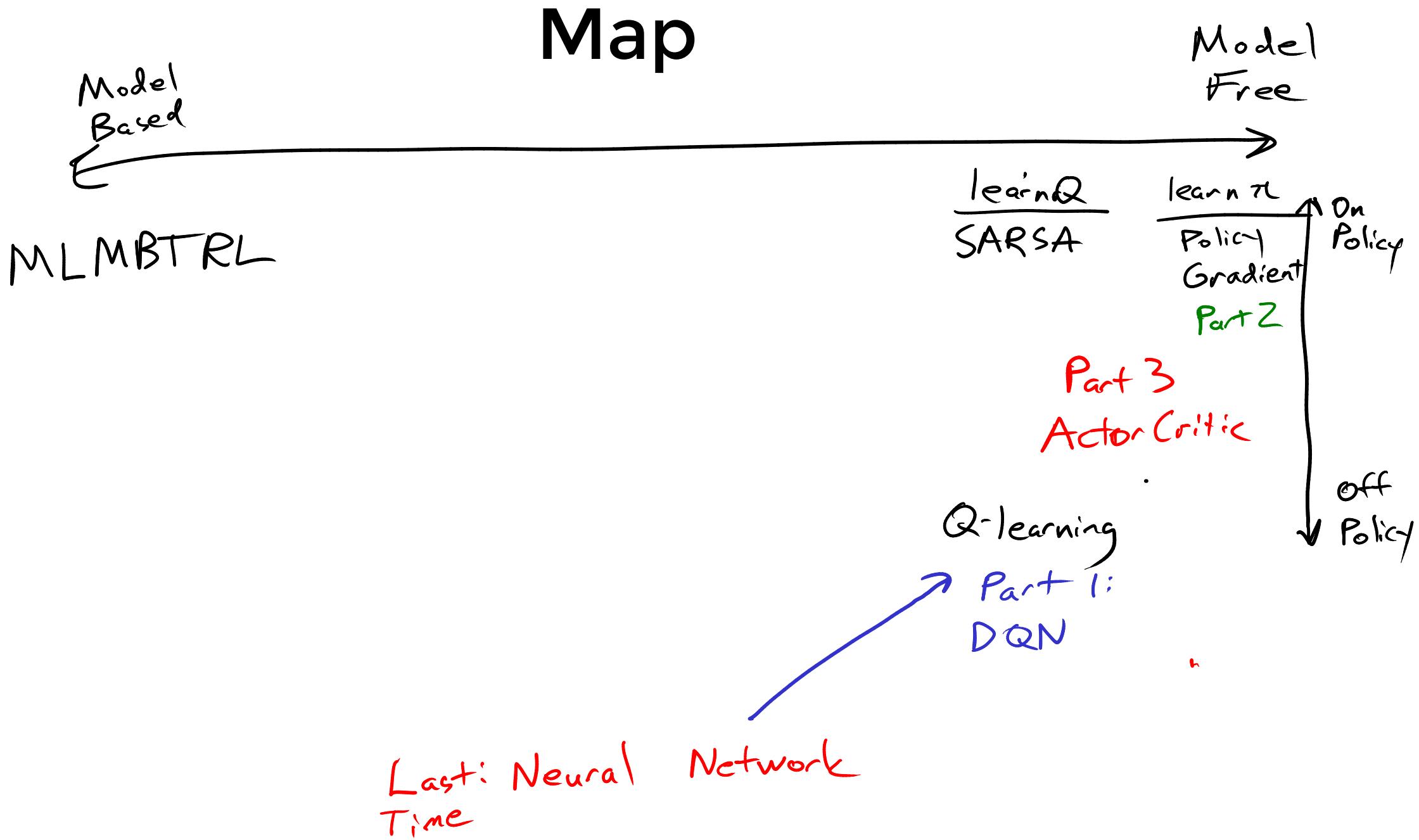


# Map



# **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

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

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

# 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:

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

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  $l$  be?

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

# 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  $l$  be?

Candidate Algorithm:

loop

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

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

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

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

$$s \leftarrow s'$$

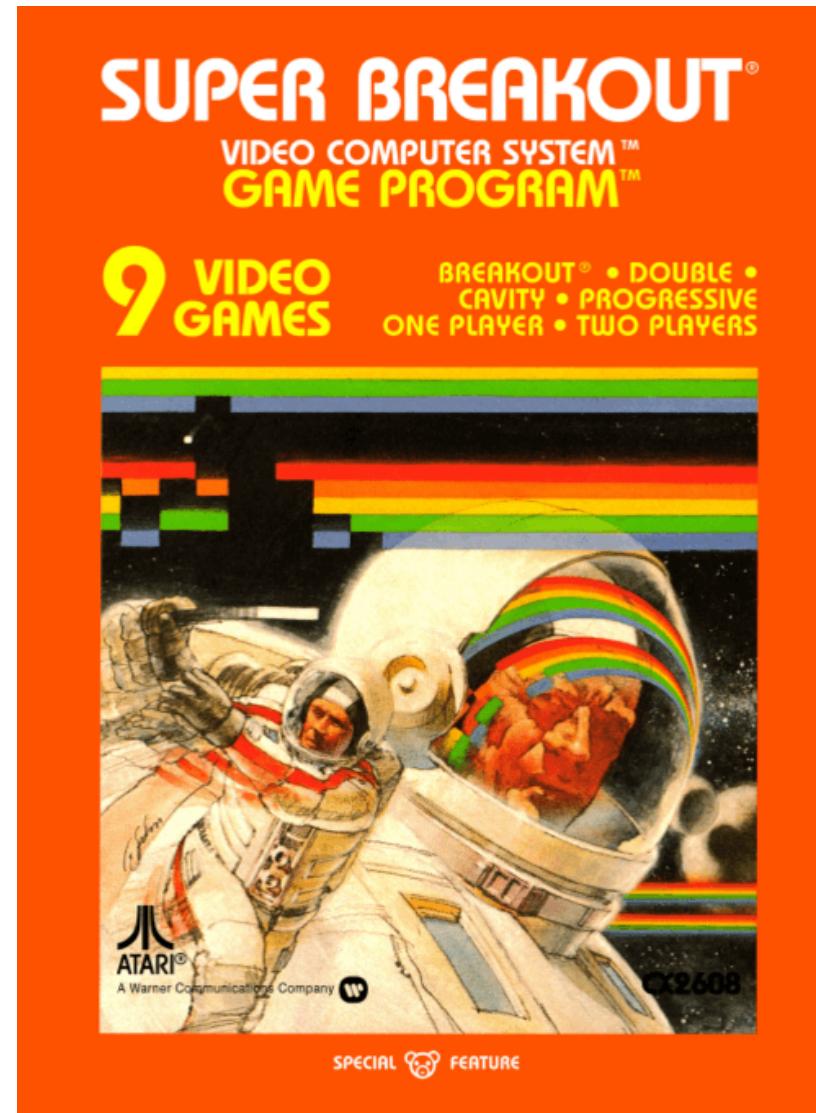


# DQN: The Atari Benchmark

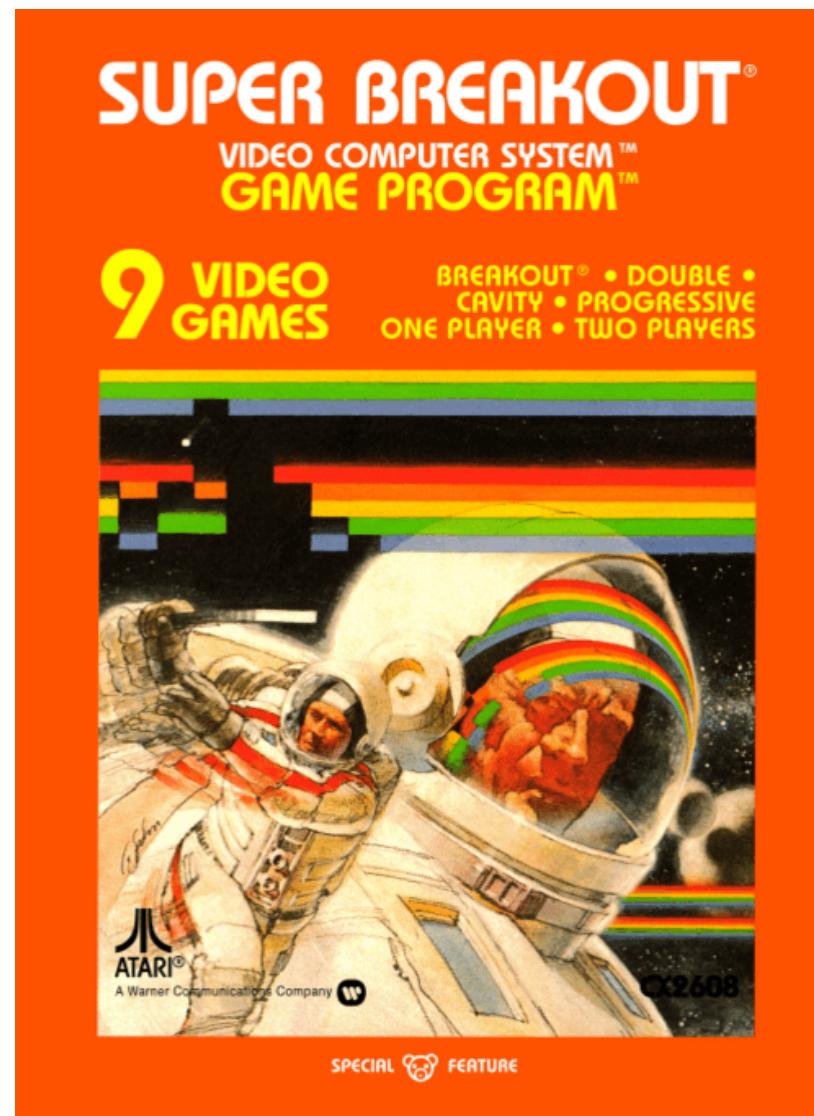
# DQN: The Atari Benchmark



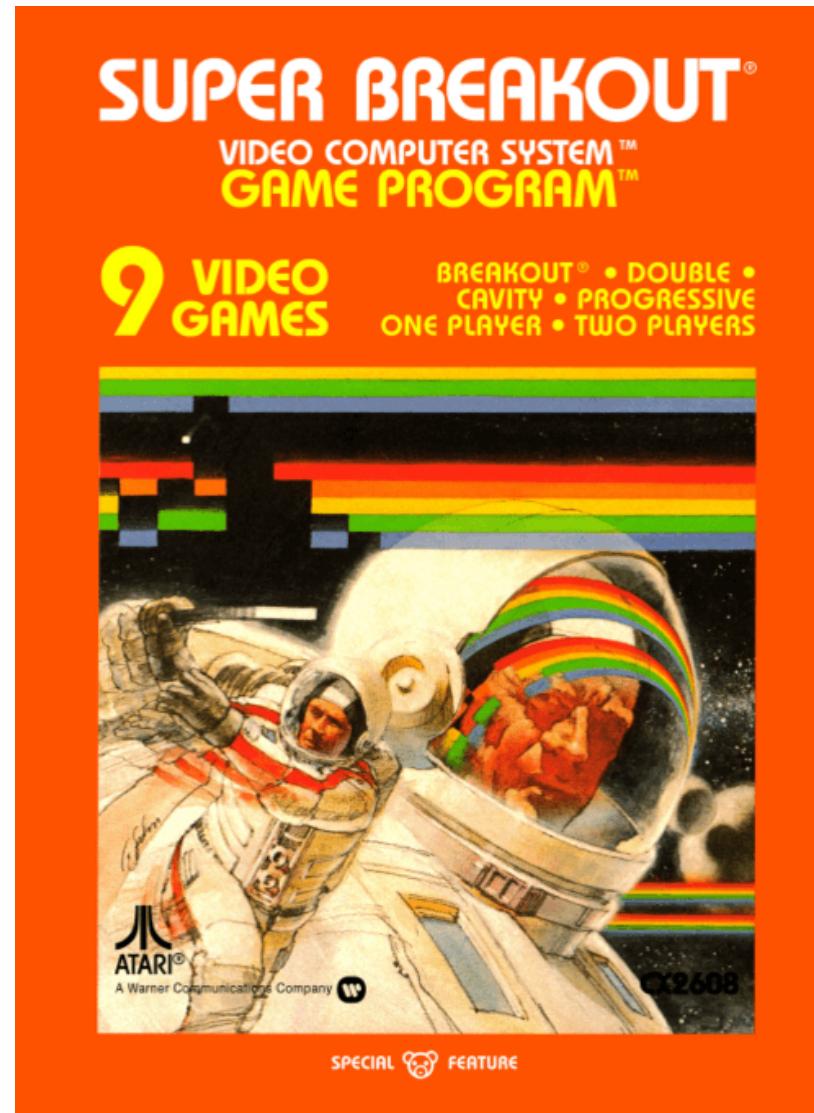
# DQN: The Atari Benchmark



# 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:

$s, a, r, s'$   
 $s', a', r', s''$

loop

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

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

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

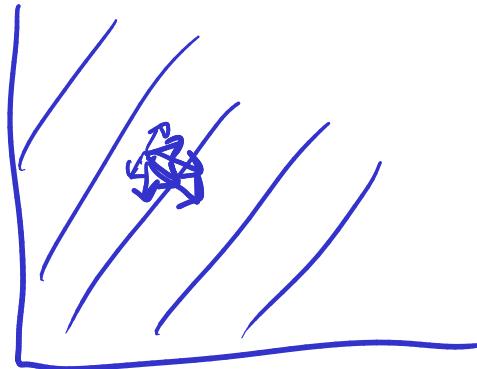
$\theta \leftarrow \theta - \alpha \nabla_{\theta} (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$

$s \leftarrow s'$

Problems:

1. Samples Highly Correlated

# DQN: Problems with Naive Approach



Candidate Algorithm:

loop

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

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

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

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

$s \leftarrow s'$

Problems:

1. Samples Highly Correlated
2. Size-1 batches

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

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

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

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

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

Problems:

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

# DQN: Problems with Naive Approach

Candidate Algorithm:

loop

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

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

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

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

$$s \leftarrow s'$$

Problems:

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

} data buffer / experience replay  
} periodically freeze target

# DQN

# DQN

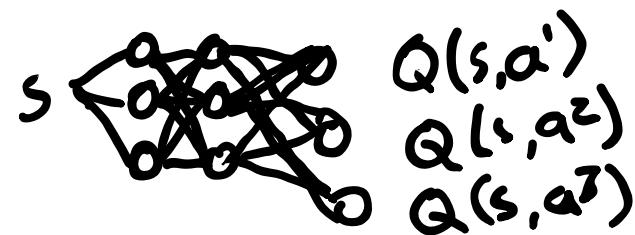
Q Network Structure:

Naive



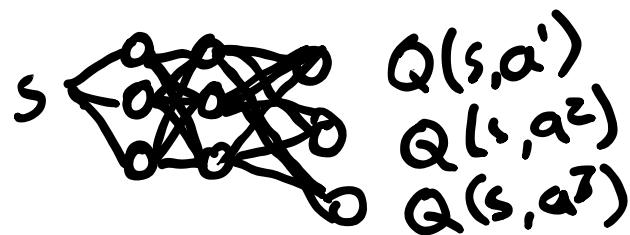
## DQN

Q Network Structure:



# DQN

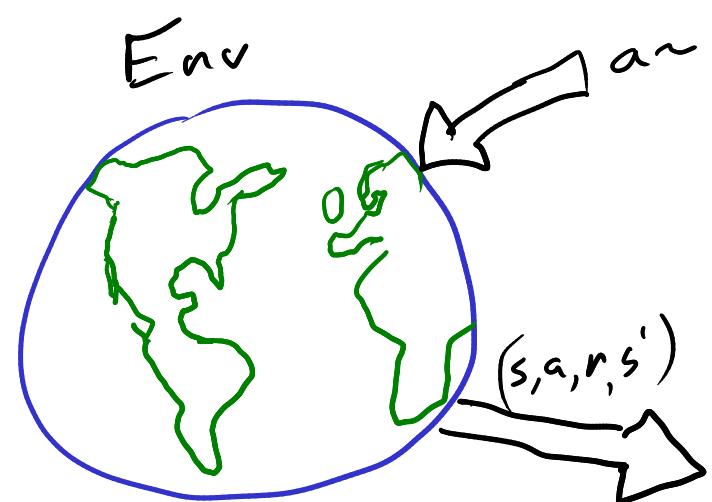
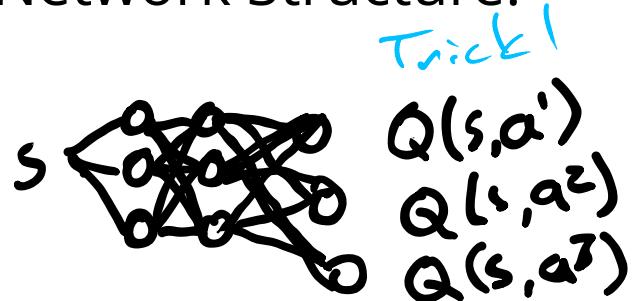
Q Network Structure:



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

# DQN

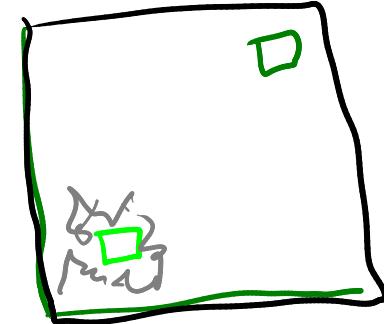
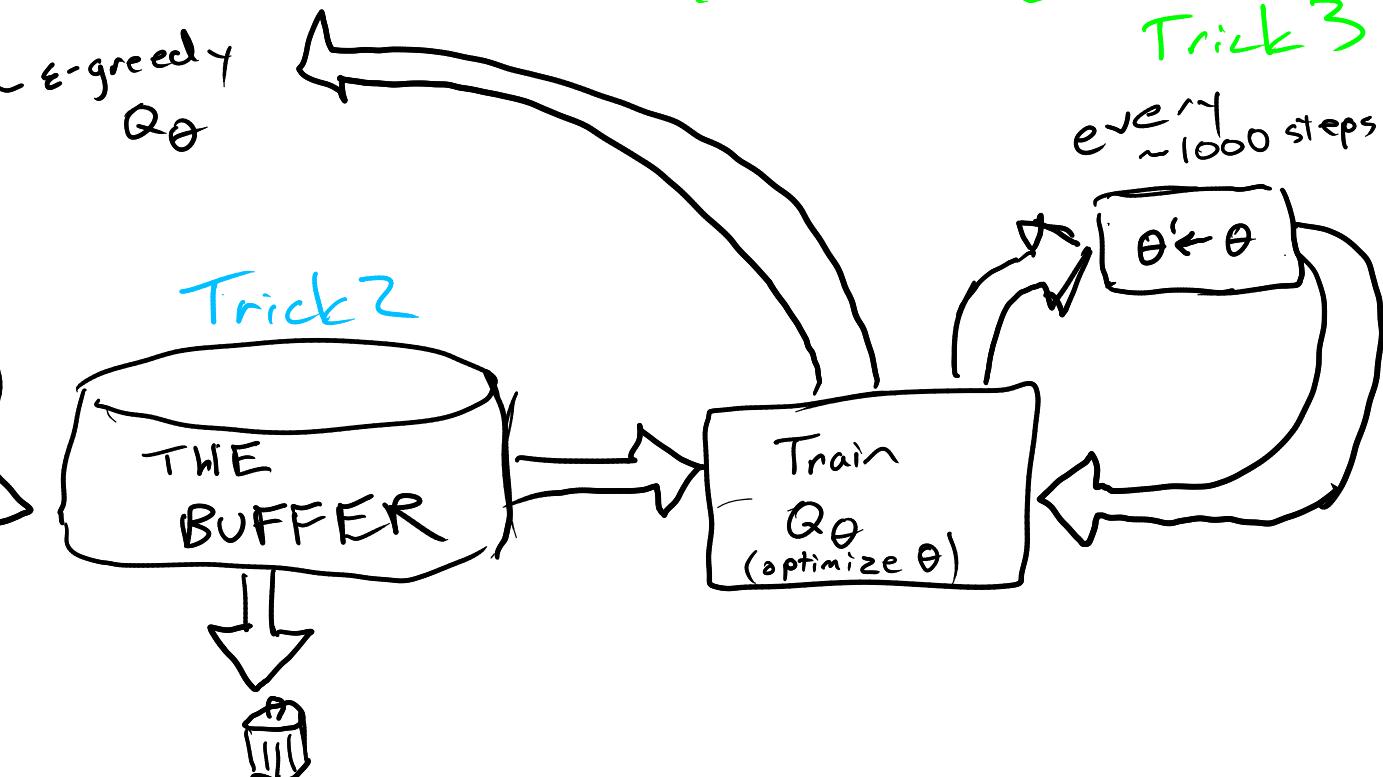
Q Network Structure:



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

Loss:

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

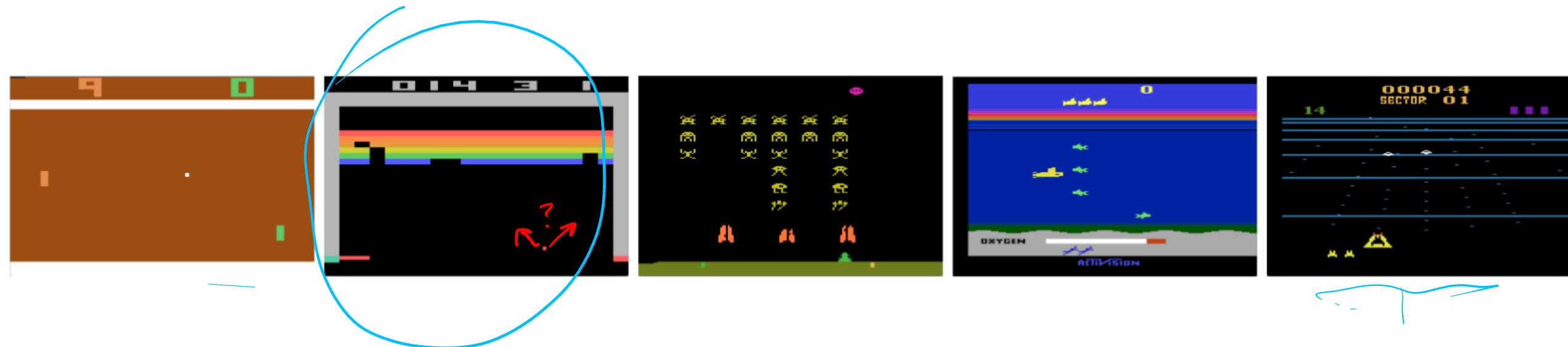


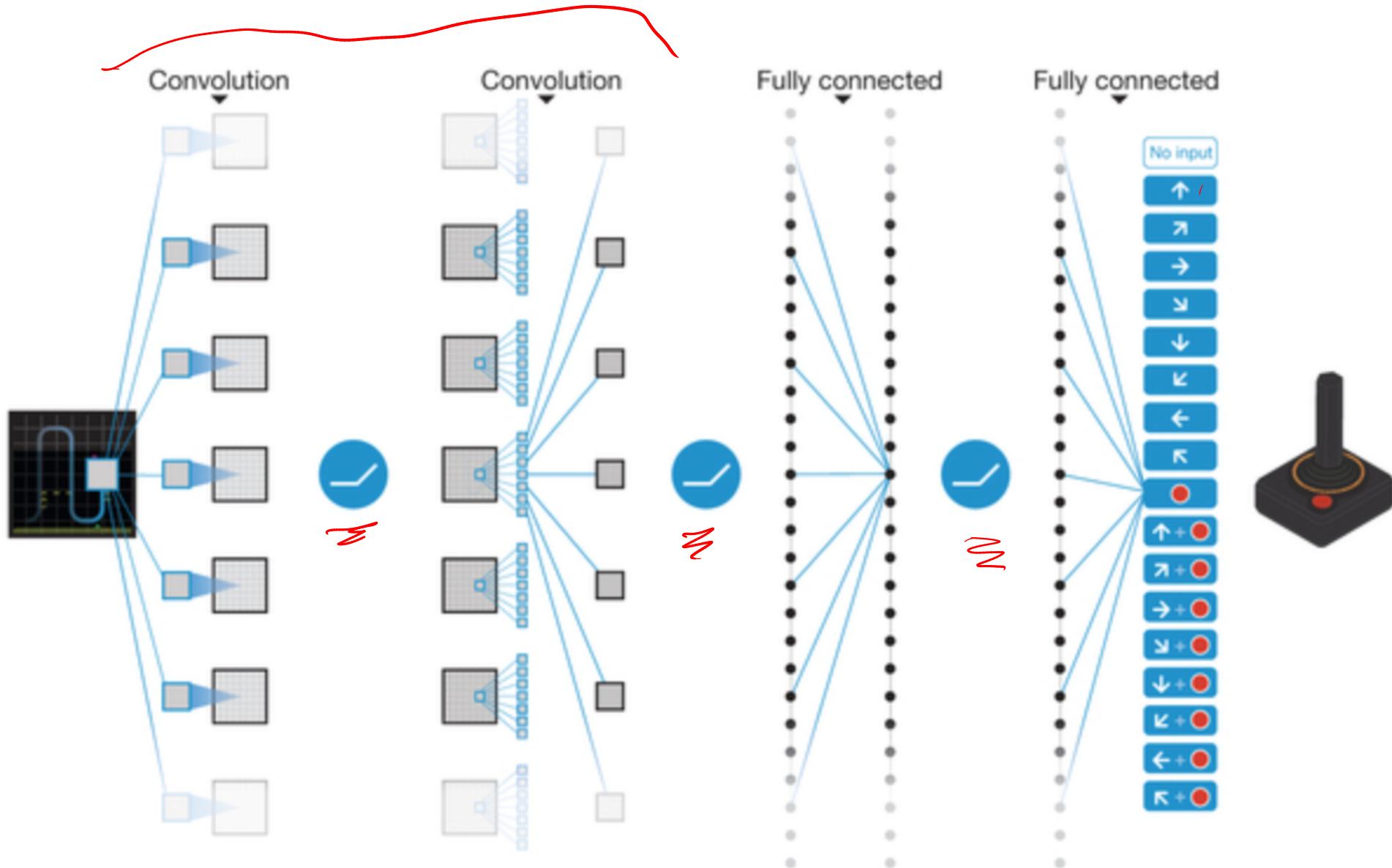
# Playing Atari with Deep Reinforcement Learning

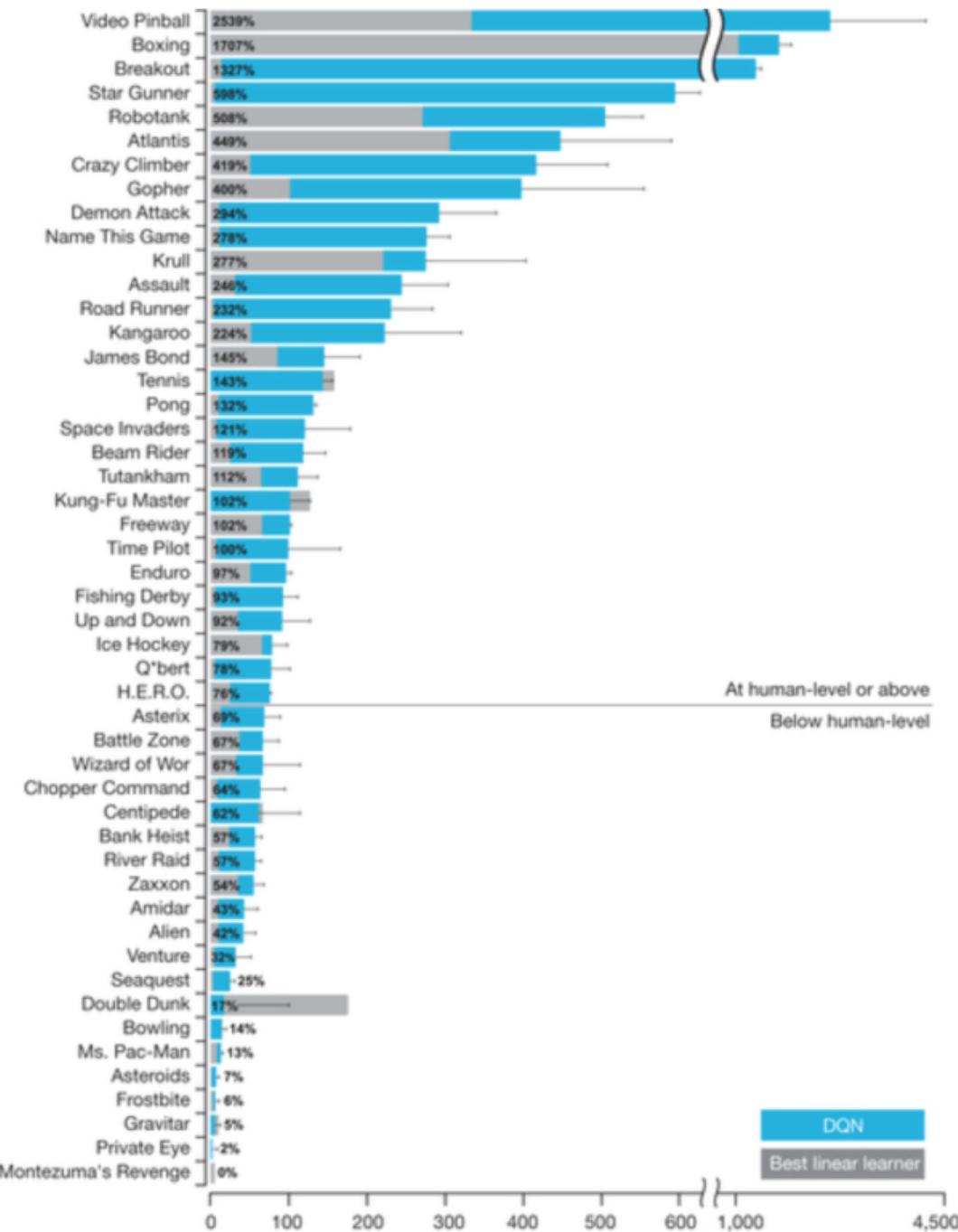
Volodymyr Mnih   Koray Kavukcuoglu   David Silver   Alex Graves   Ioannis Antonoglou

