

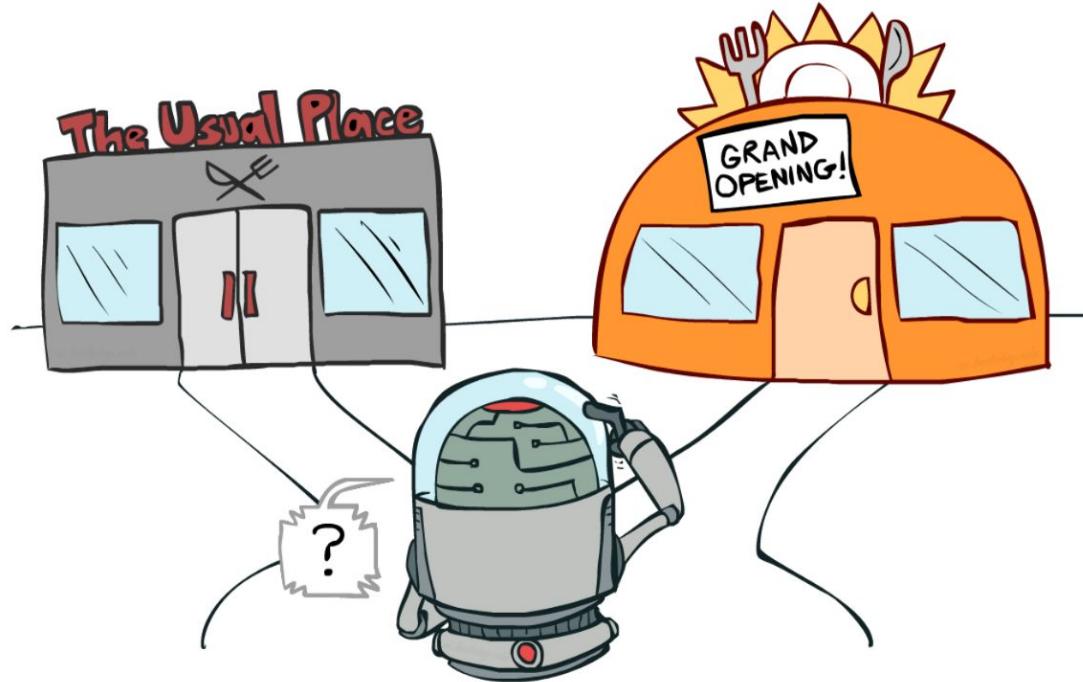
RL overview

2025.07.25
Yifei

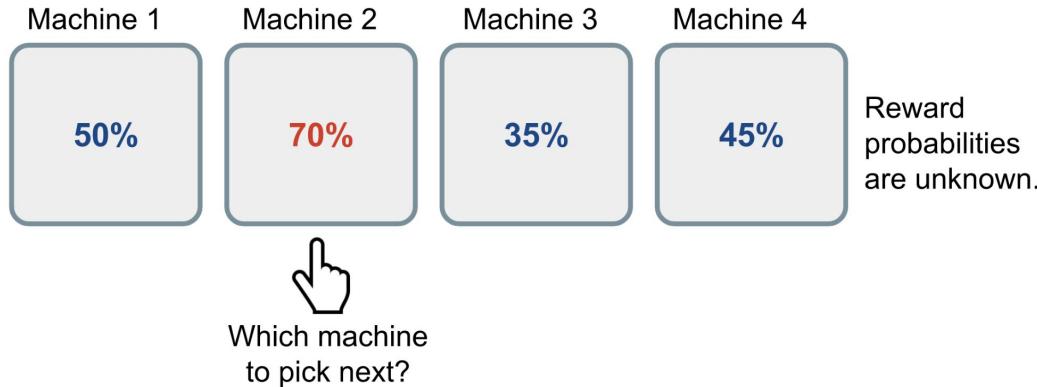
Overview

- Useful concepts in RL
 - Multi-armed bandit
 - Markov Decision Process
 - Value functions
 - Some brief intro to RL algorithm (value-based, policy gradient, actor-critic, etc)
 - Online/offline RL, on-policy/off-policy RL
- RL in LLM
 - PPO
 - RLHF
 - DPO
 - GRPO
 - AReAL (system side optimization for PPO in Reasoning LM)
- Evolutionary algorithm

The known dilemma: Explore VS Exploit



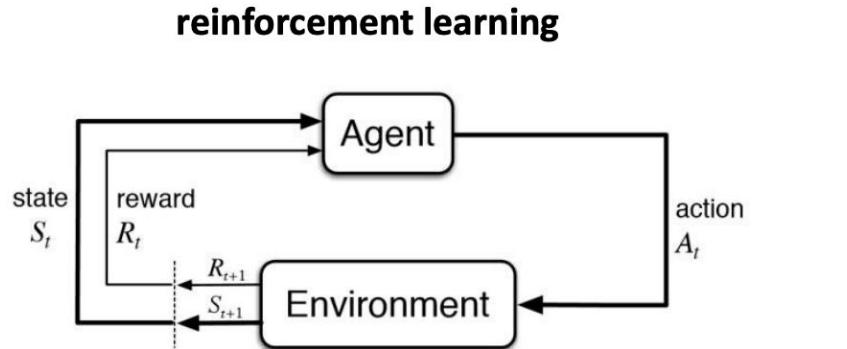
Multi-Armed bandit problem



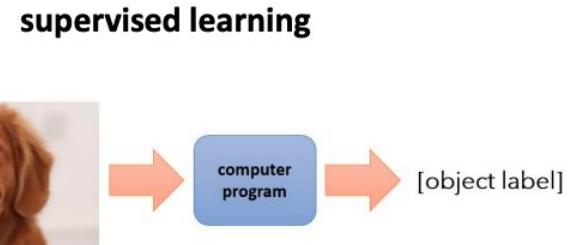
What is the best strategy to maximize the expected total reward over some time period?

Differs from Markov Decision Process because there's no state transition.

