

# DQN and Advanced Policy Gradient

# Map

# Map

# Map

Model  
Based



Model  
Free



learn Q  
SARSA

learn  $\pi$   
Policy  
Gradient

On Policy



Off Policy

ML MB TRL  
(learn  $T, R$ )

Q-learning

# Map

Model  
Based



Model  
Free



learn Q  
SARSA

learn  $\pi$   
Policy  
Gradient

On Policy



Off Policy

ML MB TRL  
(learn  $T, R$ )

Q-learning

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

# Map

Model  
Based

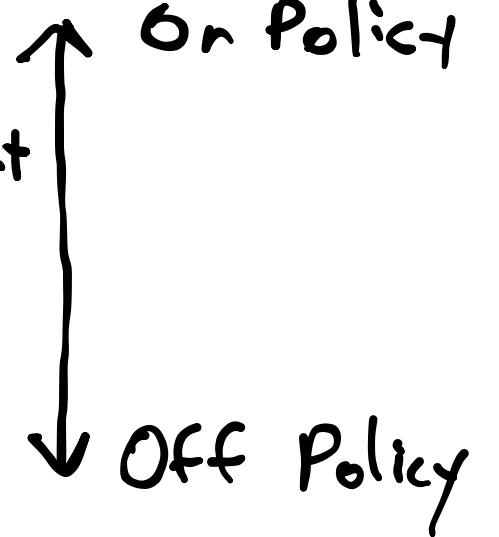


Model  
Free

learn Q  
SARSA

learn  $\pi$   
Policy  
Gradient

On Policy



ML MB TRL  
(learn  $T, R$ )

Q-learning

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

# Map

Model  
Based

Model  
Free

learn Q  
SARSA

learn  $\pi$   
Policy  
Gradient

On Policy

Off Policy

ML MB TRL  
(learn  $T, R$ )

Q-learning

Part I

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

# Map

Model  
Based



Model  
Free



learn Q  
SARSA

learn  $\pi$   
Policy Gradient  
Part 2

On Policy

Off Policy

ML MB TRL  
(learn  $T, R$ )

Q-learning

Part 1

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

# Map

Model  
Based

Model  
Free

ML MB TRL  
(learn  $T, R$ )

learn Q  
SARSA

learn  $\pi$   
Policy Gradient  
Part 2

On Policy

Off Policy

Q-learning

Part 1

Actor-Critic  
Part 3

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

# **Part I**

# **DQN**

# **Q-Learning with Neural Networks**

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$
$$\theta \leftarrow \theta - \alpha \widehat{\nabla_{\theta} l(f_{\theta}(x), y)}$$

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate  $Q$  with  $Q_{\theta}$

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate  $Q$  with  $Q_{\theta}$
- What should  $(x, y)$  be?

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

$\hat{\gamma}$        $+ d$        $\gamma$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate  $Q$  with  $Q_{\theta}$
- What should  $(x, y)$  be?  $(s, a, r, s')$
- What should  $l$  be?

$$l(s, a, r, s') = (r + \gamma \max_{a'} Q(s', a') - Q(s, a))^2$$

# Q-Learning with Neural Networks

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate  $Q$  with  $Q_{\theta}$
- What should  $(x, y)$  be?
- What should  $\underline{l}$  be?

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act!}(env, a)$$

$$s' \leftarrow \operatorname{observe}(env)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

