

REINFORCEMENT LEARNING

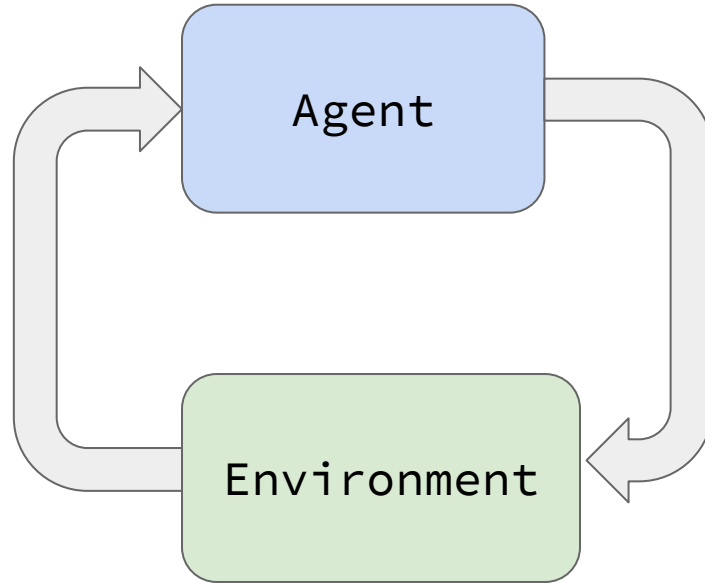
A brief overview through the A3C paper

BACK TO BASICS

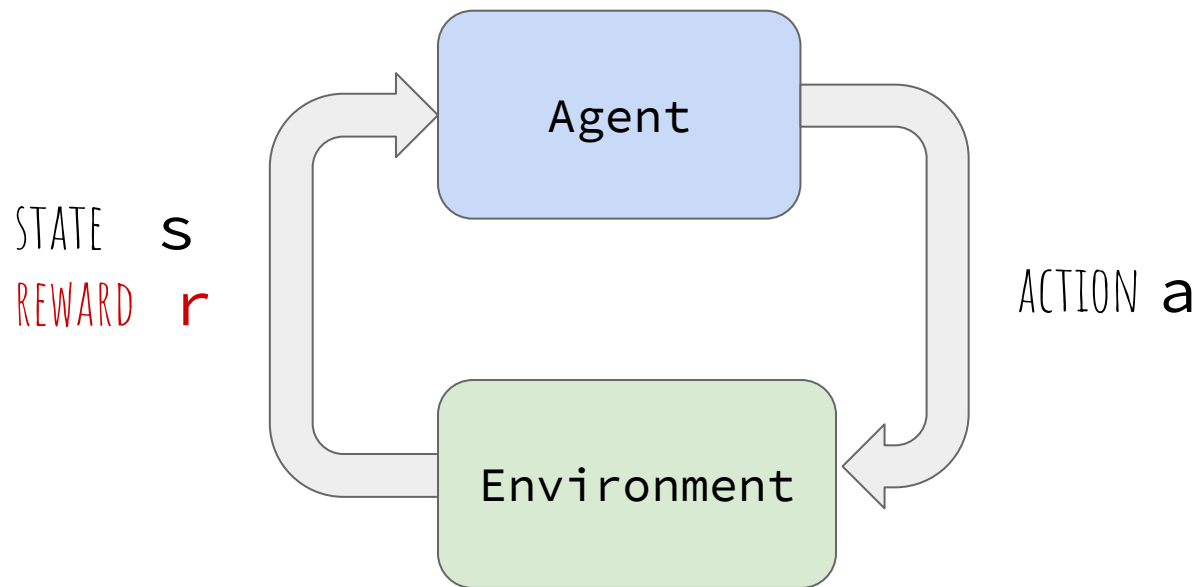
1



BACK TO BASICS



BACK TO BASICS



MARKOV DECISION PROCESSES : EXAMPLE



MARKOV DECISION PROCESSES

MARKOV DECISION PROCESSES

set of actions $a \in A$

set of states $s \in S$

transition function $T(s, a, s')$

reward function $R(s, a, s')$

start / terminal state

MARKOV DECISION PROCESSES



Andrey Markov
1856 - 1922

MARKOV DECISION PROCESSES

The stochastic process is
memoryless

MARKOV DECISION PROCESSES

The stochastic process is
memoryless

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, \dots, S_0 = s_0)$$

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

RL FRAMEWORK

RL FRAMEWORK

At each timestep t :

RL FRAMEWORK

At each timestep t :

- The agent receives a **state** s_t from S

RL FRAMEWORK

At each timestep t :

- The agent receives a **state** s_t from S
- The agent chooses **action** a_t from A , according to π

RL FRAMEWORK

At each timestep t :

- The agent receives a **state** s_t from S
- The agent chooses **action** a_t from A , according to π
- The agent receives s_{t+1} and a **reward** r_t

RL FRAMEWORK

At each timestep t :

- The agent receives a **state** s_t from S
- The agent chooses **action** a_t from A , according to π
- The agent receives s_{t+1} and a **reward** r_t

The process continues until the agent reaches a **terminal state**

RL FRAMEWORK

At each timestep t :

