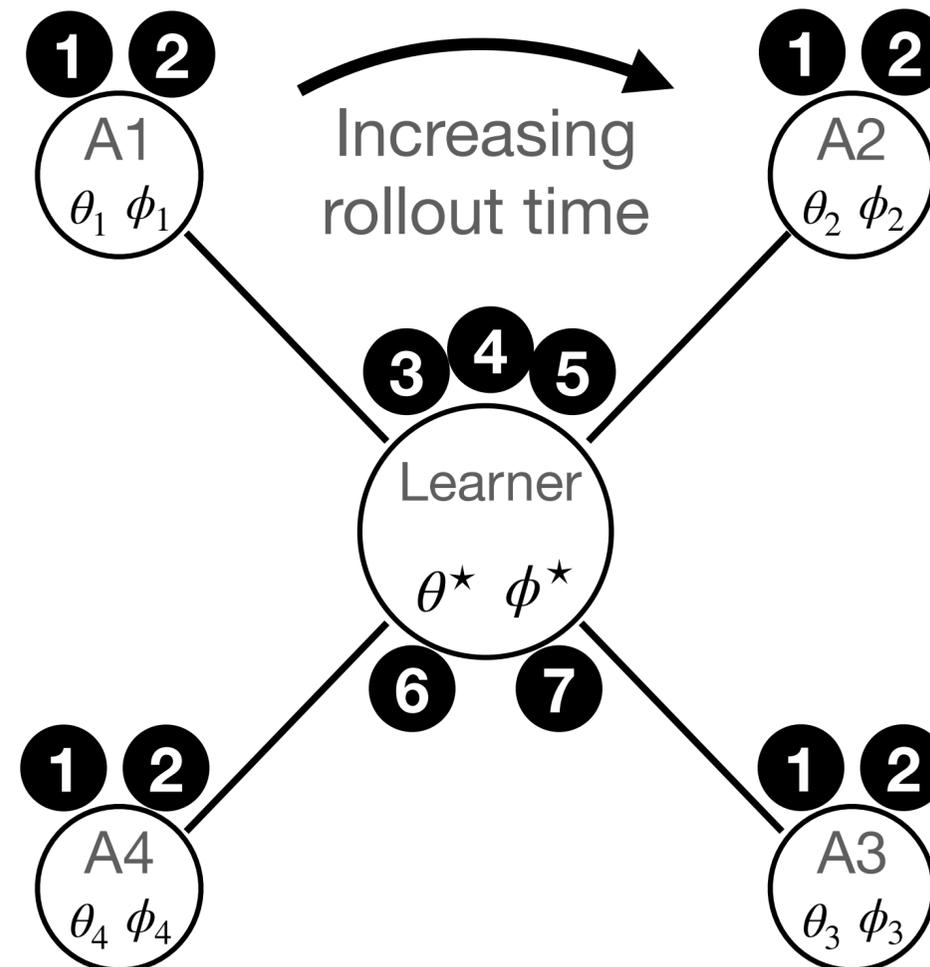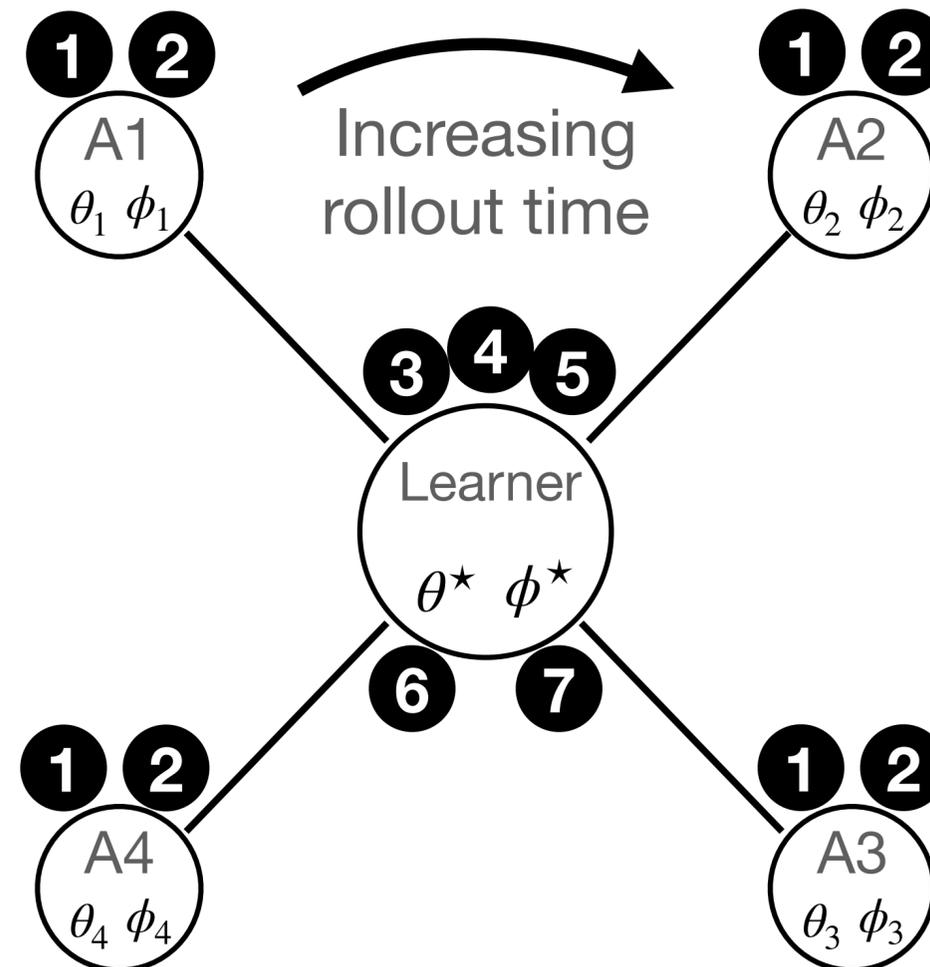**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$