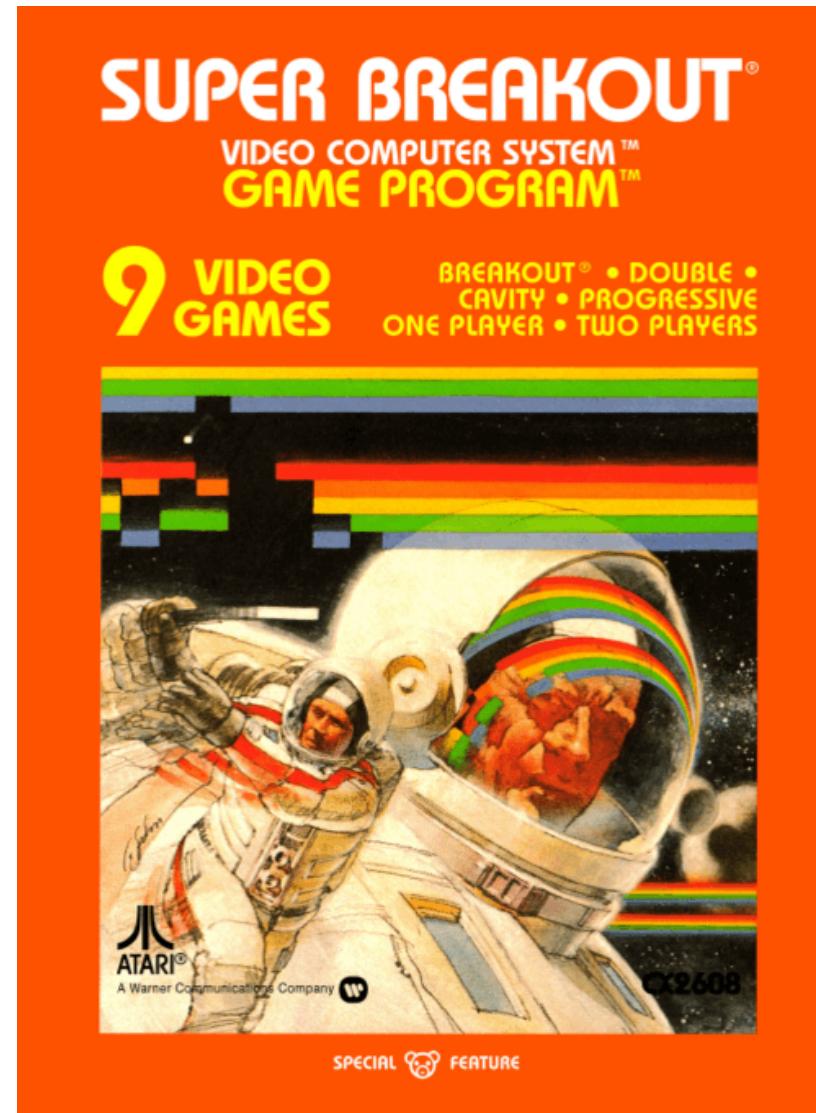
$$s \leftarrow s'$$

# DQN: The Atari Benchmark

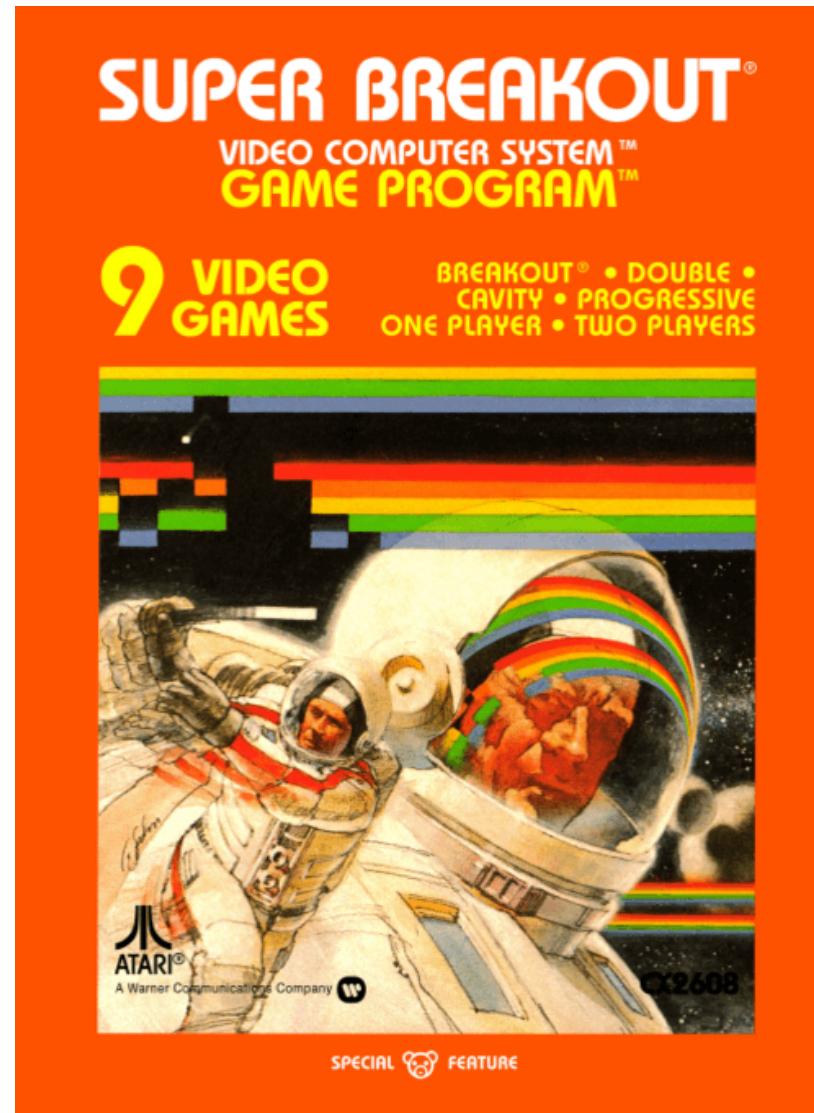
# DQN: The Atari Benchmark



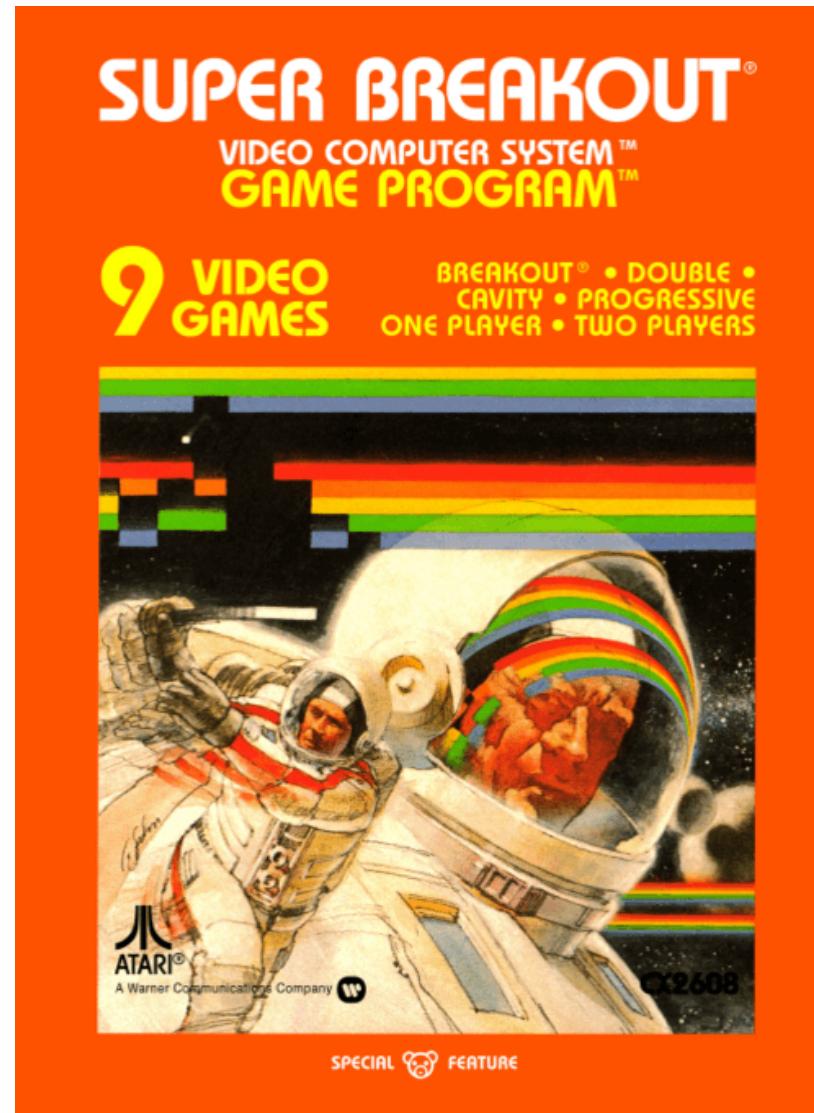
# DQN: The Atari Benchmark



# DQN: The Atari Benchmark



# DQN: The Atari Benchmark



# DQN: Problems with Naive Approach

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

$$s \leftarrow s'$$

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated
2. Size-1 batches

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated
2. Size-1 batches
3. Moving target

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated
2. Size-1 batches
3. Moving target

} data buffer/experience replay  
} periodically freeze target

# DQN

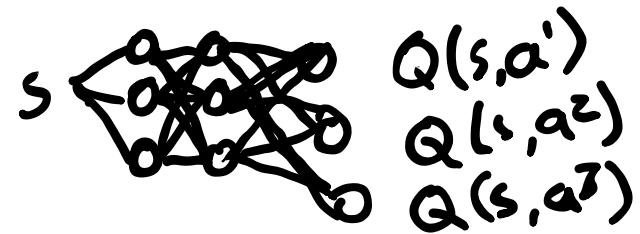
# DQN

Q Network Structure:

~~s o  
a~~  $\cancel{Q(s,a)}$

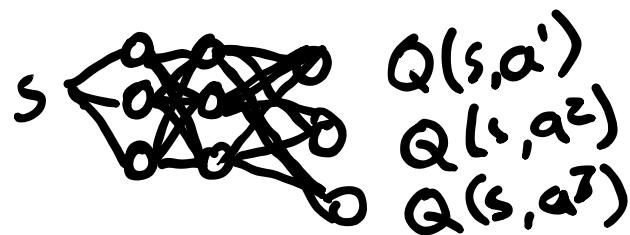
# DQN

Q Network Structure:



# DQN

Q Network Structure:

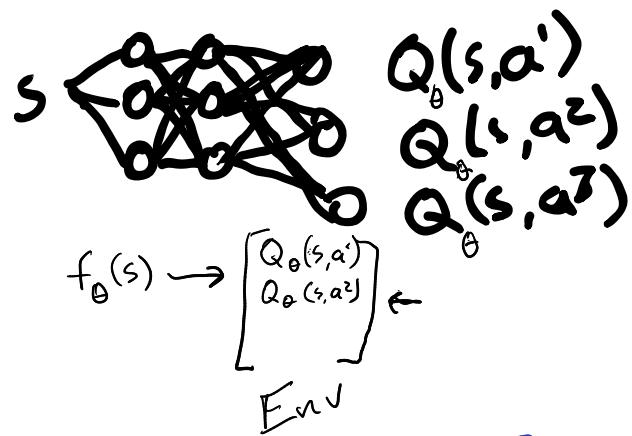


Experience Tuple:  $(s, a, r, s')$

Trick 1: network structure

# DQN

Q Network Structure:

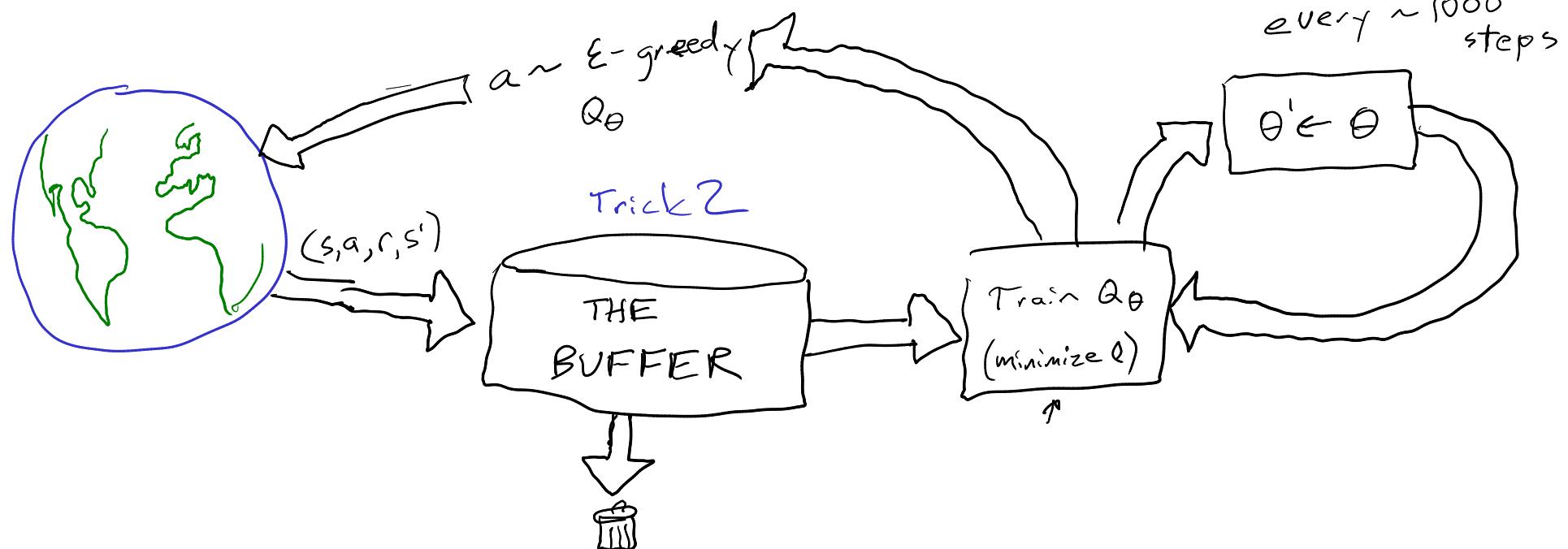


Experience Tuple:  $(s, a, r, s')$

Loss:

$$l(s, a, r, s') = \left( r + \gamma \max_{a'} Q_{\theta'}(s', a') - Q_\theta(s, a) \right)^2$$