Daan Wierstra   Martin Riedmiller

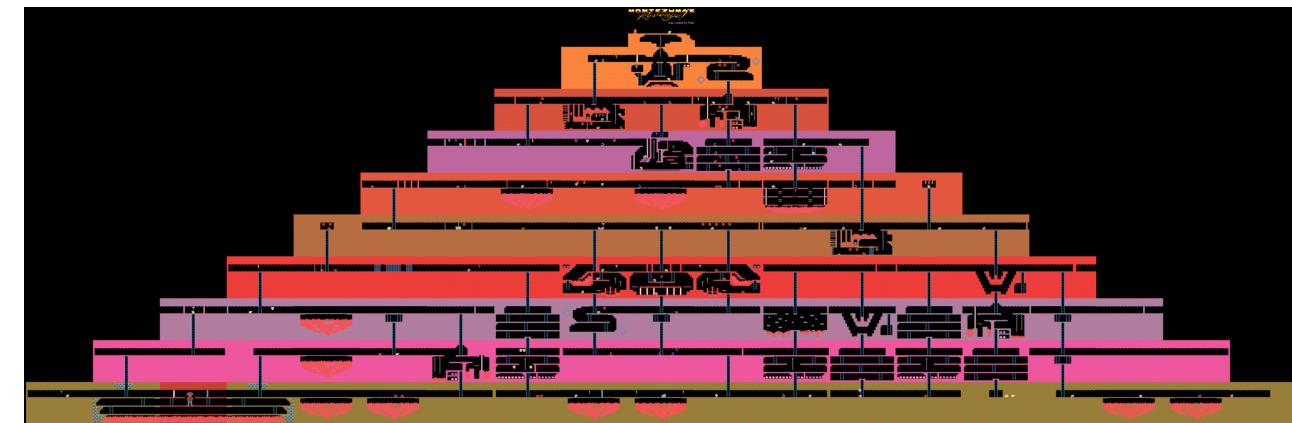
DeepMind Technologies







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



# Rainbow

# Rainbow

$$r + \gamma Q_{\theta_1}(s', \arg \max Q_{\theta_2}(s', a)) - Q_{\theta_1}(s, a)$$

- 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

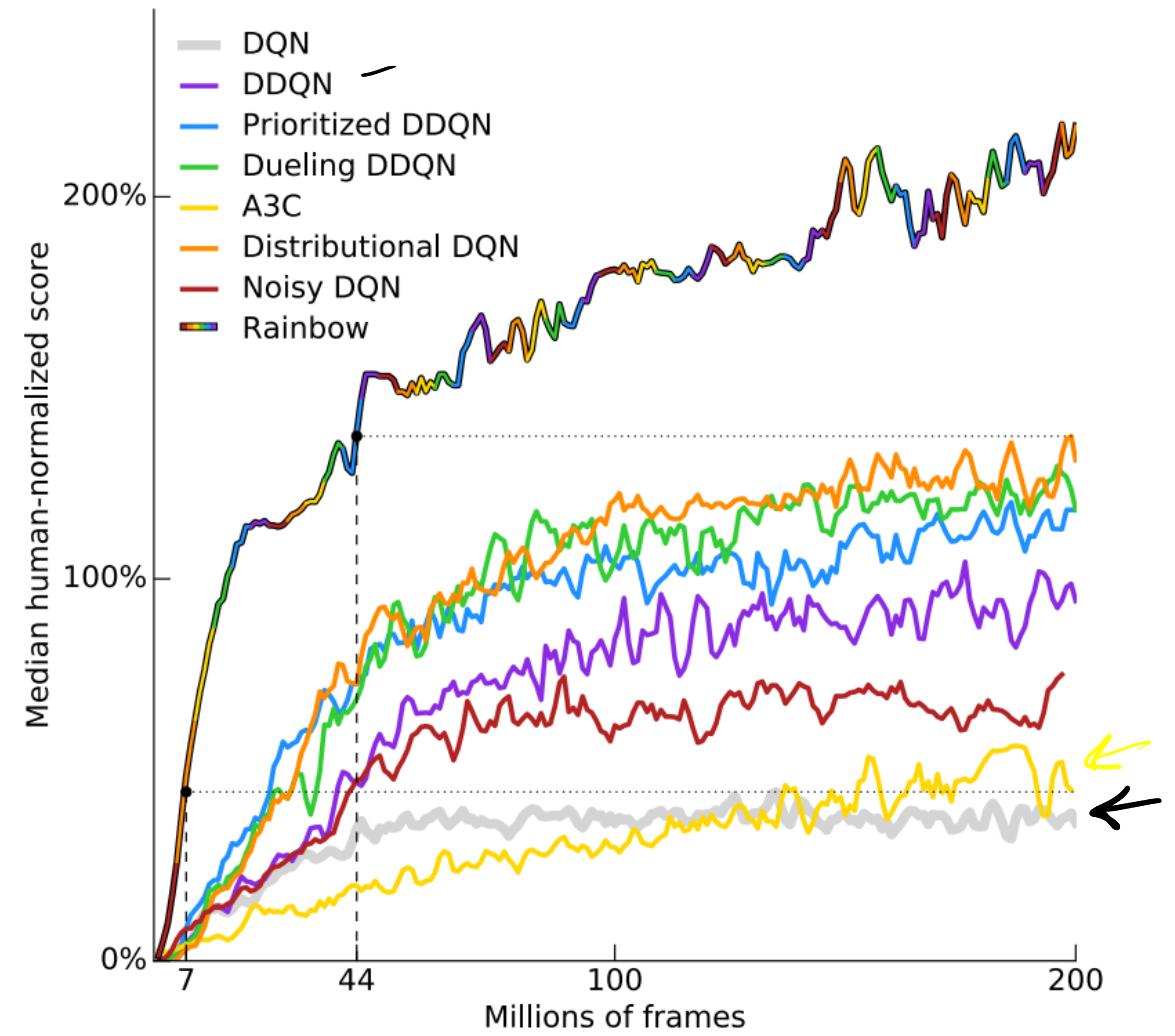
- Double Q Learning
- Prioritized Replay  
(priority proportional to last TD error)
- Dueling networks  
Value network + advantage network  
$$Q(s, a) = V(s) + A(s, a)$$
- Multi-step learning  
$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma \max Q_\theta(s_{t+n}, a') - Q_\theta(s_t,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



# Part II

# Improved Policy Gradients

# Restricted Gradient Update

# Restricted Gradient Update

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

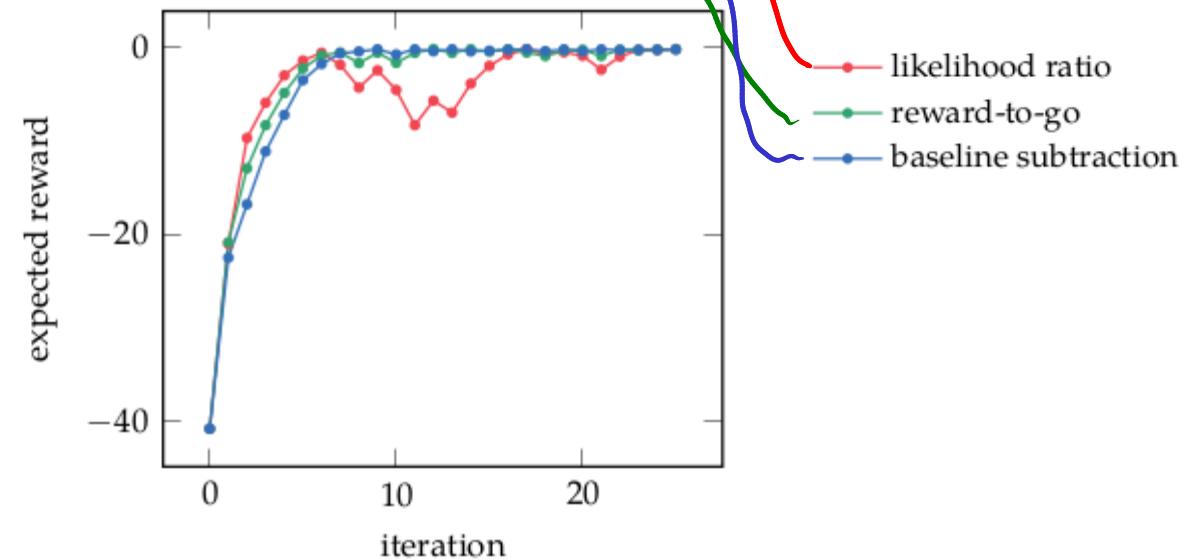
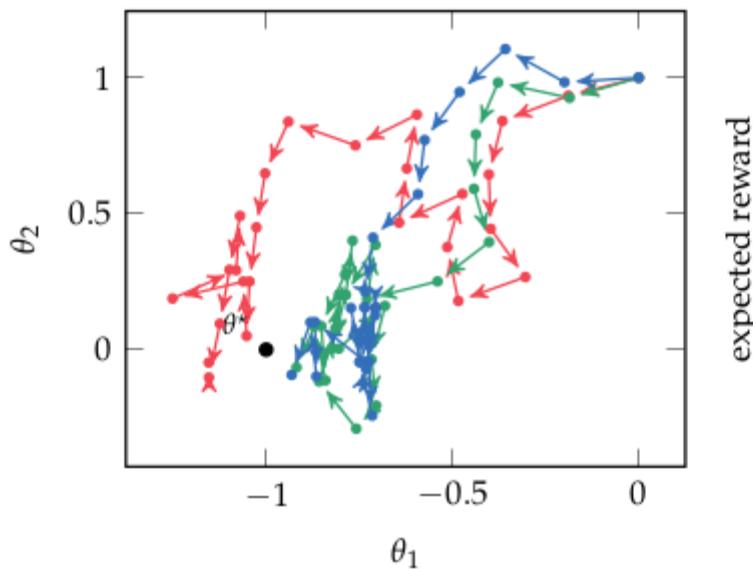
# Restricted Gradient Update