Increasing rollout time

A1 $\theta_1 \ \phi_1$

A2 $\theta_2 \ \phi_2$

Learner $\theta^\star \ \phi^\star$

A4 $\theta_4 \ \phi_4$

A3 $\theta_3 \ \phi_3$

Looks similar to A2C but…..

$$\theta_i \neq \theta^\star \qquad \phi_i \neq \phi^\star$$

We have trajectory from $\theta_i$ but we need to update $\theta^\star$

How do we do this?

# CHANGE OF VARIABLE – IMPORTANCE SAMPLING

Given two densities p and q:

$$\mathbb{E}_{x\sim p}[f(x)] = \sum_{x\in\mathcal{X}} p(x)f(x) = \sum_{x\in\mathcal{X}} q(x)\frac{p(x)}{q(x)}f(x) = \mathbb{E}_{x\sim q}\left[\frac{p(x)}{q(x)}f(x)\right]$$

Example: Single transition following $\pi_{\theta'}$, but sampled from $\pi_\theta$:

$$J(\pi_{\theta'}) = \mathbb{E}_{(s,a,r)\sim\pi_{\theta'}}[r(\cdot,s)|s] = \sum_{a\in\mathcal{A}} \pi_{\theta'}(a|s)r(a,s) = \sum_{a\in\mathcal{A}} \pi_\theta(a|s)\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)}r(a,s) = \mathbb{E}_{(s,a,r)\sim\pi_\theta}[\frac{\pi_{\theta'}(\cdot|s)}{\pi_\theta(\cdot|s)}r(\cdot,s)|s]$$

$\pi_\theta(a \mid s)$   r = r(a, s)

s

## CHANGE OF VARIABLE – IMPORTANCE SAMPLING

Given two densities p and q:

$$\mathbb{E}_{x \sim p}[f(x)] = \sum_{x \in \mathcal{X}} p(x)f(x) = \sum_{x \in \mathcal{X}} q(x)\frac{p(x)}{q(x)}f(x) = \mathbb{E}_{x \sim q}\left[\frac{p(x)}{q(x)}f(x)\right]$$