What is reinforcement learning?



input: s_t at each time step
output: a_t at each time step
data: $(s_1, a_1, r_1, \dots, s_T, a_T, r_T)$
goal: learn $\pi_\theta : s_t \rightarrow a_t$ reward
to maximize $\sum_t r_t$ policy



input: x
output: y
data: $\mathcal{D} = \{(x_i, y_i)\}$
goal: $f_\theta(x_i) \approx y_i$

someone gives this to you

- Data is not i.i.d. Previous output influences future input
- Only know if we succeed or failed (reward)

- Normally we assume the data point is i.i.d.
- With ground truth given

Definitions - Markov Decision Process

Markov decision process

$$\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$$

\mathcal{S} – state space

states $s \in \mathcal{S}$ (discrete or continuous)

\mathcal{A} – action space

actions $a \in \mathcal{A}$ (discrete or continuous)

\mathcal{T} – transition operator

r – reward function

$$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$$

$r(s_t, a_t)$ – reward

let $\mu_{t,j} = p(s_t = j)$

$$\mu_{t+1,i} = \sum_{j,k} \mathcal{T}_{i,j,k} \mu_{t,j} \xi_{t,k}$$

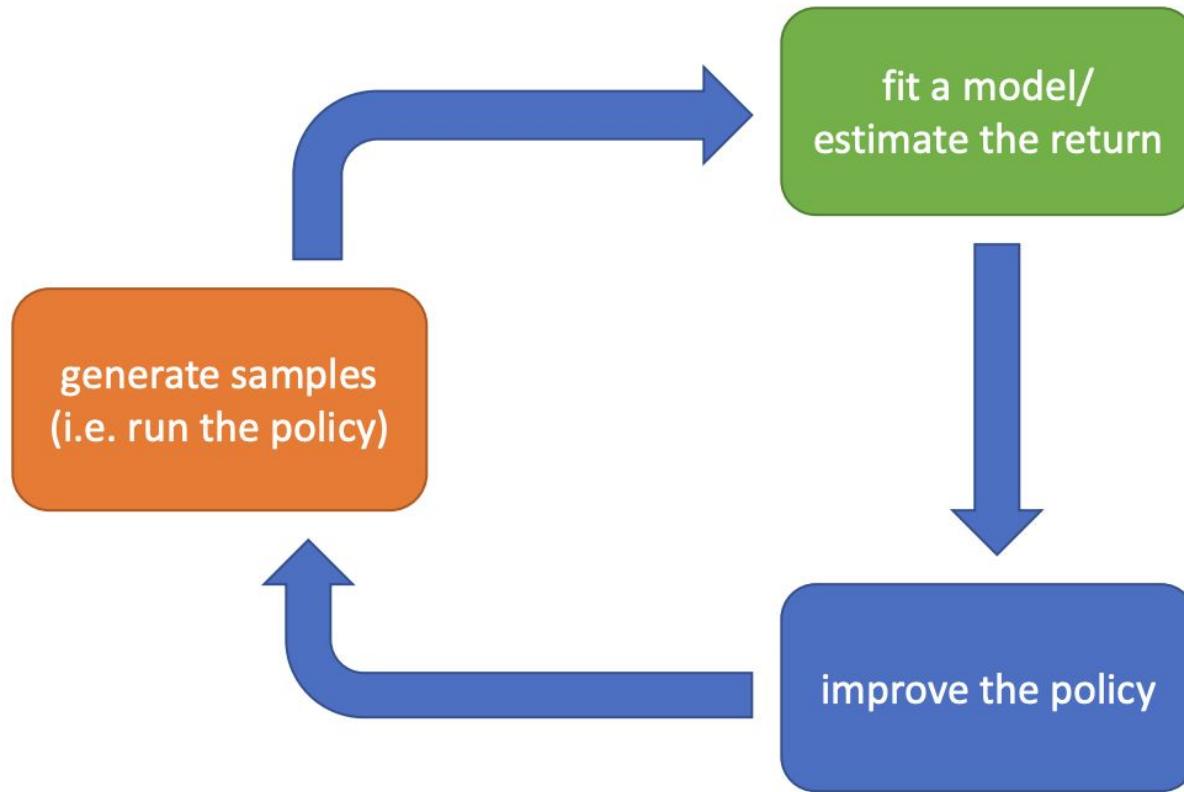
let $\xi_{t,k} = p(a_t = k)$

let $\mathcal{T}_{i,j,k} = p(s_{t+1} = i | s_t = j, a_t = k)$

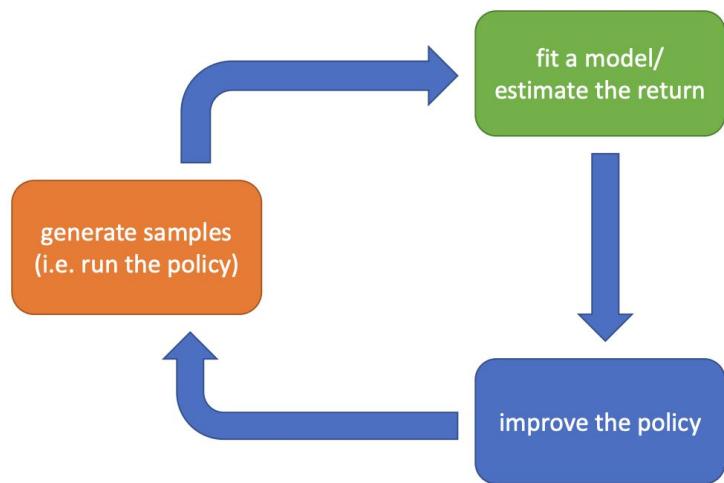
Markov property:

The probability of a certain state and its reward at time t depends on the immediately preceding state and action

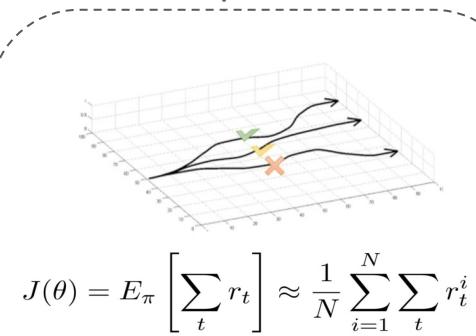
Procedure of RL algorithm



Examples



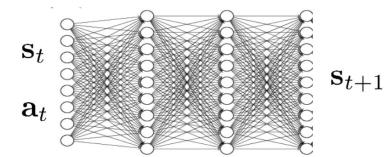
Example 1



$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

Example 2

learn f_{ϕ} such that $\mathbf{s}_{t+1} \approx f_{\phi}(\mathbf{s}_t, \mathbf{a}_t)$