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

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

# 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

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

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

1-step

$$\begin{aligned} & \underset{\theta'}{\text{maximize}} \quad V(\theta') \approx V(\theta) + \nabla V(\theta)^T (\theta' - \theta) \\ & \text{subject to} \quad g(\theta, \theta') \leq \varepsilon \end{aligned}$$

$$g(\theta, \theta') = \|\theta - \theta'\|_2^2 = \frac{1}{2} (\theta' - \theta)^T (\theta' - \theta)$$

$$\begin{aligned} \theta' &= \theta + u \sqrt{\frac{2\varepsilon}{u^T u}} = \theta + \sqrt{2\varepsilon} \frac{u}{\|u\|} \\ u &= \nabla V(\theta) \end{aligned}$$

# Natural Gradient

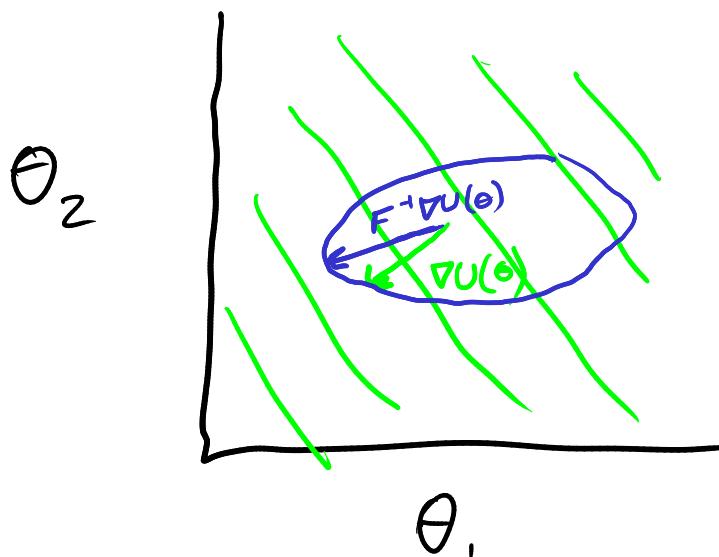
$$g(\theta, \theta') = D_{KL}(p(\tau|\theta) || p(\tau|\theta')) \leq \varepsilon$$

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

s.t.

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

$$\theta' = \theta + u \sqrt{\frac{2\varepsilon}{\nabla U(\theta)^T u}} \quad u = F_\theta^{-1} \nabla U(\theta)$$

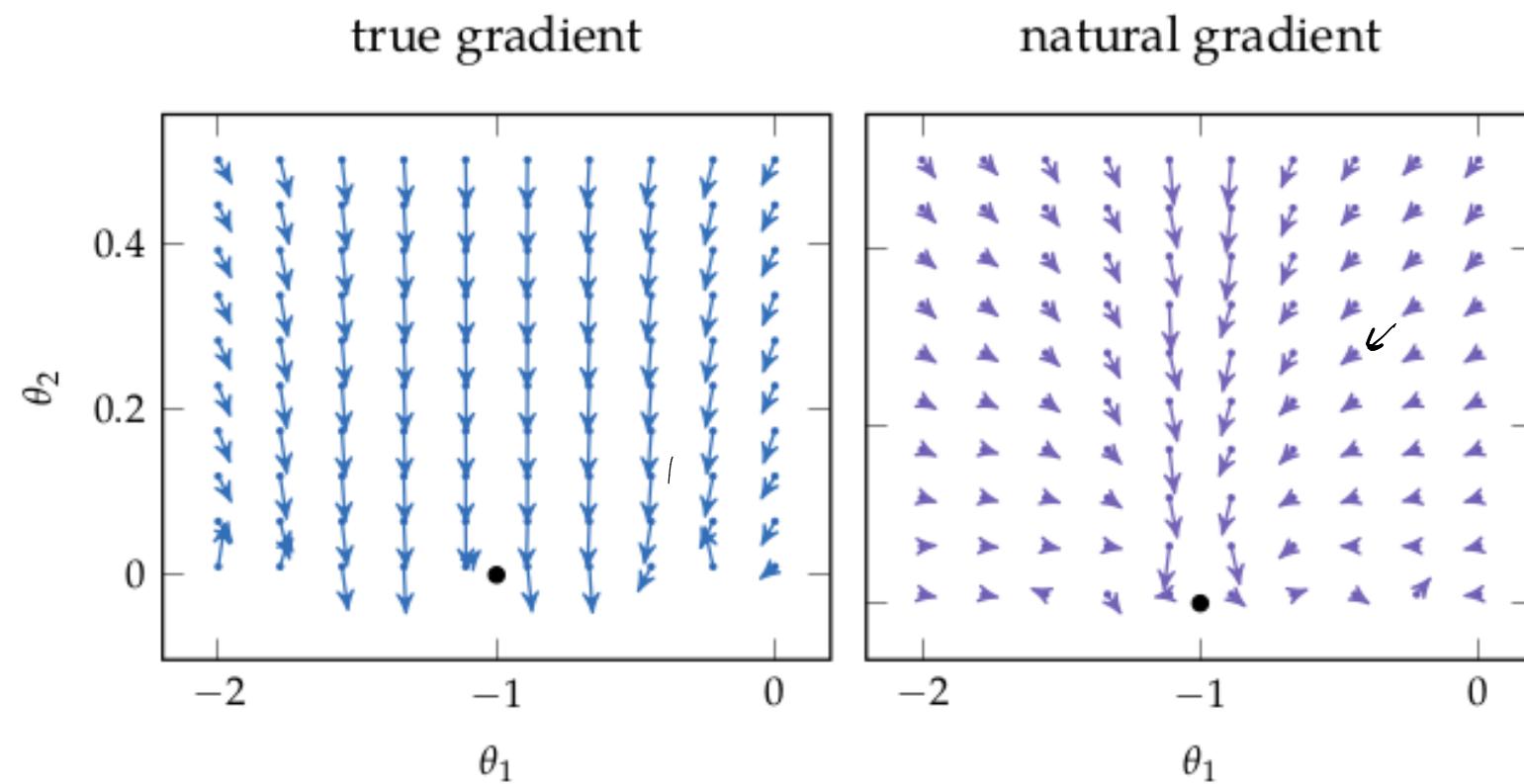


$$D_{KL}(p || q) = \int p(x) \log \frac{p(x)}{q(x)} dx \\ = - \int p(x) \log \frac{q(x)}{p(x)} dx$$

$$F_\theta = \int p(\tau|\theta) \nabla \log p(\tau|\theta) \nabla \log p(\tau|\theta)^T d\tau \\ = E_\tau [\nabla \log p(\tau|\theta) \nabla \log p(\tau|\theta)^T]$$

# Natural Gradient

# Natural Gradient



# TRPO and PPO

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

# TRPO and PPO

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

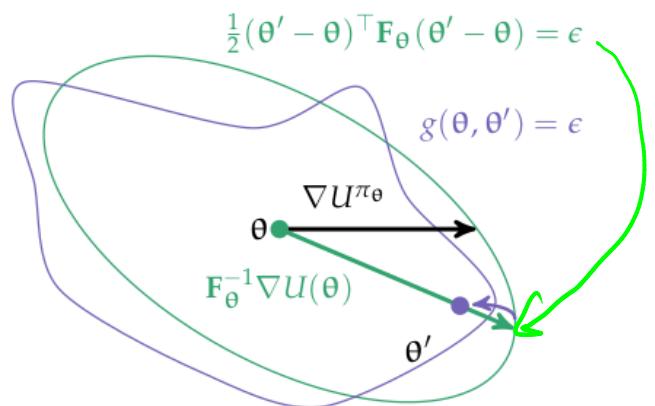
TRPO = Trust Region Policy Optimization

(Natural gradient + line search)

# TRPO and PPO

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

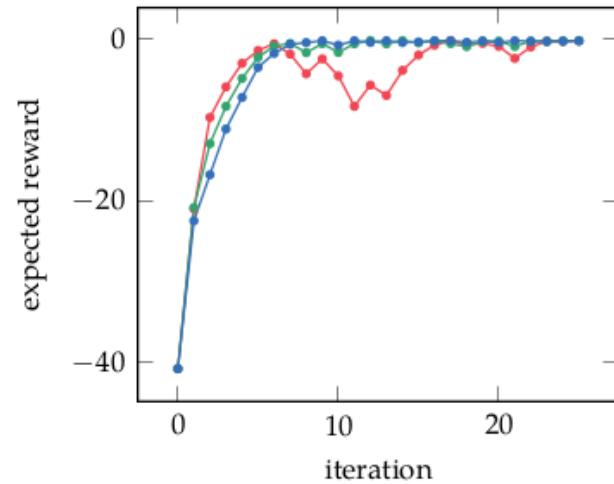
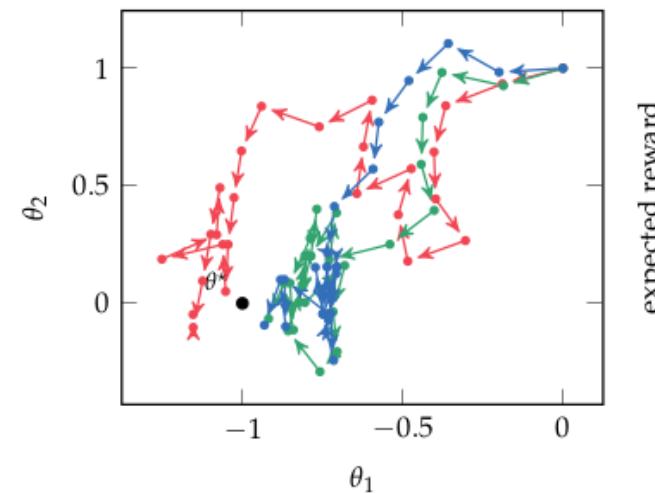
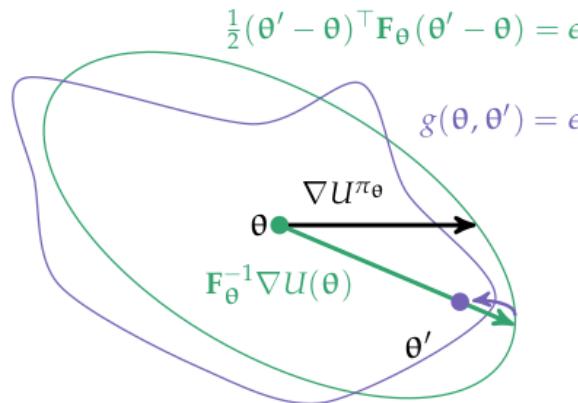
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