We can borrow
**Importance Sampling**
from TRPO and PPO
(Kahn et al., 1953)

Example: Single transition following $\pi_{\theta'}$, but sampled from $\pi_\theta$:

$$J(\pi_{\theta'}) = \mathbb{E}_{(s,a,r) \sim \pi_{\theta'}}[r(\cdot, s)|s] = \sum_{a \in \mathcal{A}} \pi_{\theta'}(a|s)r(a,s) = \sum_{a \in \mathcal{A}} \pi_\theta(a|s)\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)}r(a,s) = \mathbb{E}_{(s,a,r) \sim \pi_\theta}[\frac{\pi_{\theta'}(\cdot|s)}{\pi_\theta(\cdot|s)}r(\cdot,s)|s]$$

$\pi_\theta(a \mid s)$  ● r = r(a, s)

●

⋮

s ,  ●

Lecture 16 - Slide 24

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \,.\, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
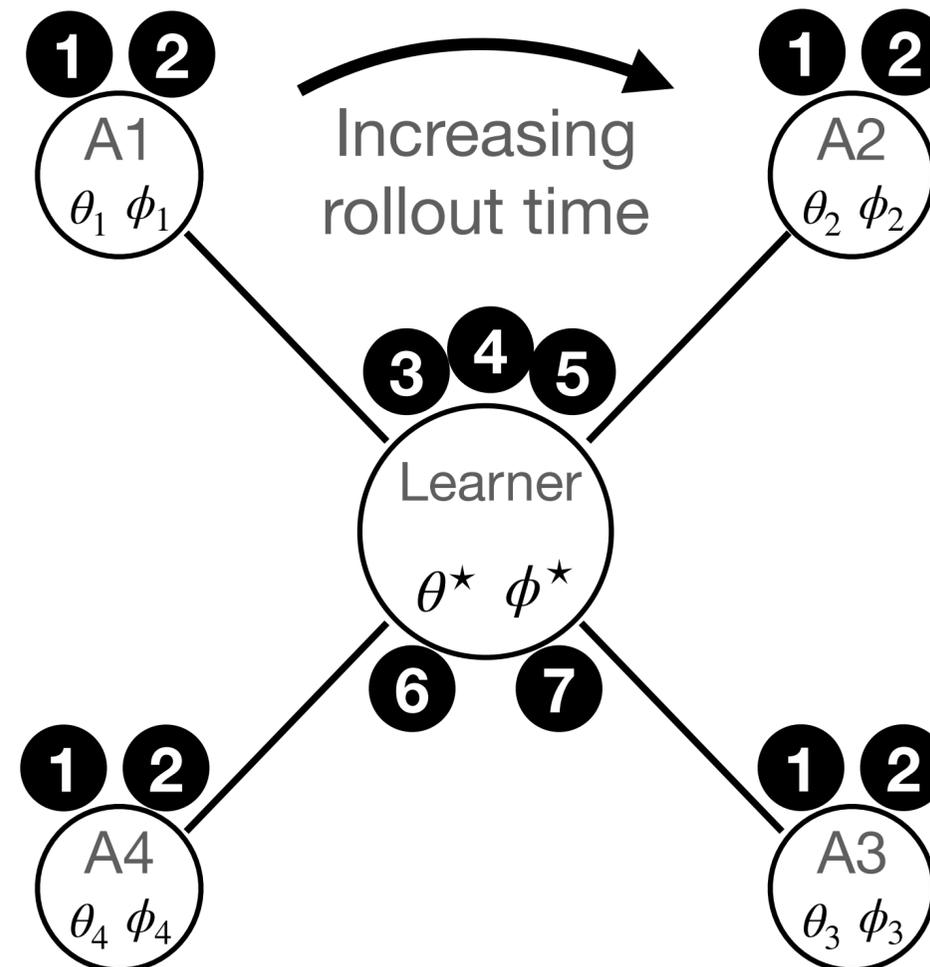
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\, \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\,.\,step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Bellman update

