# Proximal Policy Optimization (PPO)

default reinforcement learning algorithm at OpenAI
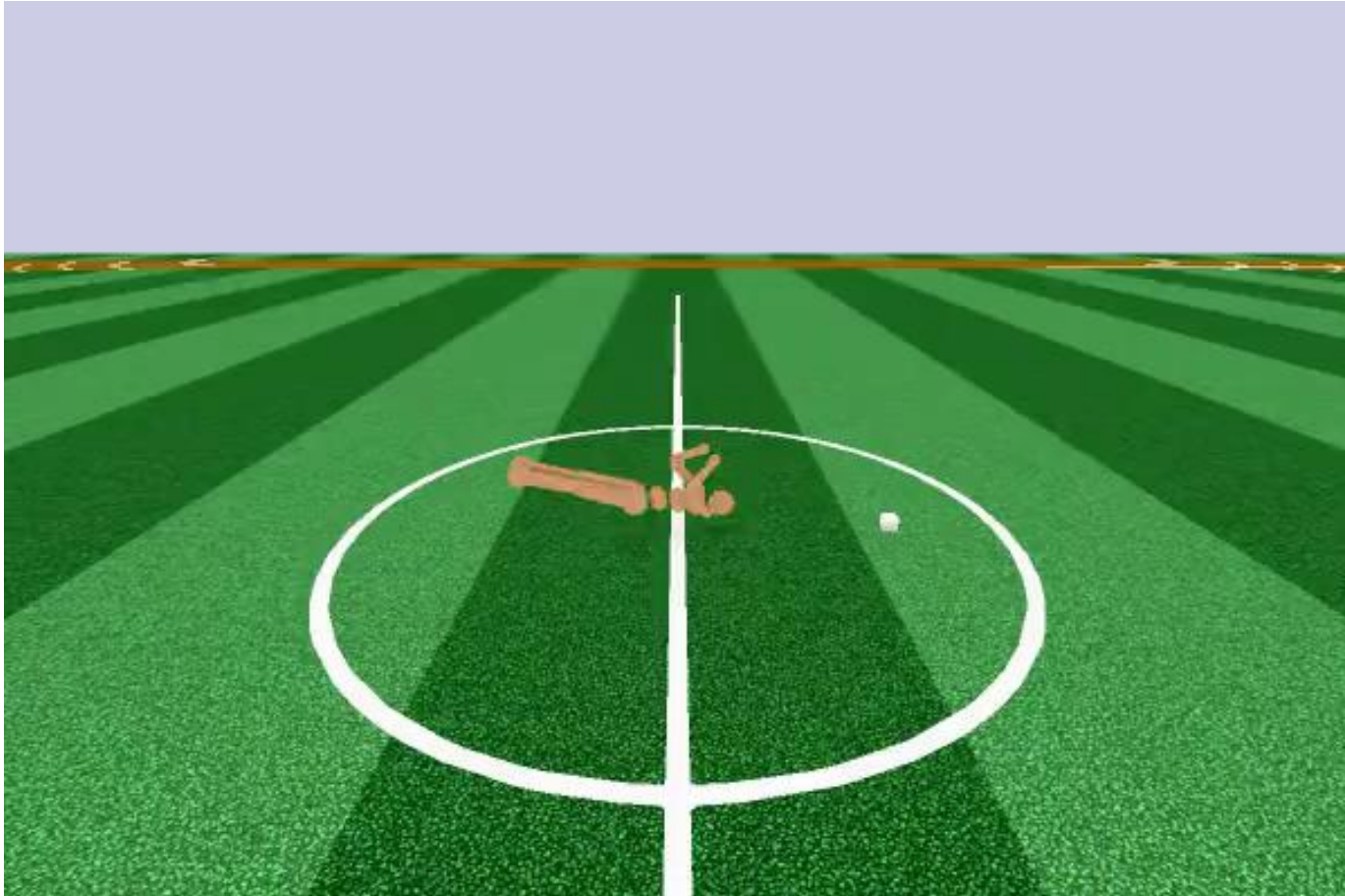
| Policy Gradient | → | On-policy → Off-policy | → | Add constraint |