# TRPO and PPO

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

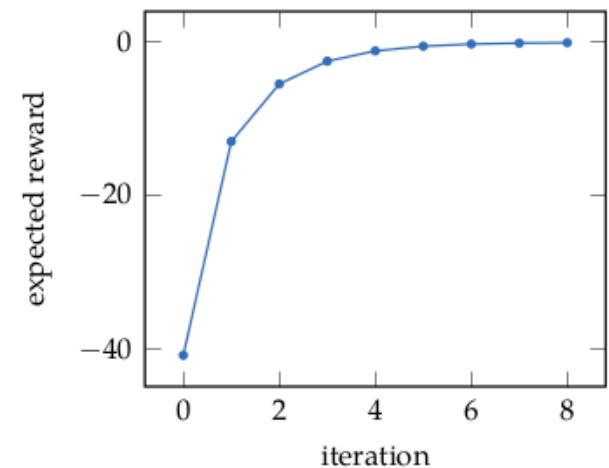
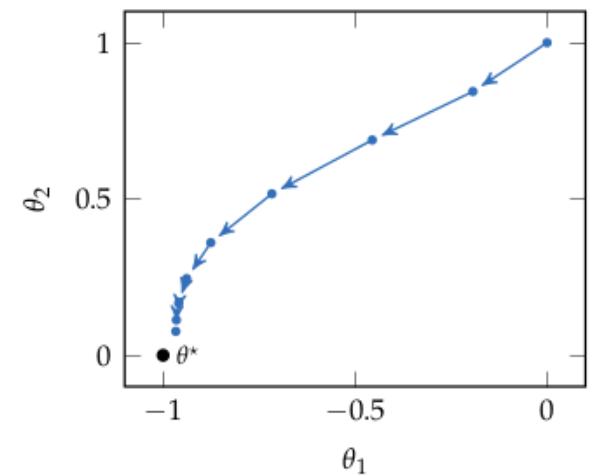
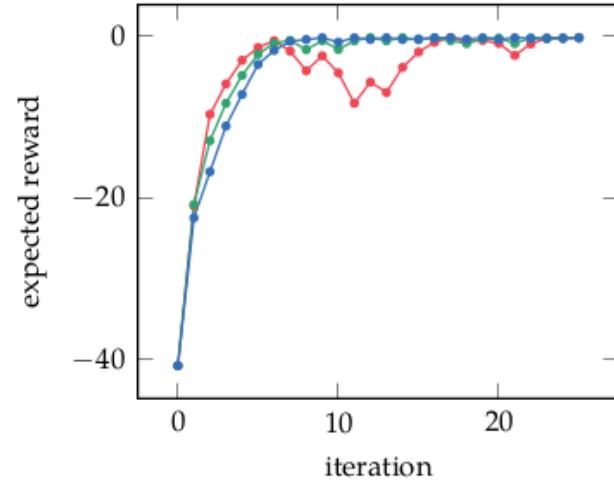
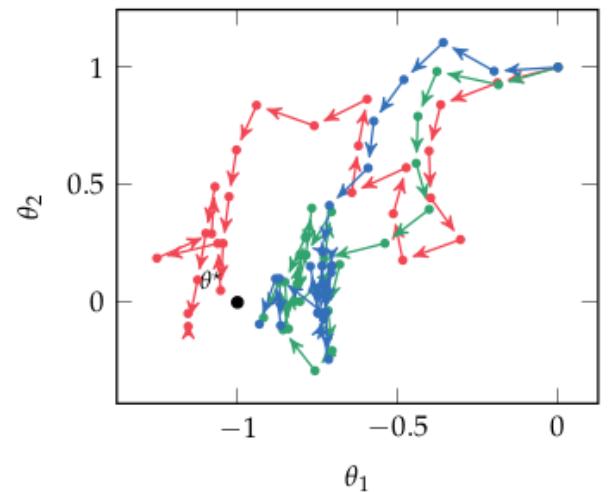
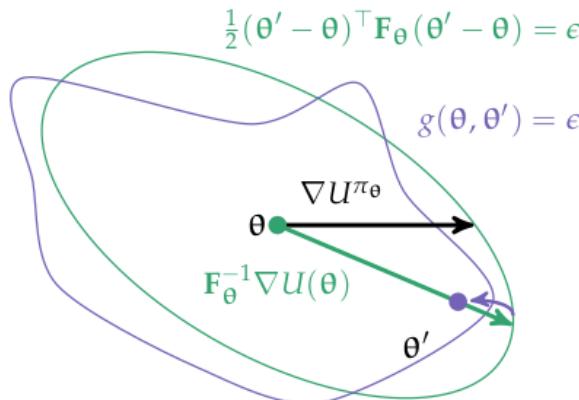
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



# TRPO and PPO

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

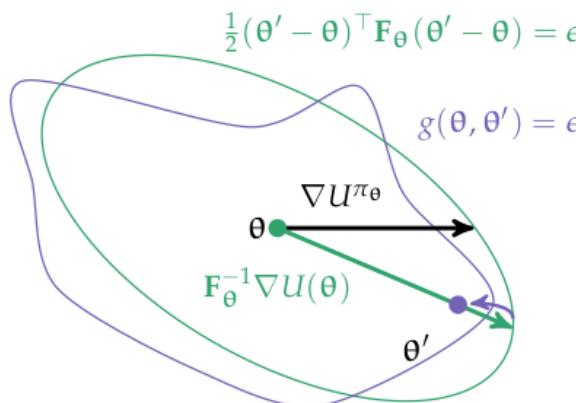
TRPO = Trust Region Policy Optimization  
(Natural gradient + line search)



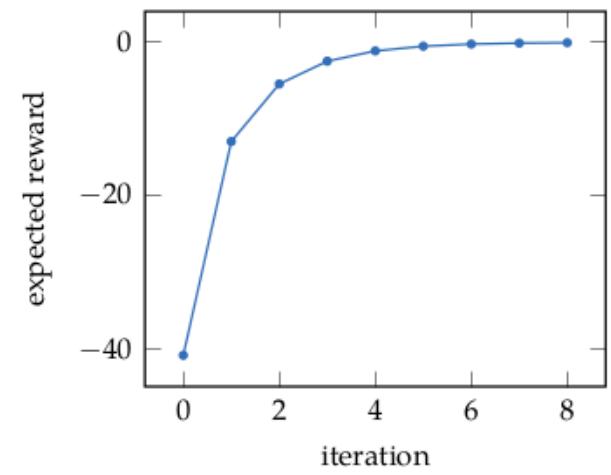
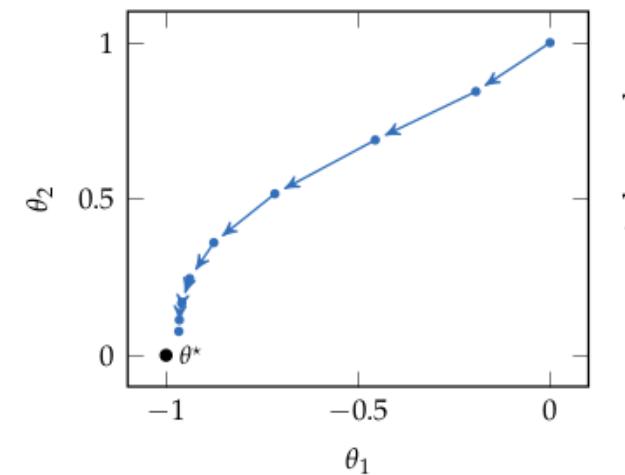
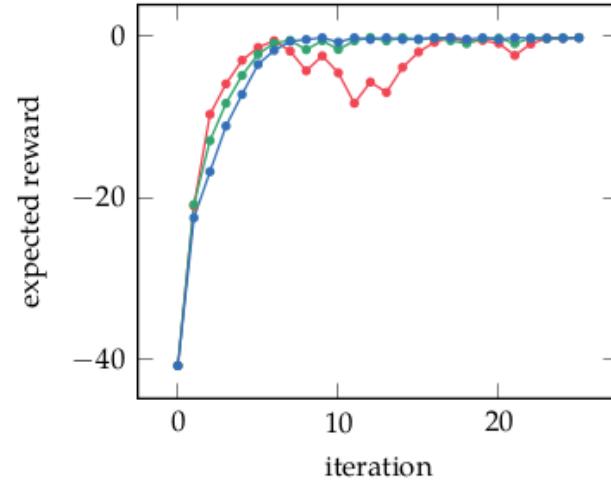
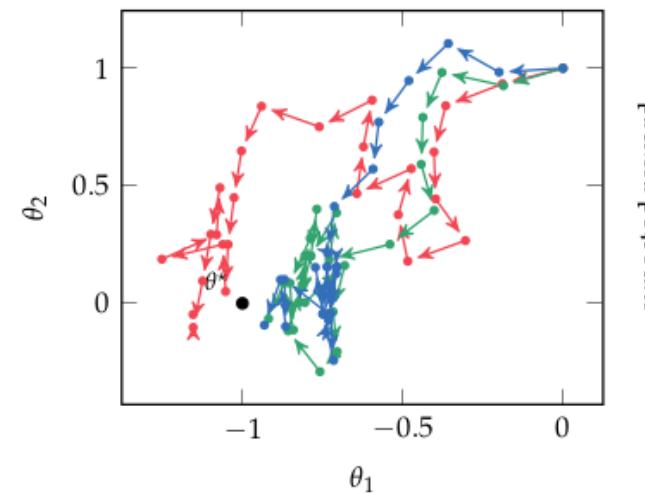
# TRPO and PPO

likelihood ratio  
reward-to-go  
baseline subtraction

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



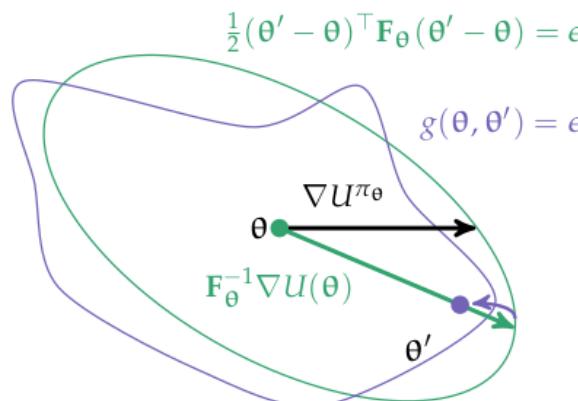
→ PPO = Proximal Policy Optimization  
(Use clamped surrogate objective to remove the need for line search)