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \,.\, step(a_t, s_t)$$

Bellman update

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

Trajectory obtained from $\theta_1$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
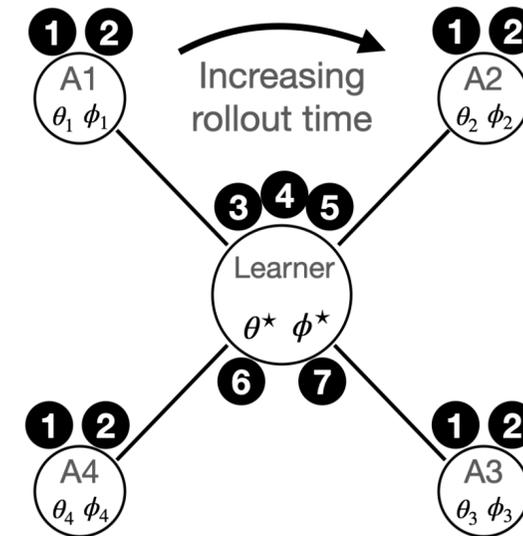
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Bellman update

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

Trajectory obtained from $\theta_1$

$$\rho_t \triangleq min(1, \frac{\pi_{\theta\star}(a_t \mid s_t)}{\pi_{\theta_1}(a_t \mid s_t)})$$

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta(\,\cdot\,|\,s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env\,.\,step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

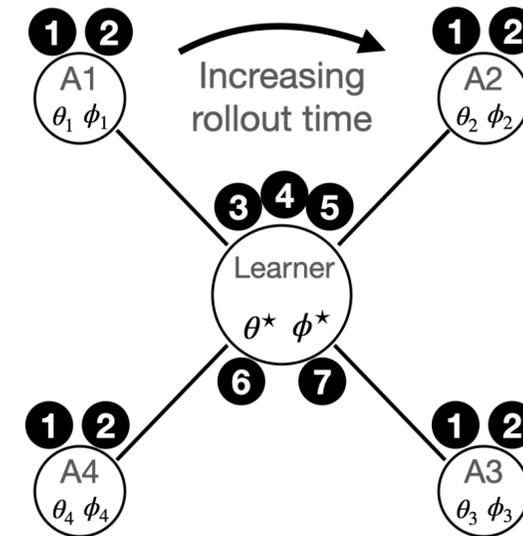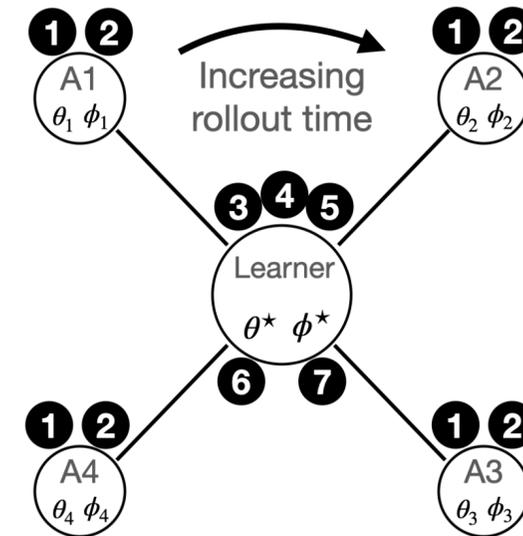**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t\,|\,s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t\,|\,s_t)$$

Bellman update

$$V_{\phi^\star}(s_t) \leftarrow V_{\phi^\star}(s_t) + \alpha_{\phi^\star} \times (r_{t+1} + \gamma V_{\phi^\star}(s_{t+1}) - V_{\phi^\star}(s_t))$$

Trajectory obtained from $\theta_1$

$$\rho_t \triangleq min(1, \frac{\pi_{\theta^\star}(a_t\,|\,s_t)}{\pi_{\theta_1}(a_t\,|\,s_t)})$$

$$V_{\phi^\star}(s_t) \leftarrow V_{\phi^\star}(s_t) + \alpha_{\phi^\star} \times \rho_t \times (r_{t+1} + \gamma V_{\phi^\star}(s_{t+1}) - V_{\phi^\star}(s_t))$$

