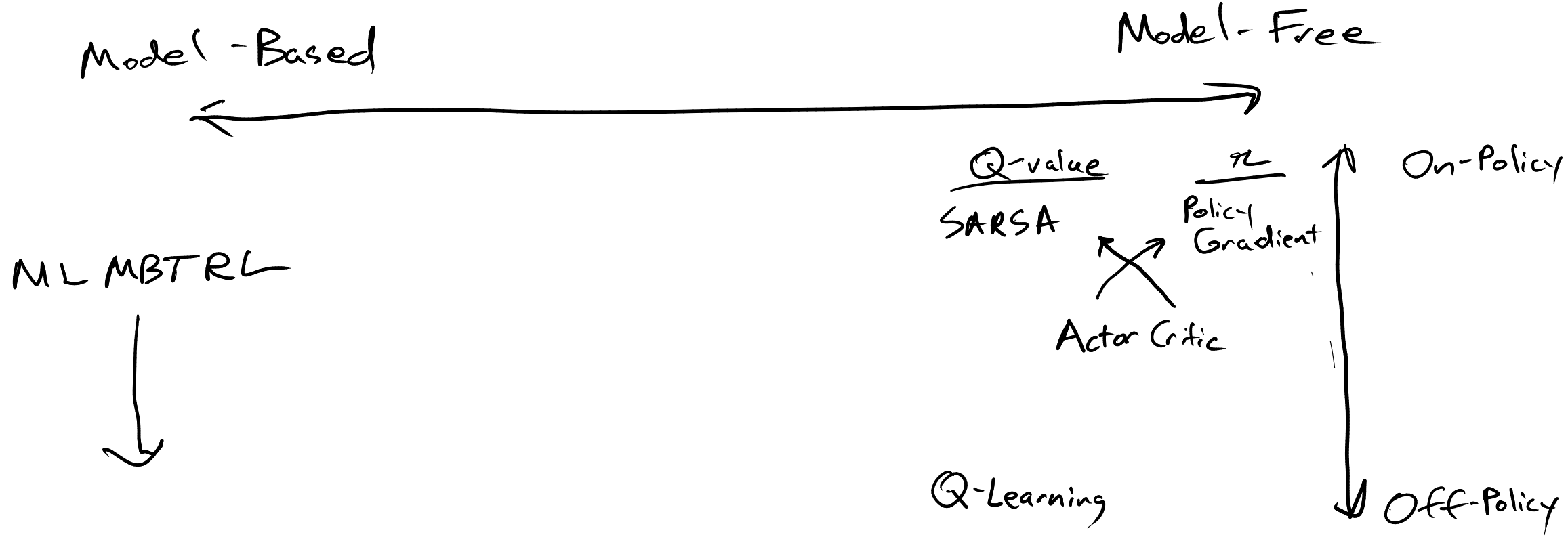


Map



1. Additional Actor-Critic
2. Advanced Exploration
3. Entropy Regularization
4. Wisdom

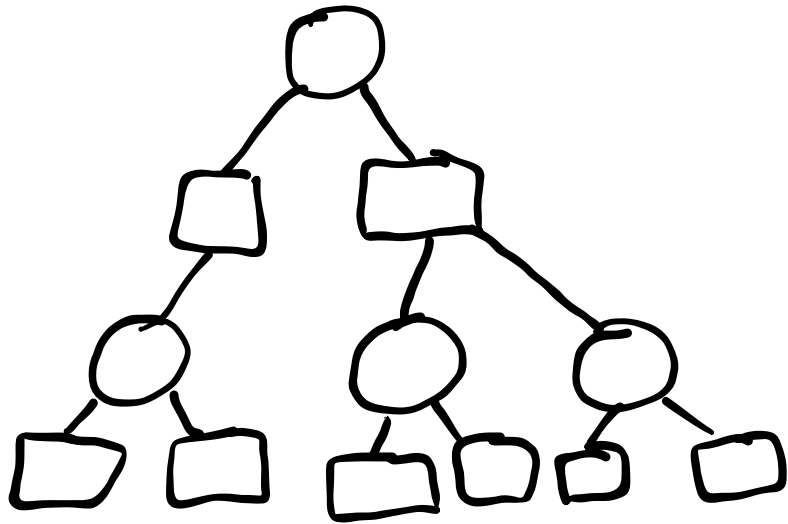
Alpha Zero: Actor Critic with MCTS

Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS
- help choose actions*
instead of rollout

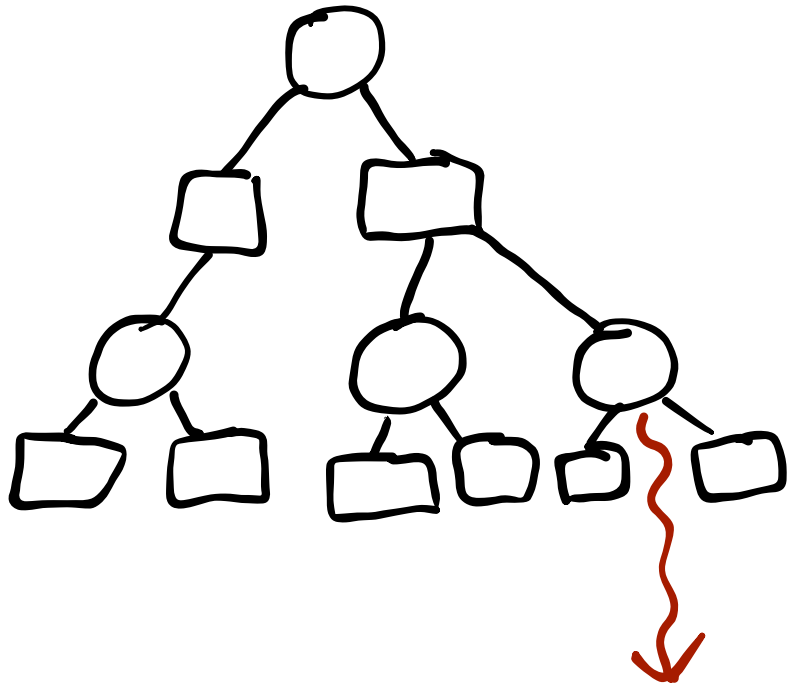
Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS



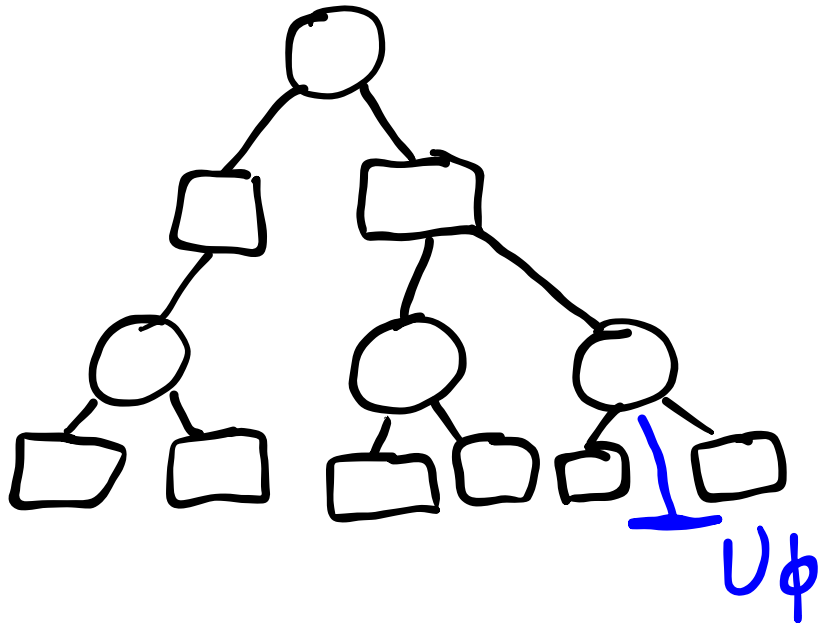
Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS



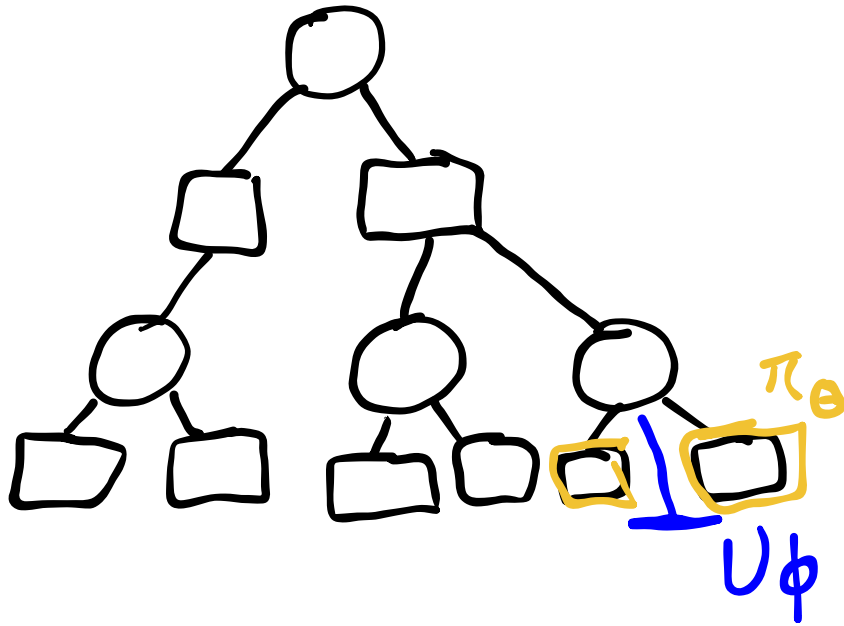
Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS



Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS

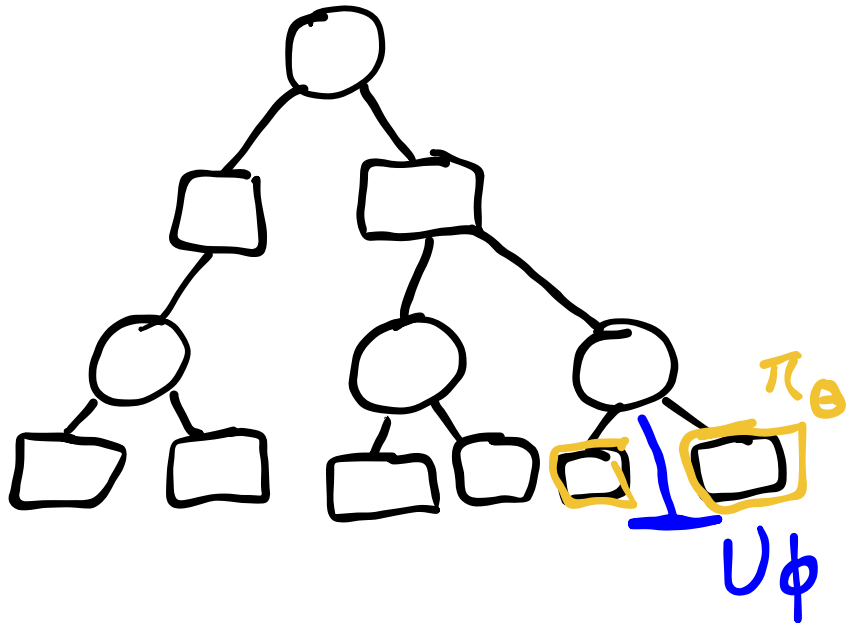


$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)} \leftarrow$$

Handwritten blue annotations: an arrow points from π_θ to the term $\pi_\theta(a | s)$, and another arrow points from the final term to the right.

Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS
2. Learn π_θ and U_ϕ from tree



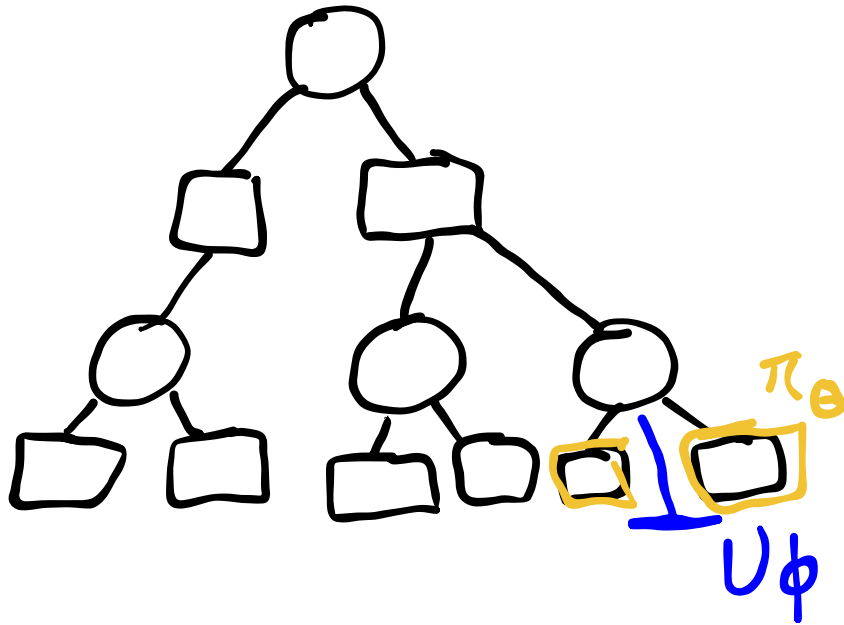
$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS
2. Learn π_θ and U_ϕ from tree

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

$$\pi_{\text{MCTS}}(a | s) \propto \underline{N(s, a)}^\eta$$



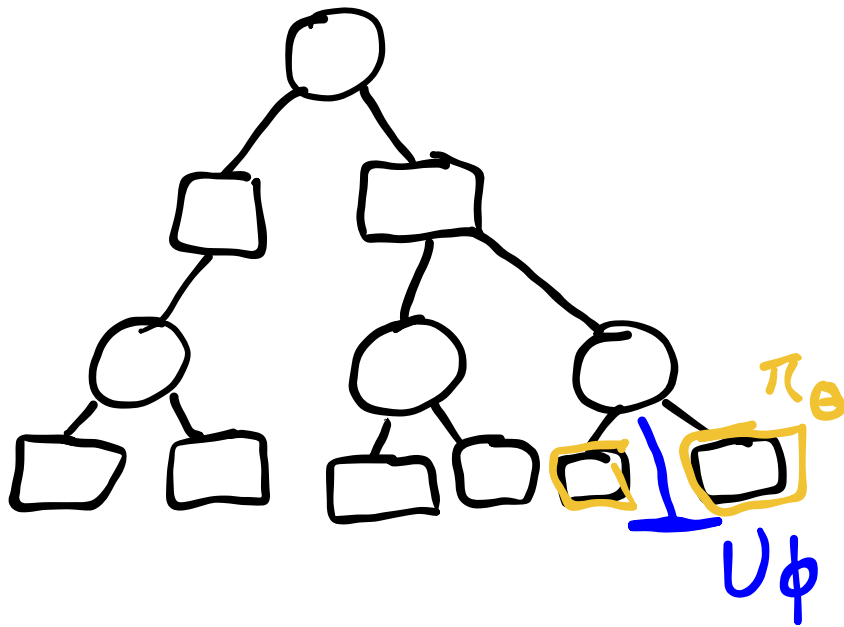
$$a = \arg \max_a Q(s, a) + c \pi_\theta(a | s) \frac{\sqrt{N(s)}}{1 + N(s, a)}$$

Alpha Zero: Actor Critic with MCTS

1. Use π_θ and U_ϕ in MCTS
2. Learn π_θ and U_ϕ from tree

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

$$\pi_{\text{MCTS}}(a | s) \propto N(s, a)^\eta$$



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

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

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

Continuous Actions: Deep Deterministic Policy Gradient

Actor Critic

$$Q_{\phi}(s,a)$$

$$\pi_{\theta}(s) = \operatorname{argmax}_a Q_{\phi}(s,a)$$

$$l(\phi) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(r + \gamma Q_{\phi}(s', \pi_{\theta}(s')) - Q_{\phi}(s,a))^2 \right]$$

$$U(\theta) = \mathbb{E}_{s \sim \mathcal{D}} [Q_{\phi}(s, \pi_{\theta}(s))]$$

$$\nabla U(\theta) = \mathbb{E}_{s \sim \mathcal{D}} [\nabla_{\theta} Q_{\phi}(s, \pi_{\theta}(s))]$$

Unstable

Is Exploration Important?

Montezuma's Revenge

Is Exploration Important?