- The agent receives a **state** s_t from S
- The agent chooses **action** a_t from A , according to π
- The agent receives s_{t+1} and a **reward** r_t

The process continues until the agent reaches a **terminal state**

The **return** is
$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

RL FRAMEWORK

At each timestep t :

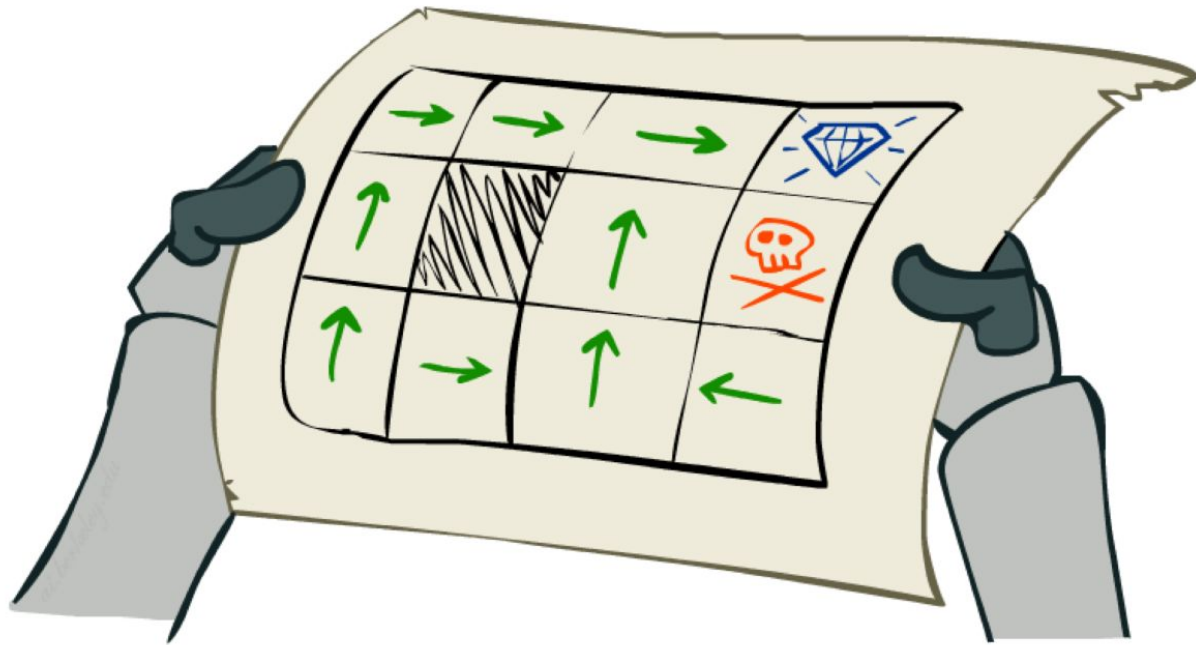
- The agent receives a **state** s_t from S
- The agent chooses **action** a_t from A , according to π
- The agent receives s_{t+1} and a **reward** r_t

The process continues until the agent reaches a **terminal state**

The **return** is
$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

The goal of the agent is to maximize the **expected return** from each state s_t

GOAL



VALUE METHODS

VALUE METHODS

State value

$$V^\pi(\mathbf{s}) = \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}]$$

VALUE METHODS

State value

$$V^\pi(\mathbf{s}) = \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}]$$

Action value

$$Q^\pi(\mathbf{s}, a) = \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}, a_t = a]$$

VALUE METHODS

State value

$$V^\pi(\mathbf{s}) = \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}]$$

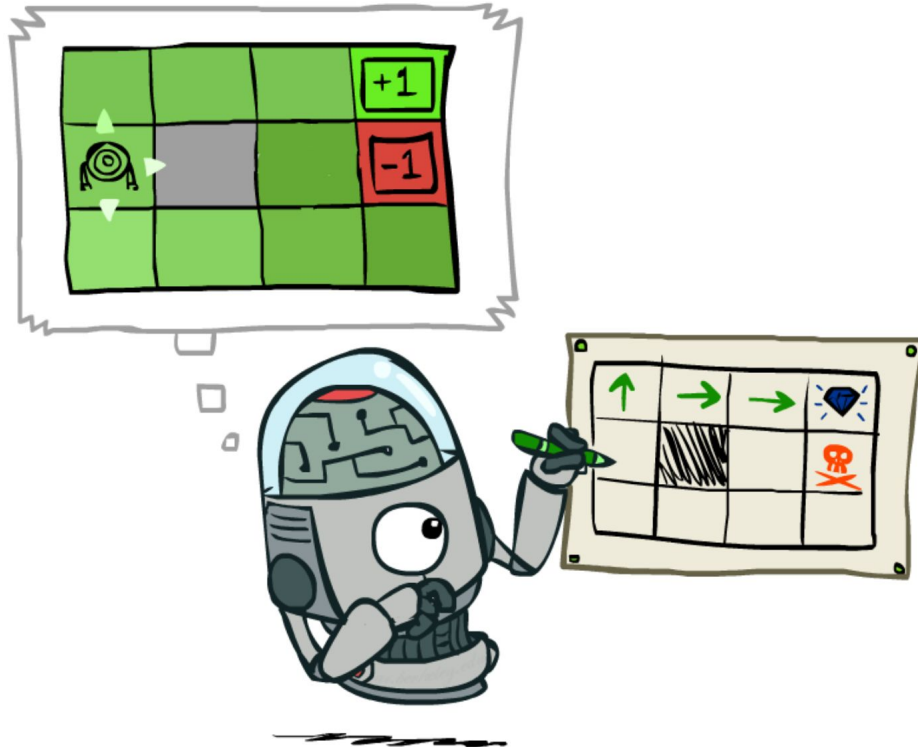
Action value

$$Q^\pi(\mathbf{s}, a) = \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}, a_t = a]$$