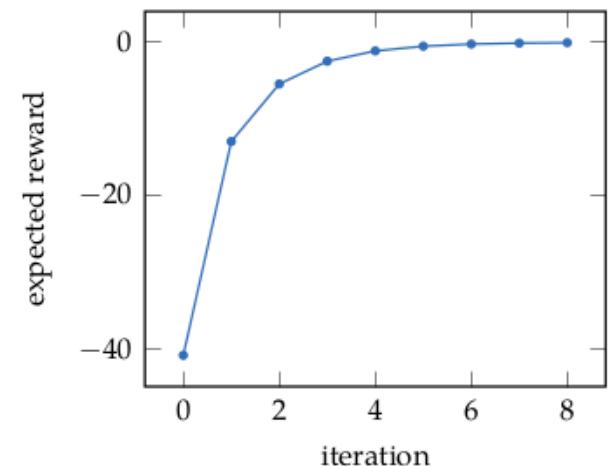
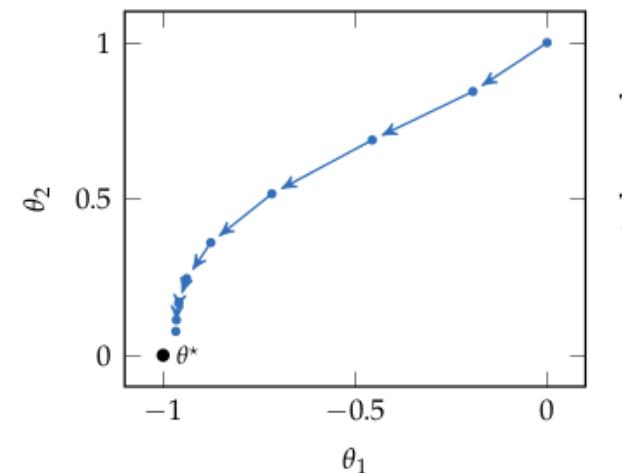
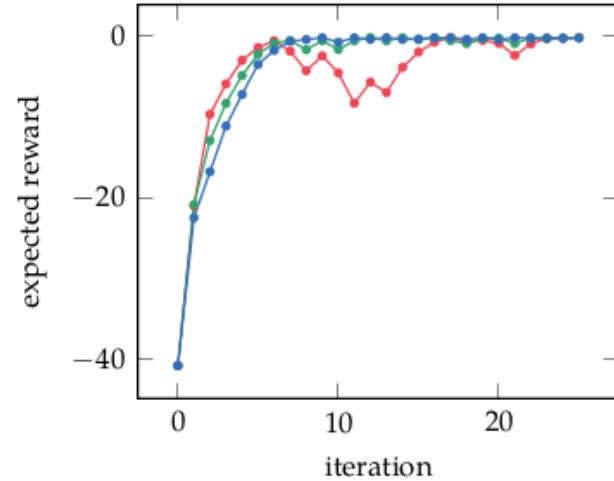
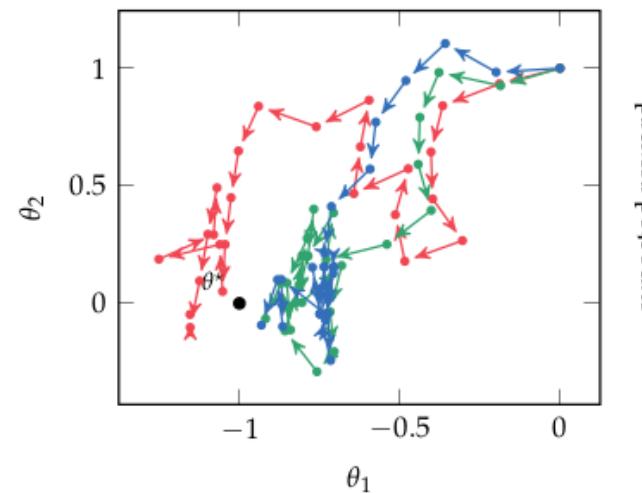
# TRPO and PPO

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

$$\nabla V(\theta) = E_c \left[ \sum_{k=0}^d \nabla_\theta \log \frac{\pi_\theta(a_k | s_k)}{\text{Actor}} \gamma^k (r_{t_0, g_0, k} - \hat{r}_{\text{base}}(s_k)) \right]$$

The diagram illustrates the Actor-Critic framework. It shows the Actor (blue) and Critic (red) components. The Actor is represented by a blue circle containing  $Q_\phi(s, a)$ . The Critic is represented by a red circle containing  $V_\phi(s)$ . A blue arrow points from the Actor's output to the Critic's input. A red arrow points from the Critic's output back to the Actor's input, labeled "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  $V$ ?

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

$\sim Q(s, a)$

# 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

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

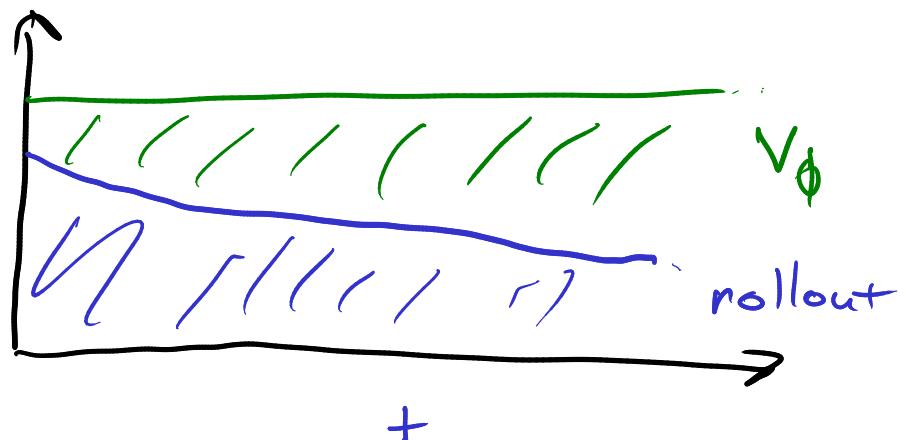
$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t - V_\phi(s_k) \quad \leftarrow \text{High Variance}$$

$$\approx \underbrace{r_k + \gamma r_{k+1} + \dots + \gamma^d r_{k+d}}_{\text{Temporal Difference}} + \underbrace{\gamma^{d+1} V_\phi(s_{k+d+1}) - V_\phi(s_k)}_{\text{Variance Reduction}}$$

let  $\delta_+ = r_+ + \gamma V_\phi(s_{++1}) - V_\phi(s_+)$

How do we choose when to stop

$$A^{\text{GAE}}(s_k, a_k) = \sum_{t=k}^{\infty} (\gamma \lambda)^{t-k} \delta_+$$