↑  
updating slower



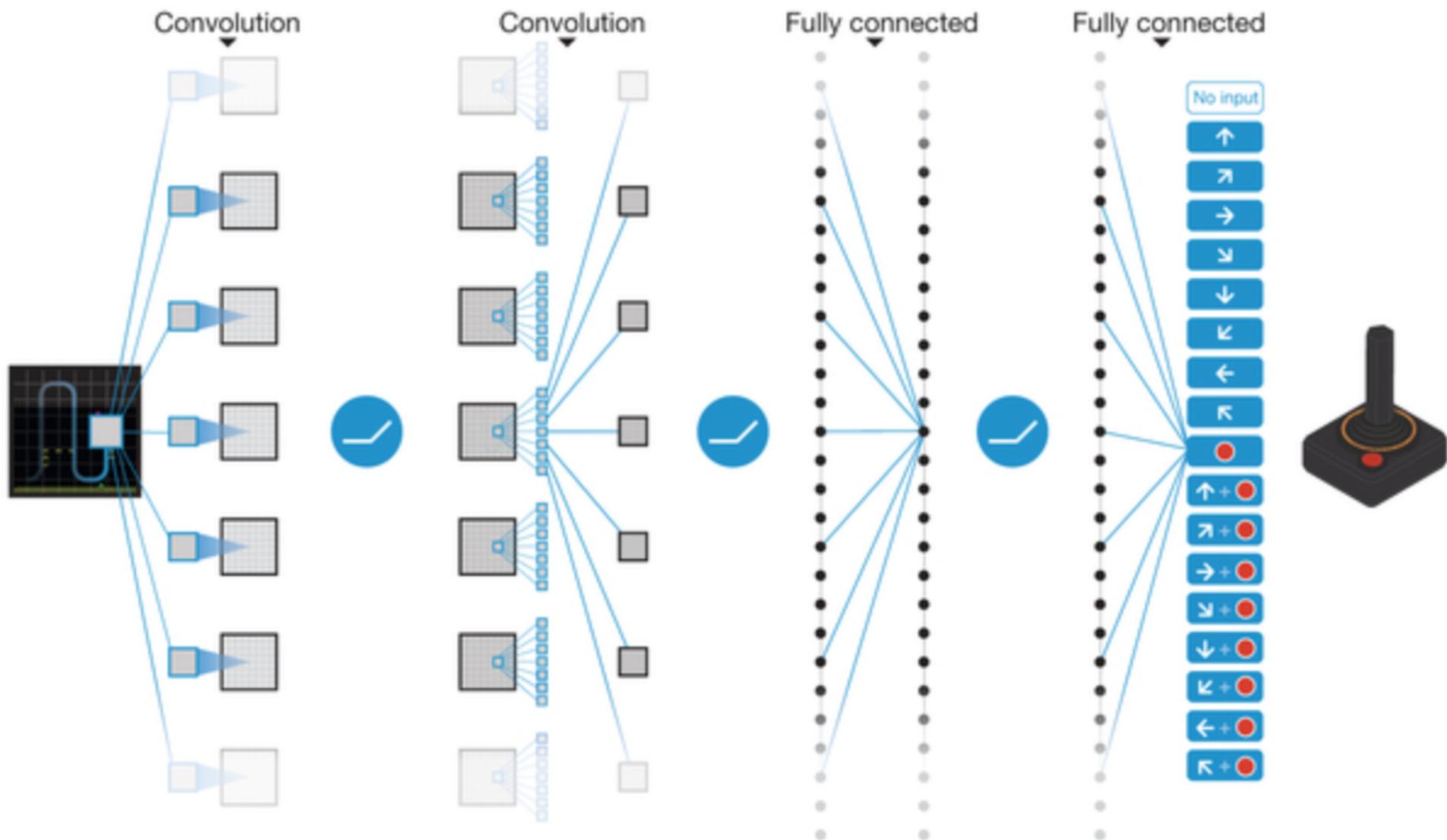
# Playing Atari with Deep Reinforcement Learning

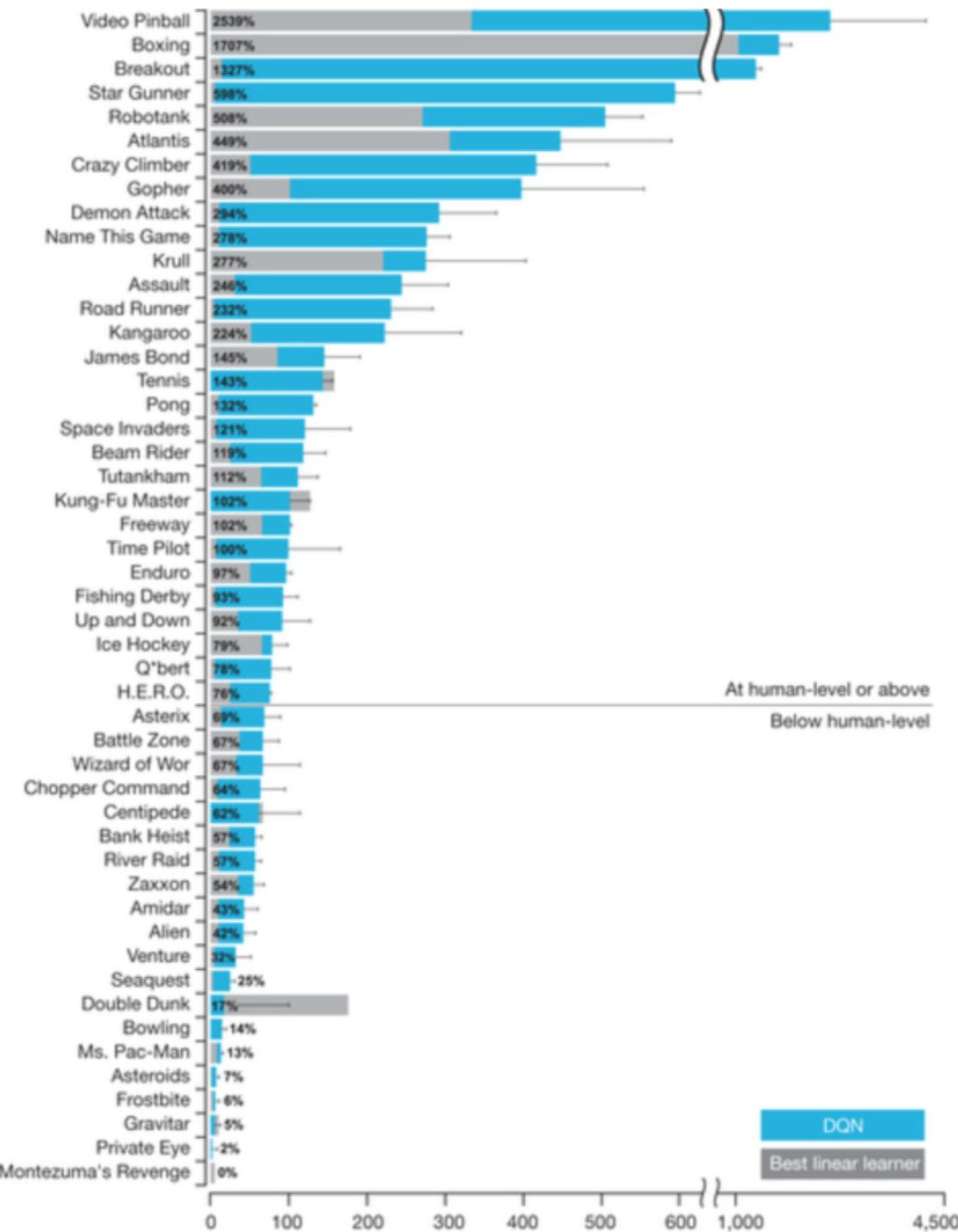
Volodymyr Mnih   Koray Kavukcuoglu   David Silver   Alex Graves   Ioannis Antonoglou

Daan Wierstra   Martin Riedmiller

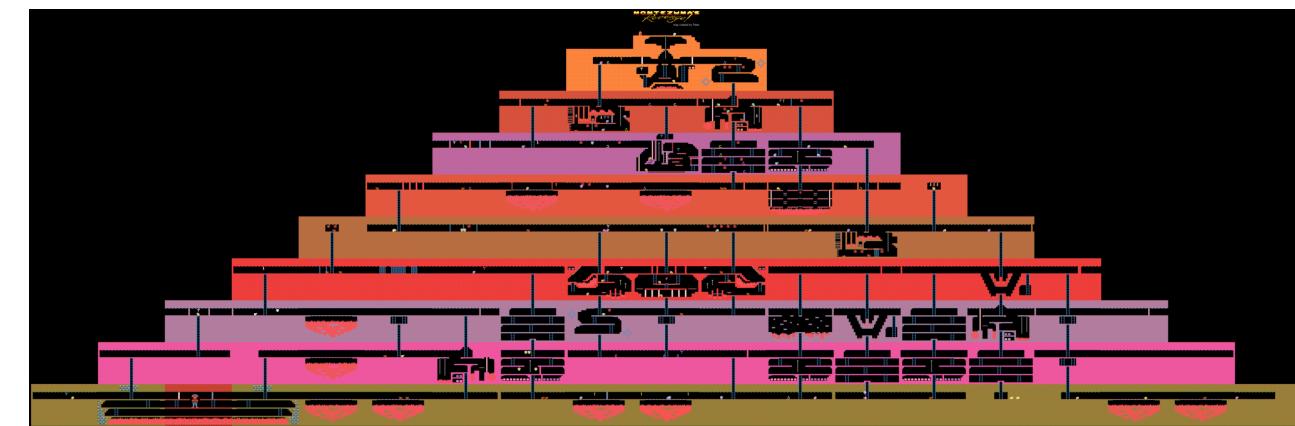
DeepMind Technologies







[https://www.youtube.com/watch?  
v=SuZVyOlgVek](https://www.youtube.com/watch?v=SuZVyOlgVek)



# Rainbow

# Rainbow

- Double Q Learning

# Rainbow

- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)

# Rainbow

- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$

# Rainbow

- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$
- Multi-step learning  
$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma \max Q_\theta(s_{t+n}, a') - Q_\theta(s_t, a_t))^2$$