Optimal value function

$$Q^*(\mathbf{s}, a) = \max_\pi Q^\pi(\mathbf{s}, a)$$

VALUE METHODS



Q-LEARNING

Q-LEARNING

$$Q^{new}(s_t, a_t) \leftarrow (1 - \alpha)Q^{old}(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q^{old}(s_{t+1}, a'))$$

APPROXIMATE SETTING

APPROXIMATE SETTING

$$Q(s, a; \theta)$$

APPROXIMATE SETTING

$$Q(s, a; \theta)$$

$$Q^*(s, a) \approx Q(s, a; \theta)$$

APPROXIMATE SETTING

$$Q(s, a; \theta)$$

$$Q^*(s, a) \approx Q(s, a; \theta)$$

$$L_i(\theta_i) = \mathbb{E}[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)]^2$$

APPROXIMATE SETTING

$$Q(s, a; \theta)$$

$$Q^*(s, a) \approx Q(s, a; \theta)$$

$$L_i(\theta_i) = \mathbb{E}[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)]^2$$

r directly affects one $Q(s, a)$

N-STEP Q-LEARNING

N-STEP Q-LEARNING

Use n-step return instead!

$$R^n = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \max_a \gamma^n Q(s_{t+n}, a)$$

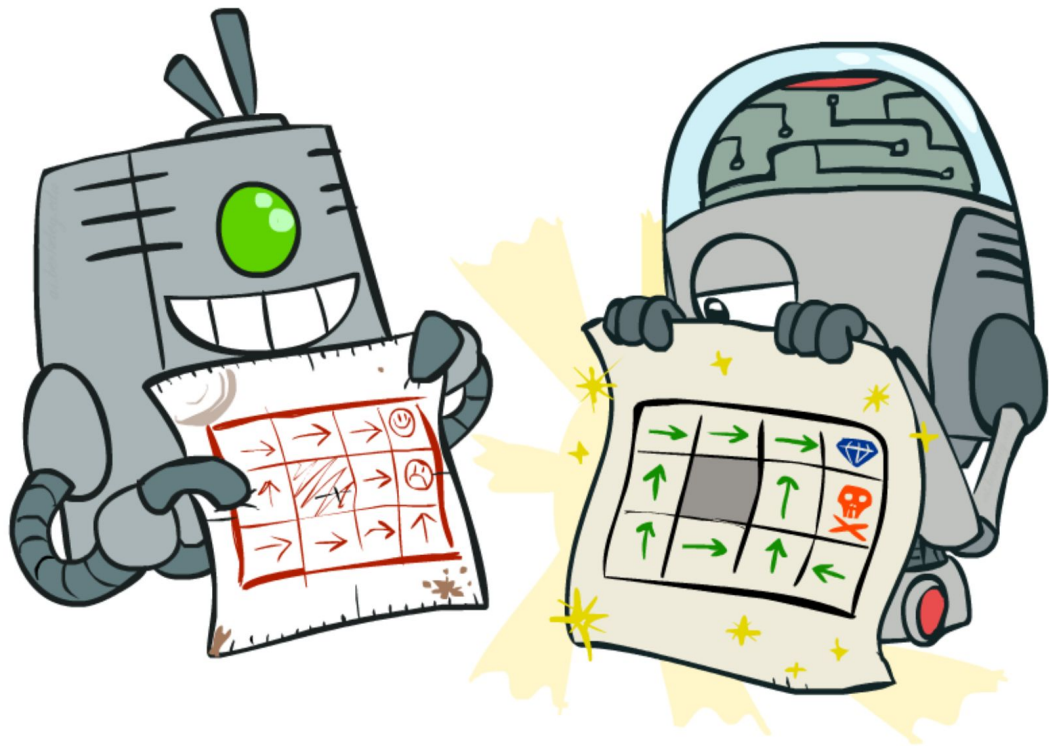
N-STEP Q-LEARNING

Use n-step return instead!

$$R^n = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \max_a \gamma^n Q(s_{t+n}, a)$$

r affects n Q(s,a)