# DeepMind

# OpenAI

# Policy Gradient (Review)

# Basic Components

You cannot control

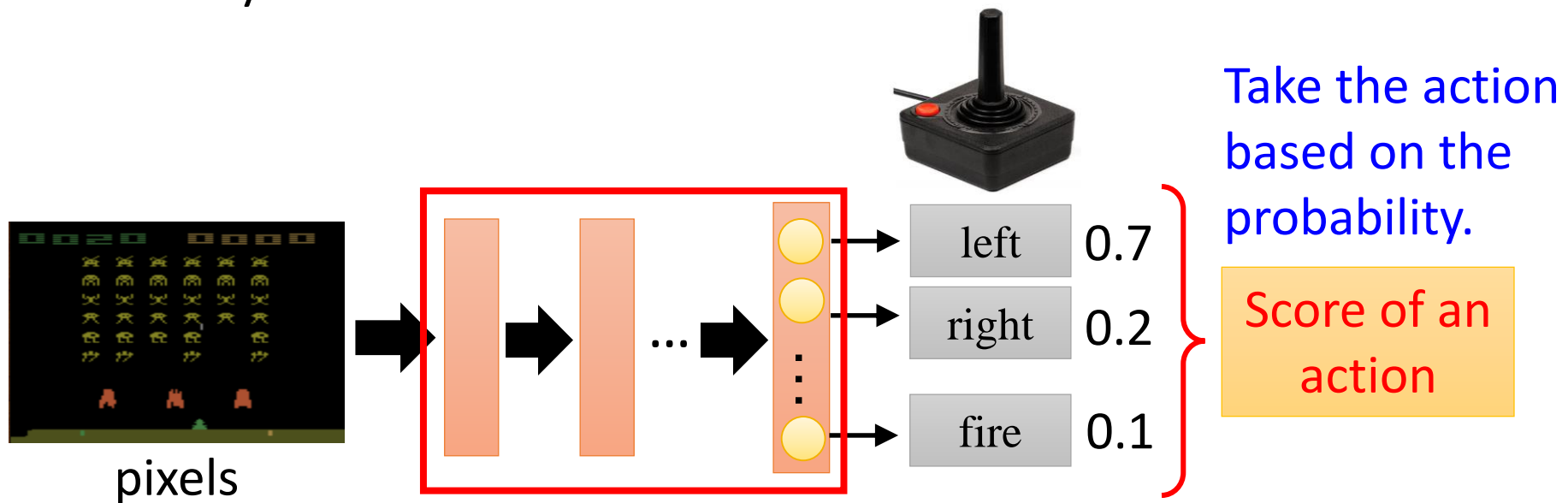|  | Actor | Env | Reward Function |
|---|---|---|---|
| Video Game |  |  | Get 20 scores when killing a monster |
| Go |  |  | The rule of GO |

# Policy of Actor

- Policy $\pi$ is a network with parameter $\theta$
  - Input: the observation of machine represented as a vector or a matrix
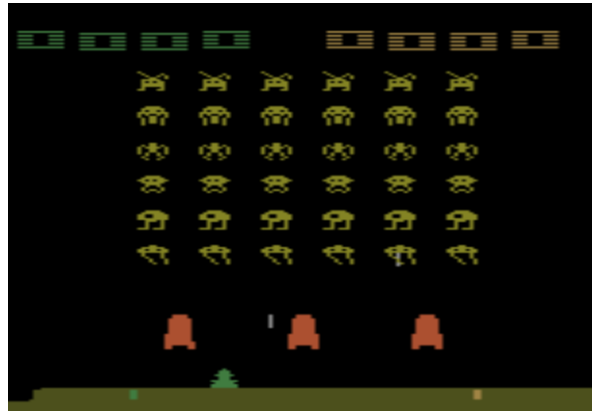  - Output: each action corresponds to a neuron in output layer



pixels

| | |
|---|---|
| left | 0.7 |
| right | 0.2 |
| fire | 0.1 |

Take the action based on the probability.

Score of an action

# Example: Playing Video Game

Start with observation $s_1$

Observation $s_2$

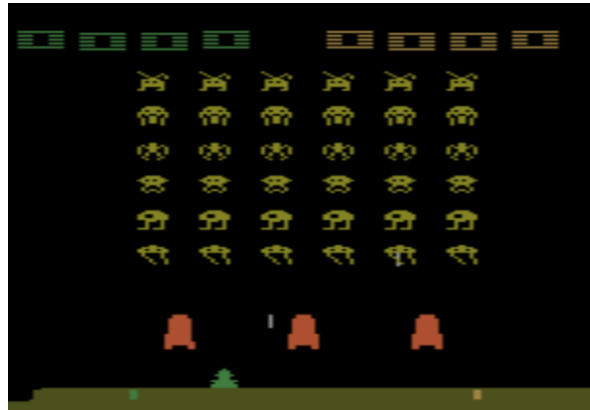Observation $s_3$

Obtain reward $r_1 = 0$

Obtain reward $r_2 = 5$

Action $a_1$: "right"

Action $a_2$: "fire"

(kill an alien)

# *Example: Playing Video Game*

Start with
observation $s_1$

Observation $s_2$

Observation $s_3$

This is an ***episode***.

Total reward:

$$R = \sum_{t=1}^{T} r_t$$

After many turns

· · · · · · · · · · · · ·➤

Game Over
(spaceship destroyed)

We want the total
reward be maximized.