# Recap

DQN

Policy Gradient → PPO

Actor - Critic

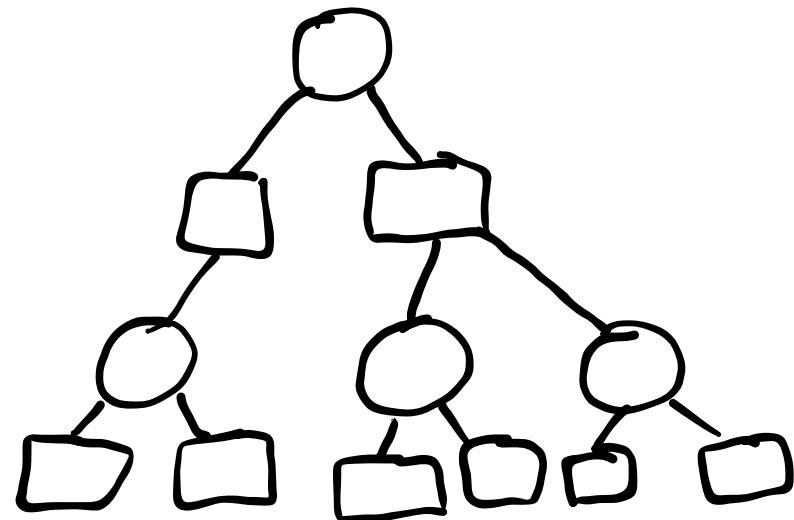
# Alpha Zero: Actor Critic with MCTS

# Alpha Zero: Actor Critic with MCTS

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

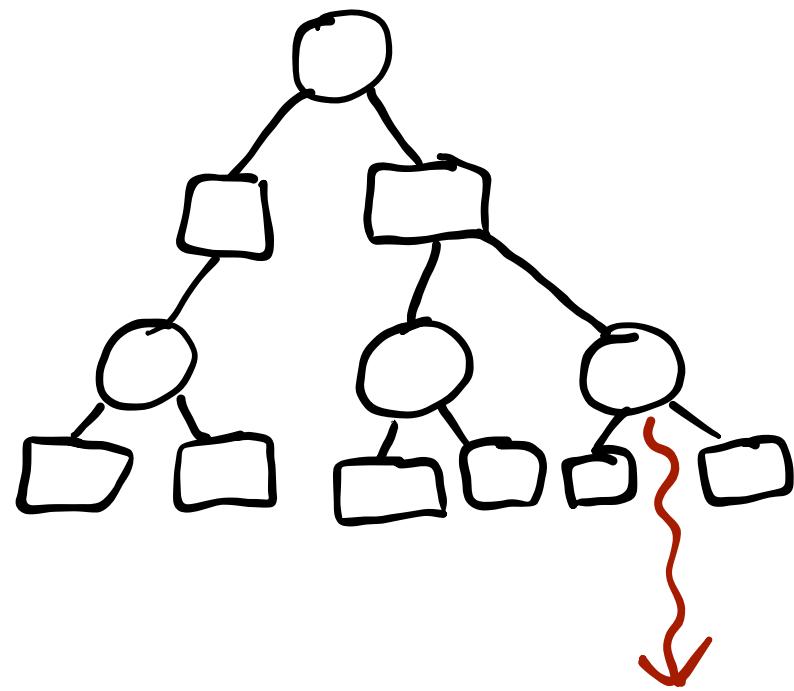
# Alpha Zero: Actor Critic with MCTS

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



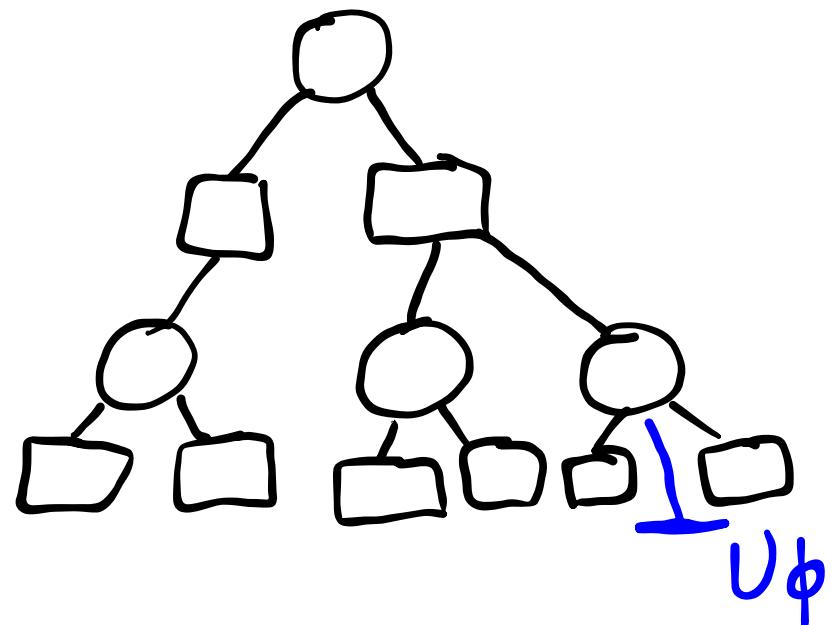
# Alpha Zero: Actor Critic with MCTS

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



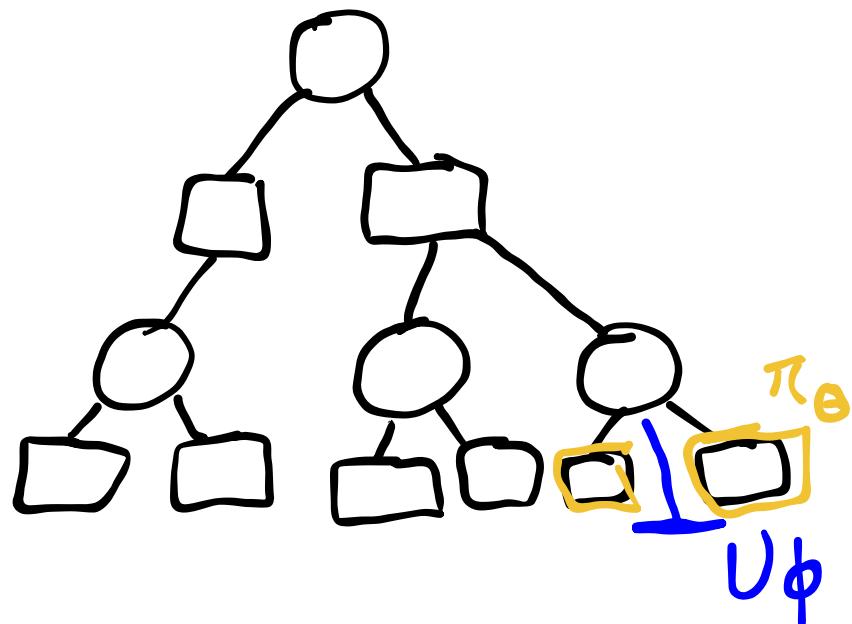
# Alpha Zero: Actor Critic with MCTS

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



# Alpha Zero: Actor Critic with MCTS

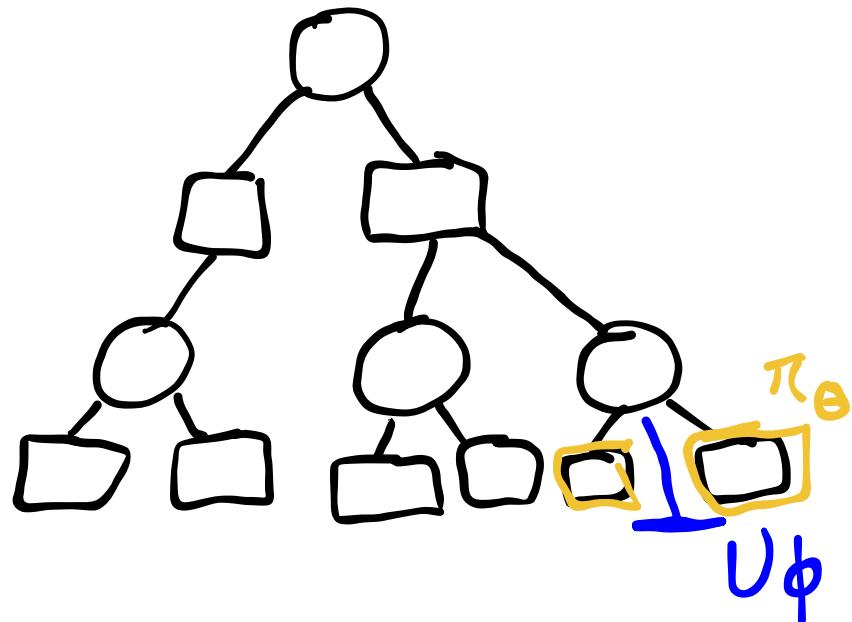
1. Use  $\pi_\theta$  and  $U_\phi$  in MCTS



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

# Alpha Zero: Actor Critic with MCTS

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

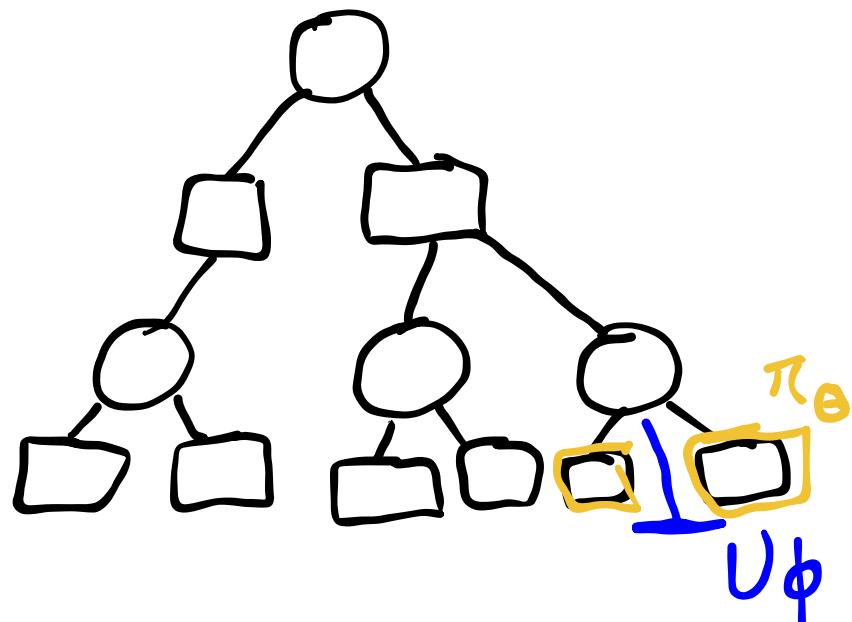


$$a = \arg \max_a Q(s, a) + c\pi_\theta(a \mid s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

# Alpha Zero: Actor Critic with MCTS

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

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

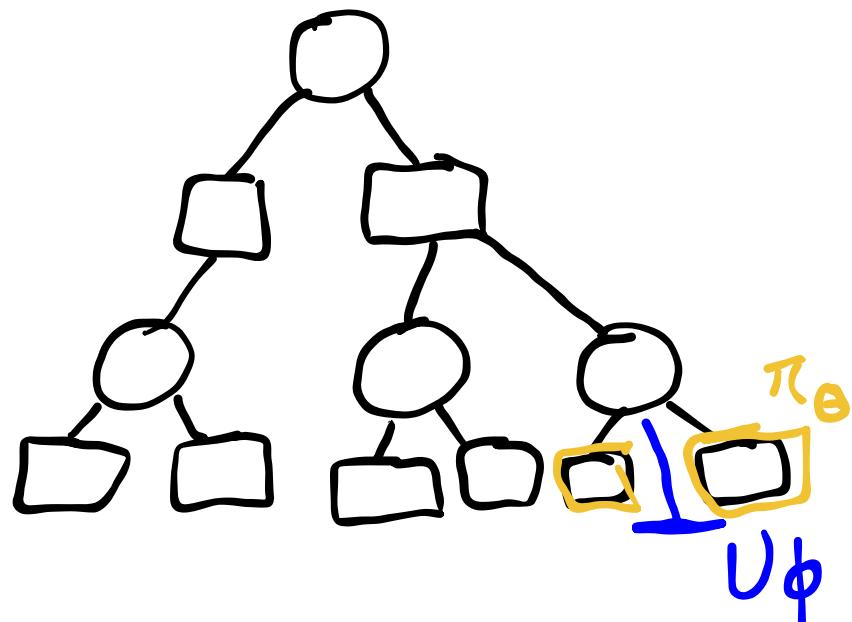


$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

# Alpha Zero: Actor Critic with MCTS

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

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



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

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

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