POLICY METHODS



POLICY METHODS

Value methods: implicit policy

$$\pi(a|s)$$

POLICY METHODS

Value methods: implicit policy

$$\pi(a|s)$$

Policy methods: explicit policy

$$\pi(a|s; \theta)$$

POLICY METHODS

Value methods: implicit policy

$$\pi(a|s)$$

Policy methods: explicit policy

$$\pi(a|s; \theta)$$

Gradient ascent on the return

$$\mathbb{E}[R_t]$$

POLICY METHODS

Value methods: implicit policy

$$\pi(a|s)$$

Policy methods: explicit policy

$$\pi(a|s; \theta)$$

Gradient ascent on the return

$$\mathbb{E}[R_t]$$

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

Variance problem

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

Variance problem

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t - b_t(s_t))$$

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

Variance problem

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t - b_t(s_t))$$

Common baseline

$$b_t(s_t) \approx V^{\pi}(s_t)$$

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

Variance problem

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t - b_t(s_t))$$

Common baseline

$$b_t(s_t) \approx V^{\pi}(s_t)$$

Advantage

$$A(a_t, s_t) = Q(a_t, s_t) - V(s_t)$$

POLICY METHODS

REINFORCE update

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) R_t$$

Variance problem

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t - b_t(s_t))$$

Common baseline

$$b_t(s_t) \approx V^{\pi}(s_t)$$

Advantage

$$A(a_t, s_t) = Q(a_t, s_t) - V(s_t)$$

ASYNCHRONOUS RL

2

Asynchronous Methods for Deep Reinforcement Learning

Volodymyr Mnih¹
Adrià Puigdomènech Badia¹
Mehdi Mirza^{1,2}
Alex Graves¹
Tim Harley¹
Timothy P. Lillicrap¹
David Silver¹
Koray Kavukcuoglu¹

VMNIEH@GOOGLE.COM
ADRIAP@GOOGLE.COM
MIRZAMOM@IRO.UMONTREAL.CA
GRAVESA@GOOGLE.COM
THARLEY@GOOGLE.COM
COUNTZERO@GOOGLE.COM
DAVIDSILVER@GOOGLE.COM
KORAYK@GOOGLE.COM

¹ Google DeepMind

² Montreal Institute for Learning Algorithms (MILA), University of Montreal

Abstract

We propose a conceptually simple and lightweight framework for deep reinforcement learning that uses asynchronous gradient descent for optimization of deep neural network controllers. We present asynchronous variants of

line RL updates are strongly correlated. By storing the agent's data in an experience replay memory, the data can be batched (Riedmiller, 2005; Schulman et al., 2015a) or randomly sampled (Mnih et al., 2013; 2015; Van Hasselt et al., 2015) from different time-steps. Aggregating over memory in this way reduces non-stationarity and decorrelates updates, but at the same time limits the methods to

MOTIVATION

MOTIVATION

RL + function approximation = ?

MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

Experience replay

MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

Experience replay



MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

Experience replay



MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

Experience replay



MOTIVATION

RL + function approximation = ?