# Rainbow

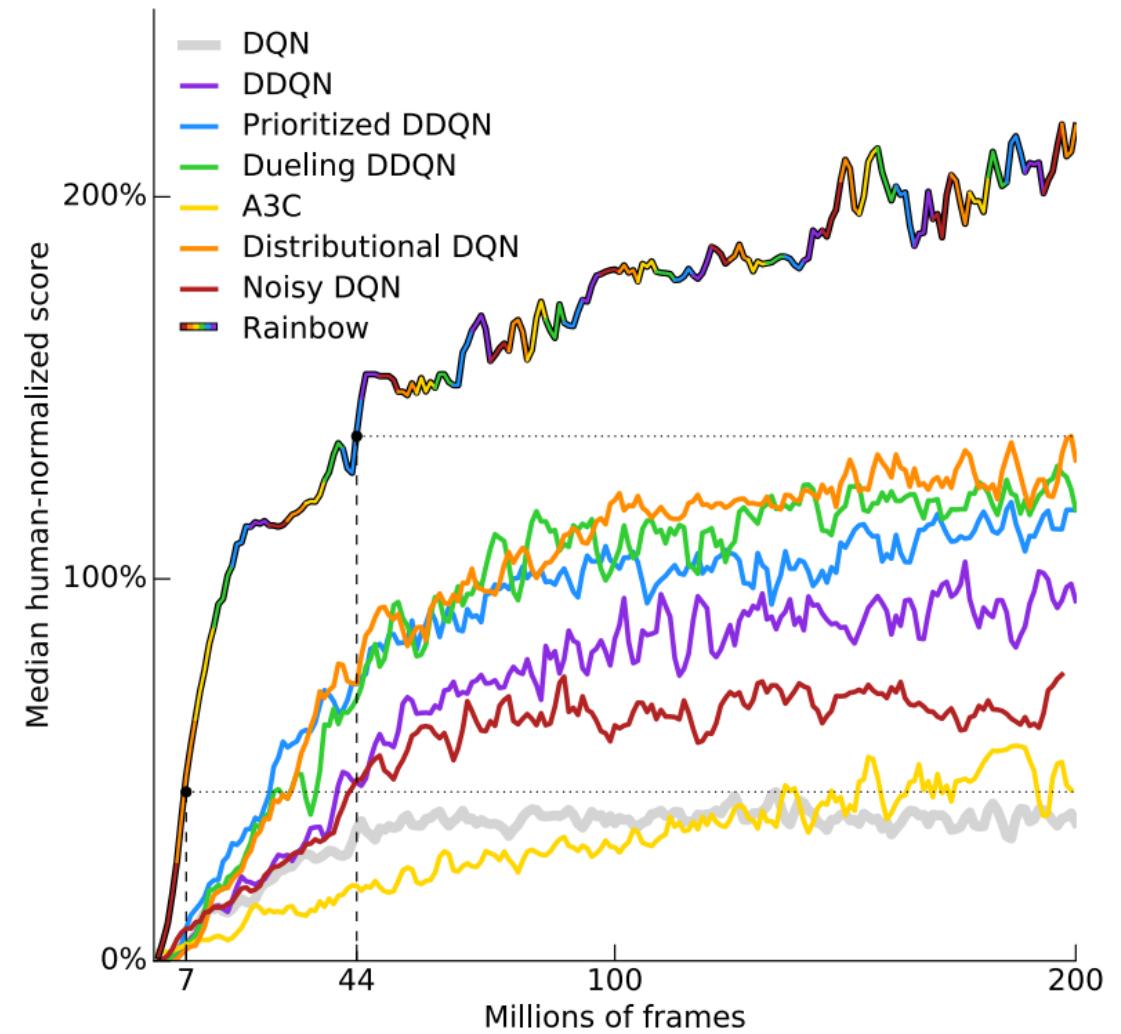
- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$
- Multi-step learning  
$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma \max Q_\theta(s_{t+n}, a') - Q_\theta(s_t, a_t))^2$$
- Distributional RL  
predict an entire distribution of values instead of just Q

# Rainbow

- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$
- Multi-step learning  
$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma \max Q_\theta(s_{t+n}, a') - Q_\theta(s_t, a_t))^2$$
- Distributional RL  
predict an entire distribution of values instead of just Q
- Noisy Nets

# Rainbow

- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$
- Multi-step learning  
$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma \max Q_\theta(s_{t+n}, a') - Q_\theta(s_t, a_t))^2$$
- Distributional RL  
predict an entire distribution of values instead of just Q
- Noisy Nets



# Actual Learning Curves

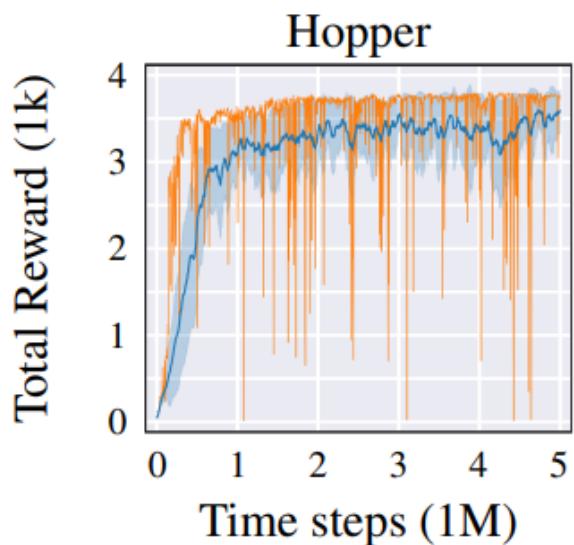


Figure 24: ■ Standard presentation ■ Single seed

The training curve as commonly presented in RL papers ■ and the training curve of a single seed ■. Both curves are from the TD3 algorithm, trained for 5M time steps. The standard presentation is to evaluate every  $N_{\text{freq}}$  steps, average scores over  $N_{\text{episodes}}$  evaluation episodes and  $N_{\text{seeds}}$  seeds, then to smooth the curve by averaging over a window of  $N_{\text{window}}$  evaluations. (In our case this corresponds to  $N_{\text{freq}} = 5000$ ,  $N_{\text{episodes}} = 10$ ,  $N_{\text{seeds}} = 10$ ,  $N_{\text{window}} = 10$ ). The learning curve of a single seed has no smoothing over seeds or evaluations ( $N_{\text{seeds}} = 1$ ,  $N_{\text{window}} = 1$ ). By averaging over many seeds and evaluations, the training curves in RL can appear deceptively smooth and stable.

Paper: [For SALE: State-Action Representation Learning for Deep Reinforcement Learning](#)

# Part II

# Improved Policy Gradients

# Restricted Gradient Update

# Restricted Gradient Update

$$\widehat{\nabla U}(\theta) = \sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k))$$

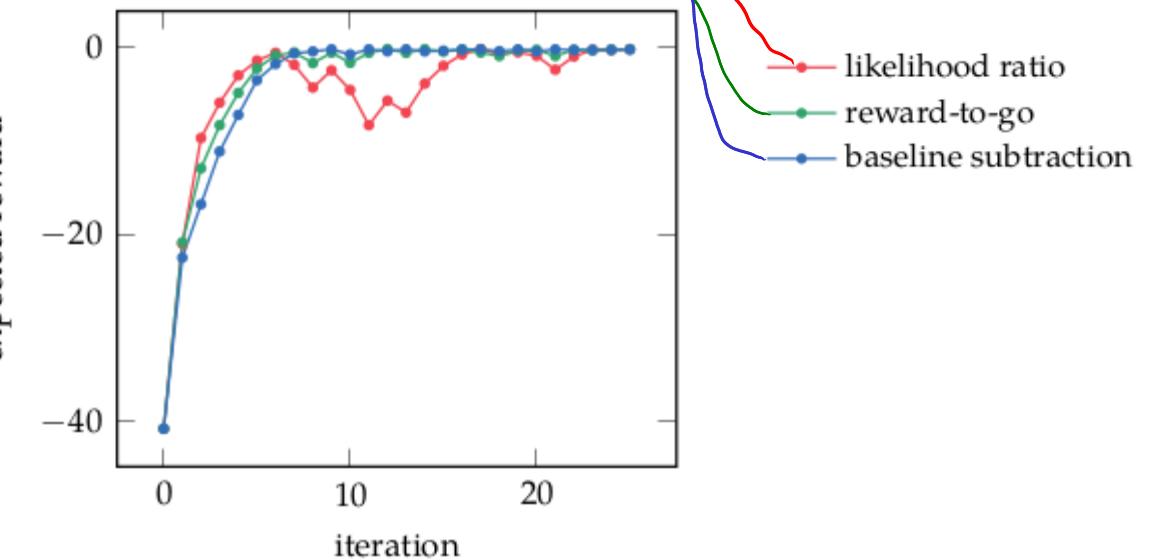
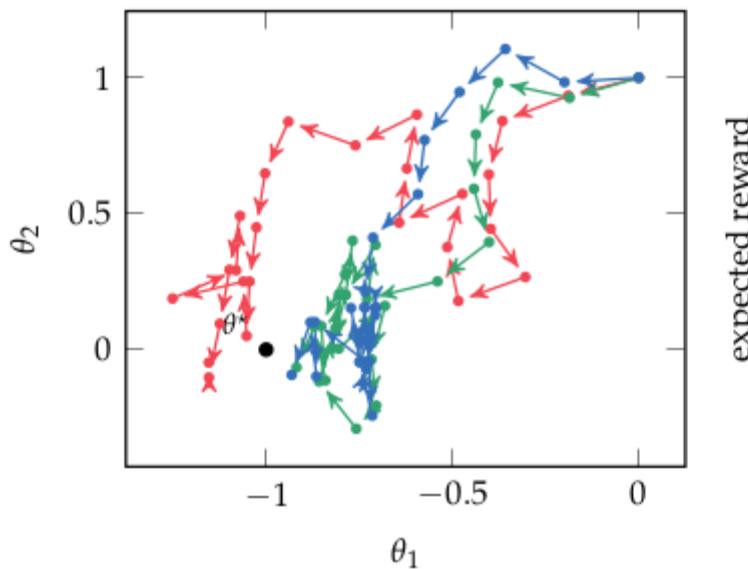
# Restricted Gradient Update

$$\widehat{\nabla U}(\theta) = \sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k))$$

$$\theta' = \theta + \alpha \widehat{\nabla U}(\theta)$$

# Restricted Gradient Update

$$\widehat{\nabla U}(\theta) = \sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k | s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k))$$
$$\theta' = \theta + \alpha \widehat{\nabla U}(\theta)$$



# Restricted Gradient Update

$$\widehat{\nabla U}(\theta) = \sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k))$$

$$\theta' = \theta + \alpha \widehat{\nabla U}(\theta)$$

maximize  $\theta'$        $U(\theta') \approx U(\theta) + \nabla U(\theta)^\top (\theta' - \theta)$

subject to       $g(\theta, \theta') \leq \varepsilon$

$$g(\theta, \theta') = \|\theta - \theta'\|_2^2 = \frac{1}{2} (\theta^\star - \theta)^\top (\theta' - \theta)$$

↗

$$\begin{aligned} \theta' &= \theta + u \sqrt{\frac{2\varepsilon}{u^\top u}} = \theta + \sqrt{2\varepsilon} \frac{u}{\|u\|} \\ u &\equiv \nabla U(\theta) \end{aligned}$$