# IMPALA | V-Trace | TD(0)



**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
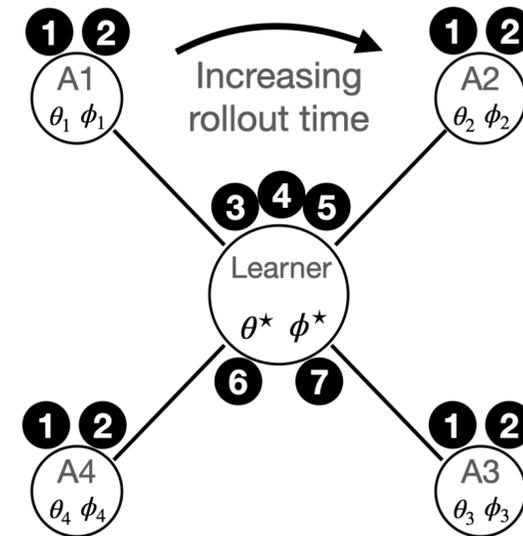
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Bellman update

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

Trajectory obtained from $\theta_1$

$$\rho_t \triangleq min(1, \frac{\pi_{\theta\star}(a_t \mid s_t)}{\pi_{\theta_1}(a_t \mid s_t)})$$

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times \rho_t \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

$$\delta_t = \rho_t \times (r_{t+1} + \gamma V_{\phi\star}(s_{t+1}) - V_{\phi\star}(s_t))$$

V-Trace correction for TD(0)

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \cdot step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

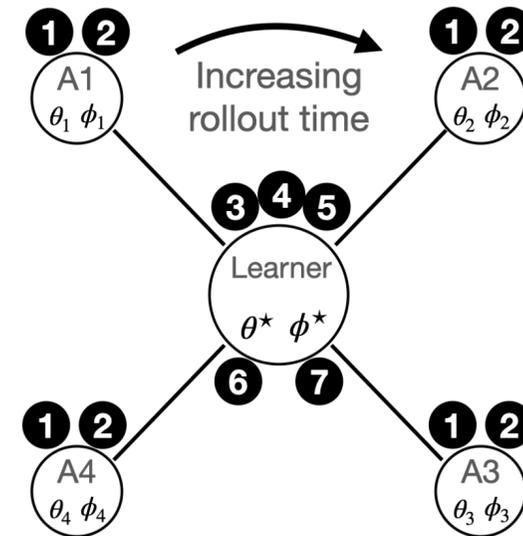**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
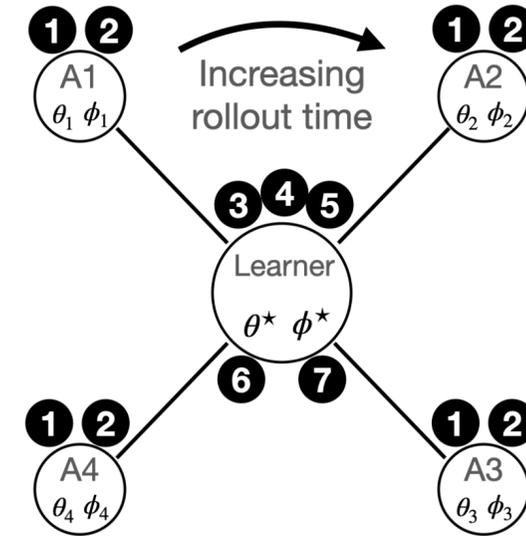
**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

$$V_{\phi^\star}(s_t) \leftarrow V_{\phi^\star}(s_t) + \alpha_{\phi^\star} \times \sum_{k=t}^{t+n} \gamma^{k-t} \times \rho_k \times (r_{k+1} + \gamma V_{\phi^\star}(s_{k+1}) - V_{\phi^\star}(s_k))$$

Naive V-Trace Bellman update for TD(n)