Strongly correlated online RL updates

Experience replay

Off-policy only



ASYNCHRONOUS RL

ASYNCHRONOUS RL

Multiple agents in parallel

ASYNCHRONOUS RL

Multiple agents in parallel

Multiple instances of the environment

ASYNCHRONOUS RL

Multiple agents in parallel

Multiple instances of the environment

On-policy: Sarsa, n-step methods, actor-critic methods

ASYNCHRONOUS RL

Multiple agents in parallel

Multiple instances of the environment

On-policy: Sarsa, n-step methods, actor-critic methods

GPU -> CPU

ASYNCHRONOUS RL FRAMEWORK

ASYNCHRONOUS RL FRAMEWORK

1 asynchronous actor-learner = 1 thread

ASYNCHRONOUS RL FRAMEWORK

1 asynchronous actor-learner = 1 thread

1 asynchronous actor-learner = 1 different exploration policy

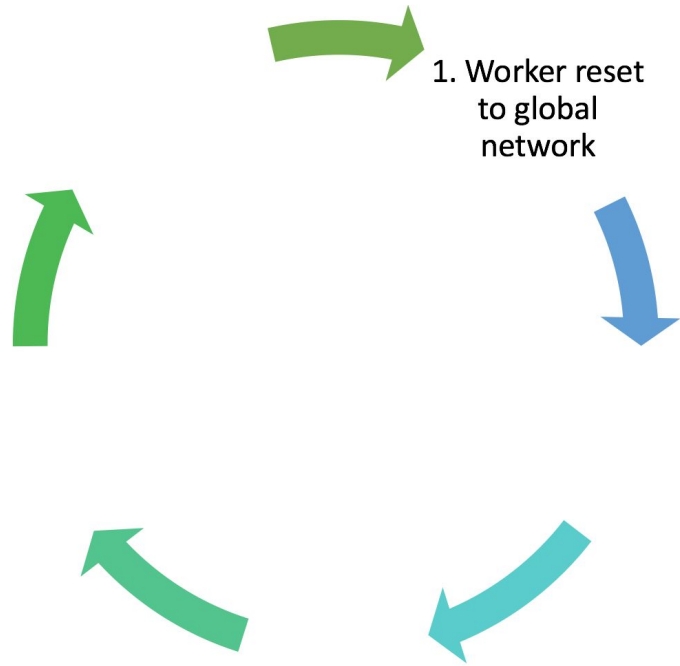
ASYNCHRONOUS RL FRAMEWORK

1 asynchronous actor-learner = 1 thread

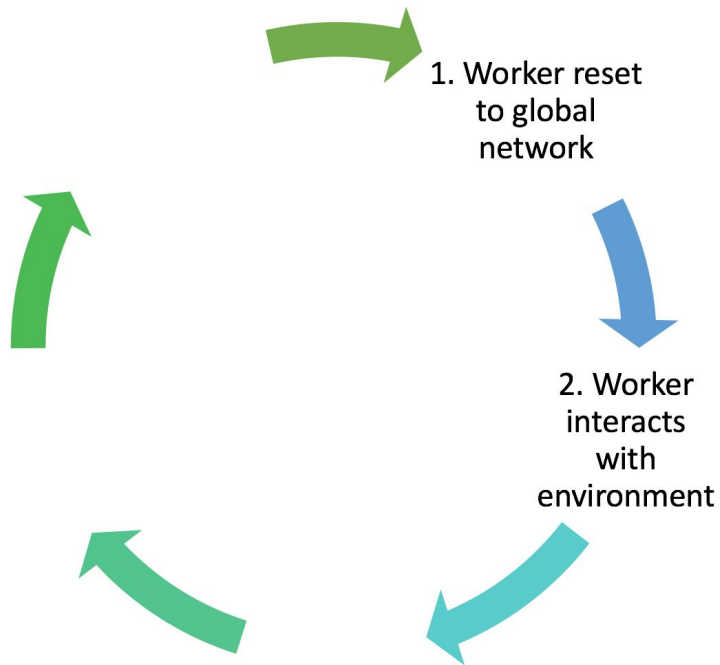
1 asynchronous actor-learner = 1 different exploration policy