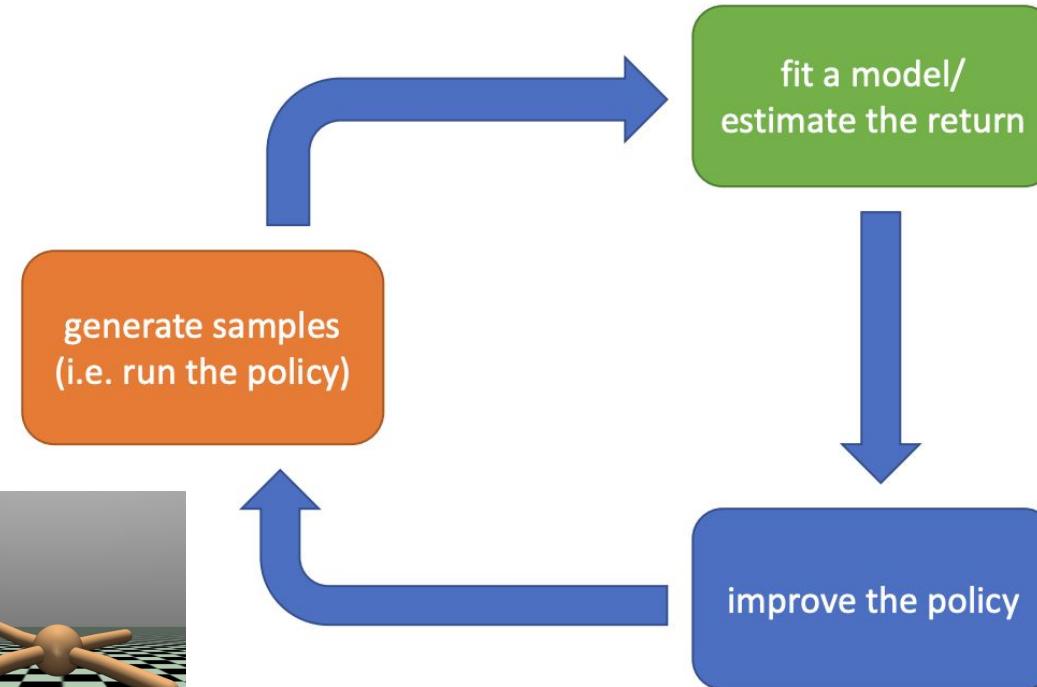
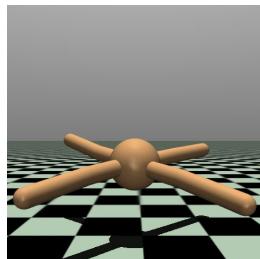


backprop through f_{ϕ} and r to train $\pi_{\theta}(\mathbf{s}_t) = \mathbf{a}_t$

Which part(s) can be expensive?

Expensive: real robot/car/language response, 1x real time

Cheap: simulation enabled, like chess/MuJoCo



Expensive: use a neural net to learn some state transition function

Cheap: naive formulation like computing the average of reward

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

backprop through f_{ϕ} and r to train $\pi_{\theta}(\mathbf{s}_t) = \mathbf{a}_t$

Estimate the return: Value functions

How good it is for an agent to perform a given action at a given state?

State-action function (Q function)

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

Total expected rewards from taking action \mathbf{a}_t in state \mathbf{s}_t

How good it is for an agent to be in a certain state?

Value function

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$V^\pi(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

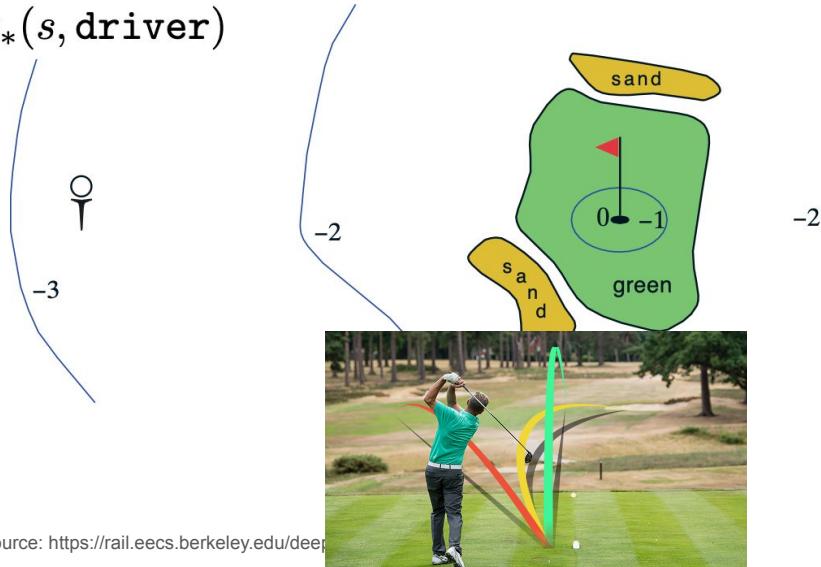
Total expected rewards from state \mathbf{s}_t

Estimate the return: Value functions

How good it is for an agent to perform a given action at a given state?

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$q_*(s, \text{driver})$

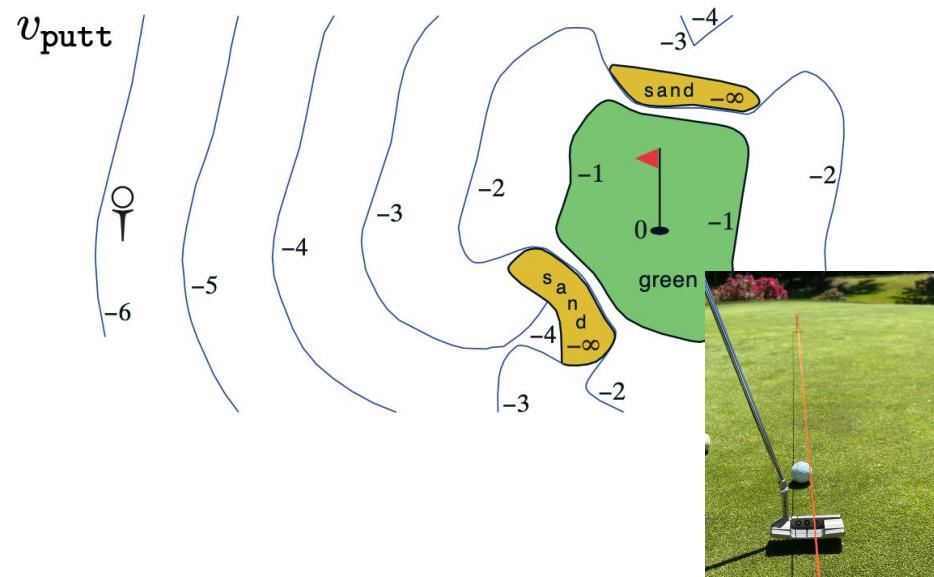


Source: <https://rail.eecs.berkeley.edu/deep-reinforcement-learning/>

How good it is for an agent to be in a certain state?