# Natural Gradient

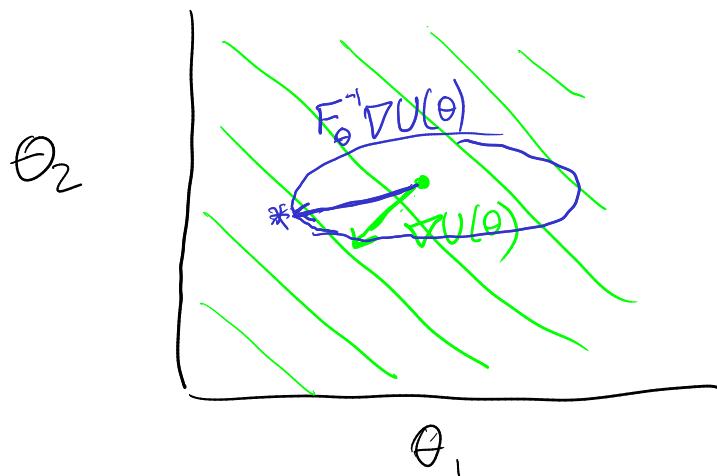
$$g(\theta, \theta') \approx D_{KL}(p(\tau|\theta) || p(\tau|\theta')) \leq \varepsilon$$

maximize  $\mathcal{U}(\theta') \approx \mathcal{U}(\theta) + \nabla \mathcal{U}(\theta)^T (\theta' - \theta)$

s.t.  $g(\theta, \theta') = \frac{1}{2} (\theta' - \theta)^T F_\theta (\theta' - \theta) \leq \varepsilon$

$$\theta' = \theta + u \sqrt{\frac{2\varepsilon}{\nabla \mathcal{U}(\theta)^T u}}$$

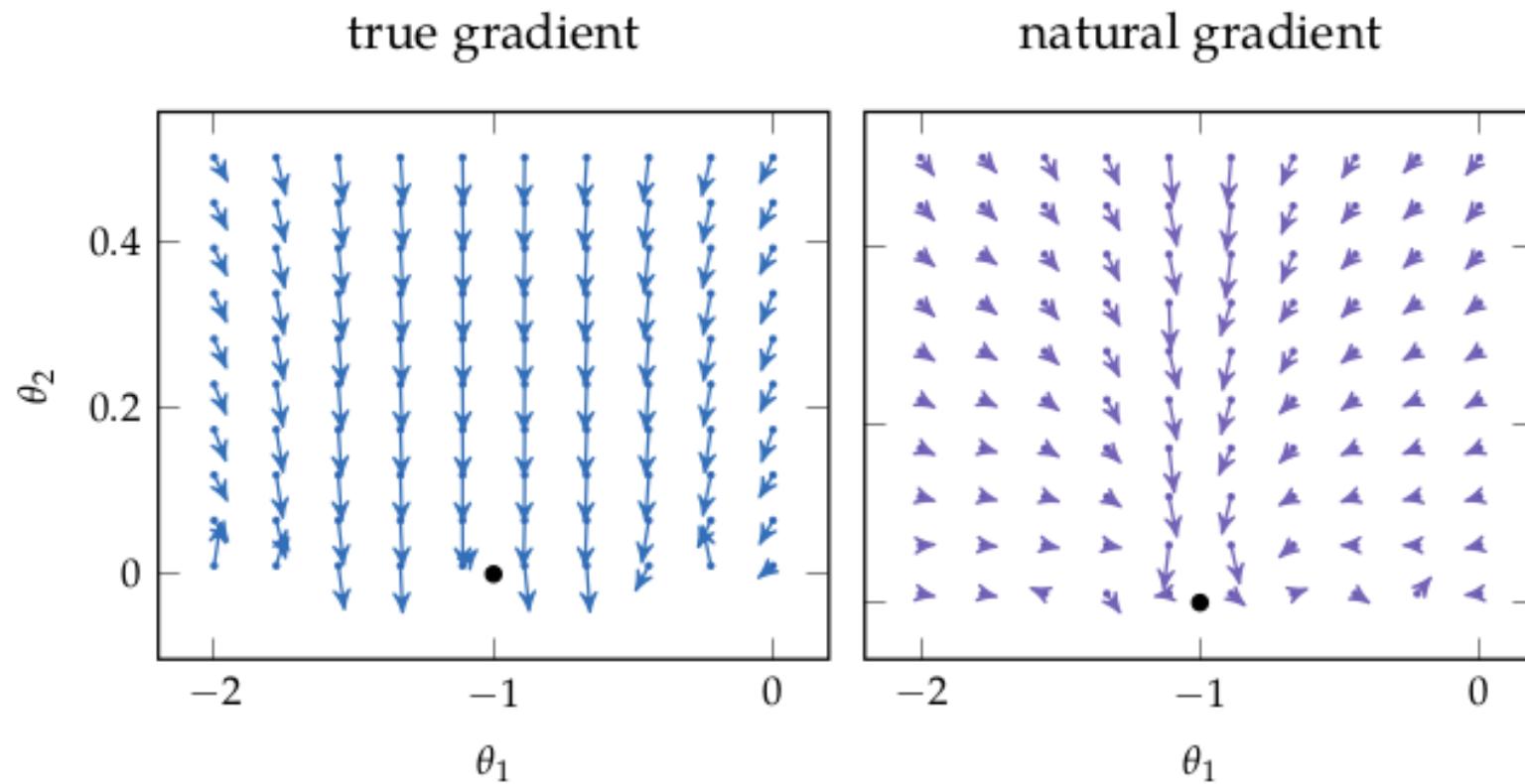
$$u = F_\theta^{-1} \nabla \mathcal{U}(\theta)$$



$$\begin{aligned} D_{KL}(p || q) &= \int p(x) \log \frac{p(x)}{q(x)} dx \\ &= - \int p(x) \log \frac{q(x)}{p(x)} dx \end{aligned}$$

$$\begin{aligned} F_\theta &= \int p(\tau|\theta) \nabla \log(p(\tau|\theta)) \nabla \log p(\tau|\theta) d\tau \\ &= E_{\tau} \left[ \underbrace{\nabla \log p(\tau|\theta)}_{\nabla \log p(\tau|\theta)} \nabla \log p(\tau|\theta)^T \right] \end{aligned}$$

# Natural Gradient



# TRPO and PPO

- likelihood ratio
- reward-to-go
- baseline subtraction

# TRPO and PPO

- likelihood ratio
- reward-to-go
- baseline subtraction

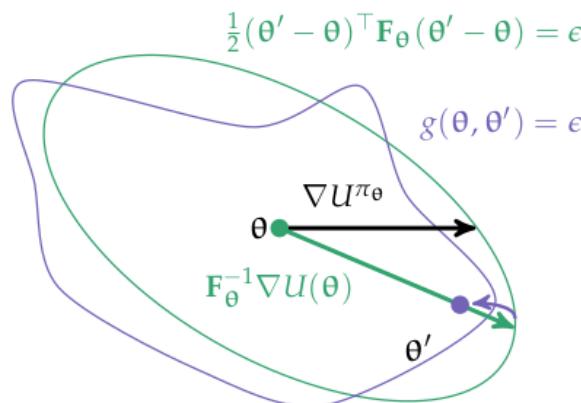
TRPO = Trust Region Policy Optimization

(Natural gradient + line search)

# TRPO and PPO

- likelihood ratio
- reward-to-go
- baseline subtraction

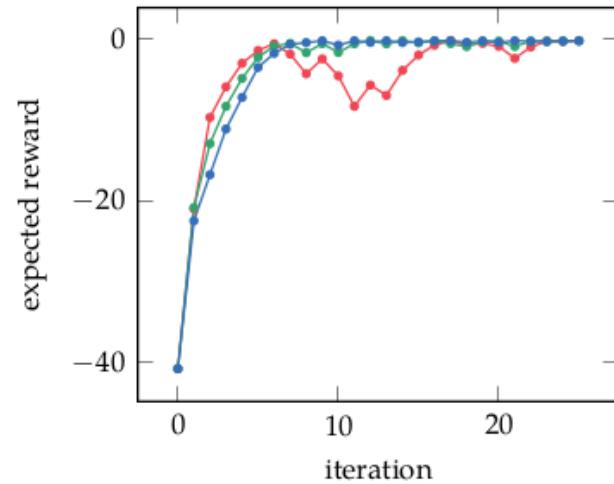
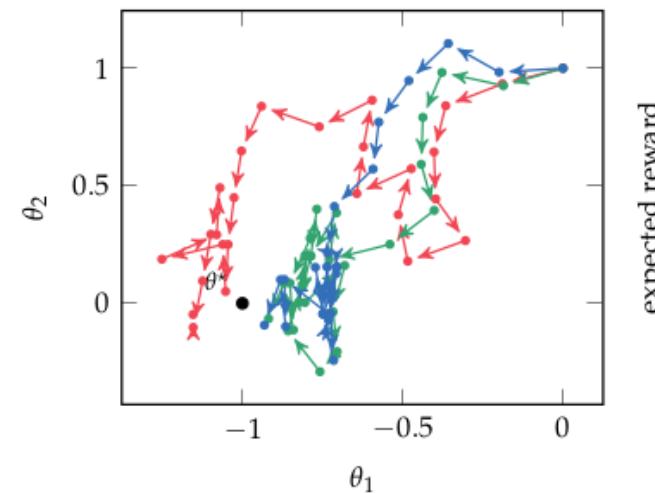
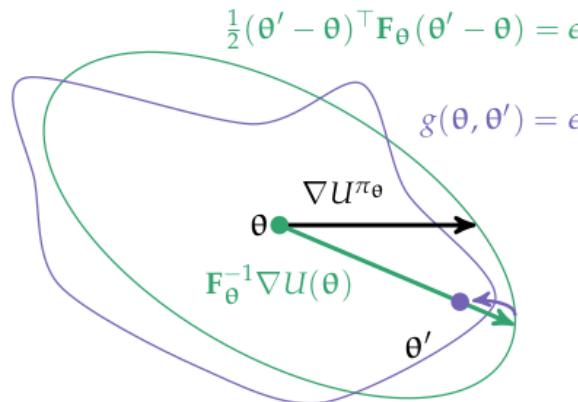
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