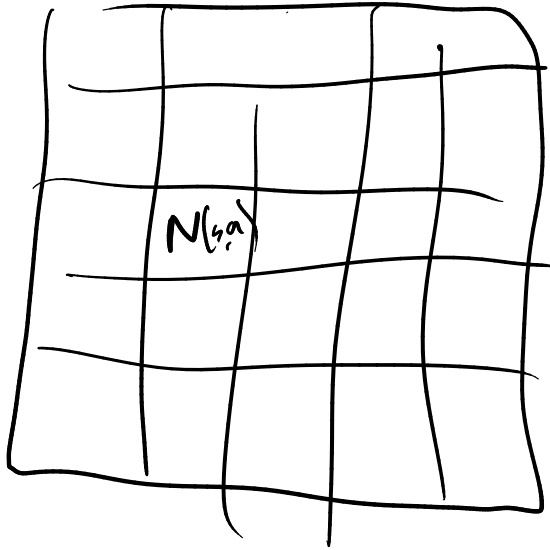
Theory

	Algorithm	Regret	Time	Space
Model-based	UCRL2 [10] ¹	at least $\tilde{\mathcal{O}}(\sqrt{H^4 S^2 AT})$	$\Omega(T S^2 A)$	$\mathcal{O}(S^2 AH)$
	Agrawal and Jia [1] ¹	at least $\tilde{\mathcal{O}}(\sqrt{H^3 S^2 AT})$		
	UCBVI [5] ²	$\tilde{\mathcal{O}}(\sqrt{H^2 SAT})$	$\tilde{\mathcal{O}}(T S^2 A)$	
	vUCQ [12] ²	$\tilde{\mathcal{O}}(\sqrt{H^2 SAT})$		
Model-free	Q-learning (ε -greedy) [14] (if 0 initialized)	$\Omega(\min\{T, \underbrace{A^{H/2}}_{\text{horizon}}\})$	$\mathcal{O}(T)$	$\mathcal{O}(SAH)$
	Delayed Q-learning [25] ³	$\tilde{\mathcal{O}}_{S,A,H}(T^{4/5})$		
	Q-learning (UCB-H)	$\tilde{\mathcal{O}}(\sqrt{H^4 SAT})$		
	Q-learning (UCB-B)	$\tilde{\mathcal{O}}(\sqrt{H^3 SAT})$		
	lower bound	$\Omega(\sqrt{H^2 SAT})$	-	-

Exploration Bonus

Exploration Bonus

- In General, $R^+(s, a) = R(s, a) + B(s, a)$
- UCB: $B(s, a) = c \sqrt{\frac{\log N(s)}{N(s, a)}}$ ←



Exploration Bonus

Example 1: Learn Pseudocount

Exploration Bonus

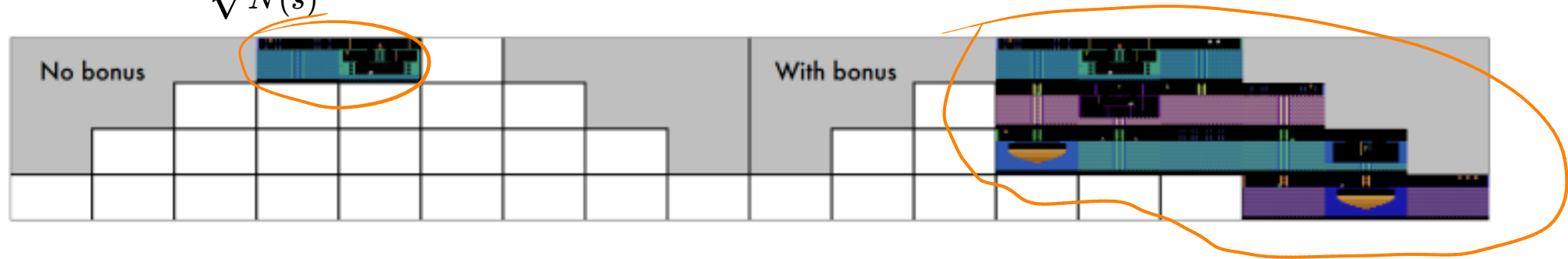
Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation

Exploration Bonus

Example 1: Learn Pseudocount

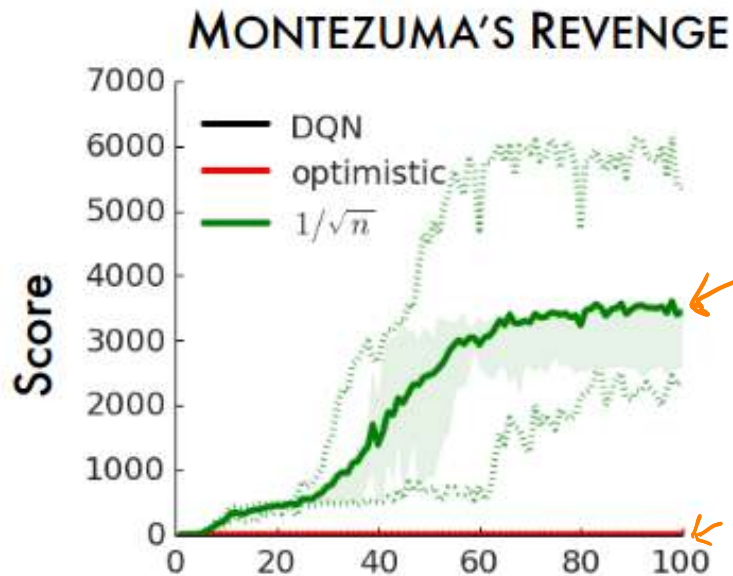
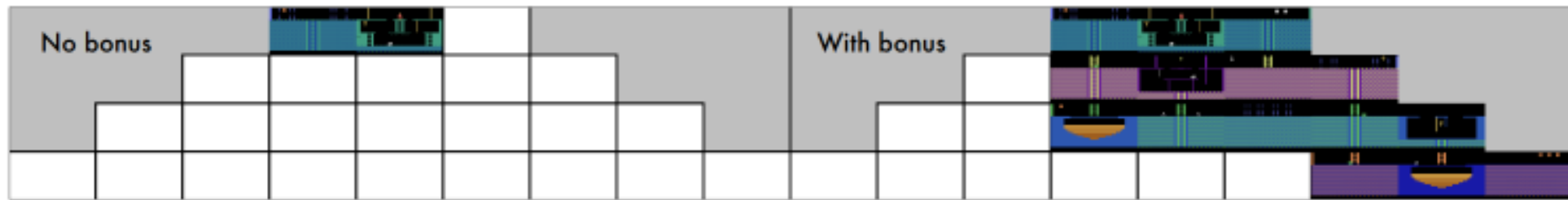
$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation



Exploration Bonus

Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation

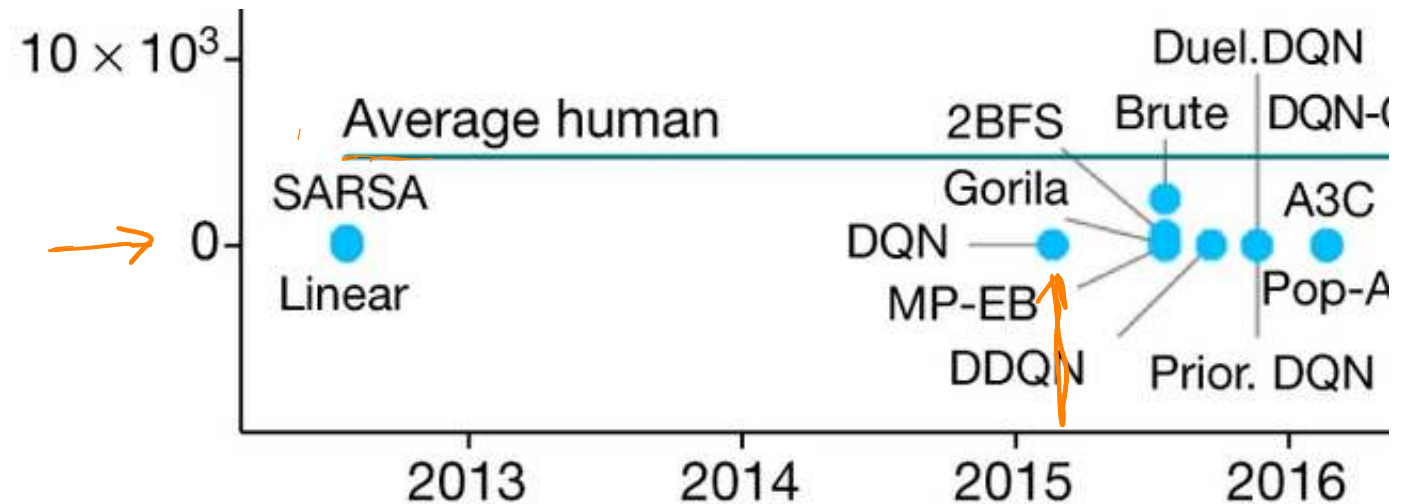
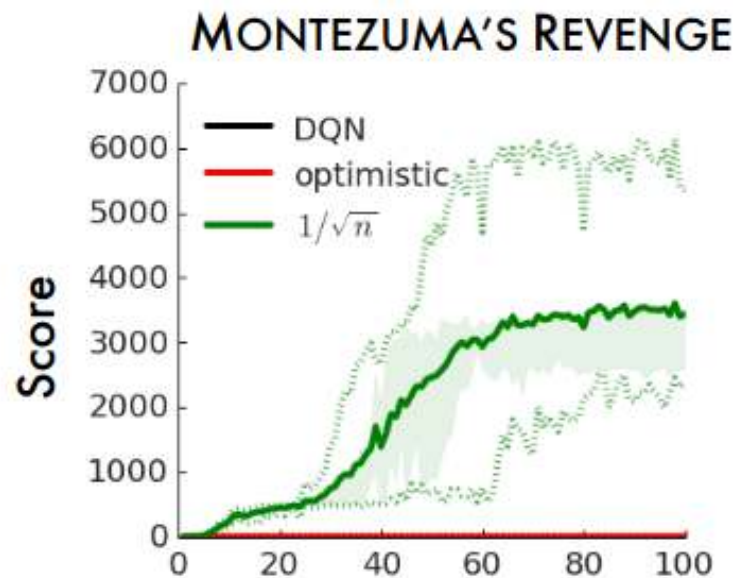
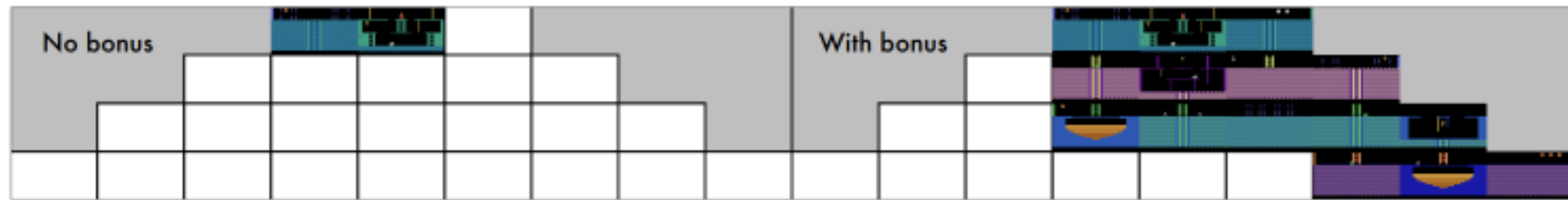


Bellemare, et al. 2016 "Unifying Count-Based Exploration..."

Exploration Bonus

Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation

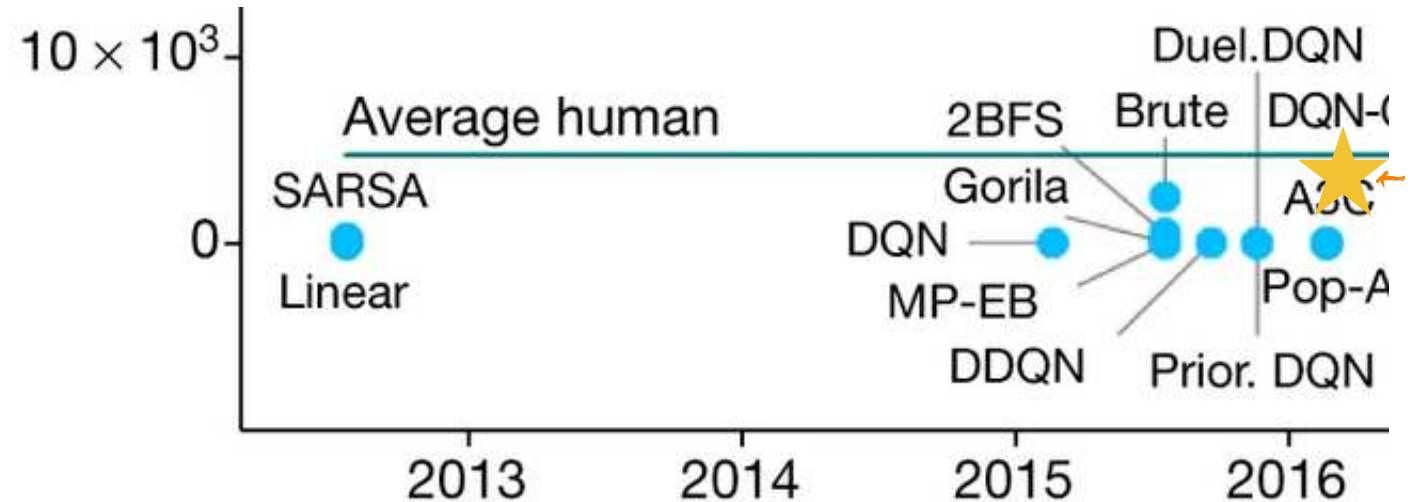
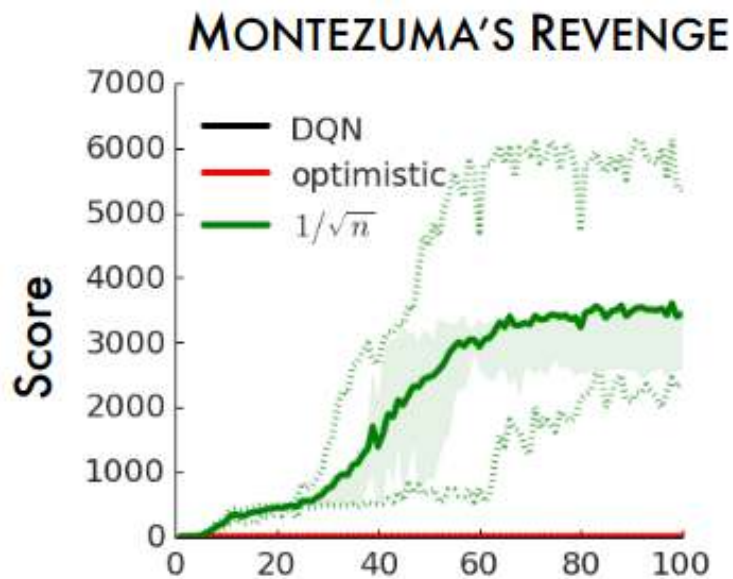
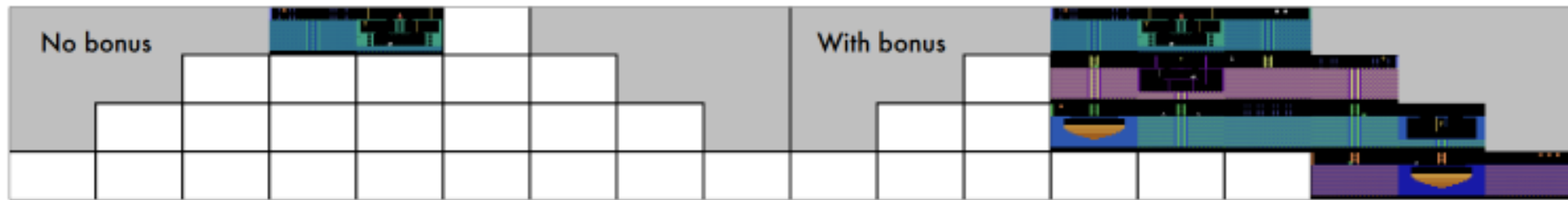


Bellemare, et al. 2016 "Unifying Count-Based Exploration..."

Exploration Bonus

Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation



Bellemare, et al. 2016 "Unifying Count-Based Exploration..."