Obtain reward $r_T$
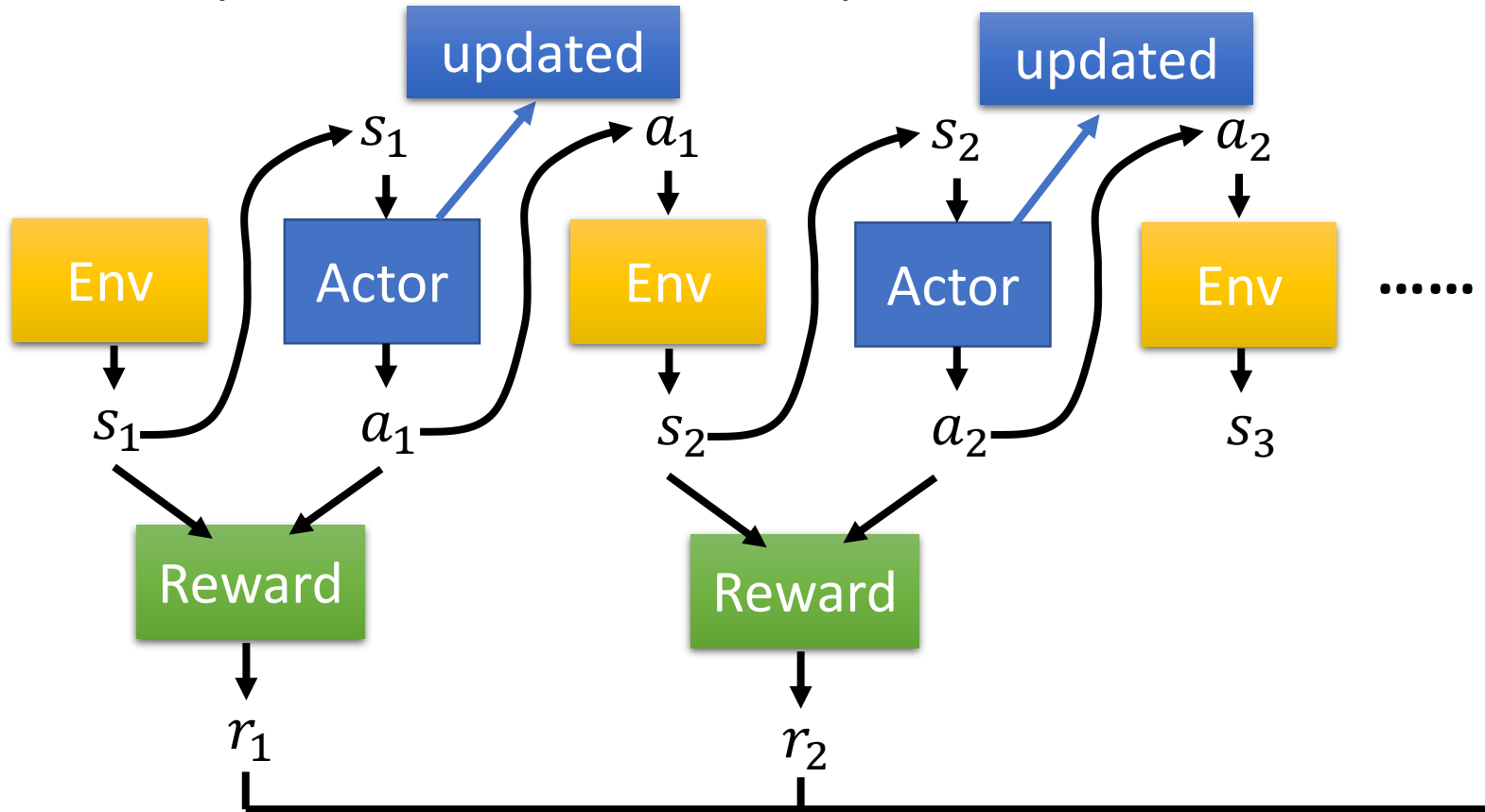
Action $a_T$

# Actor, Environment, Reward



**Trajectory** $\tau = \{s_1, a_1, s_2, a_2, \cdots, s_T, a_T\}$

$p_\theta(\tau)$

$= p(s_1)p_\theta(a_1|s_1)p(s_2|s_1,a_1)p_\theta(a_2|s_2)p(s_3|s_2,a_2)\cdots$

$= p(s_1)\prod_{t=1}^{T} p_\theta(a_t|s_t)p(s_{t+1}|s_t,a_t)$

# Actor, Environment, Reward



**Expected Reward**

$$\bar{R}_\theta = \sum_\tau R(\tau) p_\theta(\tau) = E_{\tau \sim p_\theta(\tau)}[R(\tau)]$$

$$R(\tau) = \sum_{t=1}^{T} r_t$$

# Policy Gradient

$$\bar{R}_\theta = \sum_\tau R(\tau) p_\theta(\tau) \qquad \nabla \bar{R}_\theta = ?$$

$$\nabla \bar{R}_\theta = \sum_\tau R(\tau) \nabla p_\theta(\tau) = \sum_\tau R(\tau) p_\theta(\tau) \frac{\nabla p_\theta(\tau)}{p_\theta(\tau)}$$

$R(\tau)$ do not have to be differentiable

It can even be a black box.