No replay memory

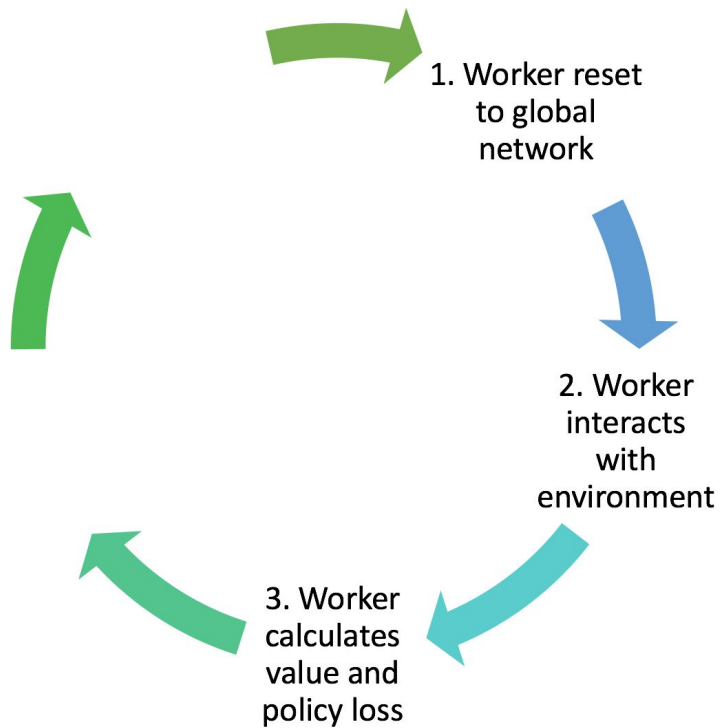
ASYNCHRONOUS RL FRAMEWORK



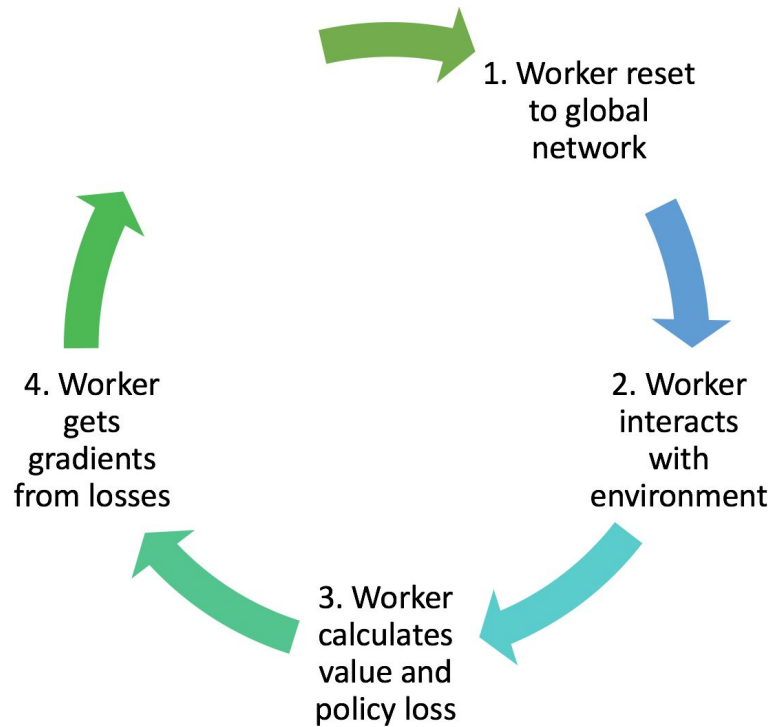
ASYNCHRONOUS RL FRAMEWORK



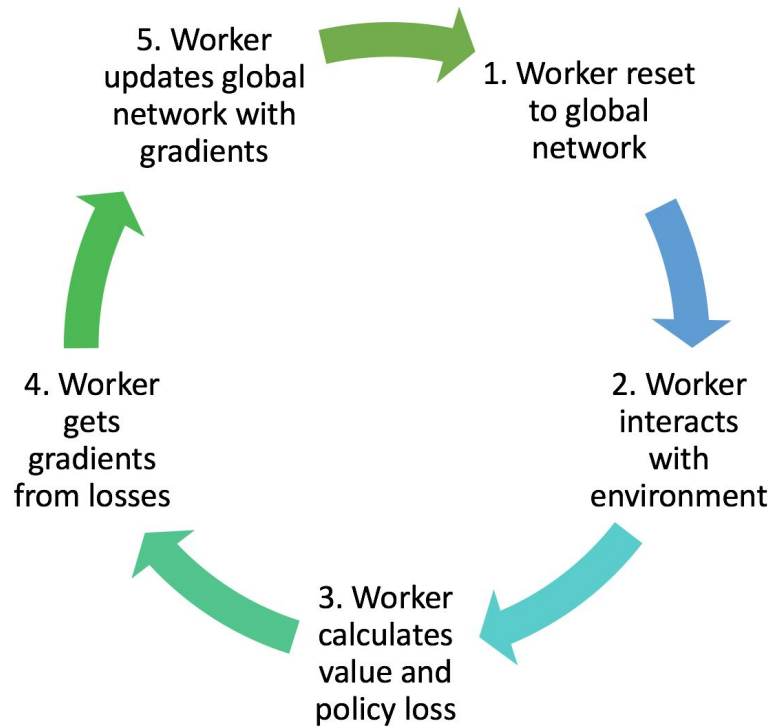
ASYNCHRONOUS RL FRAMEWORK



ASYNCHRONOUS RL FRAMEWORK



ASYNCHRONOUS RL FRAMEWORK



ASYNCHRONOUS ONE-STEP Q-LEARNING

ASYNCHRONOUS ONE-STEP Q-LEARNING

1 thread = 1 copy of the environment

ASYNCHRONOUS ONE-STEP Q-LEARNING

1 thread = 1 copy of the environment

At each step computes the gradient of the Q-learning loss

ASYNCHRONOUS ONE-STEP Q-LEARNING

1 thread = 1 copy of the environment

At each step computes the gradient of the Q-learning loss

$$L_i(\theta_i) = \mathbb{E}[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)]^2$$

ASYNCHRONOUS ONE-STEP Q-LEARNING

1 thread = 1 copy of the environment

At each step computes the gradient of the Q-learning loss

Accumulate gradients over multiple time steps

Algorithm 1 Asynchronous one-step Q-learning - pseudocode for each actor-learner thread.

// Assume global shared θ , θ^- , and counter $T = 0$.
Initialize thread step counter $t \leftarrow 0$
Initialize target network weights $\theta^- \leftarrow \theta$
Initialize network gradients $d\theta \leftarrow 0$
Get initial state s
repeat
 Take action a with ϵ -greedy policy based on $Q(s, a; \theta)$
 Receive new state s' and reward r

$$y = \begin{cases} r & \text{for terminal } s' \\ r + \gamma \max_{a'} Q(s', a'; \theta^-) & \text{for non-terminal } s' \end{cases}$$

 Accumulate gradients wrt θ : $d\theta \leftarrow d\theta + \frac{\partial(y - Q(s, a; \theta))^2}{\partial \theta}$
 $s = s'$
 $T \leftarrow T + 1$ and $t \leftarrow t + 1$
 if $T \bmod I_{target} == 0$ **then**
 Update the target network $\theta^- \leftarrow \theta$
 end if
 if $t \bmod I_{AsyncUpdate} == 0$ or s is terminal **then**
 Perform asynchronous update of θ using $d\theta$.
 Clear gradients $d\theta \leftarrow 0$.
 end if
until $T > T_{max}$

Algorithm 1 Asynchronous one-step Q-learning - pseudocode for each actor-learner thread.

// Assume global shared θ , θ^- , and counter $T = 0$.

Initialize thread step counter $t \leftarrow 0$

Initialize target network weights $\theta^- \leftarrow \theta$

Initialize network gradients $d\theta \leftarrow 0$

Get initial state s

repeat

Take action a with ϵ -greedy policy based on $Q(s, a; \theta)$

Receive new state s' and reward r

$$y = \begin{cases} r & \text{for terminal } s' \\ r + \gamma \max_{a'} Q(s', a'; \theta^-) & \text{for non-terminal } s' \end{cases}$$

Accumulate gradients wrt θ : $d\theta \leftarrow d\theta + \frac{\partial(y - Q(s, a; \theta))^2}{\partial \theta}$

$s = s'$

$T \leftarrow T + 1$ and $t \leftarrow t + 1$

if $T \bmod I_{target} == 0$ **then**

Update the target network $\theta^- \leftarrow \theta$

end if

if $t \bmod I_{AsyncUpdate} == 0$ or s is terminal **then**

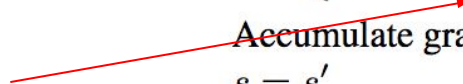
Perform asynchronous update of θ using $d\theta$.