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

v_{putt}



Types of RL algorithms

Value-based: estimate value function or Q-function of the optimal policy (no explicit policy)

Policy gradients: directly differentiate this objective:

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

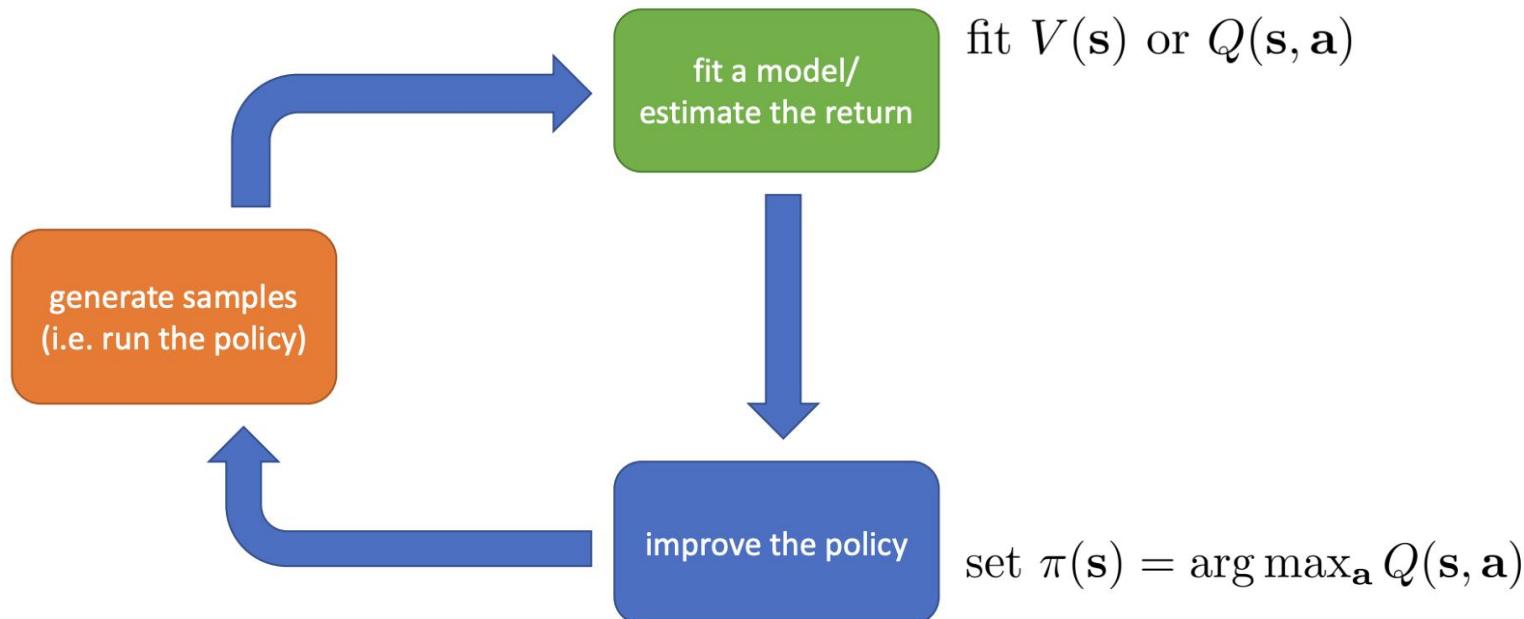
Actor-critic: estimate value function or Q-function of the current policy, use it to improve policy

Model-based RL: estimate the transition model (given current state and action, what is the next state)

- Use it for planning (no policy), e.g. optimal control, MCTS
- Use it to improve a policy
- use model to learn value function, e.g. dynamic programming

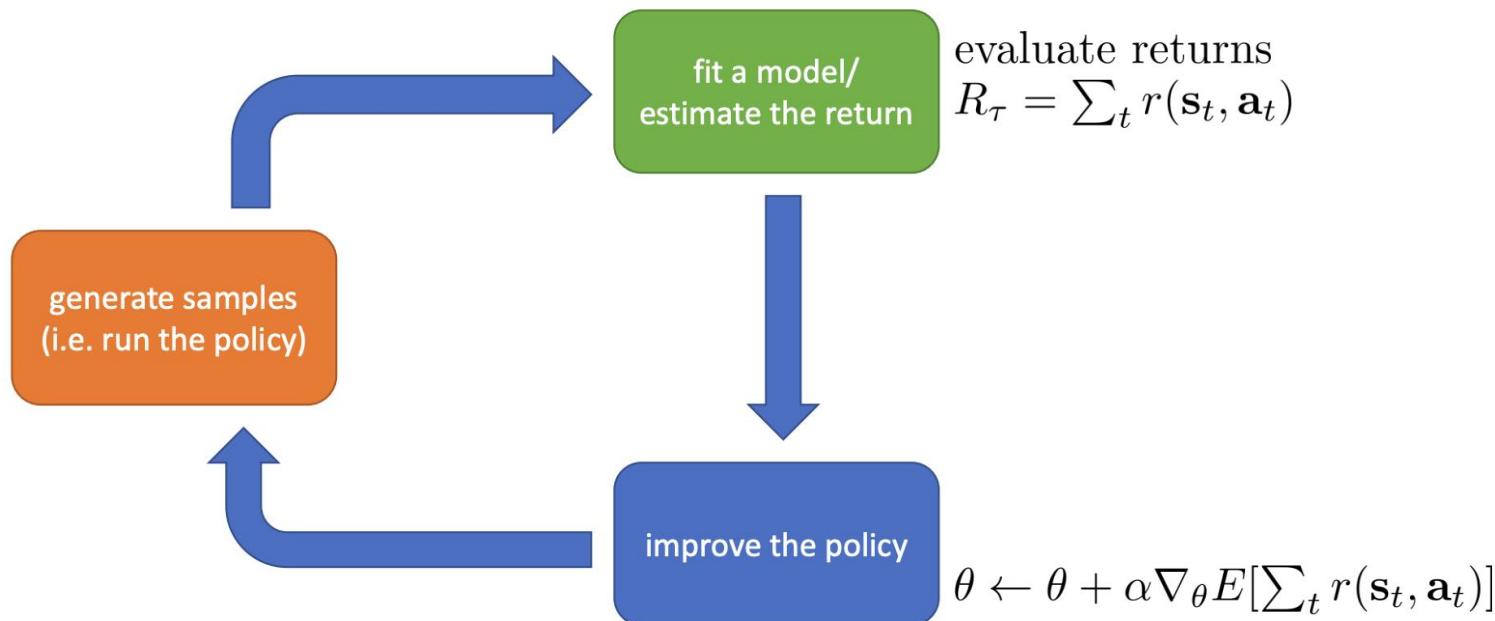
Value-based algorithms

Can we learn the value function or Q function and pick the most promising action?