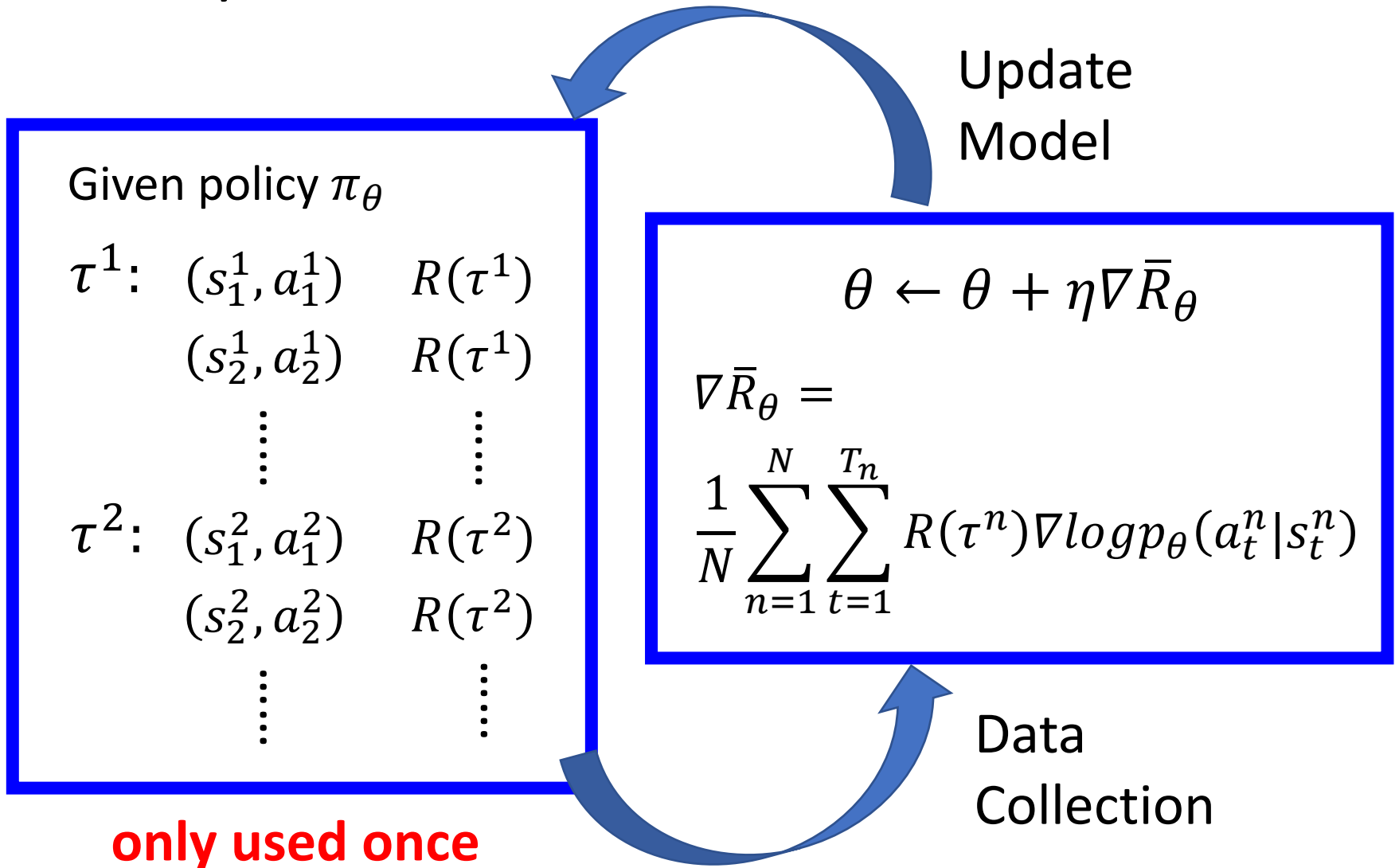$$= \sum_\tau R(\tau) p_\theta(\tau) \nabla log p_\theta(\tau)$$

$$\nabla f(x) = f(x) \nabla log f(x)$$

$$= E_{\tau \sim p_\theta(\tau)}[R(\tau) \nabla log p_\theta(\tau)] \approx \frac{1}{N} \sum_{n=1}^{N} R(\tau^n) \nabla log p_\theta(\tau^n)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} R(\tau^n) \nabla log p_\theta(a_t^n | s_t^n)$$

# Policy Gradient

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)}[R(\tau)\nabla log p_\theta(\tau)]$$

Given policy $\pi_\theta$

$\tau^1:$    $(s_1^1, a_1^1)$    $R(\tau^1)$

        $(s_2^1, a_2^1)$    $R(\tau^1)$

       $\vdots$         $\vdots$

$\tau^2:$    $(s_1^2, a_1^2)$    $R(\tau^2)$

        $(s_2^2, a_2^2)$    $R(\tau^2)$

       $\vdots$         $\vdots$

**only used once**

Update Model

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta = \frac{1}{N}\sum_{n=1}^{N}\sum_{t=1}^{T_n} R(\tau^n)\nabla log p_\theta(a_t^n | s_t^n)$$
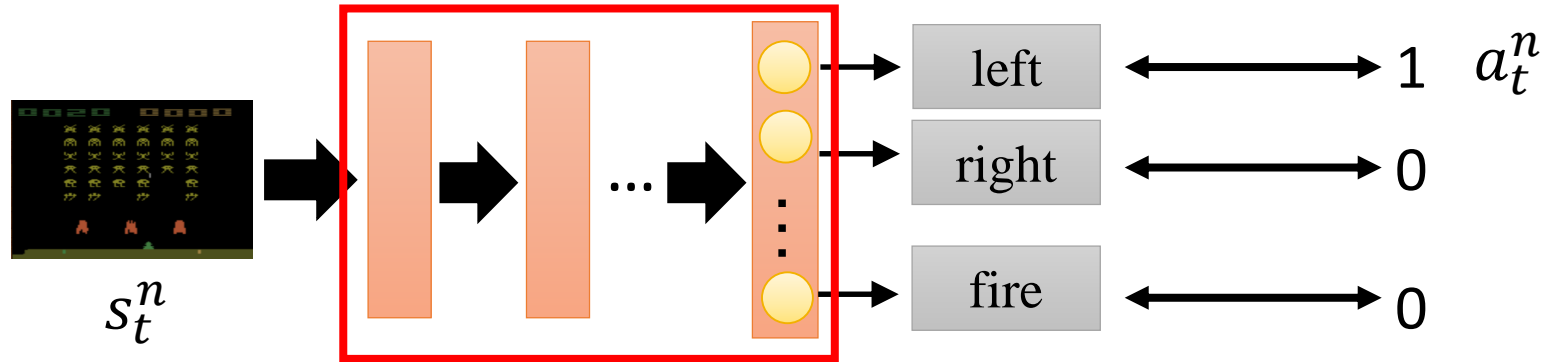
Data Collection