Exploration Bonus

^

—
,

Exploration Bonus

Example 2: Learn a function of the state and action

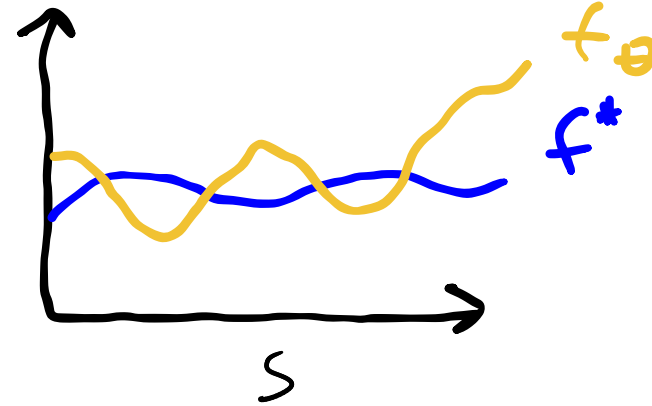
^

—
/

Exploration Bonus

Example 2: Learn a function of the state and action

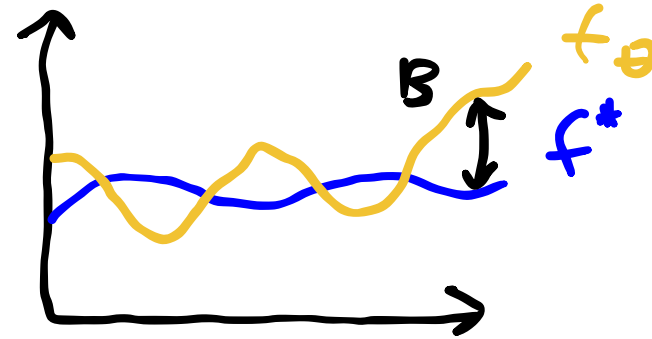
$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$



Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

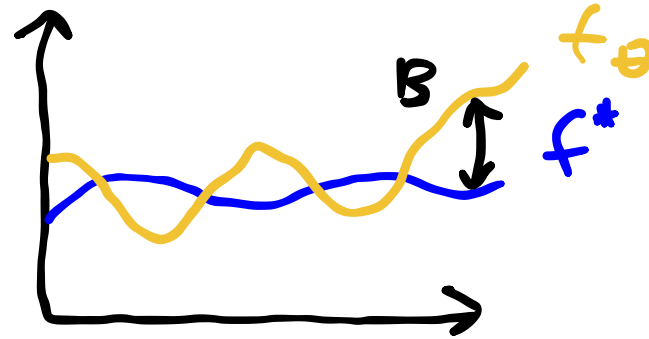


Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

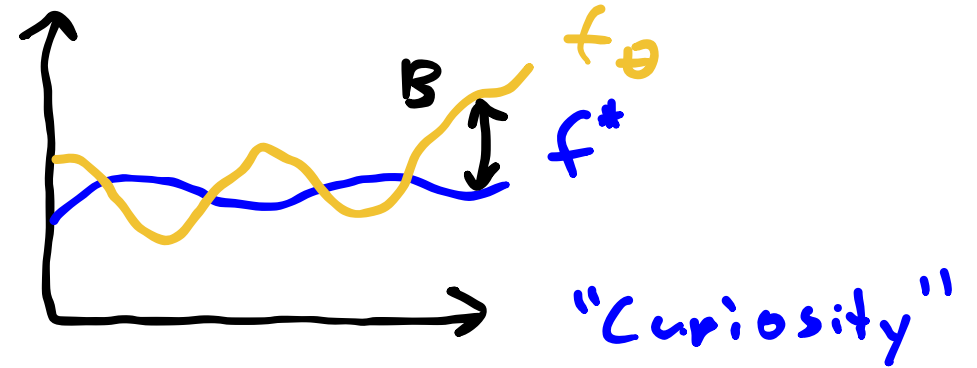


Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?



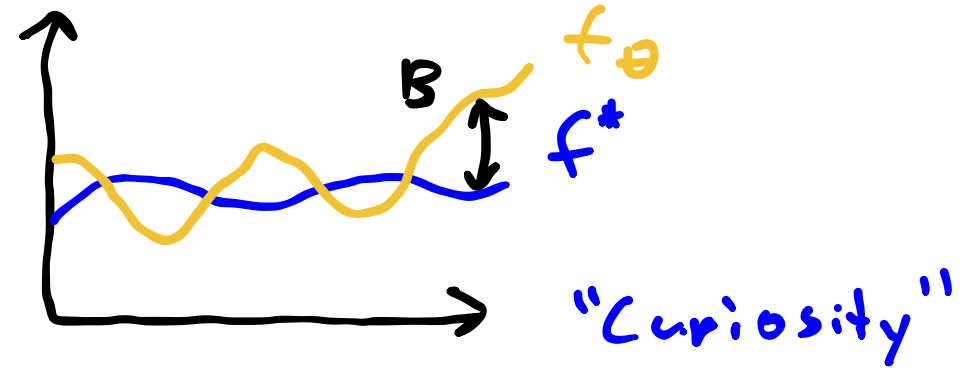
Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)



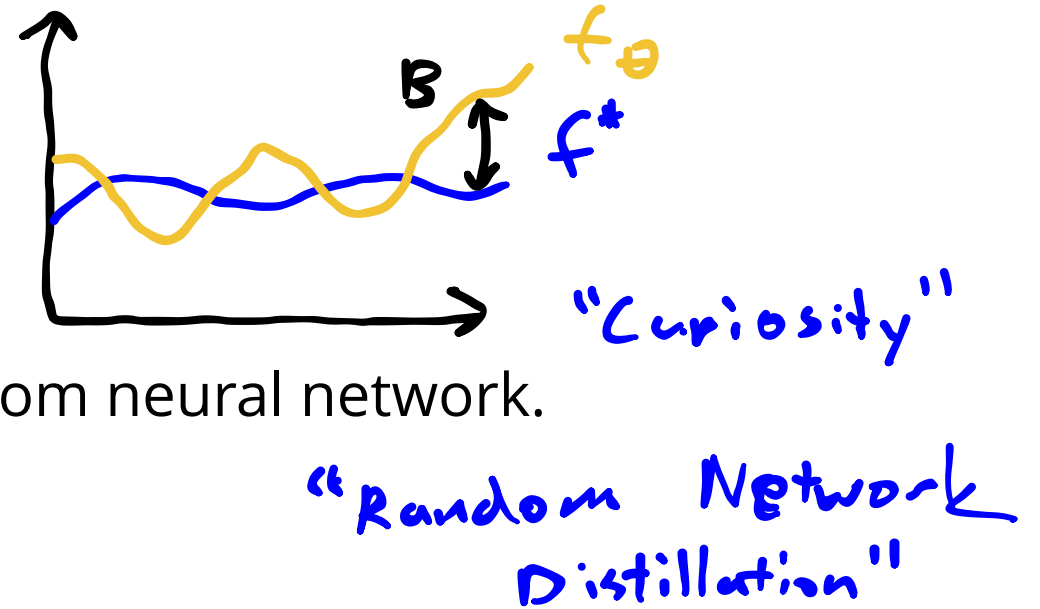
Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)
- $f^*(s, a) = f_\phi(s, a)$ where f_ϕ is a random neural network.



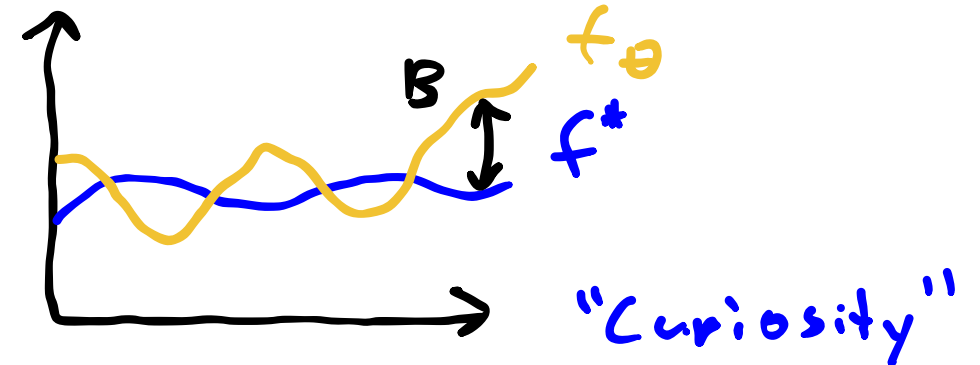
Exploration Bonus

Example 2: Learn a function of the state and action

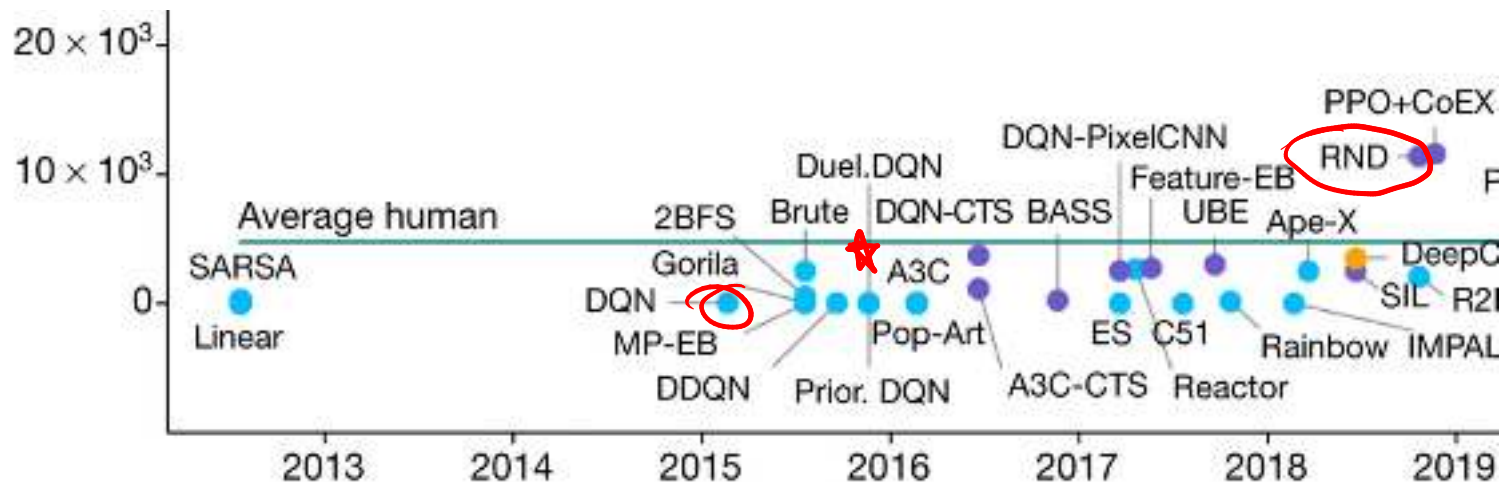
$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)
- $f^*(s, a) = f_\phi(s, a)$ where f_ϕ is a random neural network.



"Random Network Distillation"



Exploration Bonus

Exploration Bonus

Example 3: Thompson Sampling

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q
2. Sample Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q
2. Sample \hat{Q} from ~~Q~~ dist.
3. Act according to \hat{Q}

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q ← Hard
2. Sample Q
3. Act according to Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q ← Hard
 2. Sample Q
 3. Act according to Q
- Bootstrapping with multiple Q networks

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q ← Hard
2. Sample Q
3. Act according to Q

- Bootstrapping with multiple Q networks
- Dropout

Exploration Bonus

Exploration Bonus

Example 4: Go-Explore

Exploration Bonus

Example 4: Go-Explore

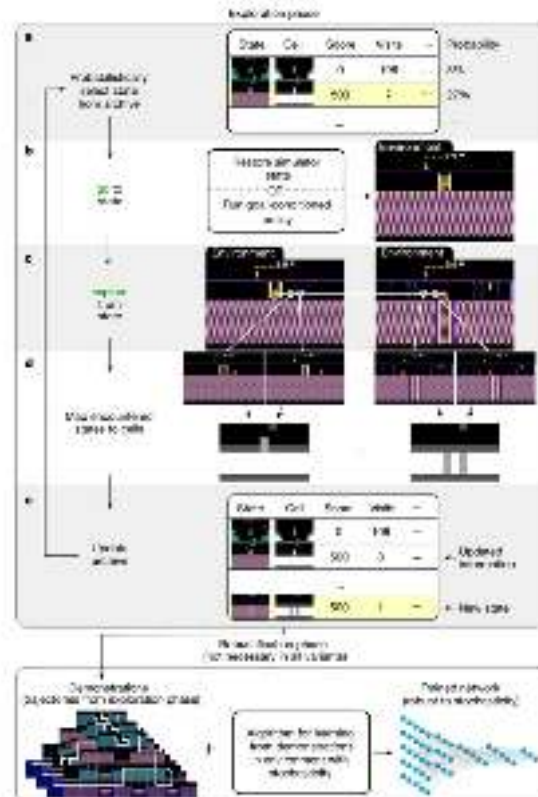
"First return, then explore"

Exploration Bonus

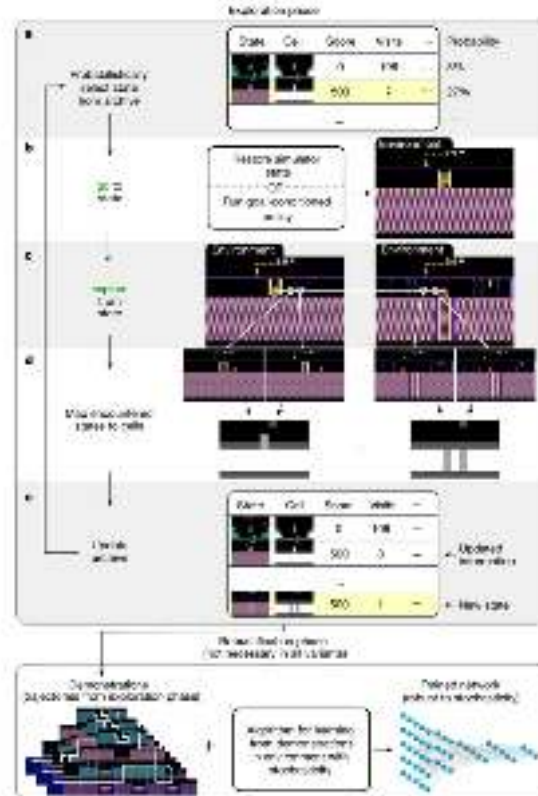
Example 4: Go-Explore

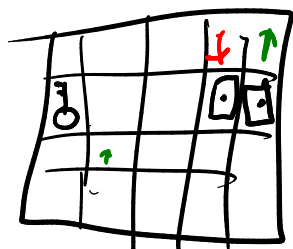
"First return, then explore"

Fig. 1: Overview of Go-Explore.



Example 4: Go-Explore

Fig. 1. Curves of α vs. β for α -explosive.



Exploration Bonus

Example 4: Go-Explore

"First return, then explore"