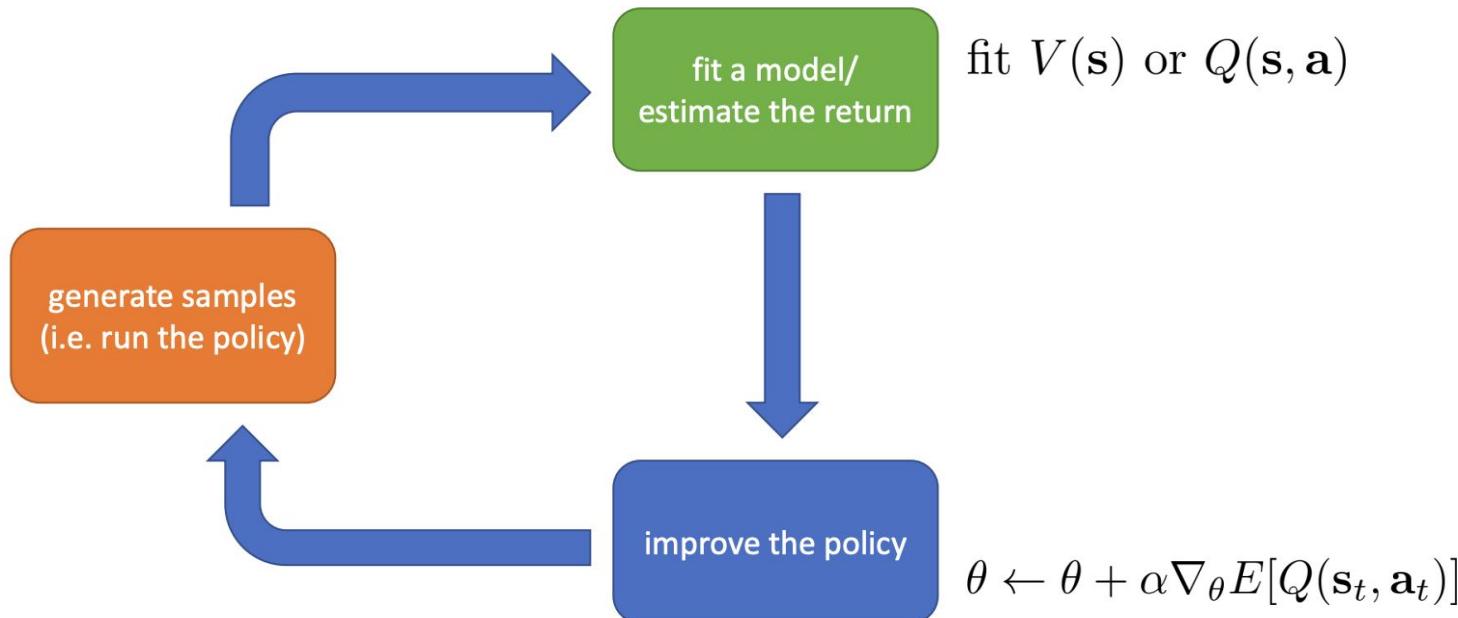
Policy Gradient methods

Can we directly learn the policy on how to act?



Actor-Critic

Can we both learn the policy and an estimation of the reward?

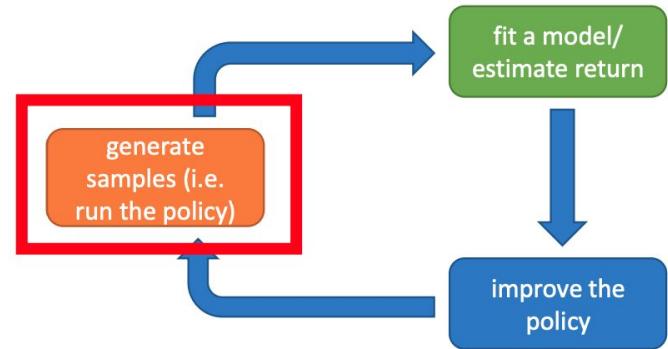


On-policy RL VS off-policy RL

Sampling can be very time consuming especially for on-policy algorithms.

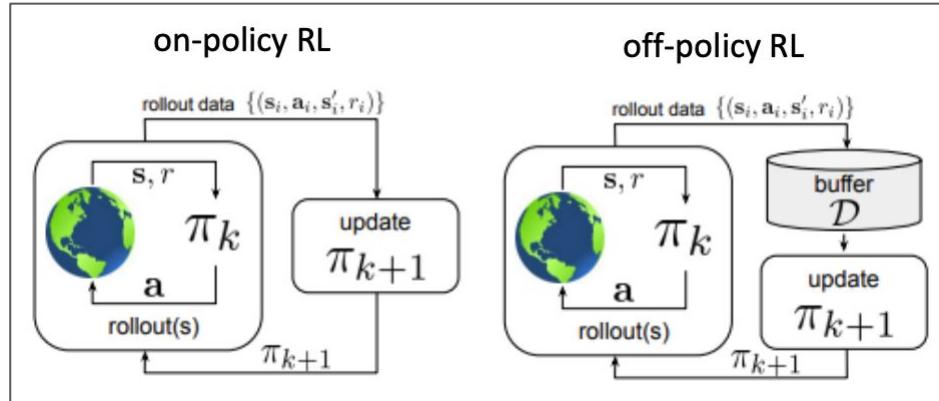
On-policy: each time the policy is changed, even a little bit, we need to generate new samples

Off-policy: able to improve the policy without generating new samples from that policy

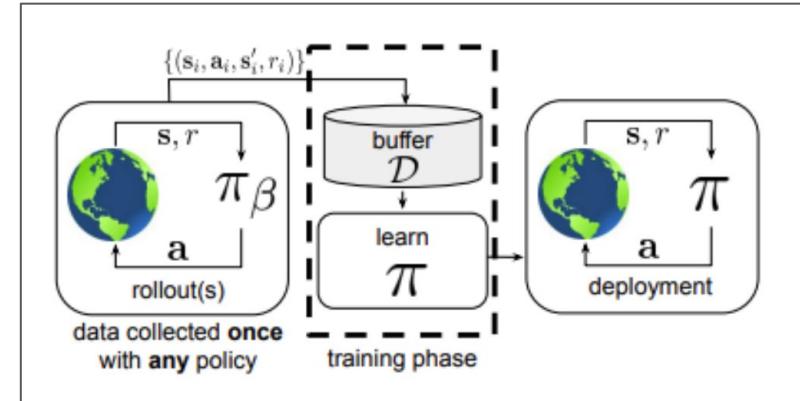


Online RL VS Offline RL

Given a lot of data and no interaction with env, can we learn the best policy from it?

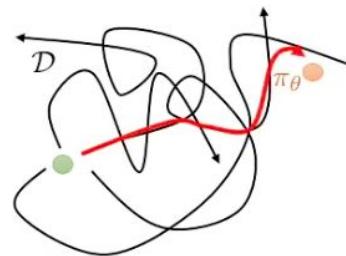


Online RL



Offline RL

It's like ‘micro-scale’ stitching small good trajectories from a lot of suboptimal trajectories