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

# *Implementation*

$$\nabla \bar{R}_\theta = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} R(\tau^n) \nabla log p_\theta(a_t^n | s_t^n)$$

Consider as classification problem

$$s_t^n \quad a_t^n \quad R(\tau^n)$$



$s_t^n$

left $\longleftrightarrow$ 1 $a_t^n$

right $\longleftrightarrow$ 0

fire $\longleftrightarrow$ 0

$$\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} log p_\theta(a_t^n | s_t^n) \quad \xrightarrow{\text{TF, pyTorch ...}} \quad \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} \nabla log p_\theta(a_t^n | s_t^n)$$

$$\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} \underline{R(\tau^n)} log p_\theta(a_t^n | s_t^n) \quad \longrightarrow \quad \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} \underline{R(\tau^n)} \nabla log p_\theta(a_t^n | s_t^n)$$
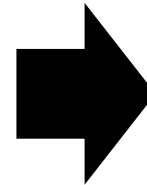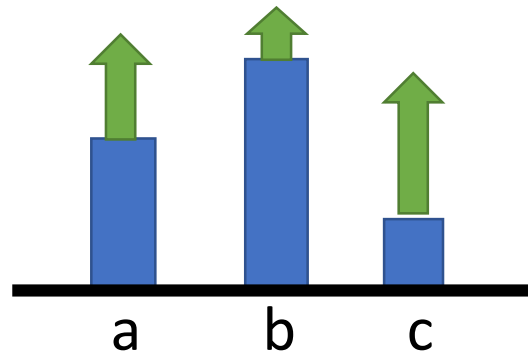
# Tip 1: Add a Baseline

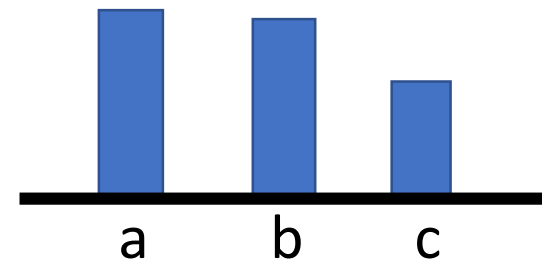$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

It is possible that $R(\tau^n)$ is always positive.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla log p_\theta(a_t^n | s_t^n) \qquad b \approx E[R(\tau)]$$
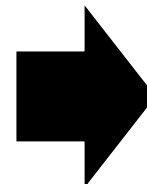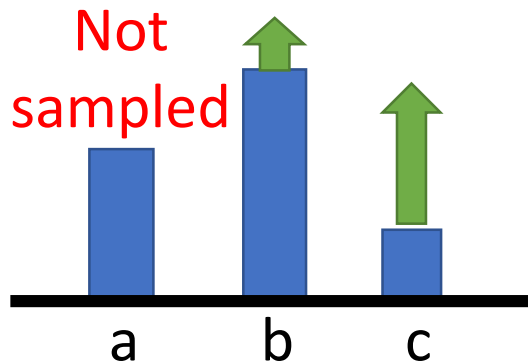
Ideal case

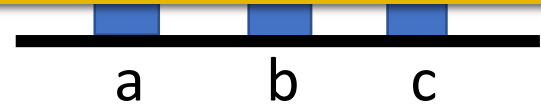It is probability …

Sampling
……

Not sampled

The probability of the actions not sampled will decrease.

a b c

# Tip 2: Assign Suitable Credit

考慮某個時間點以後到終點的reward

$\times 3$    $\times -2$    $\times -2$      $\times -7$    $\times -2$    $\times -2$

0+-2      -2

$(s_a, a_1)$    $(s_b, a_2)$    $(s_c, a_3)$      $(s_a, a_2)$    $(s_b, a_2)$    $(s_c, a_3)$

+5      +0      -2      -5      +0      -2

$R = +3$        $R = -7$

discount factor: +5 -> reward = +5+0*0.99+(-2)*(0.99)^2 使得越遠越沒有關係

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} (\cancel{R(\tau^n)} - b) \nabla log p_\theta(a_t^n | s_t^n)$$

$$\sum_{t'=t}^{T_n} r_{t'}^n$$

# Tip 2: Assign Suitable Credit

相對於同一個state的其他actions有多好

Advantage Function $A^\theta(s_t, a_t)$

How good it is if we take $a_t$ other than other actions at $s_t$.

Estimated by "*critic*" (later)

Can be state-dependent

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} (R(t^n) - b) \nabla log p_\theta(a_t^n | s_t^n)$$

$$\sum_{t'=t}^{T_n} r_{t'}^n \rightarrow \sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n$$

Add discount factor $\gamma < 1$

# From on-policy to off-policy

当跟環境互動的成本很高的時候，不要更新以後就把原來的data丟掉

Using the experience more than once

# On-policy v.s. Off-policy

- On-policy: The agent learned and the agent interacting with the environment is the same.

- Off-policy: The agent learned and the agent interacting with the environment is different.


阿光下棋


佐為下棋、阿光在旁邊看

# On-policy → Off-policy

reward 的期望值

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla log p_\theta(\tau)]$$

- Use $\pi_\theta$ to collect data. When $\theta$ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train $\theta$. $\theta'$ is fixed, so we can re-use the sample data.