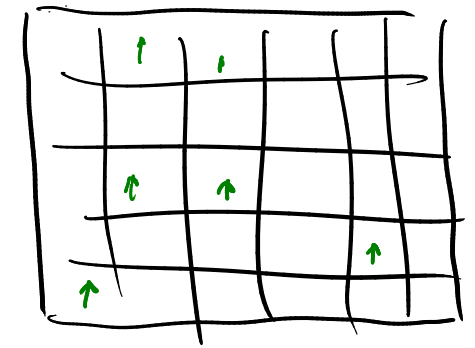
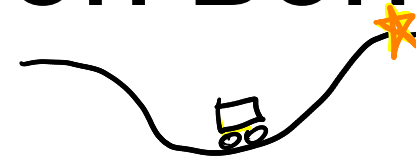
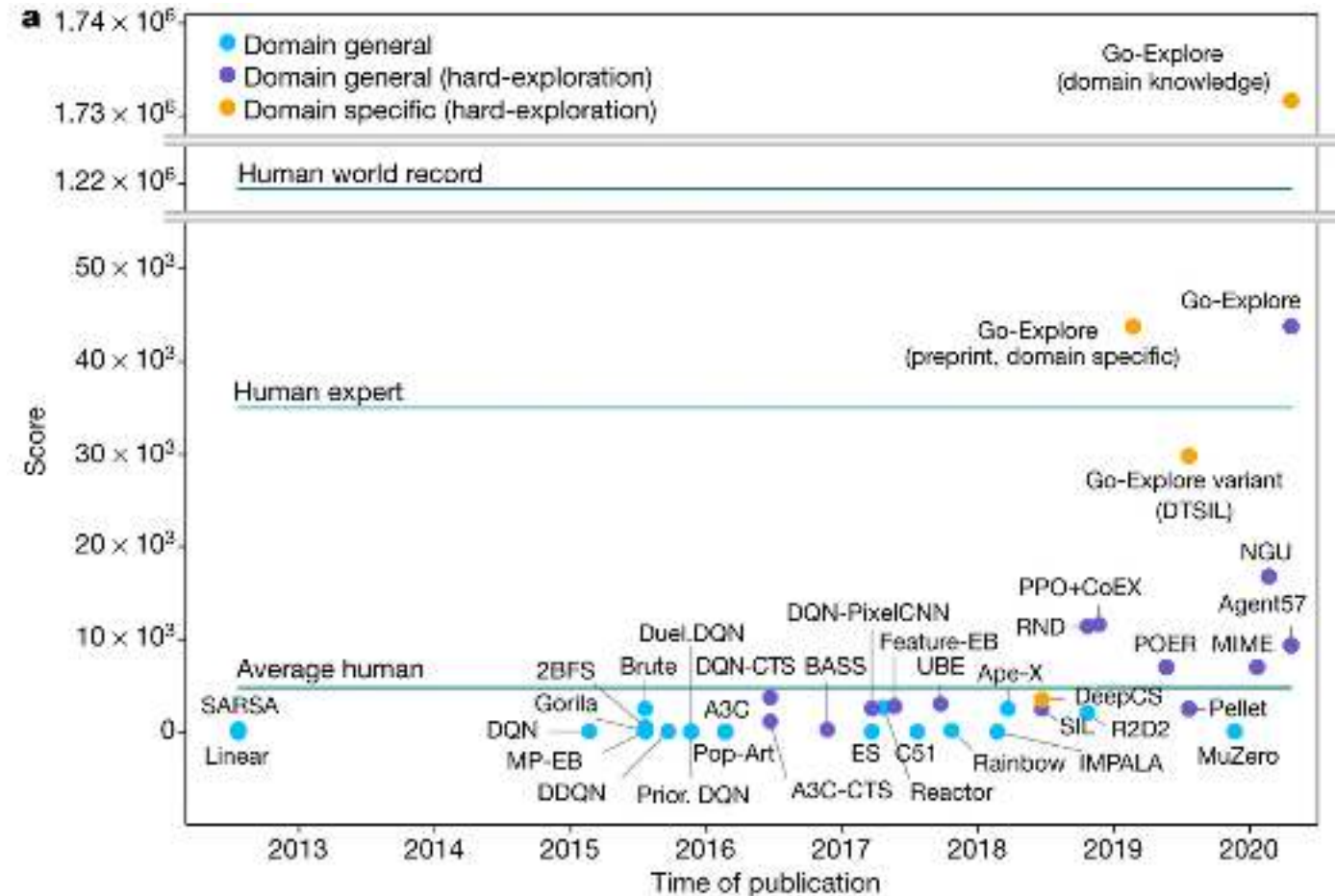
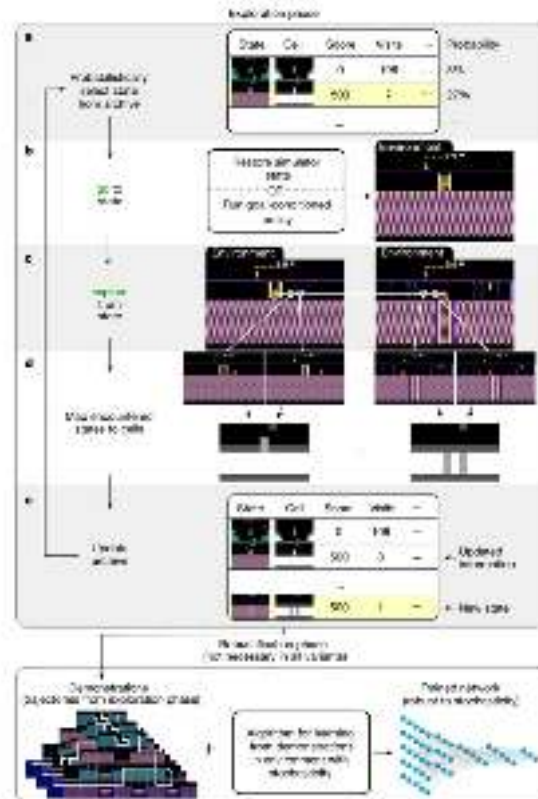
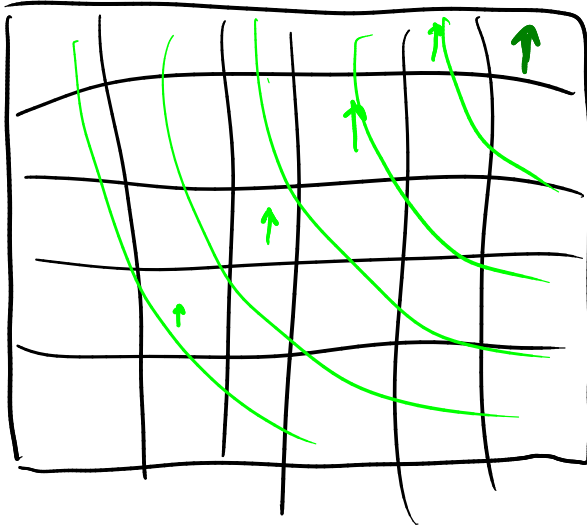


Fig. 1: Overview of Go-Explore.



(Uber AI Labs)

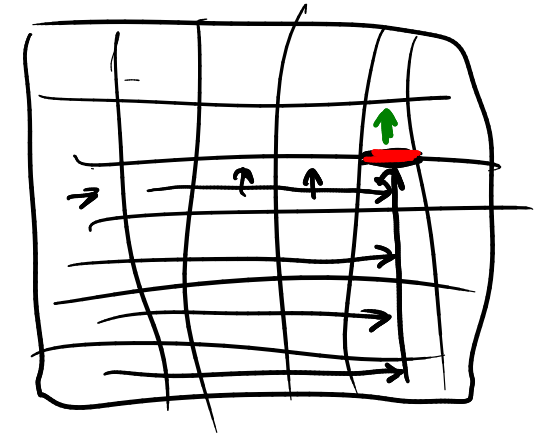
Soft Actor Critic: Entropy Regularization



Intrinsic

Soft Actor Critic: Entropy Regularization

$$U(\pi) = E[\sum \gamma^t r_t]$$
$$U(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \underbrace{\alpha \mathcal{H}(\pi(\cdot | s_t))}_{\text{entropy}}) \right]$$



Soft Actor Critic: Entropy Regularization

$$U(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

$$V(s_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(s_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | s_t)]$$

Soft Actor Critic: Entropy Regularization

$$U(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

$$\underline{V(s_t)} = \mathbb{E}_{\mathbf{a}_t \sim \pi} [\underline{Q(s_t, \mathbf{a}_t)} - \underline{\log \pi(\mathbf{a}_t | s_t)}]$$

$$\underline{\mathcal{T}^\pi Q(s_t, \mathbf{a}_t)} \triangleq \underline{r(s_t, \mathbf{a}_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V(s_{t+1})]}$$

iterative
policy
evaluation

Soft Actor Critic: Entropy Regularization

$$\rightarrow U(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

soft policy iteration

soft policy eval

$$V(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)]$$
$$\mathcal{T}^{\pi} Q(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V(\mathbf{s}_{t+1})]$$
$$\pi_{\text{new}} = \arg \min_{\pi' \in \Pi} D_{\text{KL}} \left(\pi'(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot))}{Z^{\pi_{\text{old}}}(\mathbf{s}_t)} \right)$$

Soft Actor Critic

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.

for each iteration **do**

for each environment step **do**

$$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$$

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$$

end for

for each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$$

$$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i) \text{ for } i \in \{1, 2\}$$

$$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$$

$$\bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}$$

end for

end for

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.

for each iteration **do**

for each environment step **do**

$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)$

$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$

end for

for each gradient step **do**

$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$

$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$

$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$

$\bar{\psi} \leftarrow \tau\psi + (1 - \tau)\bar{\psi}$

end for

end for

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[\frac{1}{2} (V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)])^2 \right]$$

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.

for each iteration **do**

for each environment step **do**

$$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$$

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$$

end for

for each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$$

$$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i) \text{ for } i \in \{1, 2\}$$

$$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$$

$$\bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}$$

end for

end for

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[\frac{1}{2} \left(V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)] \right)^2 \right]$$

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.

for each iteration **do**

for each environment step **do**

$$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)$$

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$$

$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$$

end for

for each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$$

$$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i) \text{ for } i \in \{1, 2\}$$

$$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$$

$$\bar{\psi} \leftarrow \tau\psi + (1 - \tau)\bar{\psi}$$

end for

end for

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[\frac{1}{2} \left(V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)] \right)^2 \right]$$

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

$$\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V_{\bar{\psi}}(\mathbf{s}_{t+1})]$$

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.

for each iteration **do** $V_\psi \quad Q_\theta \quad \pi_\phi$

for each environment step **do**

$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$

$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$

end for

for each gradient step **do**

$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$

$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$

$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$

$\bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}$

end for

end for

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[\frac{1}{2} \left(V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)] \right)^2 \right]$$

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

$$\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V_{\bar{\psi}}(\mathbf{s}_{t+1})]$$

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[D_{\text{KL}} \left(\pi_\phi(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q_\theta(\mathbf{s}_t, \cdot))}{Z_\theta(\mathbf{s}_t)} \right) \right]$$

Soft Actor Critic

Advantages:

Soft Actor Critic

Advantages:

- Stable

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies
- Exploration

Soft Actor Critic

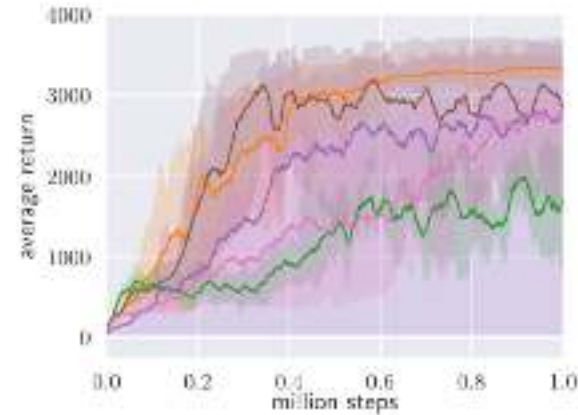
Advantages:

- Stable
- Learns many near-optimal policies
- Exploration
- Insensitivity to hyperparameters
- Off-Policy

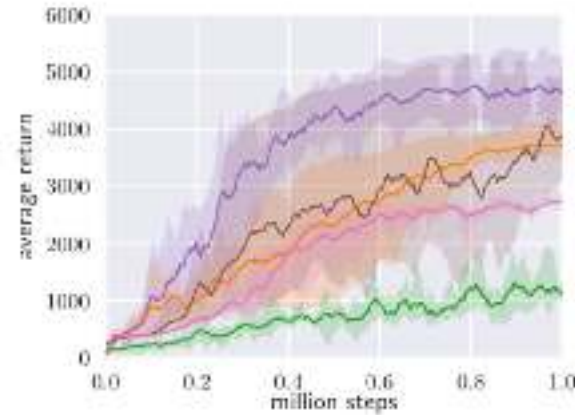
Soft Actor Critic

Advantages:

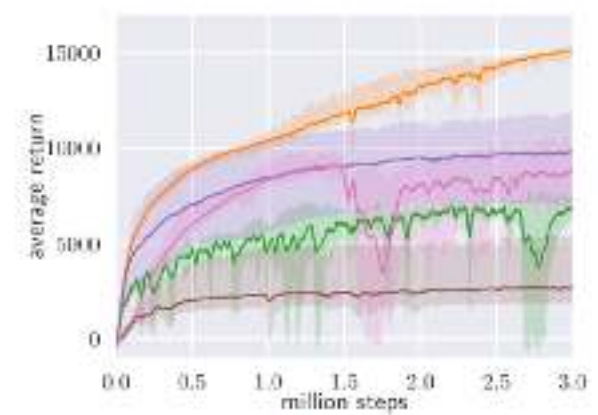
- Stable
- Learns
- Explora
- Insensit
- Off-Poli



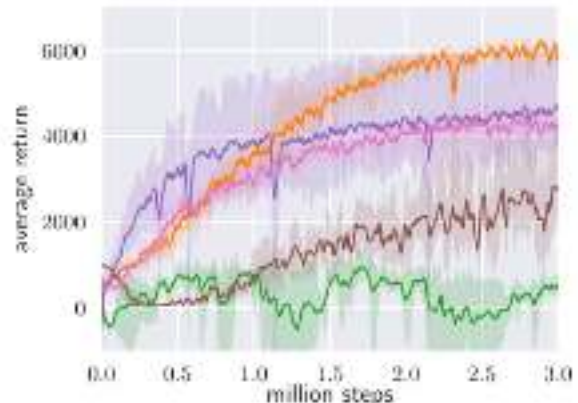
(a) Hopper-v1



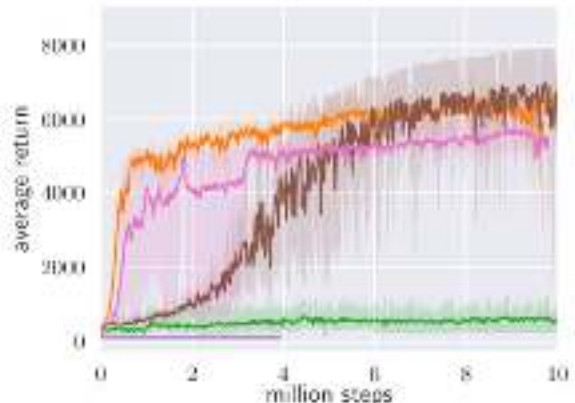
(b) Walker2d-v1



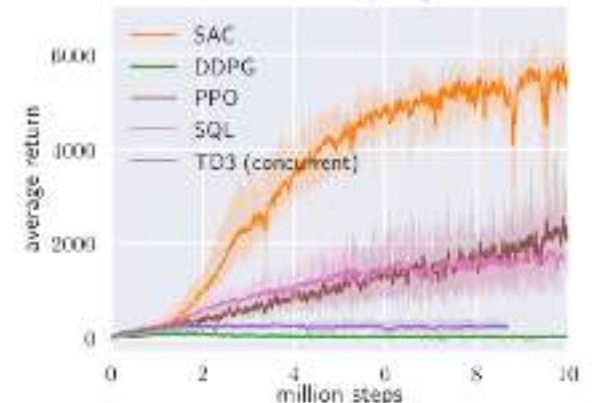
(c) HalfCheetah-v1



(d) Ant-v1



(e) Humanoid-v1



(f) Humanoid (rllab)

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies
- Exploration
- Insensitivity to hyperparameters
- Off-Policy

Disadvantages

- Sensitive to α Solution = Entropy
constraint and adjust α

Soft Actor Critic

Advantages:

Disadvantages:

- Stable
- Less
- Exploration
- Inference
- Off-policy

Algorithm 1 Soft Actor-Critic

Input: θ_1, θ_2, ϕ

$\theta_1 \leftarrow \theta_1, \theta_2 \leftarrow \theta_2$

$\mathcal{D} \leftarrow \emptyset$

for each iteration **do**

for each environment step **do**

$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$

$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$

end for

for each gradient step **do**

$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$

$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$

$\alpha \leftarrow \alpha - \lambda \hat{\nabla}_\alpha J(\alpha)$

$\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$ for $i \in \{1, 2\}$

end for

end for

Output: θ_1, θ_2, ϕ

▷ Initial parameters

▷ Initialize target network weights

▷ Initialize an empty replay pool

▷ Sample action from the policy

▷ Sample transition from the environment

▷ Store the transition in the replay pool

▷ Update the Q-function parameters

▷ Update policy weights

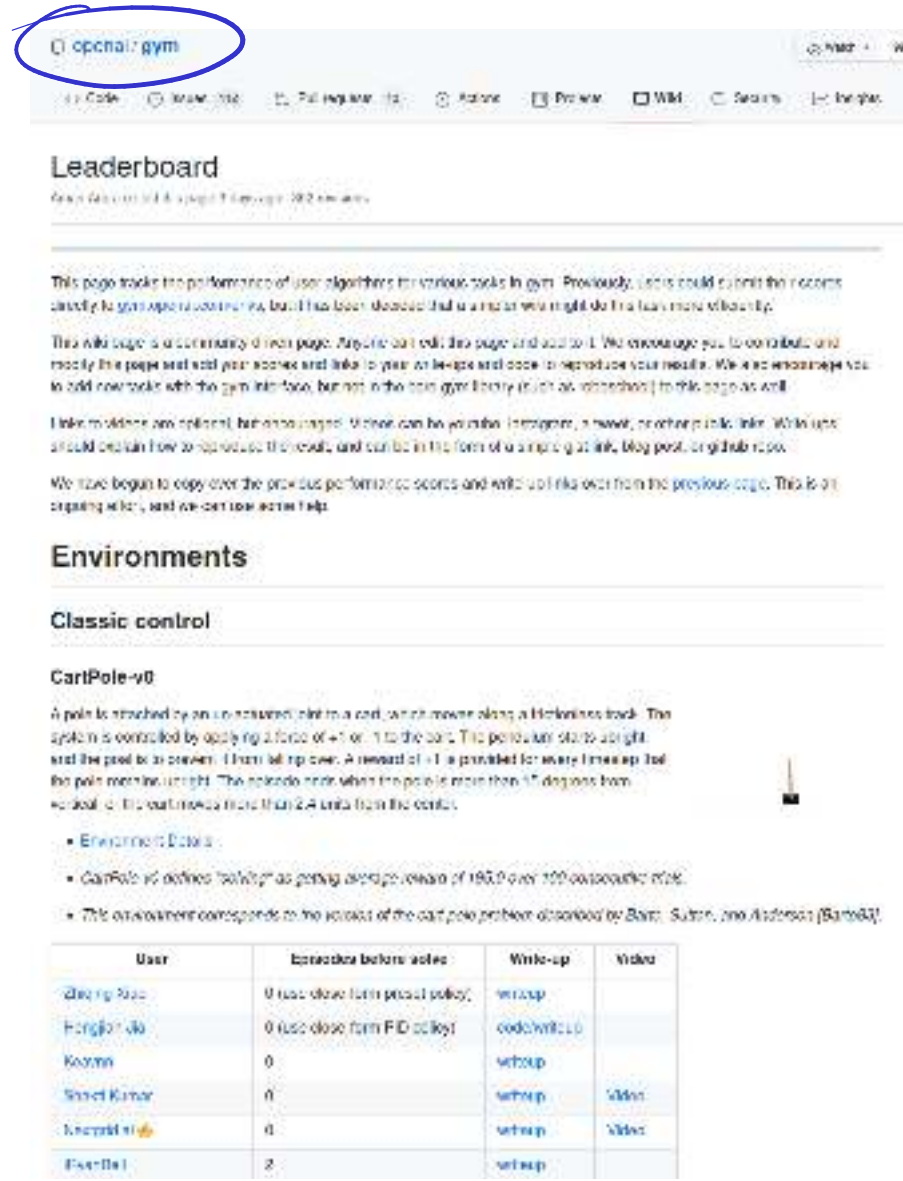
▷ Adjust temperature

▷ Update target network weights

▷ Optimized parameters

Wisdom

Deep RL: The Dream



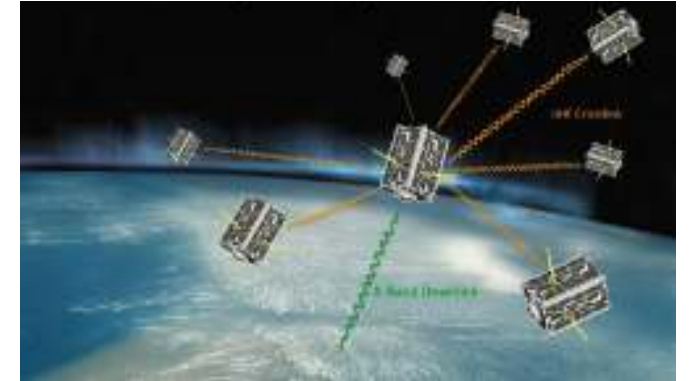
The screenshot shows the OpenAI Gym website. The URL bar shows 'openai: gym' circled in blue. The page has a navigation bar with links: Code, About, Full regular (beta), Actions, Preview, Wiki, Security, and the flag icon. The main heading is 'Leaderboard' with a subtitle 'Rank 1st out of 8 (page 1) (page 100 new items)'. Below this is a paragraph explaining the leaderboard's purpose and a community-driven page. It then lists 'Environments' and 'Classic control'. Under 'Classic control', it features 'CartPole-v0' with a description of the game and a small image of the cart. Below the description are three bullet points: 'Environment Details', 'CartPole-v0 defines `best_ep` as getting average reward of 190.0 over 100 episodes max.', and 'This environment corresponds to the version of the cart pole problem described by Bart, Sutton, and Anderson (Barto83)'. At the bottom is a table with columns: User, Episodes before solved, Wrote-up, and Video.

User	Episodes before solved	Wrote-up	Video
ding-rui	0 (use close form policy)	setup	
Ergebnis-0	0 (use close form PID policy)	code/wiki	
Kozmin	0	setup	
Sherif-Kumar	0	setup	Video
Kucukpilot	0	setup	Video
Esentel	2	setup	

Using Deep RL for your problem

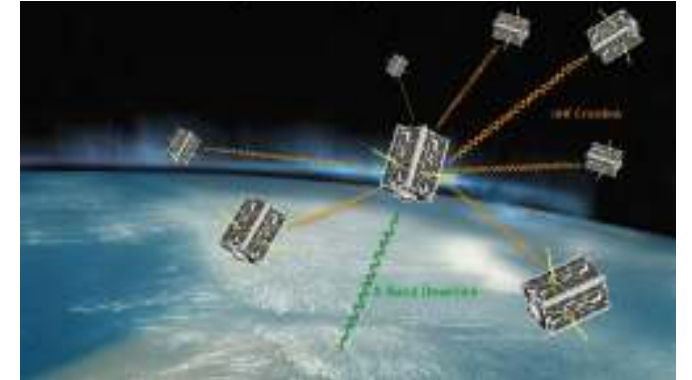
Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)



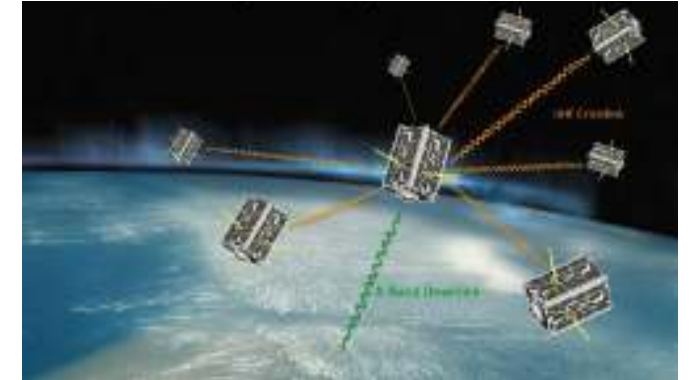
Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics



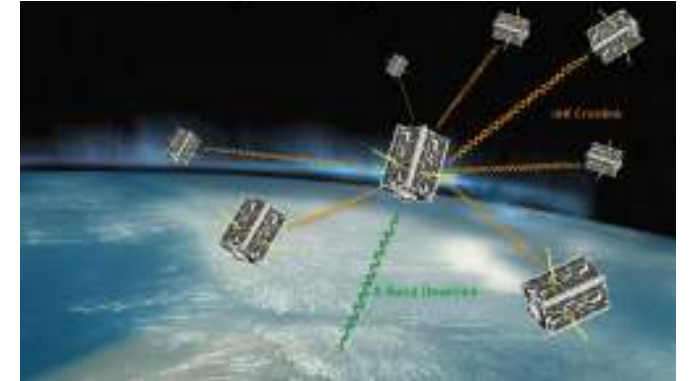
Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems



Using Deep RL for your problem

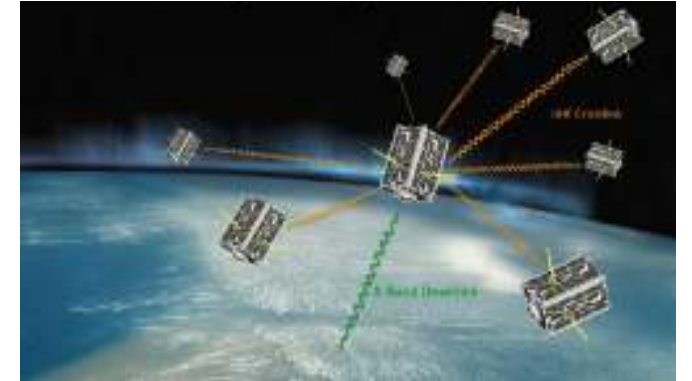
1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines



 [openai / baselines](#)

Using Deep RL for your problem

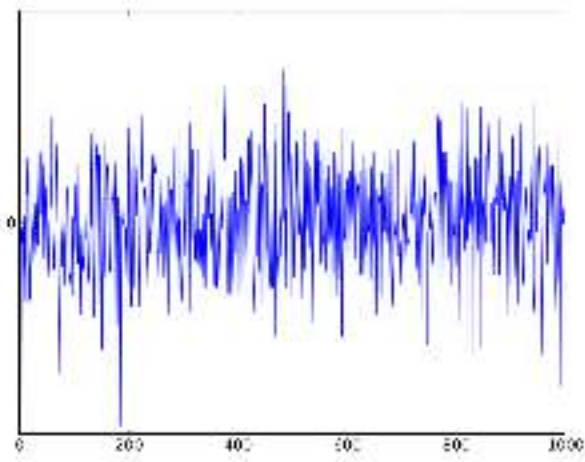
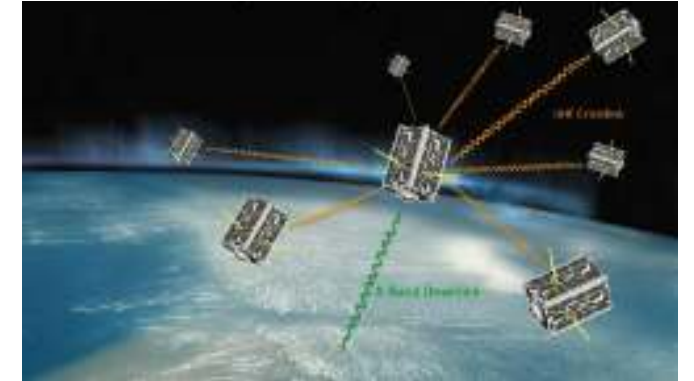
1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



 [openai / baselines](#)

Using Deep RL for your problem

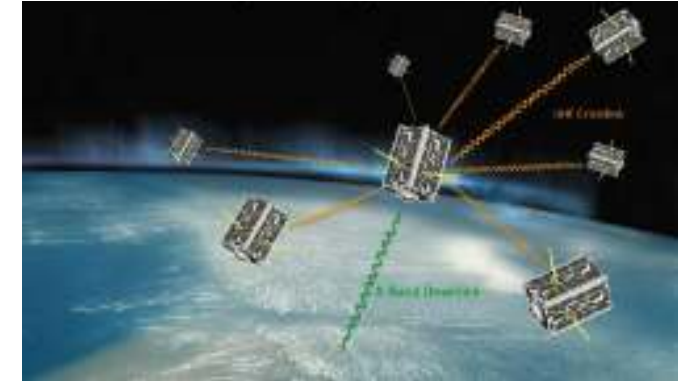
1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



 [openai / baselines](#)

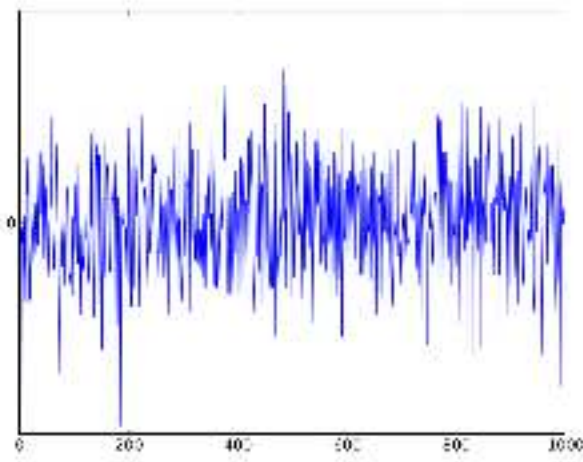
Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



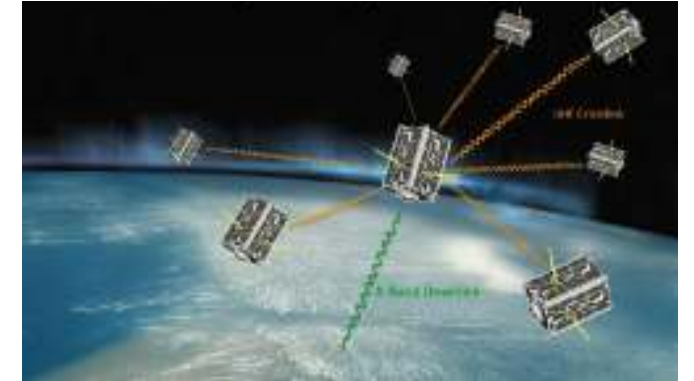
 [openai / baselines](#)

Why not?

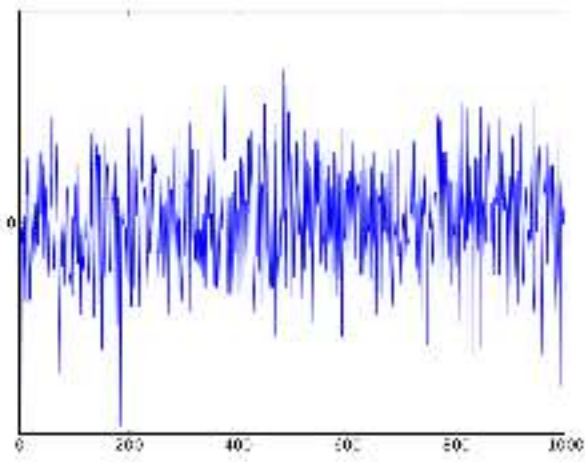


Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



 openai / **baselines**

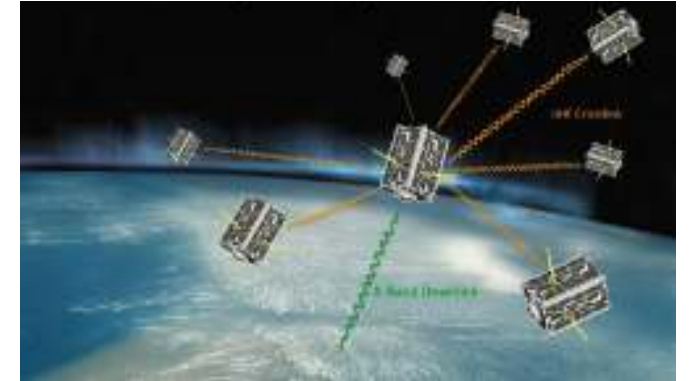


Why not?

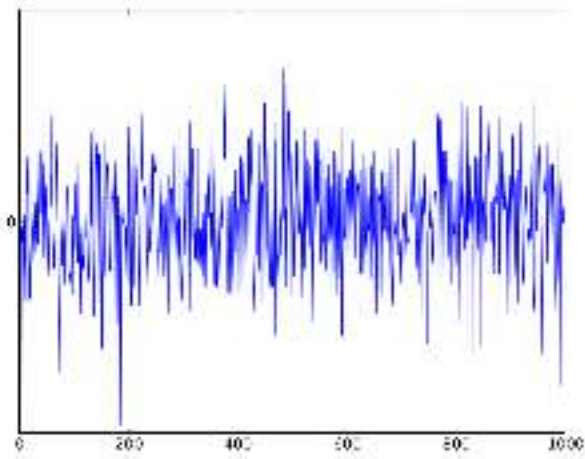
- Hyperparameters?

Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



 openai / **baselines**



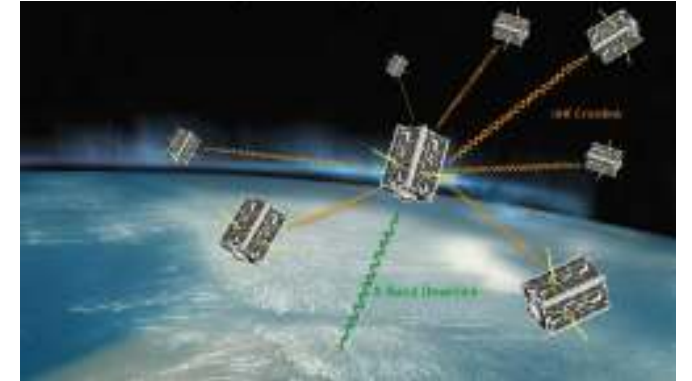
Why not?

- Hyperparameters?
- Reward scaling?

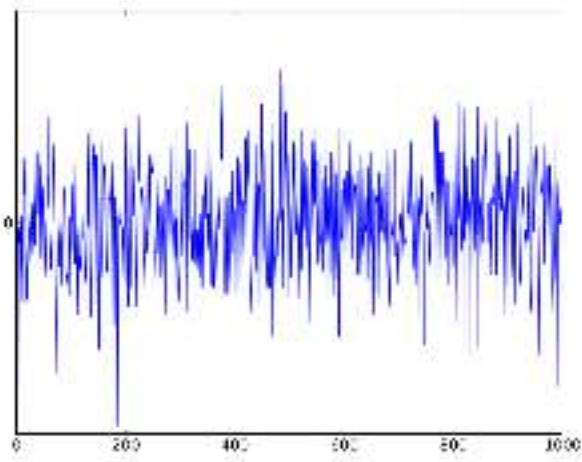
Deep RL that matters

Using Deep RL for your problem

1. Some interesting problem (smallsat swarm)
2. Spend weeks theorizing about the exact-right cost function and dynamics
3. Decide RL can solve all of your problems
4. Fire up open-ai baselines
5. Does it work??