# TRPO and PPO

likelihood ratio  
reward-to-go  
baseline subtraction

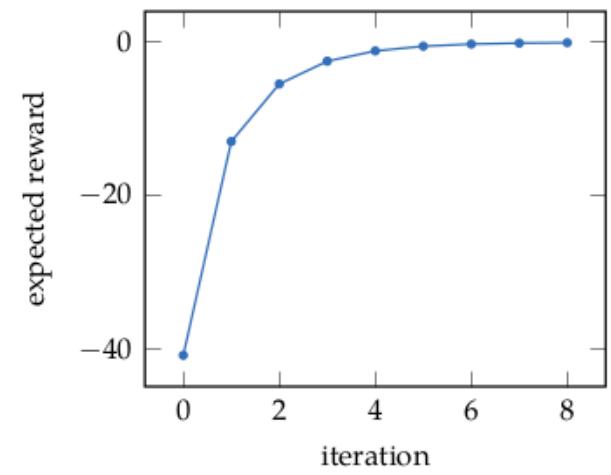
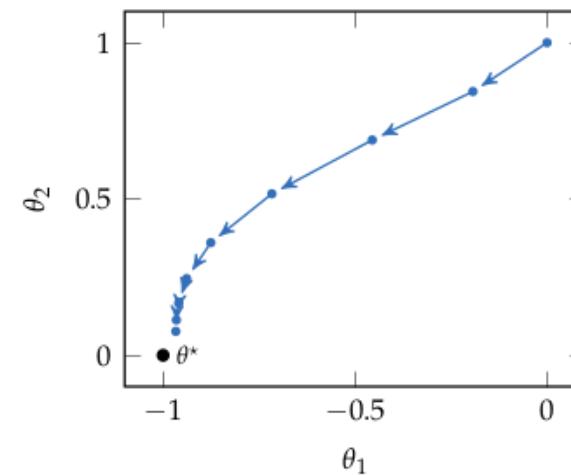
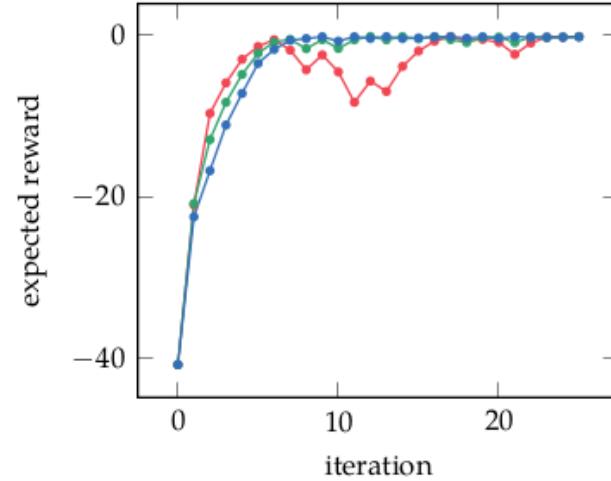
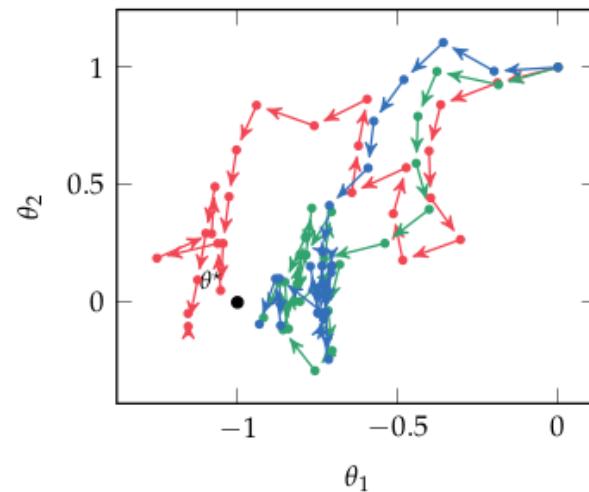
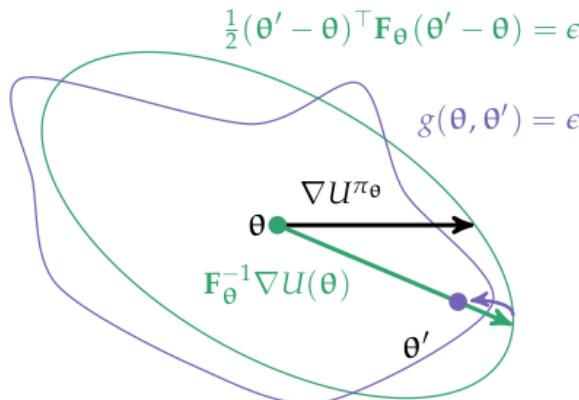
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



# TRPO and PPO

- likelihood ratio
- reward-to-go
- baseline subtraction

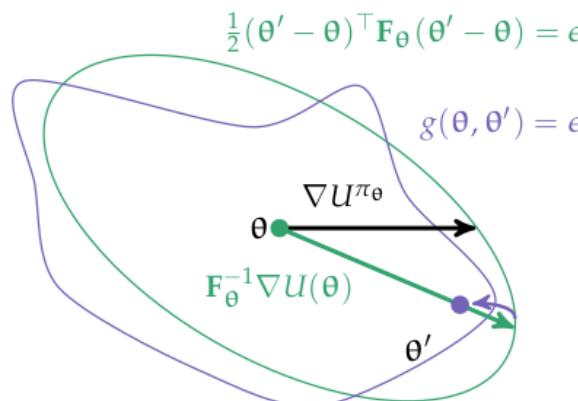
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



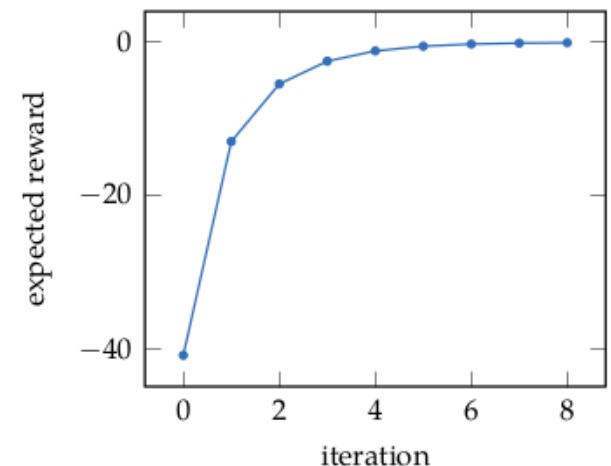
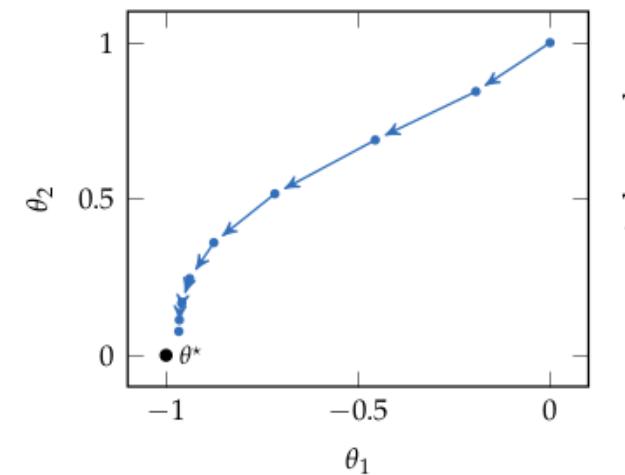
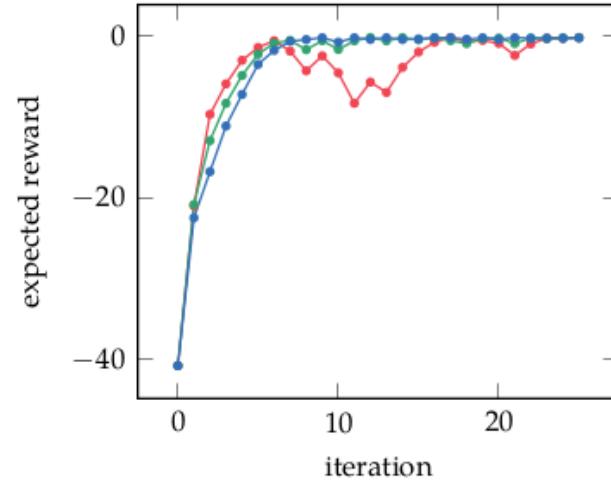
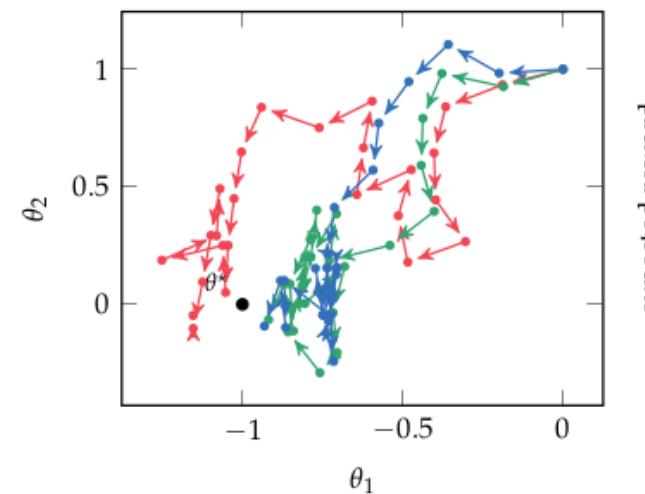
# TRPO and PPO

likelihood ratio  
reward-to-go  
baseline subtraction

TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



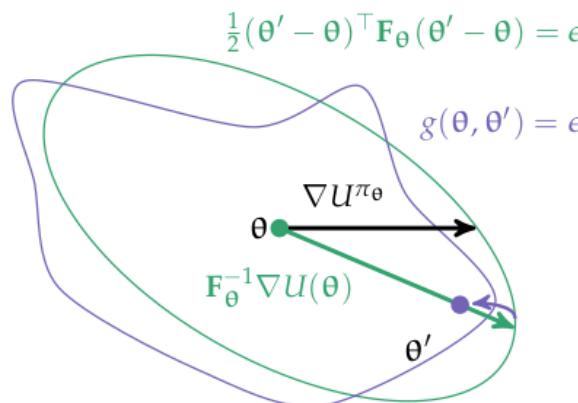
PPO = Proximal Policy Optimization  
(Use clamped surrogate objective to remove the need for line search)