## *Importance Sampling*

$x^i$ is sampled from $p(x)$

$$E_{x \sim p}[f(x)] \approx \frac{1}{N} \sum_{i=1}^{N} f(x^i)$$

We only have $x^i$ sampled from $q(x)$ 看別人學習 (p&q are both distribution)

修正項

從q sample x

$$= \int f(x)p(x)dx = \int f(x)\frac{p(x)}{q(x)}q(x)dx = E_{x \sim q}\left[f(x)\frac{p(x)}{q(x)}\right]$$

Importance weight

# Issue of Importance Sampling

$$E_{x \sim p}[f(x)] = E_{x \sim q}[f(x) \frac{p(x)}{q(x)}]$$

$$Var_{x \sim p}[f(x)] \qquad Var_{x \sim q}[f(x) \frac{p(x)}{q(x)}]$$

$$VAR[X]$$
$$= E[X^2] - (E[X])^2$$

mean同 var不一定相同>>希望相同 >>使得分布可以相像

$$Var_{x \sim p}[f(x)] = E_{x \sim p}[f(x)^2] - \left(E_{x \sim p}[f(x)]\right)^2$$

$$Var_{x \sim q}[f(x) \frac{p(x)}{q(x)}] = E_{x \sim q}\left[\left(f(x) \frac{p(x)}{q(x)}\right)^2\right] - \left(E_{x \sim q}\left[f(x) \frac{p(x)}{q(x)}\right]\right)^2$$

$$= E_{x \sim p}\left[f(x)^2 \frac{p(x)}{q(x)}\right] - \left(E_{x \sim p}[f(x)]\right)^2$$
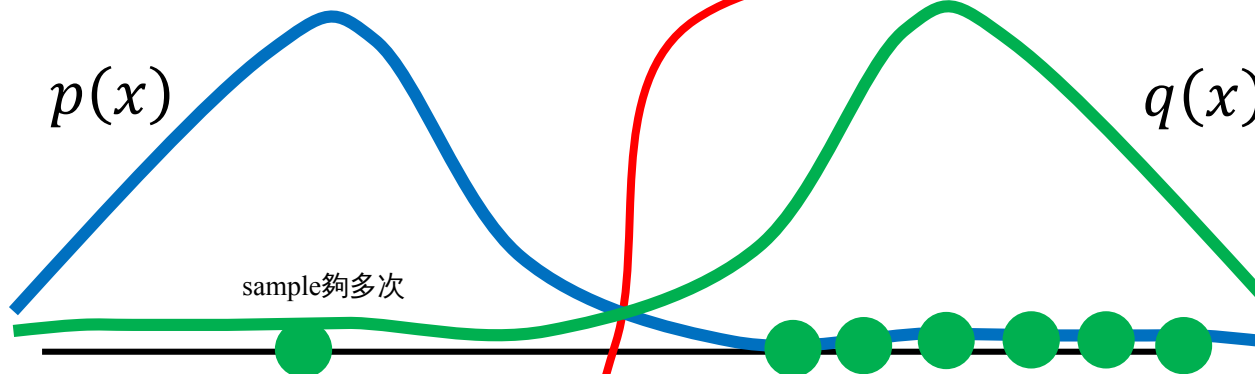
# Issue of Importance Sampling

使得q sample出來的可以接近p

$$E_{x \sim p}[f(x)] = E_{x \sim q}[f(x)\frac{p(x)}{q(x)}]$$

負值

大

$E_{x \sim p}[f(x)]$ is negative

$f(x)$

theta's reward

$p(x)$

$q(x)$

sample夠多次

Very large weight

q sample出來的f(x)負很大 可以抹去右方很多正項
因此可以讓兩者的期望值還是接近的

$E_{x \sim p}[f(x)]$ is ~~positive?~~
negative

# On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)}[R(\tau)\nabla logp_\theta(\tau)]$$

- Use $\pi_\theta$ to collect data. When $\theta$ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train $\theta$. $\theta'$ is fixed, so we can re-use the sample data.

$$\nabla \bar{R}_\theta = E_{\tau \sim p_{\theta'}(\tau)}\left[\frac{p_\theta(\tau)}{p_{\theta'}(\tau)}R(\tau)\nabla logp_\theta(\tau)\right]$$

跟環境互動　update parameters

之前sample的data還是可以拿來train

- Sample the data from $\theta'$.
- Use the data to train $\theta$ many times.

***Importance Sampling***    $E_{x \sim p}[f(x)] = E_{x \sim q}[f(x)\frac{p(x)}{q(x)}]$

# On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x)\nabla log f(x)$$

$$= E_{(s_t,a_t)\sim\pi_\theta}[A^\theta(s_t,a_t)\nabla log p_\theta(a_t^n|s_t^n)]$$

好的or壞的
advantage function

因為是theta prime在跟環境互動

$$A^{\theta'}(s_t,a_t)$$

This term is from sampled data.

$$= E_{(s_t,a_t)\sim\pi_{\theta'}}[\frac{P_\theta(s_t,a_t)}{P_{\theta'}(s_t,a_t)}A^\theta(s_t,a_t)\nabla log p_\theta(a_t^n|s_t^n)]$$

nn

會看到的state是差不多的 省略

$$= E_{(s_t,a_t)\sim\pi_{\theta'}}[\frac{p_\theta(a_t|s_t)}{p_{\theta'}(a_t|s_t)}\frac{p_\theta(s_t)}{p_{\theta'}(s_t)}A^\theta(s_t,a_t)\nabla log p_\theta(a_t^n|s_t^n)]$$

reward(以期望值表示)

$$J^{\theta'}(\theta) = E_{(s_t,a_t)\sim\pi_{\theta'}}\left[\frac{p_\theta(a_t|s_t)}{p_{\theta'}(a_t|s_t)}A^{\theta'}(s_t,a_t)\right]$$

什麼情況下theta prime不能再被用?