Clear gradients $d\theta \leftarrow 0$.

end if

until $T > T_{max}$

Sarsa



ASYNCHRONOUS N-STEP Q-LEARNING

ASYNCHRONOUS N-STEP Q-LEARNING

Use forward view instead of backward view

ASYNCHRONOUS N-STEP Q-LEARNING

Use forward view instead of backward view

Easier for training NNs with momentum

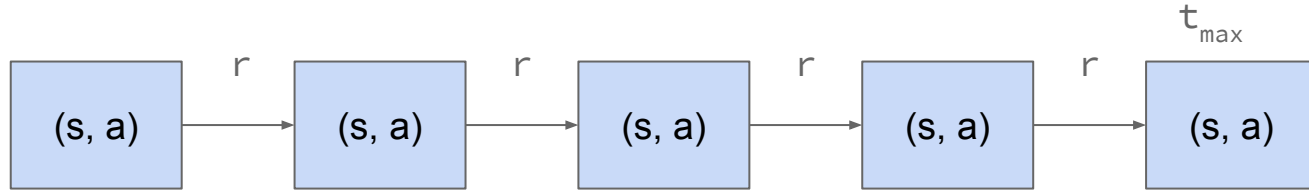
ASYNCHRONOUS N-STEP Q-LEARNING

Use forward view instead of backward view

Easier for training NNs with momentum

First select actions up to t_{\max} or until terminal state

ASYNCHRONOUS N-STEP Q-LEARNING



ASYNCHRONOUS ADVANTAGE ACTOR CRITIC (A3C)

ASYNCHRONOUS ADVANTAGE ACTOR CRITIC (A3C)

Maintains $\pi(a_t, s_t; \theta)$ and $V(s_t; \theta_v)$

ASYNCHRONOUS ADVANTAGE ACTOR CRITIC (A3C)

Maintains $\pi(a_t, s_t; \theta)$ and $V(s_t; \theta_v)$

Also operates in forward view (update after every t_{\max})

ASYNCHRONOUS ADVANTAGE ACTOR CRITIC (A3C)

Maintains $\pi(a_t, s_t; \theta)$ and $V(s_t; \theta_v)$

Also operates in forward view (update after every t_{\max})

Update $\nabla_{\theta'} \log \pi(a_t | s_t; \theta') A(s_t, a_t; \theta, \theta_v)$

ASYNCHRONOUS ADVANTAGE ACTOR CRITIC (A3C)

Maintains $\pi(a_t, s_t; \theta)$ and $V(s_t; \theta_v)$

Also operates in forward view (update after every t_{\max})

Update $\nabla_{\theta'} \log \pi(a_t | s_t; \theta') A(s_t, a_t; \theta, \theta_v)$

Advantage (k up to t_{\max}) $\sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v)$

A3C DETAILS

A3C DETAILS

θ and θ_v come from the same network with two heads

A3C DETAILS

θ and θ_v come from the same network with two heads

Add entropy regularization term

RESULTS

3

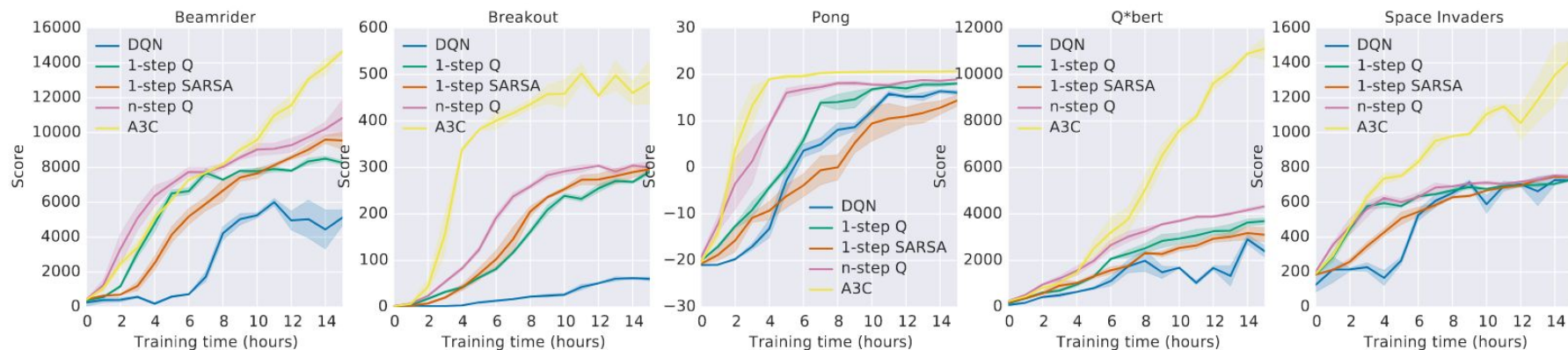
BENCHMARKS

Arcade Learning Environment (Bellemare et al., 2012)

TORCS 3D Racing simulator (Wyman et al., 2013)