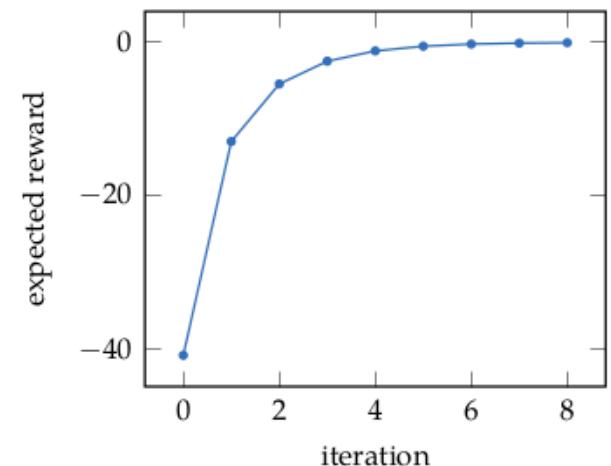
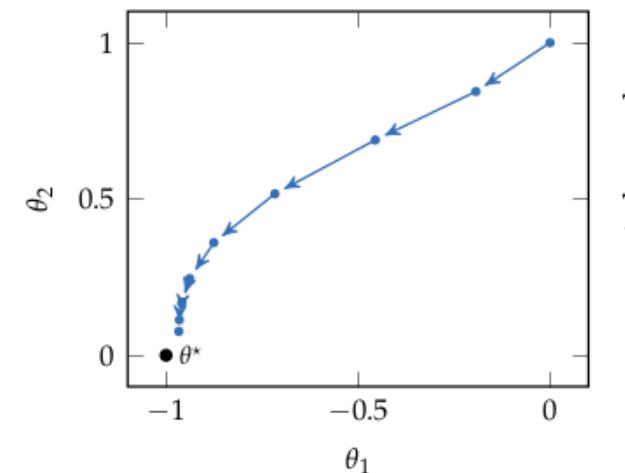
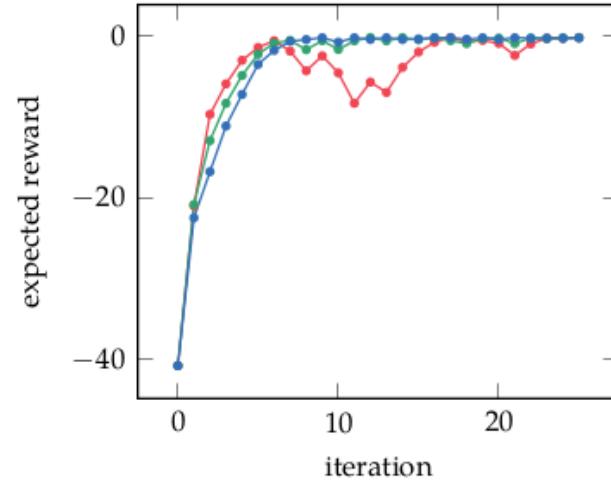
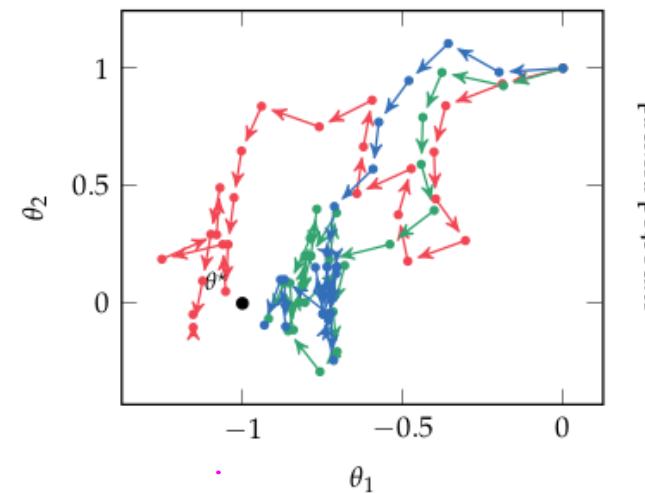
# TRPO and PPO

likelihood ratio  
reward-to-go  
baseline subtraction

TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



PPO = Proximal Policy Optimization  
(Use clamped surrogate objective to remove the need for line search)



# **Part III**

## **Actor-Critic**

# Actor-Critic

# Actor-Critic

Which should we learn?  $A$ ,  $Q$ , or  $V$ ?

# Actor-Critic

Which should we learn?  $A$ ,  $Q$ , or  $V$ ?

$$\nabla U(\theta) = E_{\tau} \left[ \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

# Actor-Critic

Which should we learn?  $A$ ,  $Q$ , or  $V$ ?

$$\nabla U(\theta) = E_{\tau} \left[ \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

*temporal difference residual*

-

# Actor-Critic

Which should we learn?  $A$ ,  $Q$ , or  $V$ ?

$$\nabla U(\theta) = E_{\tau} \left[ \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

*temporal difference residual*

$$l(\phi) = E \left[ (V_{\phi}(s) - V^{\pi_{\theta}}(s))^2 \right]$$

.

# Actor-Critic

Which should we learn?  $A$ ,  $Q$ , or  $V$ ?

$$\nabla U(\theta) = E_{\tau} \left[ \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

*temporal difference residual*

$$l(\phi) = E \left[ (V_{\phi}(s) - V^{\pi_{\theta}}(s))^2 \right]$$

*estimate with reward to go from sims*

# Generalized Advantage Estimation

# Recap

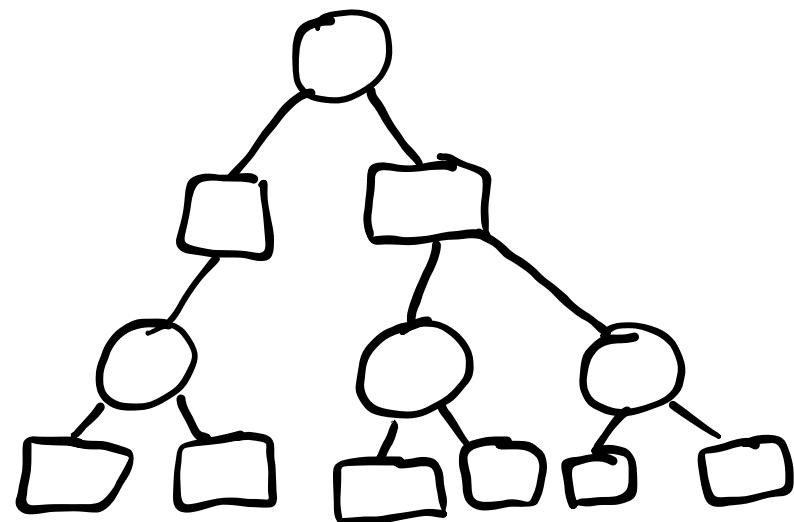
# Alpha Zero: Actor Critic with MCTS

# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS

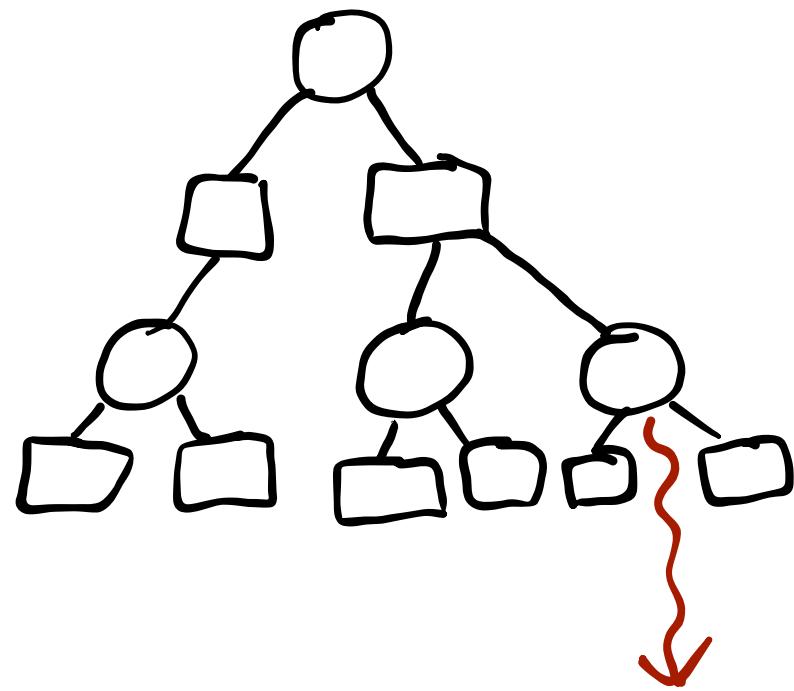
# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS