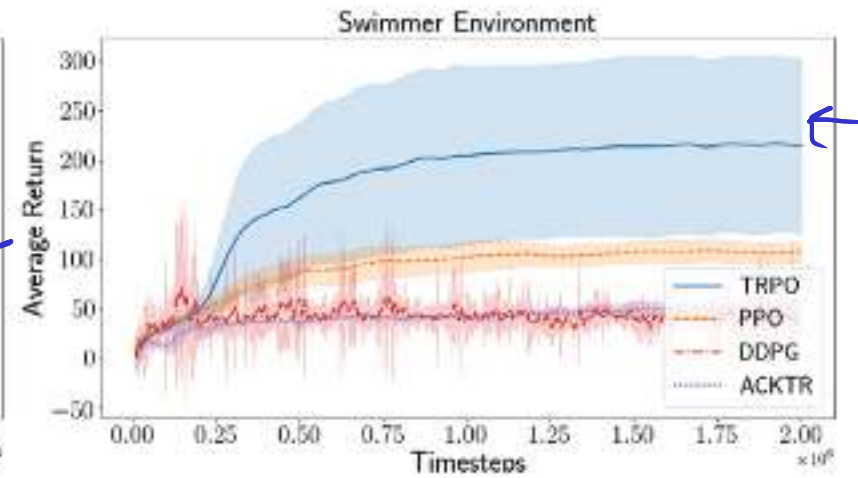
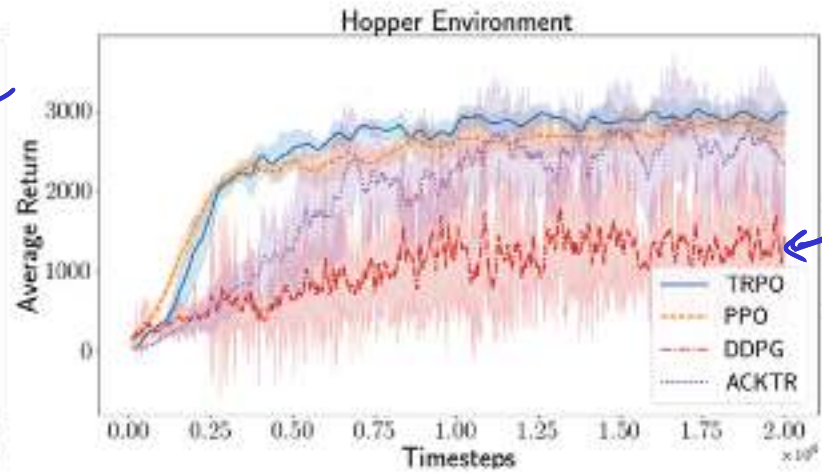
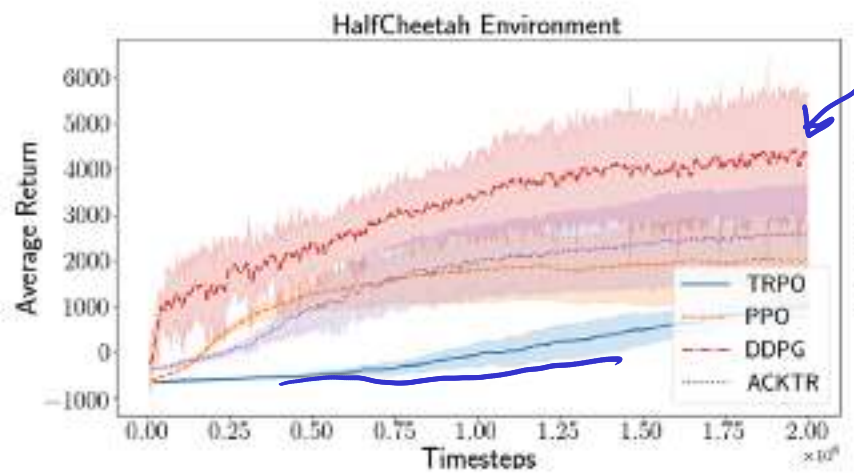
 openai / **baselines**



Why not?

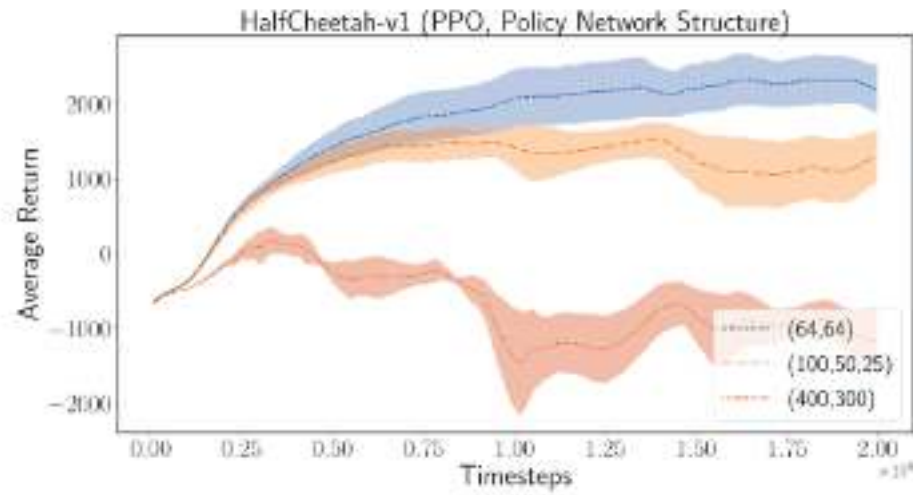
- Hyperparameters?
- Reward scaling?
- Not enough training time????

Algorithms

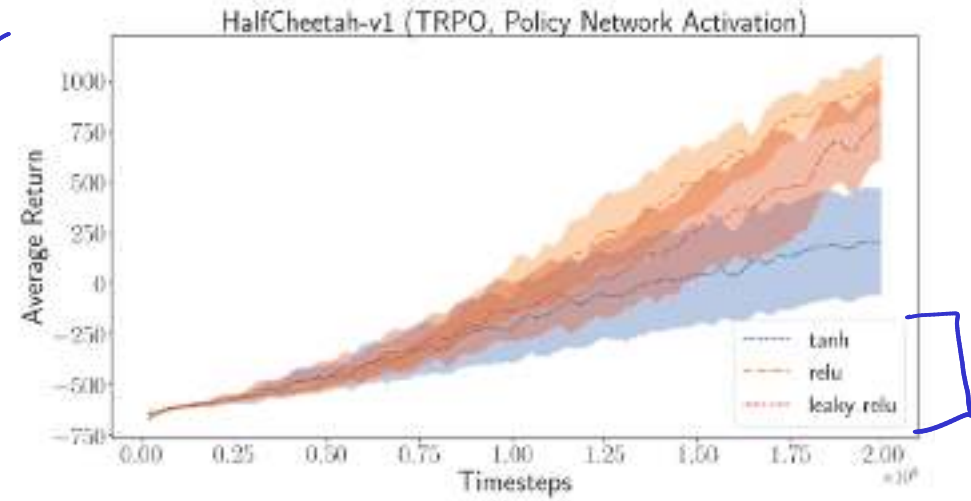
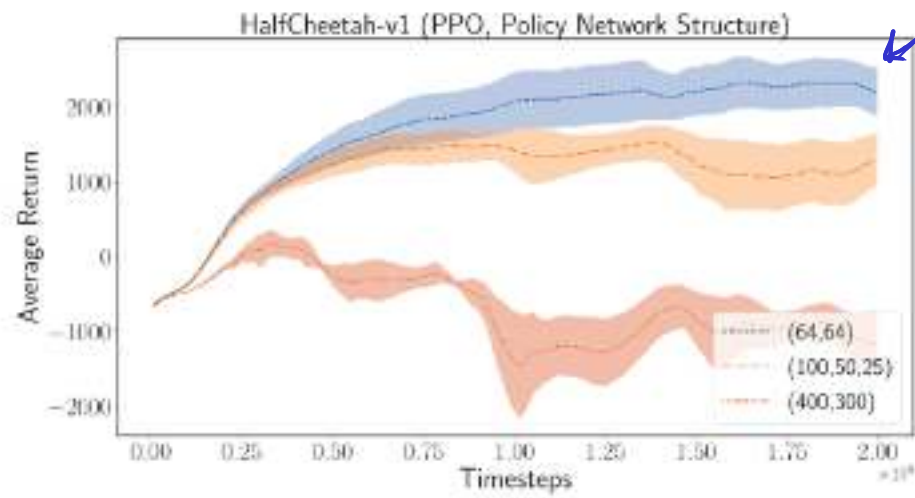


Policy Network Architecture

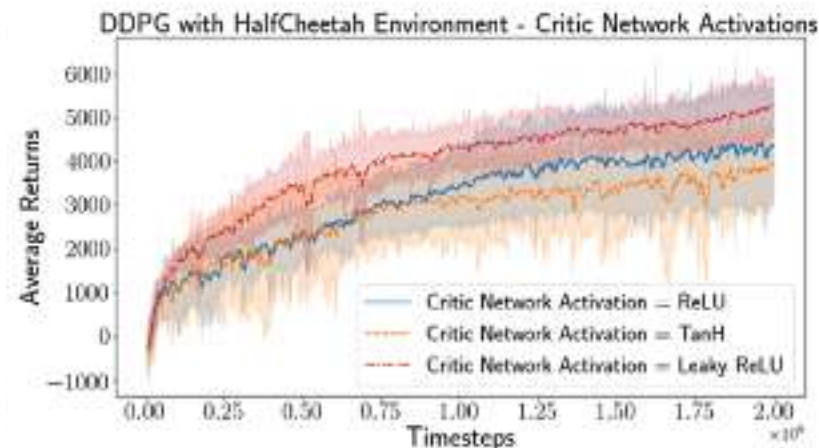
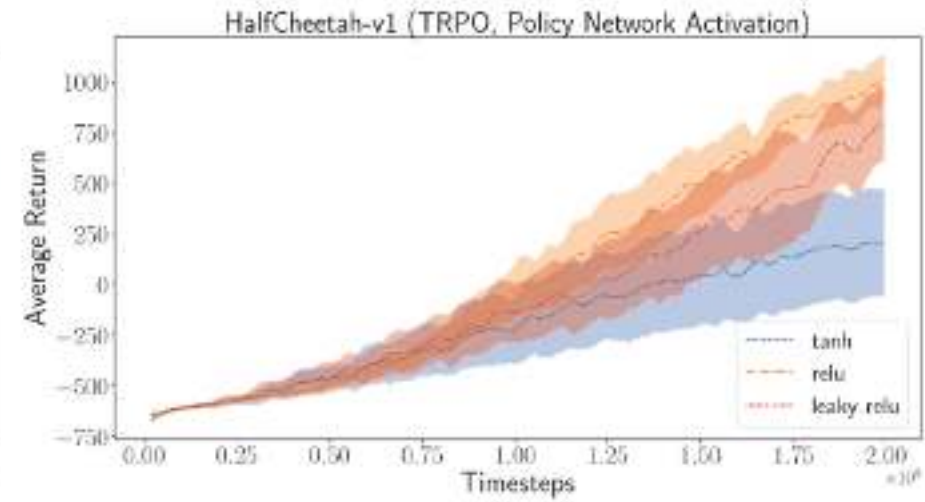
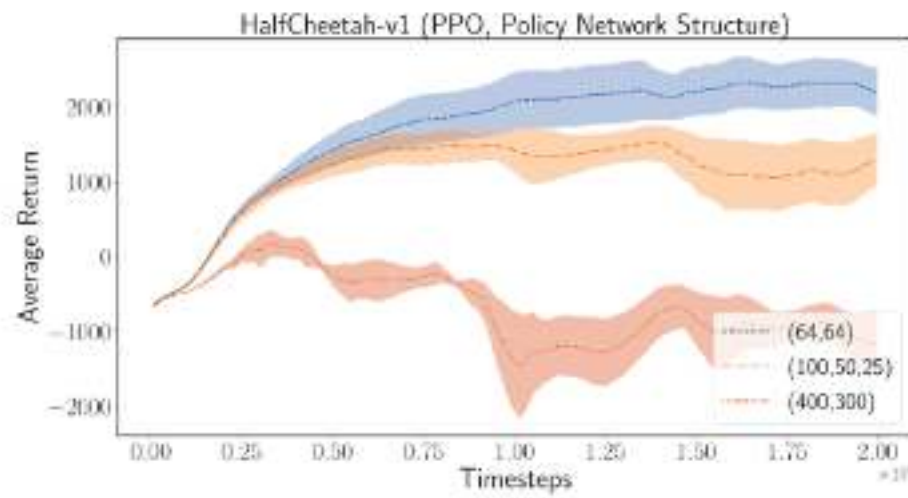
Policy Network Architecture



Policy Network Architecture



Policy Network Architecture



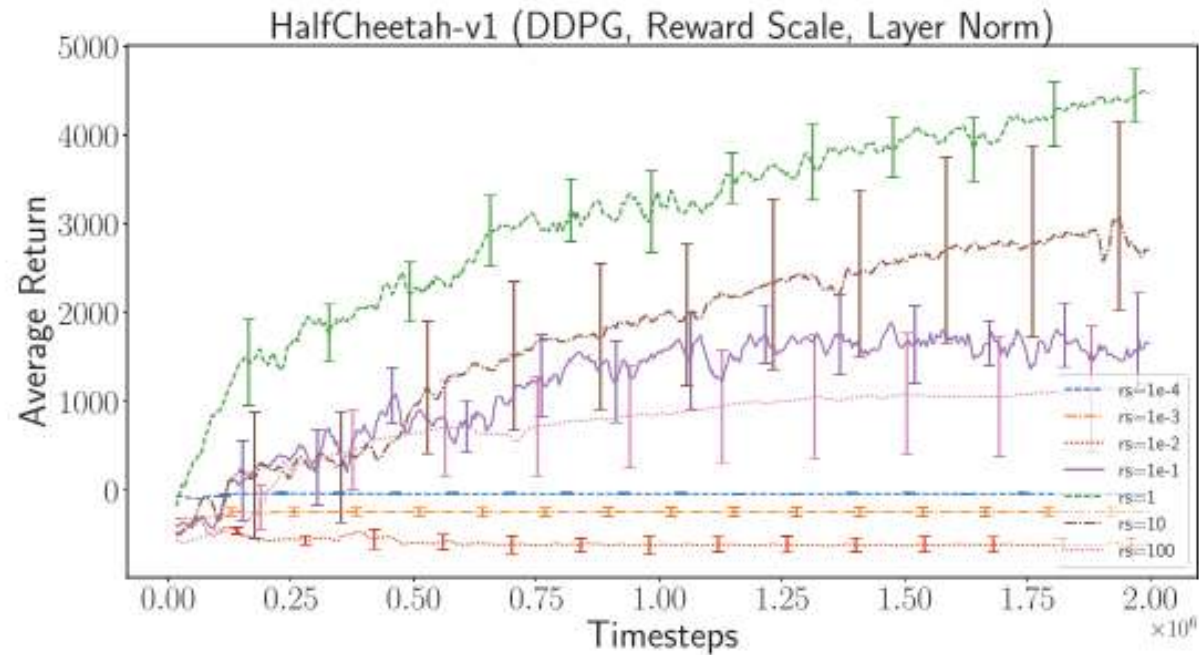
Reward Rescaling

Reward Rescaling

"simply multiplying the rewards generated from an environment by some scalar"