## When to stop?

前提假設就是p和q不能差太多

# Add Constraint

穩紮穩打，步步為營

# PPO / TRPO

$\theta$ cannot be very different from $\theta'$

Constraint on behavior not parameters

L1,L2都是在計算參數之間的距離
但RL中action具有隨機性，參數的改變不能直接代表action的變化

KL: 兩分布的差異

## *Proximal Policy Optimization (PPO)*

超參數: 這個KL有多重要?
NN train出來的
有多像?(越像越好)

$$\nabla f(x) = f(x) \nabla log f(x)$$

把限制放在objective function中去更新gradient

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KL(\theta, \theta')$$

objective function
behavior distance, not parameters的距離
同樣的state output出來的動作相近(action 是一個distribution~~~discrete or continuous 都可以視為)

$$J^{\theta'}(\theta) = E_{(s_t,a_t) \sim \pi_{\theta'}} \left[ \frac{p_\theta(a_t|s_t)}{p_{\theta'}(a_t|s_t)} A^{\theta'}(s_t, a_t) \right]$$

ppo直接把kl放到objective function裡面考慮
假如KL很大 J出來就是0 代表不能更新 代表這個theta 組合太爛

## *TRPO (Trust Region Policy Optimization)*

PPO的前身 把KL額外拿出來 兩邊分開要一起考慮比較複雜

$$J_{TRPO}^{\theta'}(\theta) = E_{(s_t,a_t) \sim \pi_{\theta'}} \left[ \frac{p_\theta(a_t|s_t)}{p_{\theta'}(a_t|s_t)} A^{\theta'}(s_t, a_t) \right]$$

額外的限制

$$KL(\theta, \theta') < \delta$$

# PPO algorithm

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)} A^{\theta^k}(s_t, a_t)$$

- Initial policy parameters $\theta^0$

- In each iteration
  - Using $\theta^k$ to interact with the environment to collect $\{s_t, a_t\}$ and compute advantage $A^{\theta^k}(s_t, a_t)$
  - Find $\theta$ optimizing $J_{PPO}(\theta)$

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

Update parameters several times

- If $KL(\theta, \theta^k) > KL_{max}^{自己訂}$, increase $\beta$
- If $KL(\theta, \theta^k) < KL_{min}$, decrease $\beta$
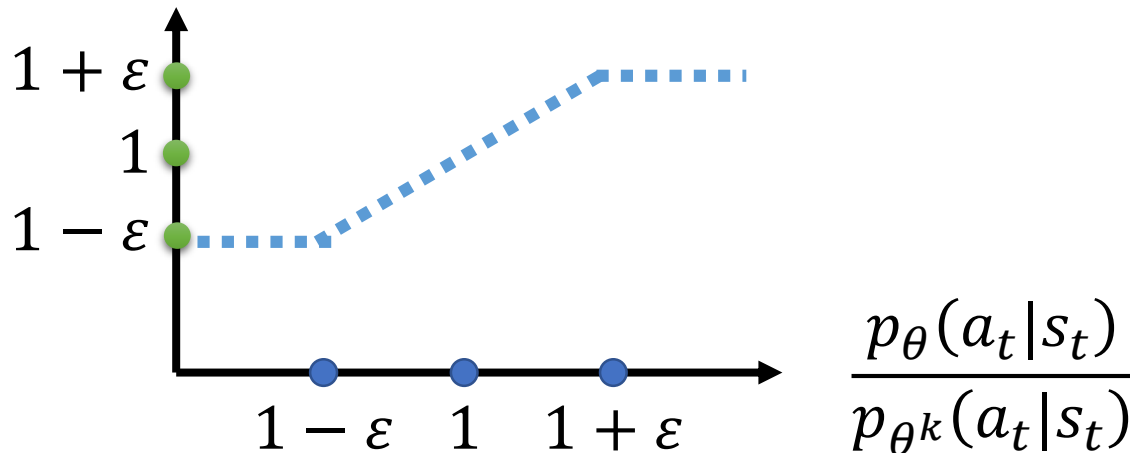
Adaptive KL Penalty

## *PPO algorithm*

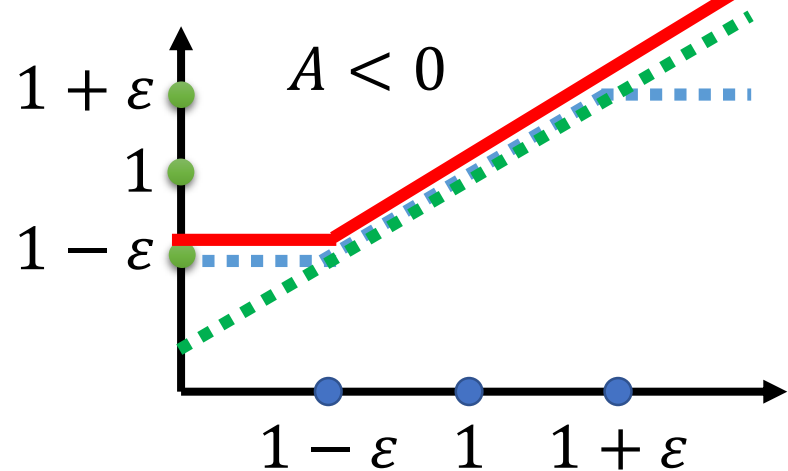$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)} A^{\theta^k}(s_t, a_t)$$