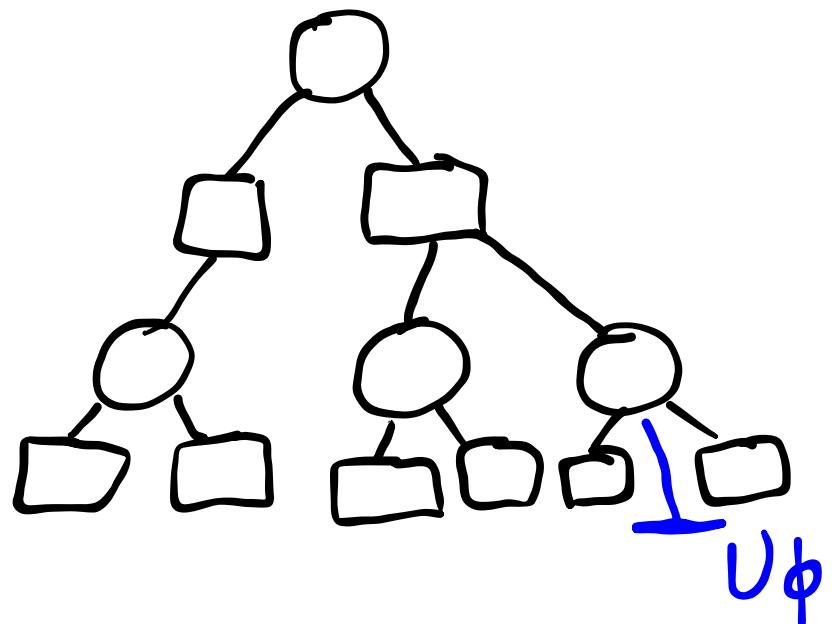
# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS



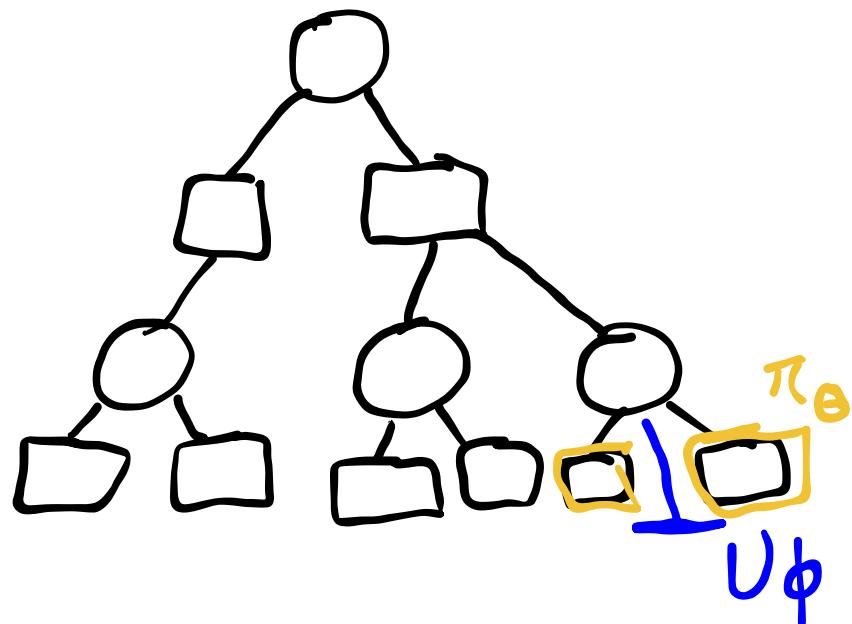
# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS



# Alpha Zero: Actor Critic with MCTS

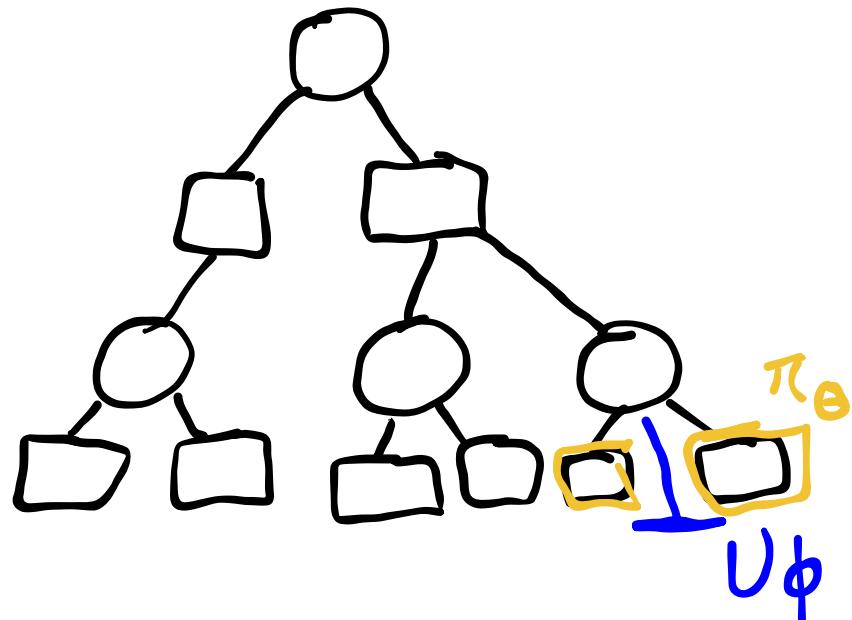
1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS



$$a = \arg \max_a Q(s, a) + c\pi_\theta(a \mid s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS
2. Learn  $\pi_\theta$  and  $U_\phi$  from tree

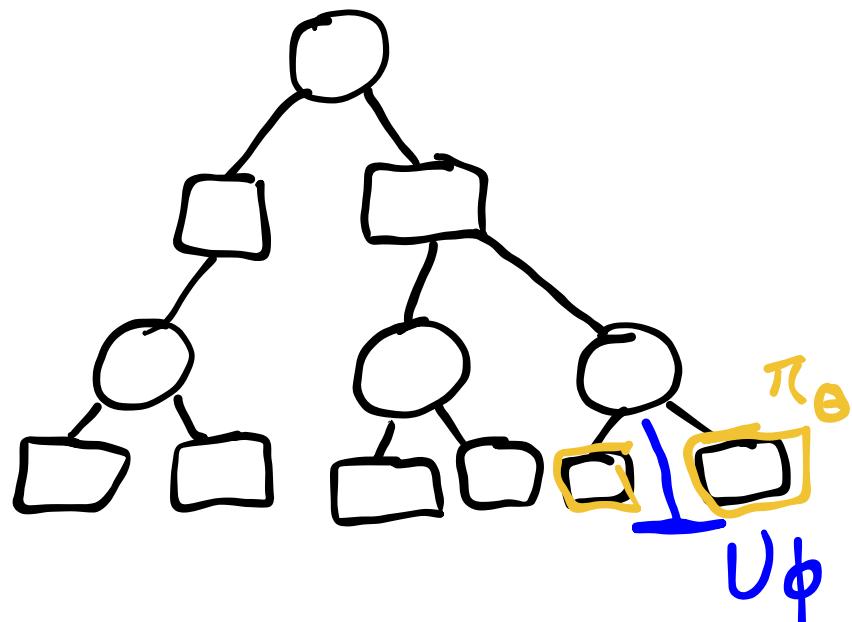


$$a = \arg \max_a Q(s, a) + c\pi_\theta(a \mid s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS
2. Learn  $\pi_\theta$  and  $U_\phi$  from tree

$$\ell(\boldsymbol{\theta}) = -\mathbb{E}_s \left[ \sum_a \pi_{\text{MCTS}}(a | s) \log \pi_\theta(a | s) \right]$$
$$\pi_{\text{MCTS}}(a | s) \propto N(s, a)^\eta$$

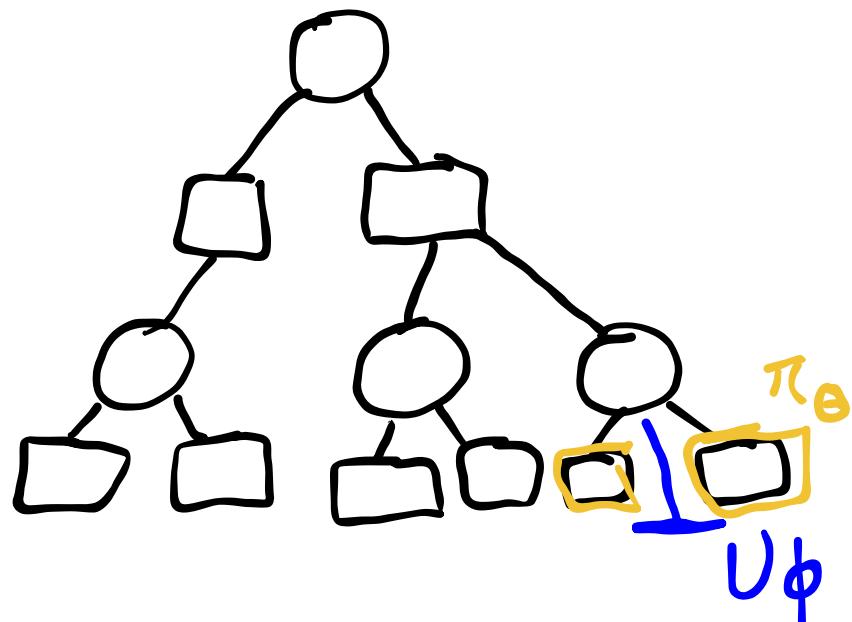


$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

# Alpha Zero: Actor Critic with MCTS

1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS
2. Learn  $\pi_\theta$  and  $U_\phi$  from tree

$$\ell(\boldsymbol{\theta}) = -\mathbb{E}_s \left[ \sum_a \pi_{\text{MCTS}}(a | s) \log \pi_\theta(a | s) \right]$$
$$\pi_{\text{MCTS}}(a | s) \propto N(s, a)^\eta$$



$$\ell(\boldsymbol{\Phi}) = \frac{1}{2} \mathbb{E}_s \left[ (U_\Phi(s) - U_{\text{MCTS}}(s))^2 \right]$$

$$U_{\text{MCTS}}(s) = \max_a Q(s, a)$$

$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$