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times \sum_{k=t}^{t+n} \gamma^{k-t} \times \rho_k \times (r_{k+1} + \gamma V_{\phi\star}(s_{k+1}) - V_{\phi\star}(s_k))$$

Naive V-Trace Bellman update for TD(n)

$$V_{\phi\star}(s_0) \leftarrow V_{\phi\star}(s_0) + \alpha_{\phi\star} \times (\rho_0(r_1 + \gamma V_{\phi\star}(s_1) - V_{\phi\star}(s_0))$$
$$+\rho_1(r_2 + \gamma V_{\phi\star}(s_2) - V_{\phi\star}(s_1)))$$

① S

② G

③ C

④ V

⑤ P

⑥ U

⑦ Update Policy parameters $\theta$

$V_{\phi\star}(s_k))$

$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$

## Example

Bellman:

$$V(S0) = \boxed{p0*(r0 + (p1*r1 + p2*r2))}$$
$$= p0*r0 + p0*p1*r1 + p0*p2*r2$$

Trajectories:
$$V(S0) = \boxed{p0*p1*(r0+r1) + p0*p2*(r0+r2)}$$
$$= p0*p1*r0 + p0*p1*r1 + p0*p2*r0 + p0*p2*r2$$
$$= p0*r0*\underbrace{(p1+p2)}_{=1} + p0*p1*r1 + p0*p2*r2$$
$$= p0*r0 + p0*p1*r1 + p0*p2*r2$$



Unrolled MDP for H = 2, policy $\pi$

$p0 = 1$
$p1 + p2 = 1$

Lecture 7 - Slide 24

Intuition: rooted DAG $\rightarrow$ set of paths starting at root
Contribution of each edge to expected reward at *S0* is distributed over trajectories that contain it

① S

② G

$V_{\phi\star}(s_k))$

③ C

Example

**Bellman:** $V(S0) = \boxed{p0*(r0 + (p1*r1 + p2*r2))}$

$= p0*r0 + \boxed{p0*p1*r1} + p0*p2*r2$

**Trajectories:** $V(S0) = \boxed{p0*p1*(r0+r1) + p0*p2*(r0+r2)}$

$= p0*p1*r0 + p0*p1*r1 + p0*p2*r0 + p0*p2*r2$

$= p0*r0*\underbrace{(p1+p2)}_{=1} + p0*p1*r1 + p0*p2*r2$

$= p0*r0 + p0*p1*r1 + p0*p2*r2$

④ V

⑤ P

⑥ U



Unrolled MDP for H = 2, policy $\pi$

$p0 = 1$

$p1 + p2 = 1$

Reward contribution proportional to probability of its *Path*
We need to weight future samples (t + n) from t

Lecture 7 - Slide 24

Intuition: rooted DAG → set of paths starting at root
Contribution of each edge to expected reward at *S0* is distributed over trajectories that contain it

⑦ Update Policy parameters $\theta$

$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

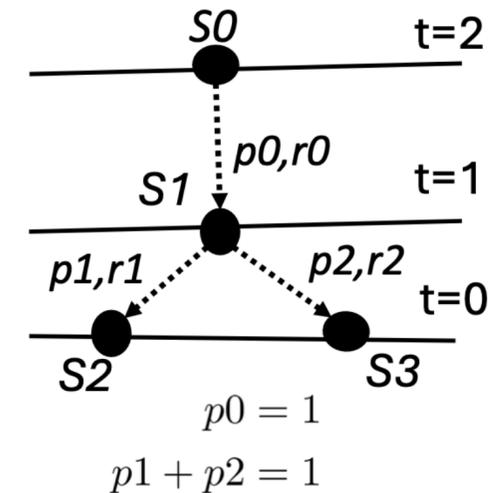$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env.step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

$$V_{\phi\star}(s_0) \leftarrow V_{\phi\star}(s_0) + \alpha_{\phi\star} \times (\rho_0(r_1 + \gamma V_{\phi\star}(s_1) - V_{\phi\star}(s_0))$$
$$+ \quad \rho_1(r_2 + \gamma V_{\phi\star}(s_2) - V_{\phi\star}(s_1)))$$