## *PPO2 algorithm*

$$J_{PPO2}^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)}$$

$$clip\left(\frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)}, 1 - \varepsilon, 1 + \varepsilon\right) A^{\theta^k}(s_t, a_t)$$
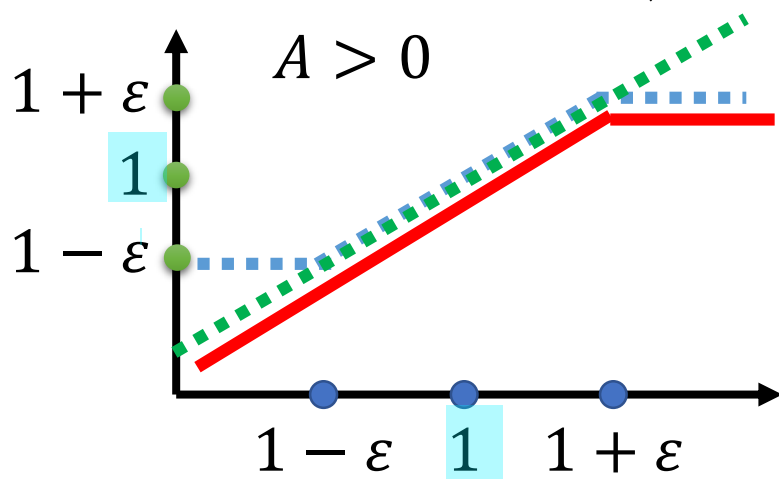
## *PPO algorithm*

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)} A^{\theta^k}(s_t, a_t)$$

## *PPO2 algorithm*

$$J_{PPO2}^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} min\left( \frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)} A^{\theta^k}(s_t, a_t), \right.$$

$$\left. clip\left( \frac{p_\theta(a_t|s_t)}{p_{\theta^k}(a_t|s_t)}, 1-\varepsilon, 1+\varepsilon \right) A^{\theta^k}(s_t, a_t) \right)$$
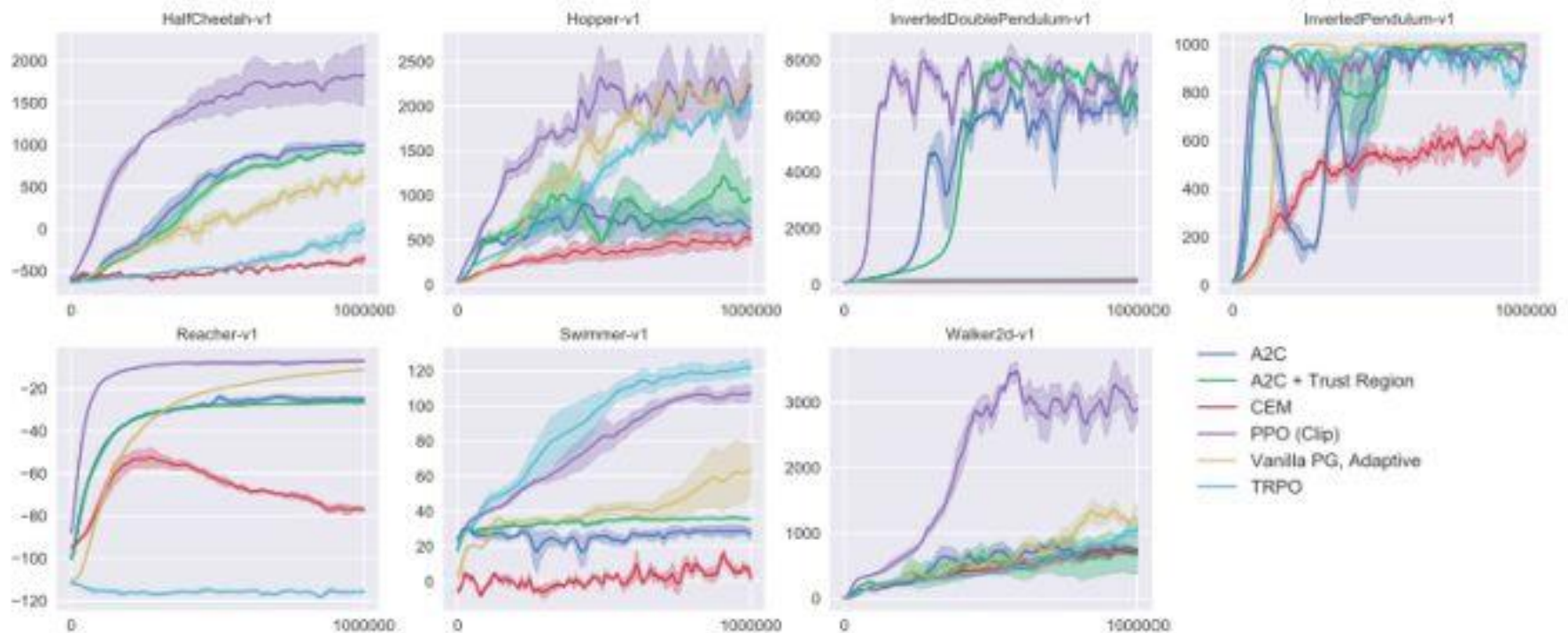
# Experimental Results



Figure 3: Comparison of several algorithms on several MuJoCo environments, training for one million timesteps.