Summary of RL's terminologies

Environment: everything the agent interacts with

Agent: The learner that makes decision based on its policy

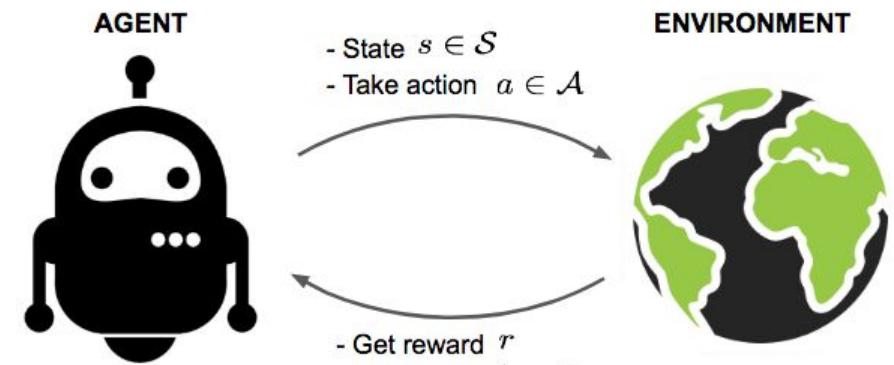
Policy: a strategy to act upon given state

State: encodes current situation

Action: decision made in a given state

Reward: feedback for an action

Value: expected cumulative reward



Goal: By interacting with environment, agent learns a good strategy to maximize future rewards

RL for LLM

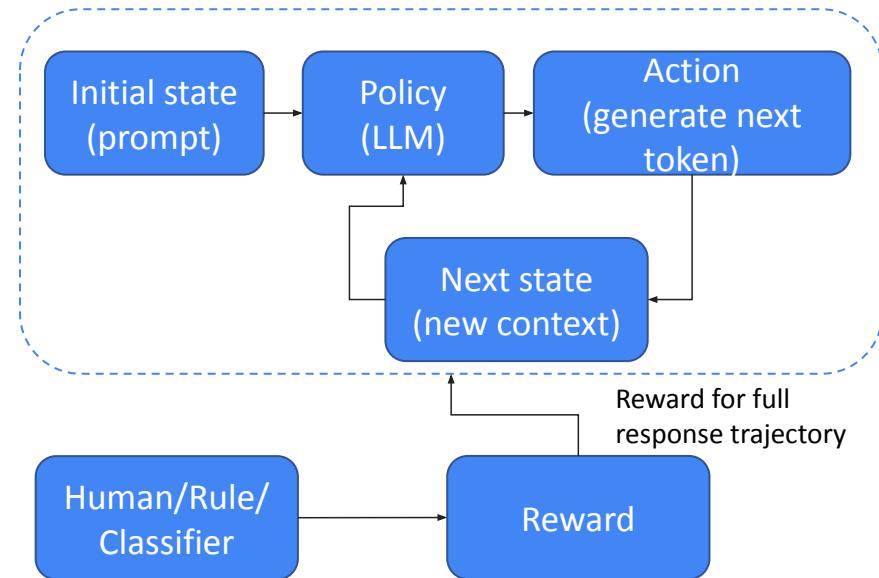
Policy: Language model

State: current context/prompt

Action: next token prediction

Reward: feedback from
Human/Rule/Classifier

Trajectory/Rollout: a sequence of
tokens generated by LLM



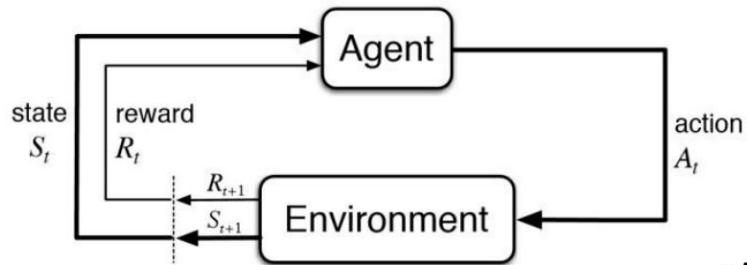
Goal: LLM learns to generate outputs that align with certain reward feedback, can help LLM align with human preference or improve reasoning ability

Proximal Policy Optimization

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.

Recap

reinforcement learning



input: s_t at each time step

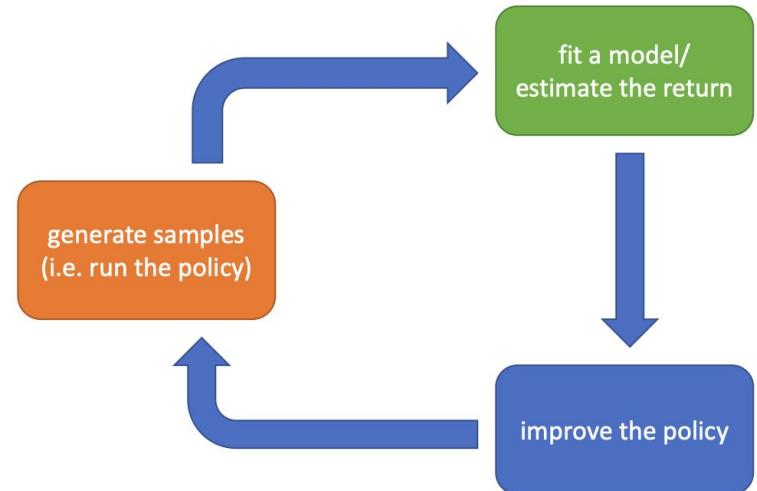
output: a_t at each time step

data: $(s_1, a_1, r_1, \dots, s_T, a_T, r_T)$

goal: learn $\pi_\theta : s_t \rightarrow a_t$

to maximize $\sum_t r_t$

pick your
own actions



Policy Gradient methods

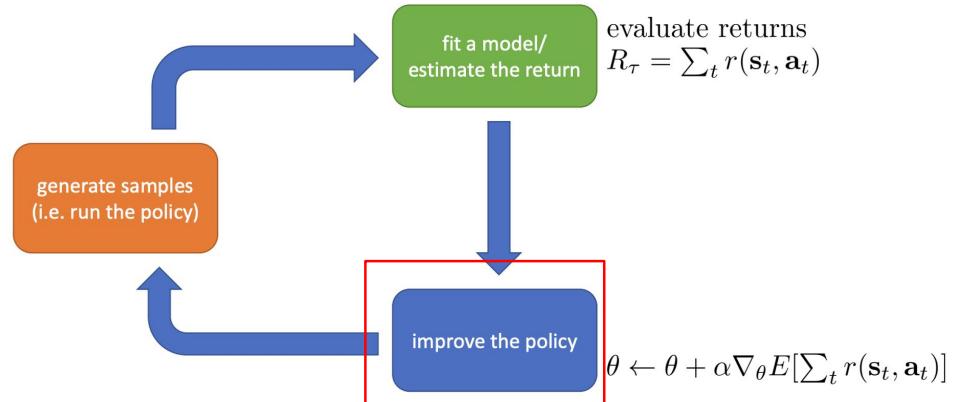
Can we directly learn the policy on how to act?