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \, \cdot \, | \, s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env \, . \, step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \, | \, s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \, | \, s_t)$$

$$V_{\phi^\star}(s_0) \leftarrow V_{\phi^\star}(s_0) + \alpha_{\phi^\star} \times (\rho_0(r_1 + \gamma V_{\phi^\star}(s_1) - V_{\phi^\star}(s_0))$$
$$+ \rho_0 \rho_1 (r_2 + \gamma V_{\phi^\star}(s_2) - V_{\phi^\star}(s_1)))$$

# IMPALA | V-Trace | TD(n)

**1** Sample an action with current policy for $s_t$

$$a_t \sim \pi_\theta( \cdot \mid s_t)$$

**2** Get trajectory for sampled action $a_t$

$$s_t, a_t, s_{t+1}, r_{t+1} \leftarrow env . step(a_t, s_t)$$

**3** Compute TD-Error $\delta_t$

$$\delta_t = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

**4** Value Gradient $\nabla_\phi V_\phi(s_t)$

**5** Policy Gradient $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$

**6** Update Value parameters $\phi$

$$\phi \leftarrow \phi + \alpha_\phi \delta_t \nabla_\phi V_\phi(s_t)$$

**7** Update Policy parameters $\theta$

$$\theta \leftarrow \theta + \alpha_\theta \delta_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

$$V_{\phi\star}(s_0) \leftarrow V_{\phi\star}(s_0) + \alpha_{\phi\star} \times (\rho_0(r_1 + \gamma V_{\phi\star}(s_1) - V_{\phi\star}(s_0))$$
$$+ \rho_0 \rho_1(r_2 + \gamma V_{\phi\star}(s_2) - V_{\phi\star}(s_1)))$$

$$V_{\phi\star}(s_t) \leftarrow V_{\phi\star}(s_t) + \alpha_{\phi\star} \times \sum_{k=t}^{t+n} \gamma^{k-t} \times (\prod_{i=t}^{k-1} \rho_i) \times \rho_k \times (r_{k+1} + \gamma V_{\phi\star}(s_{k+1}) - V_{\phi\star}(s_k))$$

Final V-Trace correction for TD(n)

# IMPALA | Evaluation

| Architecture | CPUs | GPUs[1] | FPS[2] | |
|---|---|---|---|---|
| **Single-Machine** | | | Task 1 | Task 2 |
| A3C 32 workers | 64 | 0 | 6.5K | 9K |
| Batched A2C (sync step) | 48 | 0 | 9K | 5K |
| Batched A2C (sync step) | 48 | 1 | 13K | 5.5K |
| Batched A2C (sync traj.) | 48 | 0 | 16K | 17.5K |
| Batched A2C (dyn. batch) | 48 | 1 | 16K | 13K |
| IMPALA 48 actors | 48 | 0 | 17K | 20.5K |
| IMPALA (dyn. batch) 48 actors[3] | 48 | 1 | 21K | 24K |
| **Distributed** | | | | |
| A3C | 200 | 0 | 46K | 50K |
| IMPALA | 150 | 1 | 80K | |
| IMPALA (optimised) | 375 | 1 | 200K | |
| IMPALA (optimised) batch 128 | 500 | 1 | 250K | |

[1] Nvidia P100 [2] In frames/sec (4 times the agent steps due to action repeat). [3] Limited by amount of rendering possible on a single machine.