A3C only: MuJoCo (Todorov, 2015)

ATARI 2600



DQN: Nvidia K40 GPU

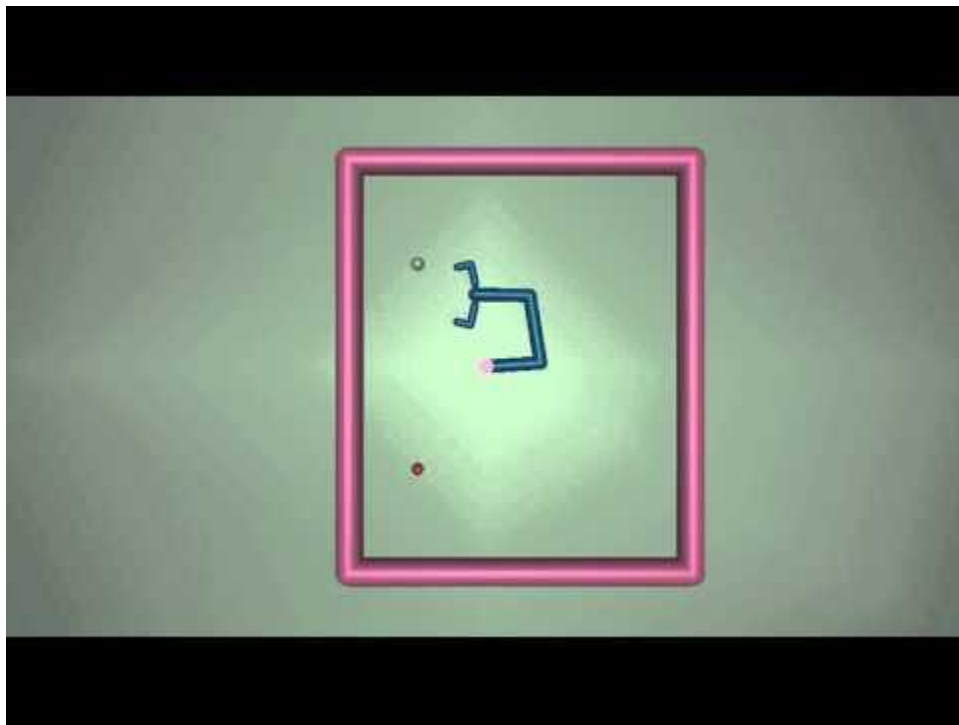
Asynchronous: CPU (16 cores)

ATARI 2600

Method	Training Time	Mean	Median
DQN	8 days on GPU	121.9%	47.5%
Gorila	4 days, 100 machines	215.2%	71.3%
D-DQN	8 days on GPU	332.9%	110.9%
Dueling D-DQN	8 days on GPU	343.8%	117.1%
Prioritized DQN	8 days on GPU	463.6%	127.6%
A3C, FF	1 day on CPU	344.1%	68.2%
A3C, FF	4 days on CPU	496.8%	116.6%
A3C, LSTM	4 days on CPU	623.0%	112.6%

Table 1. Mean and median human-normalized scores on 57 Atari games using the human starts evaluation metric. Supplementary Table SS3 shows the raw scores for all games.

MUJOCO



REFERENCES AND CREDITS

Illustrations from Dan Klein and Pieter Abbeel for CS188
Intro to AI at UC Berkeley <http://ai.berkeley.edu>

REFERENCES AND CREDITS

Bellemare, Marc G, Naddaf, Yavar, Veness, Joel, and Bowling, Michael. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 2012.

Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Graves, Alex, Antonoglou, Ioannis, Wierstra, Daan, and Riedmiller, Martin. Playing atari with deep reinforcement learning. In *NIPS Deep Learning Workshop*. 2013.

Degrís, Thomas, Pilarski, Patrick M, and Sutton, Richard S. Model-free reinforcement learning with continuous action in practice. In *American Control Conference (ACC)*, 2012, pp. 2177–2182. IEEE, 2012.

Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Rusu, Andrei A., Veness, Joel, Bellemare, Marc G., Graves, Alex, Riedmiller, Martin, Fidjeland, Andreas K., Ostrovski, Georg, Petersen, Stig, Beattie, Charles, Sadik, Amir, Antonoglou, Ioannis, King, Helen, Kumaran, Dharshan, Wierstra, Daan, Legg, Shane, and Hassabis, Demis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 02 2015. URL <http://dx.doi.org/10.1038/nature14236>.

REFERENCES AND CREDITS

Riedmiller, Martin. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *Machine Learning: ECML 2005*, pp. 317–328. Springer Berlin Heidelberg, 2005.

Schulman, John, Levine, Sergey, Moritz, Philipp, Jordan, Michael I, and Abbeel, Pieter. Trust region policy optimization. In *International Conference on Machine Learning (ICML)*, 2015a.

Sutton, R. and Barto, A. *Reinforcement Learning: an Introduction*. MIT Press, 1998.

Todorov, E. *MuJoCo: Modeling, Simulation and Visualization of Multi-Joint Dynamics with Contact (ed 1.0)*. Roboti Publishing, 2015.

Williams, R.J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, 1992.

QUESTIONS

4