REINFORCE algorithm:

- 1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run the policy)
- 2. $\nabla_\theta J(\theta) \approx \sum_i \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i) \left(\sum_{t'=t}^T r(\mathbf{s}_{t'}^i, \mathbf{a}_{t'}^i) \right) \right)$
- 3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t}) \hat{Q}_{i,t}^\pi$$

“reward to go”



Improving the policy gradient

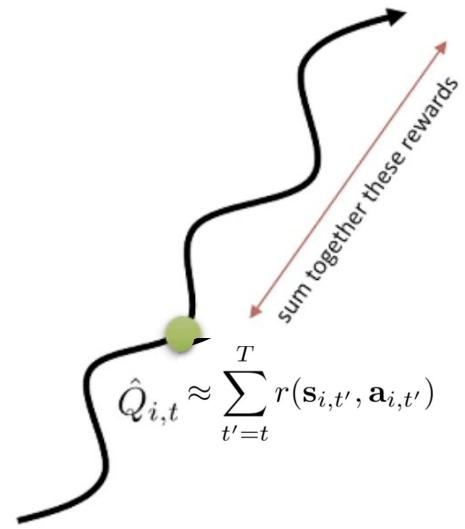
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\underbrace{\sum_{t'=1}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})}_{\text{"reward to go"}} \right)$$

“reward to go”

$$\hat{Q}_{i,t}$$

$\hat{Q}_{i,t}$: estimate of expected reward if we take action $\mathbf{a}_{i,t}$ in state $\mathbf{s}_{i,t}$

can we get a better estimate?



Improving the policy gradient

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\underbrace{\sum_{t'=1}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})}_{\text{"reward to go"}} \right)$$

“reward to go”

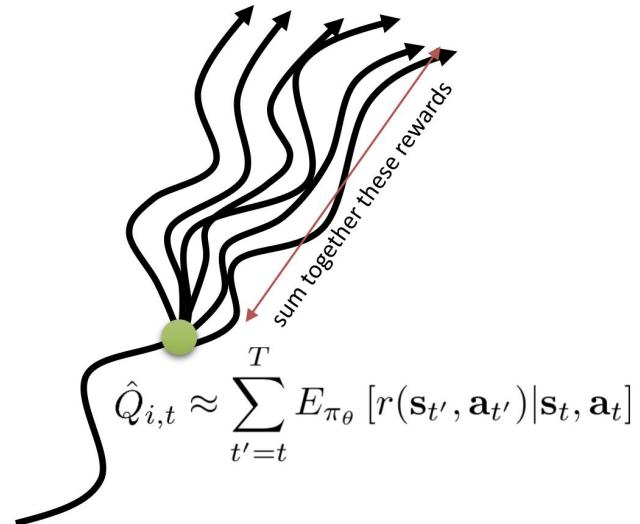
$$\hat{Q}_{i,t}$$

$\hat{Q}_{i,t}$: estimate of expected reward if we take action $\mathbf{a}_{i,t}$ in state $\mathbf{s}_{i,t}$

can we get a better estimate?

$Q(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: true *expected* reward-to-go

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$



Can we have a baseline of expected reward?

$Q(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: true *expected* reward-to-go

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) (Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - b)$$

$$b_t = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

We can take the average Q value as baseline

A better baseline can be a function that depends on state

$$V(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)]$$

Can we have a baseline of expected reward?

$Q(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: true *expected* reward-to-go

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \frac{Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - V(\mathbf{s}_{i,t})}{Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - b}$$

This is called the
Advantage function

$$b_t = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

We can take the average Q value as
baseline

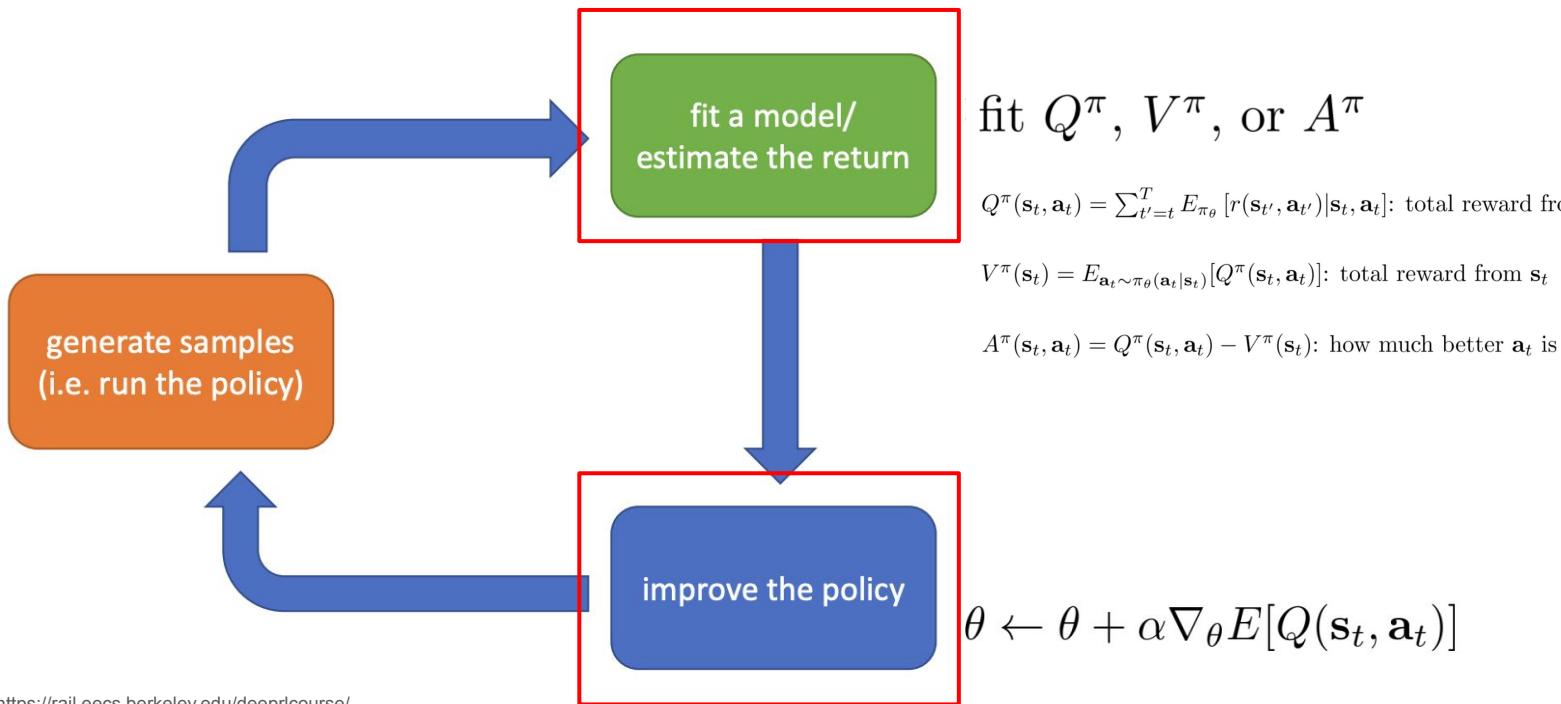
A better baseline can be a function that depends on state

$$V(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)]$$

Intuition:
**We increase the probability
of actions that is better
than average and decrease
the probability of actions
that is worse than average**

Actor-Critic

Can we both learn the policy and an estimation of the reward?