*Table 1.* Throughput on `seekavoid_arena_01` (task 1) and `rooms_keys_doors_puzzle` (task 2) with the shallow model in Figure 3. The latter has variable length episodes and slow restarts. Batched A2C and IMPALA use batch size 32 if not otherwise mentioned.
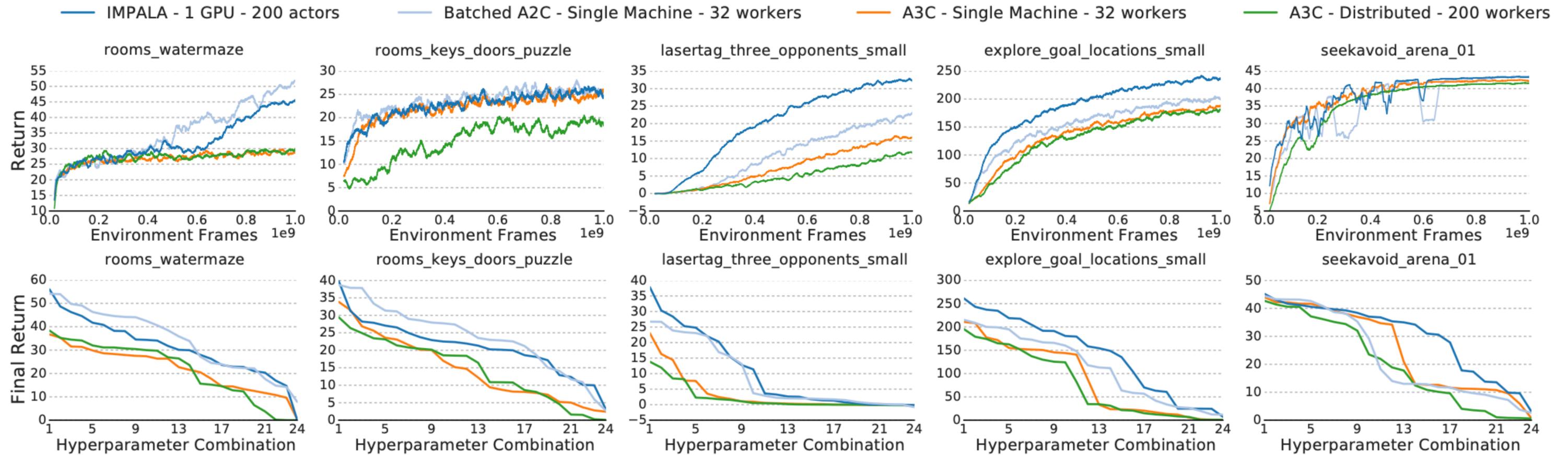
## Throughput Evaluation

| | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 |
|---|---|---|---|---|---|
| **Without Replay** | | | | | |
| V-trace | 46.8 | 32.9 | **31.3** | **229.2** | **43.8** |
| 1-Step | **51.8** | **35.9** | 25.4 | 215.8 | 43.7 |
| $\varepsilon$-correction | 44.2 | 27.3 | 4.3 | 107.7 | 41.5 |
| No-correction | 40.3 | 29.1 | 5.0 | 94.9 | 16.1 |
| **With Replay** | | | | | |
| V-trace | 47.1 | **35.8** | **34.5** | **250.8** | **46.9** |
| 1-Step | **54.7** | 34.4 | 26.4 | 204.8 | 41.6 |
| $\varepsilon$-correction | 30.4 | 30.2 | 3.9 | 101.5 | 37.6 |
| No-correction | 35.0 | 21.1 | 2.8 | 85.0 | 11.2 |

Tasks: `rooms_watermaze`, `rooms_keys_doors_puzzle`, `lasertag_three_opponents_small`, `explore_goal_locations_small`, `seekavoid_arena_01`

*Table 2.* Average final return over 3 best hyperparameters for different off-policy correction methods on 5 DeepMind Lab tasks. When the lag in policy is negligible both V-trace and 1-step importance sampling perform similarly well and better than $\varepsilon$-correction/No-correction. However, when the lag increases due to use of experience replay, V-trace performs better than all other methods in 4 out 5 tasks.

## Policy Lag Correction Evaluation

Figure 4. **Top Row:** Single task training on 5 DeepMind Lab tasks. Each curve is the mean of the best 3 runs based on final return. IMPALA achieves better performance than A3C. **Bottom Row:** Stability across hyperparameter combinations sorted by the final performance across different hyperparameter combinations. IMPALA is consistently more stable than A3C.

Stability Evaluation