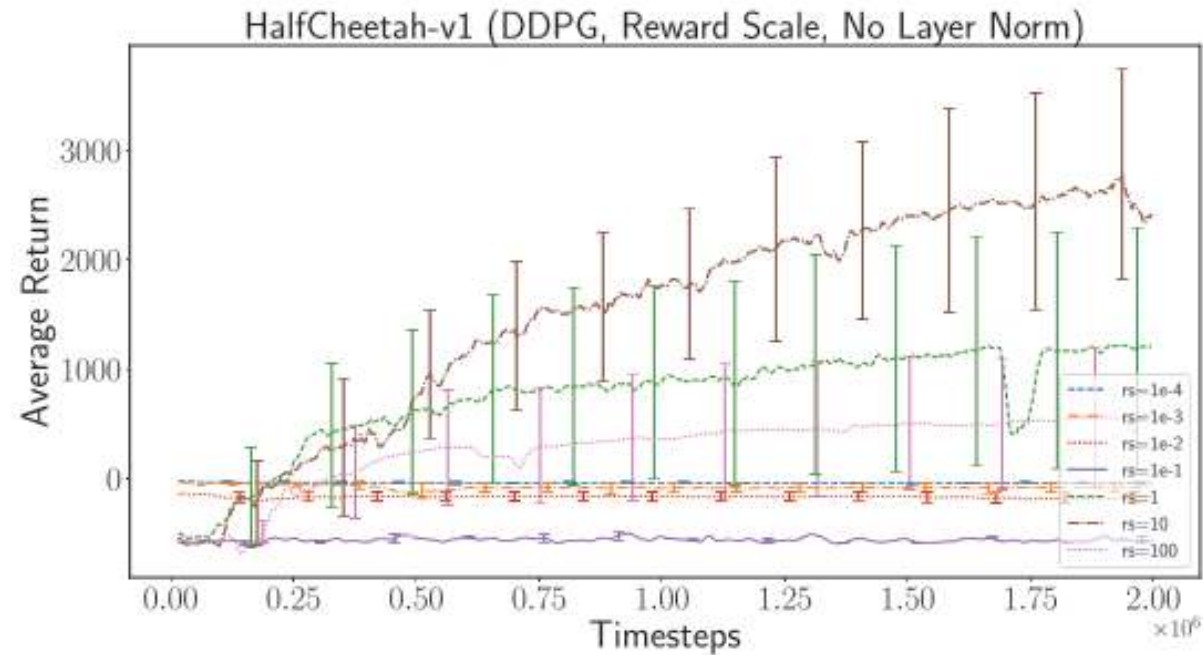
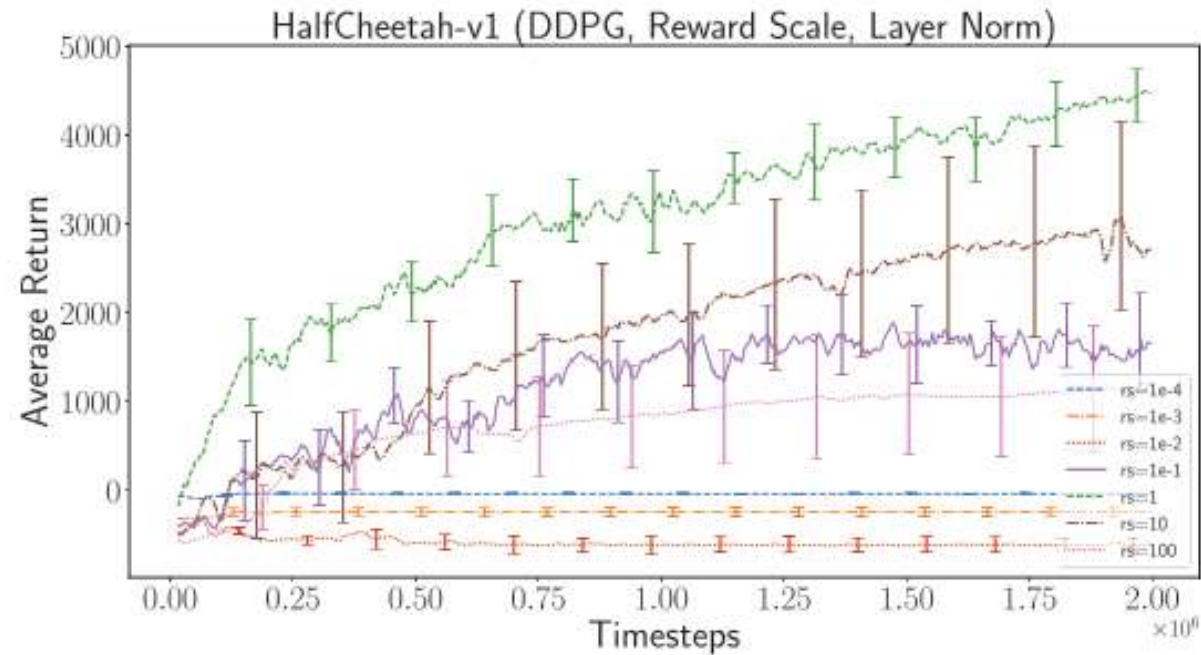
Reward Rescaling

"simply multiplying the rewards generated from an environment by some scalar"



Reward Rescaling

"simply multiplying the rewards generated from an environment by some scalar"



Statistical Significance

Statistical Significance

"Unfortunately, in recent reported results, it is not uncommon for the top-N trials to be selected from among several trials (Wu et al. 2017; Mnih et al. 2016)"

Statistical Significance

"Unfortunately, in recent reported results, it is not uncommon for the top-N trials to be selected from among several trials (Wu et al. 2017; Mnih et al. 2016)"

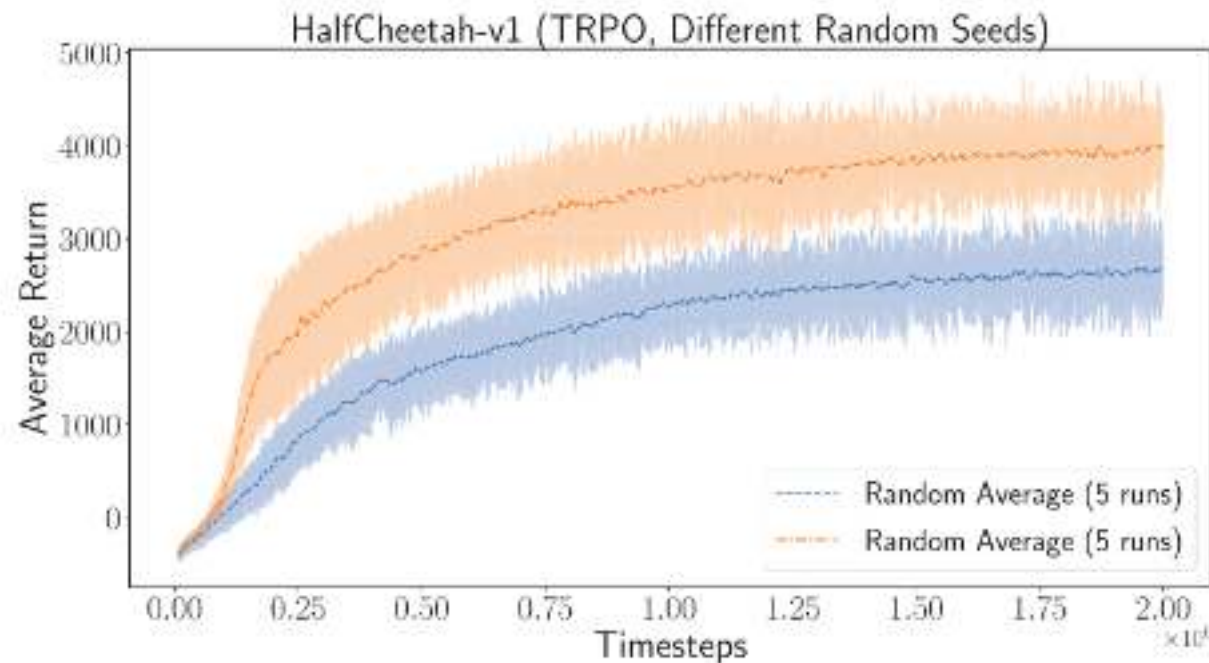
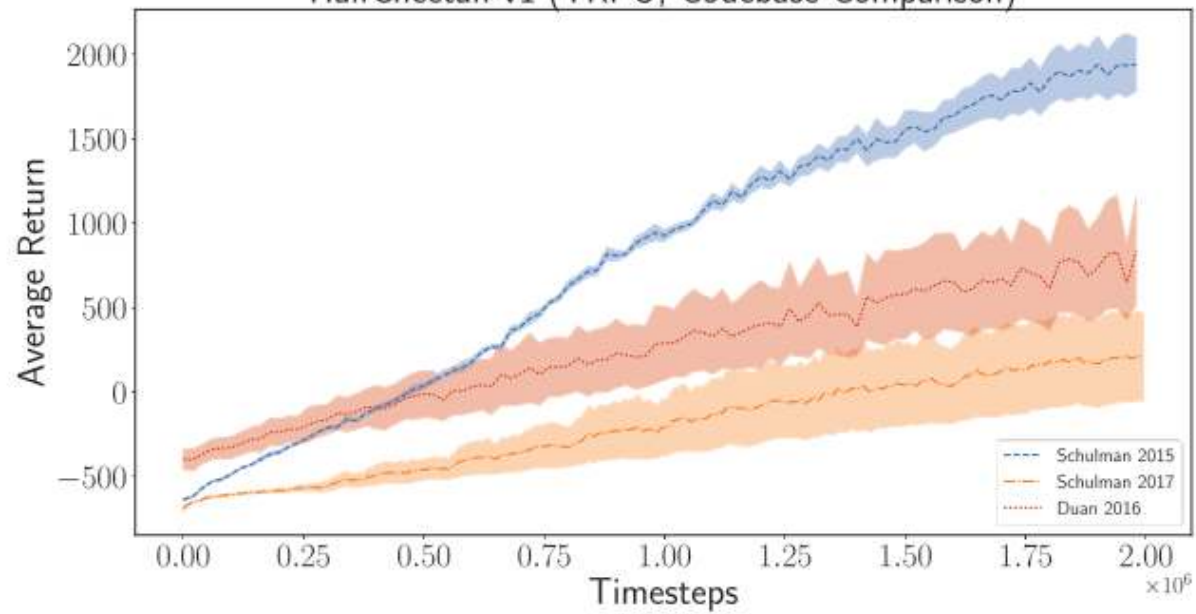


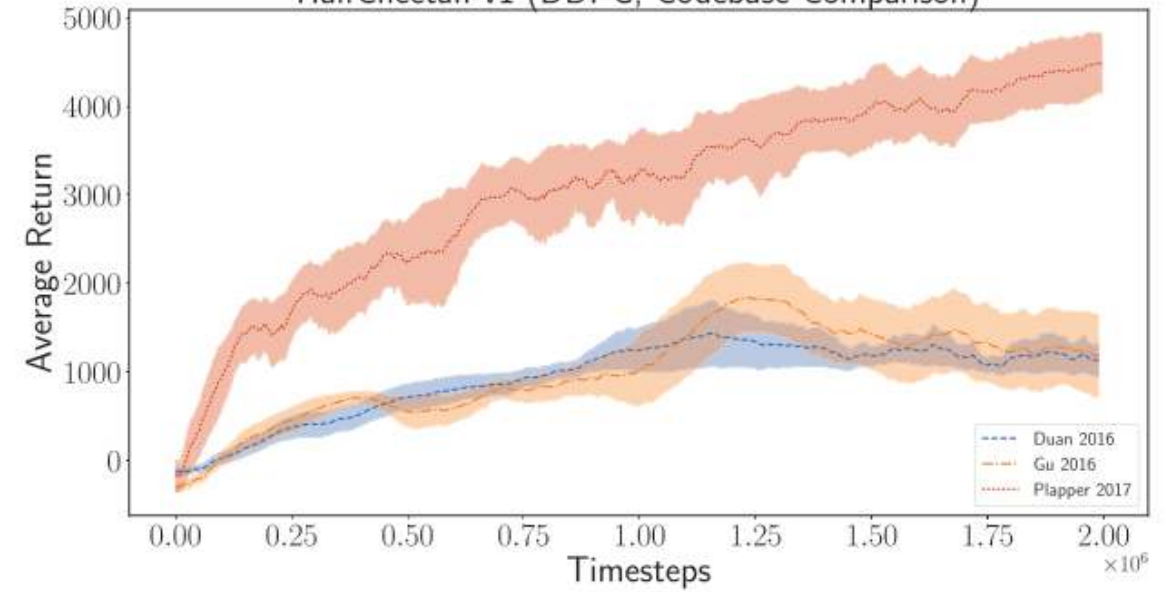
Figure 5: TRPO on HalfCheetah-v1 using the same hyperparameter configurations averaged over two sets of 5 different random seeds each. The average 2-sample t -test across entire training distribution resulted in $t = -9.0916$, $p = 0.0016$.

Codebases

HalfCheetah-v1 (TRPO, Codebase Comparison)



HalfCheetah-v1 (DDPG, Codebase Comparison)



How to choose an RL Algorithm

How to choose an RL Algorithm

(According to Sergey Levine)

How to choose an RL Algorithm

(According to Sergey Levine)



How to choose an RL Algorithm

(According to Sergey Levine)



Sample
Efficiency

How to choose an RL Algorithm

(According to Sergey Levine)



Sample
Efficiency

Ease of Use
/ Stability

How to choose an RL Algorithm

(According to Sergey Levine)



Sample
Efficiency

Ease of Use
/ Stability

← fewer samples



How to choose an RL Algorithm

(According to Sergey Levine)



Sample
Efficiency

Ease of Use
/ Stability

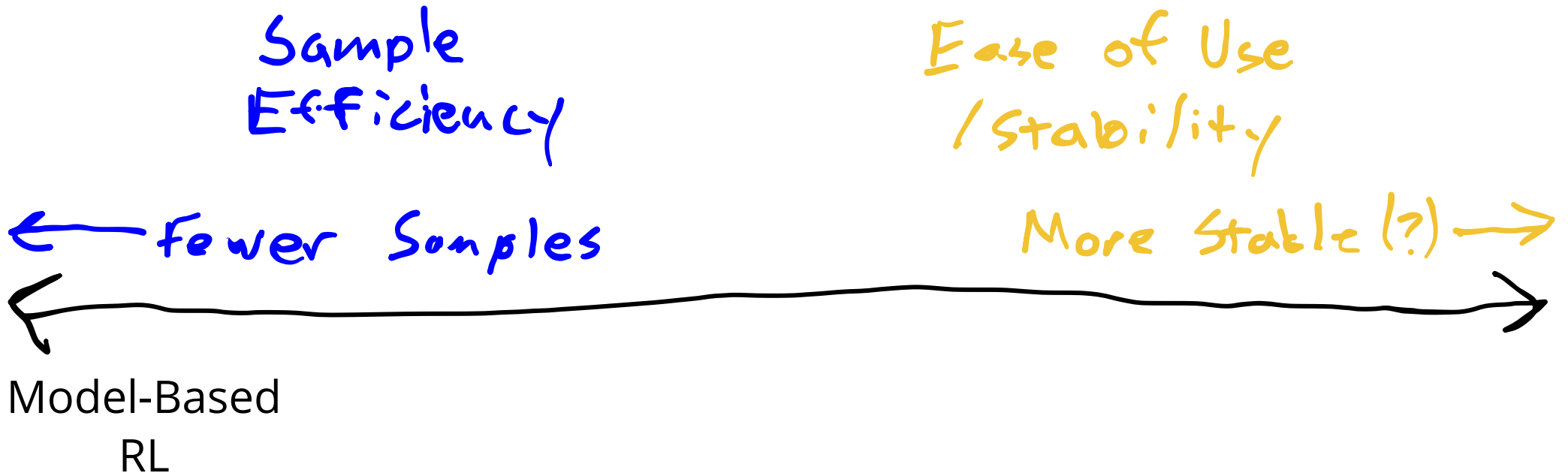
← fewer Samples

More Stable (?) →



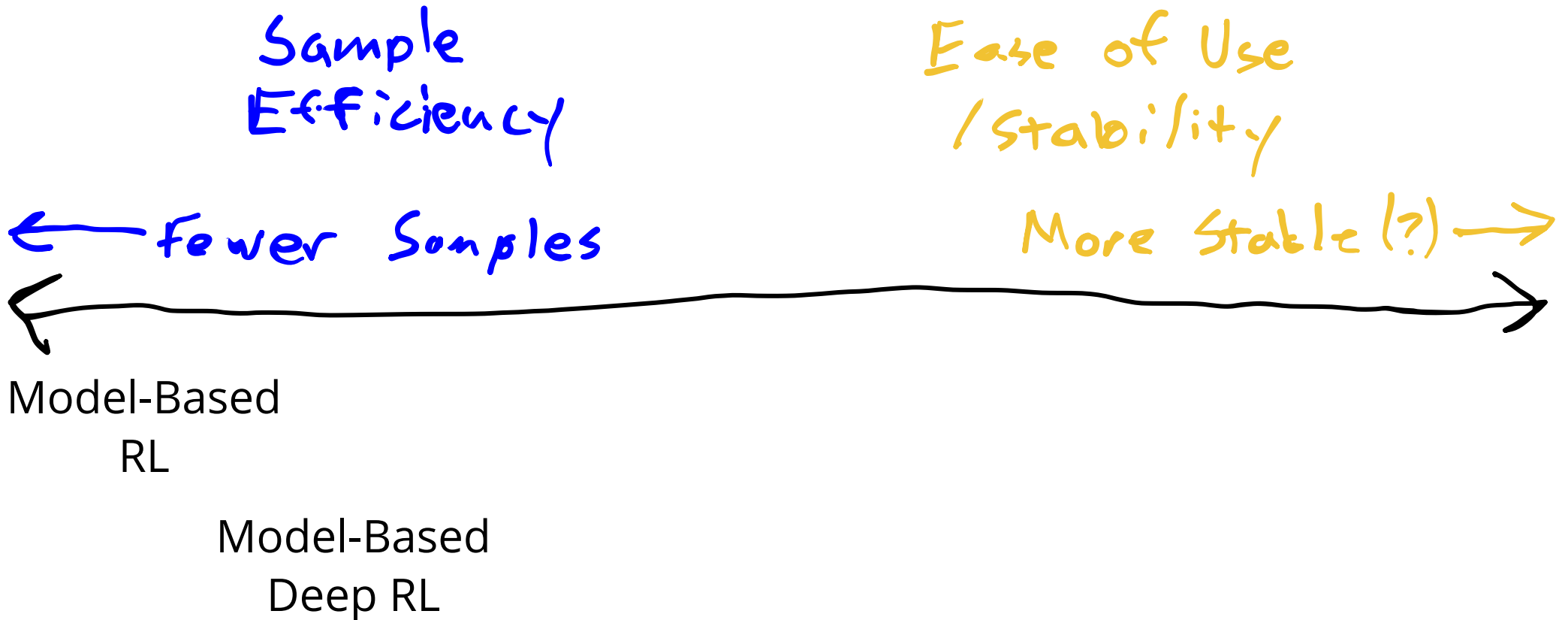
How to choose an RL Algorithm

(According to Sergey Levine)



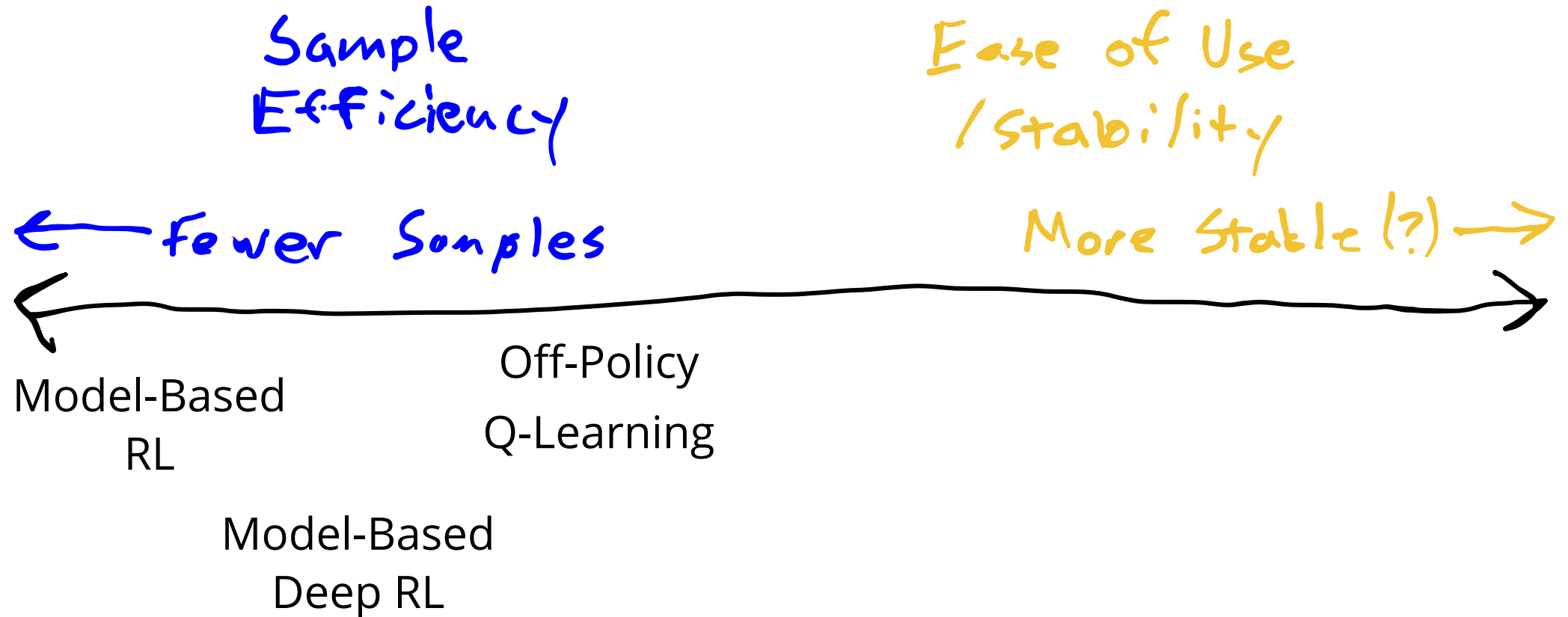
How to choose an RL Algorithm

(According to Sergey Levine)



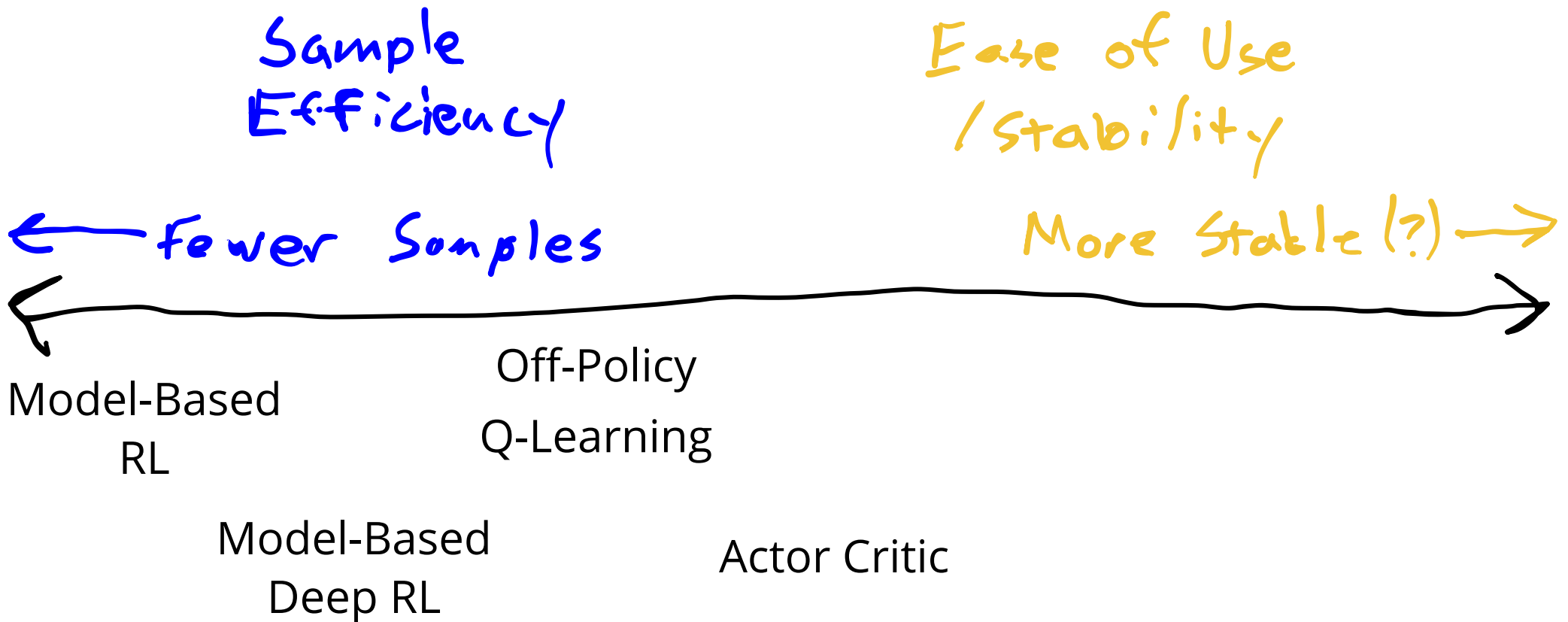
How to choose an RL Algorithm

(According to Sergey Levine)



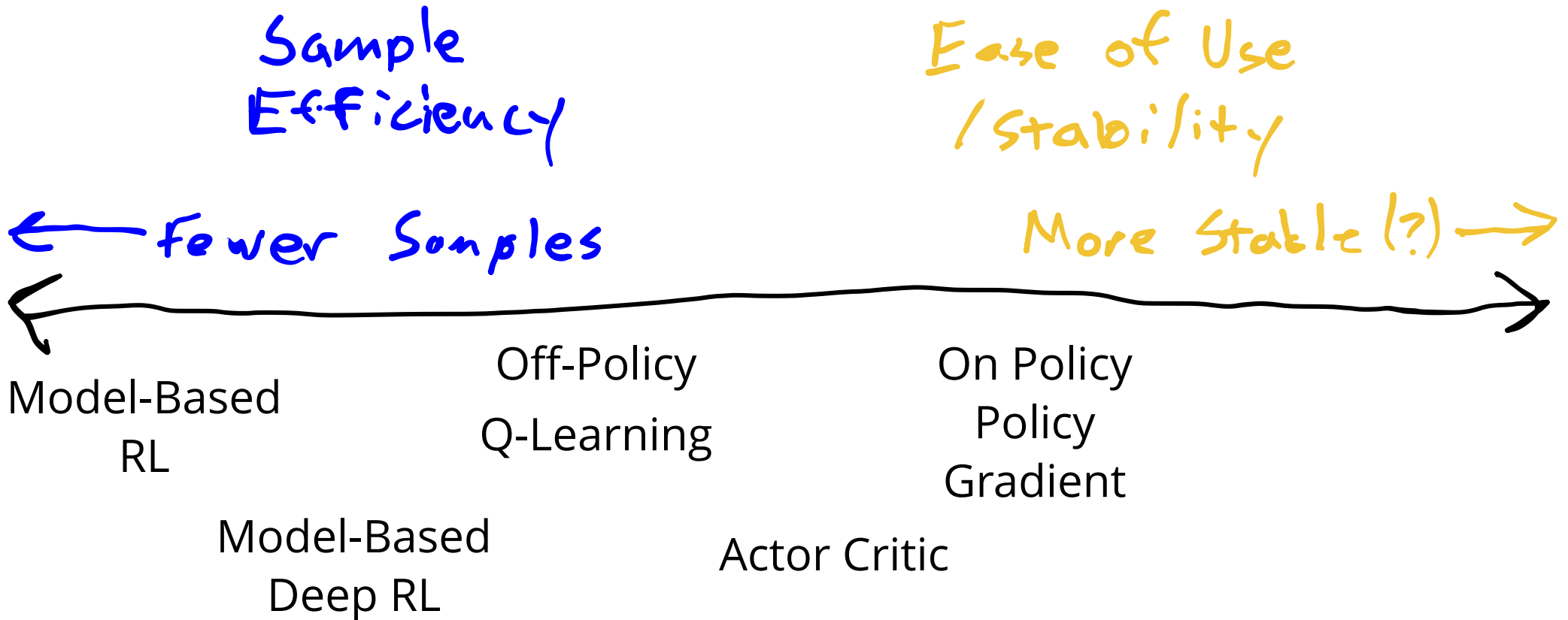
How to choose an RL Algorithm

(According to Sergey Levine)



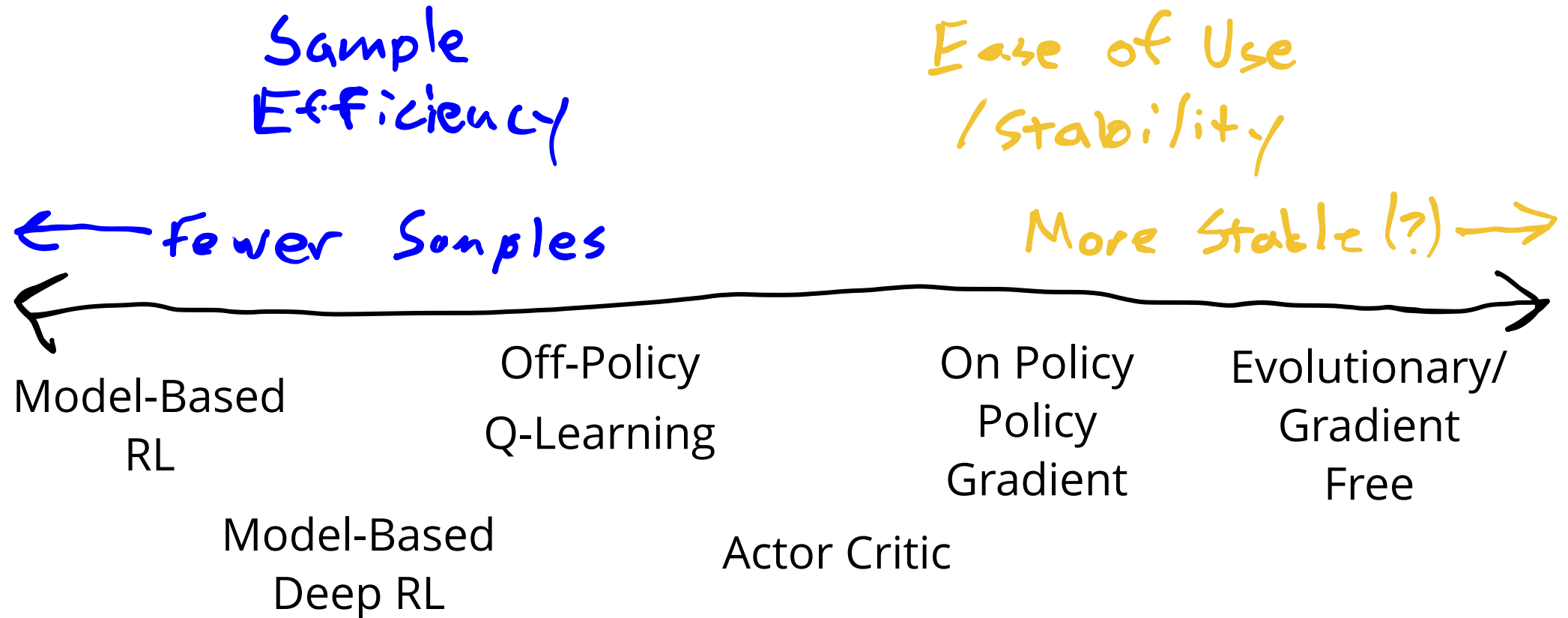
How to choose an RL Algorithm

(According to Sergey Levine)



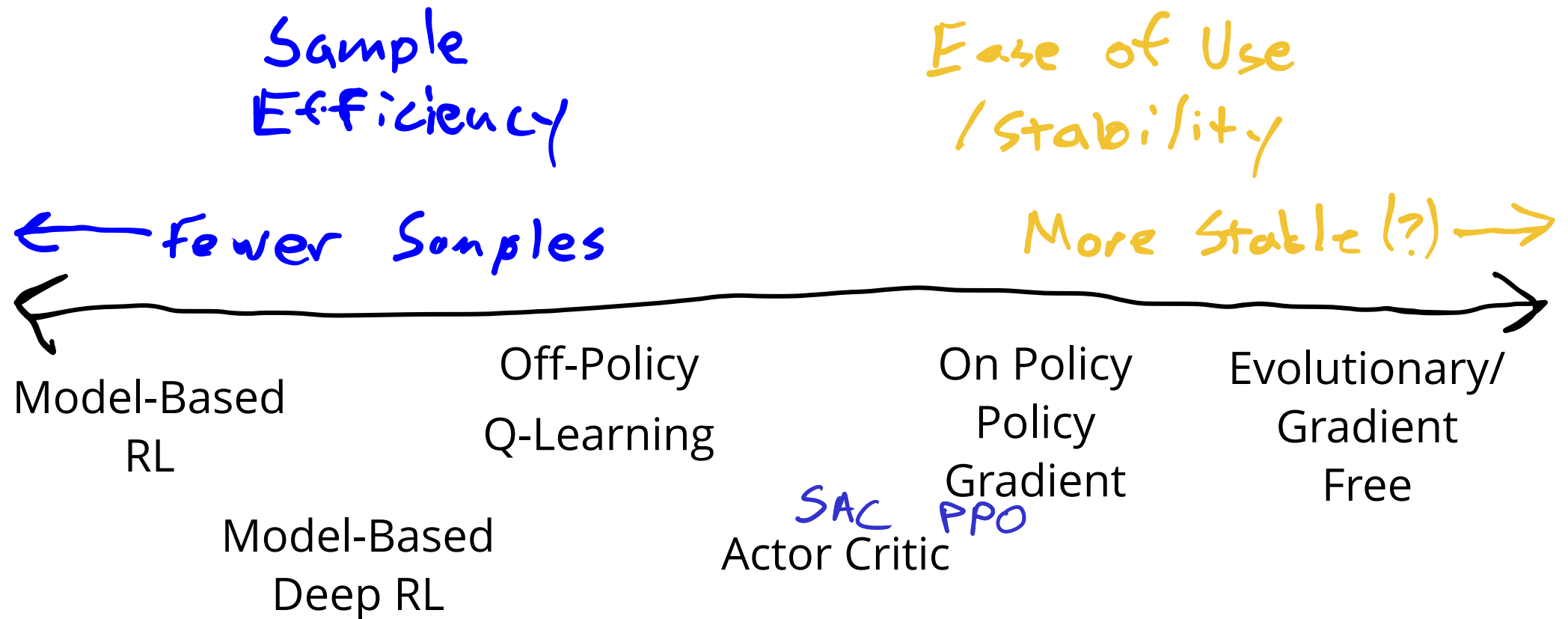
How to choose an RL Algorithm

(According to Sergey Levine)



How to choose an RL Algorithm

(According to Sergey Levine)



(Most people use SAC or PPO)

Where Does RL Work?

Where Does RL Work?

- Cooling servers



Where Does RL Work?

- Cooling servers
- Winning at Go