State and state-action functions

$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: total reward from taking \mathbf{a}_t in \mathbf{s}_t

$V^\pi(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$: total reward from \mathbf{s}_t

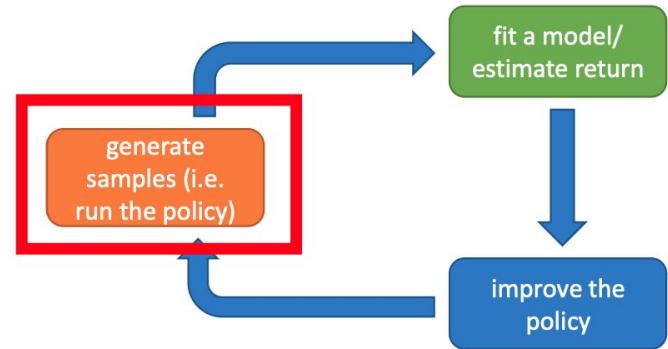
$A^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t)$: how much better \mathbf{a}_t is

On-policy RL VS off-policy RL

Sampling can be very time consuming especially for on-policy algorithms.

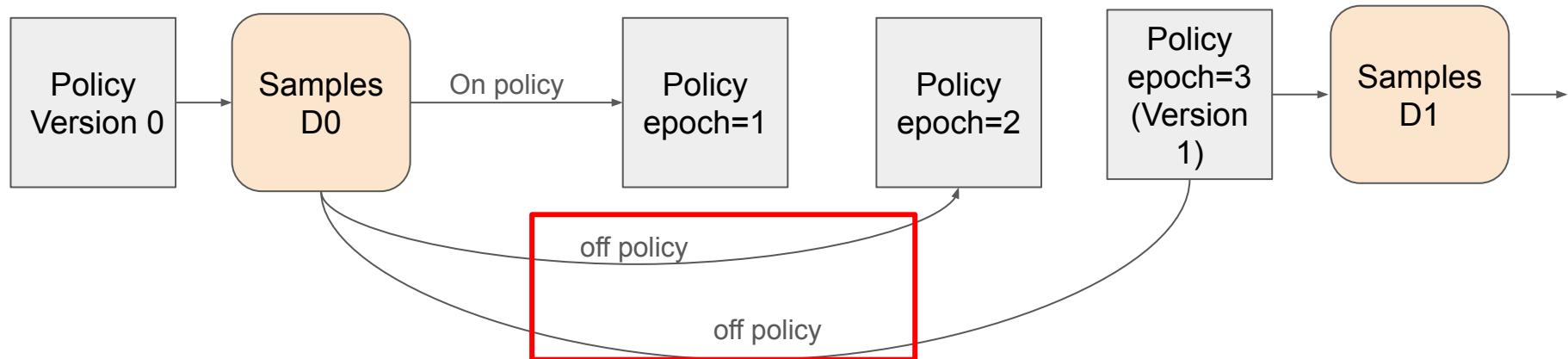
On-policy: each time the policy is changed, even a little bit, we need to generate new samples

Off-policy: able to improve the policy without strictly using the samples generated from that policy



An illustration

Assume we set epoch_num=3 and batch_size=entire_train_set



Could lead to large policy updates -> need constraints!

Problem? Data inefficiency!

- On-policy method: for each new policy, we need to generate completely new trajectories
- The trajectory is thrown out after just one gradient update
- Performing multiple steps of optimization on the loss with the same trajectory is appealing but can lead to destructively large policy updates

Can we add a constraint on the updated policy so that it does not diverge too much from the policy that generated the trajectory?

Policy Gradient methods

Objective for Policy
Gradient methods

Estimator of Advantage: how
much better (or worse) an action a_t
at is compared to the policy's
average expected value at state : s_t

$$L^{PG}(\theta) = \hat{\mathbb{E}}_t \left[\log \pi_\theta(a_t | s_t) \hat{A}_t \right]$$

Empirical average over a finite
batch of samples

Commonly used
gradient estimator

$$\hat{g} = \hat{\mathbb{E}}_t \left[\nabla_\theta \log \pi_\theta(a_t | s_t) \hat{A}_t \right]$$

Trust Region Policy Optimization (TRPO)

Hard constraint version:

$$\begin{aligned} & \underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] \\ & \text{subject to} \quad \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta \end{aligned}$$

Importance ratio from the
importance sampling

Penalty version:

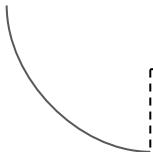
$$\underset{\theta}{\text{maximize}} \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)] \right]$$

TRPO uses a **hard constraint** rather than a penalty because it is hard to choose a single value of β that performs well across different problems—or even within a single problem, where the characteristics change over the course of learning

Proximal Policy Optimization

Motivation:

- The hard constraint helps in the training process from TRPO.
- Maybe the constraint is not a strict constraint: Does it matter if we only break the constraint just a few times?
- What if we treat it as a “soft” constraint?
 - Add proximal value to objective function

- 
- 1. Adaptive KL penalty
 - 2. Clipped Surrogate

PPO with adaptive KL Penalty

$$L^{KL PEN}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)] \right]$$

Hard to pick β value \rightarrow use adaptive β

Compute $d = \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)]]$

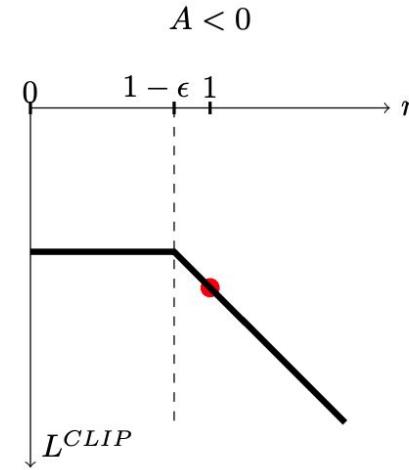
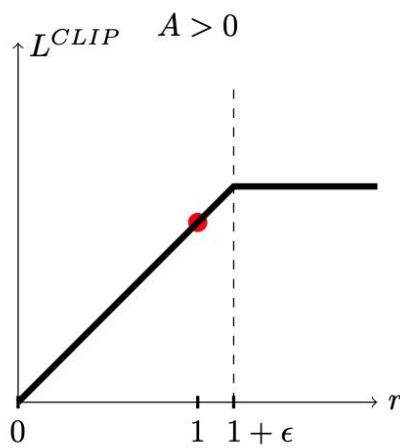
- If $d < d_{\text{targ}}/1.5$, $\beta \leftarrow \beta/2$
- If $d > d_{\text{targ}} \times 1.5$, $\beta \leftarrow \beta \times 2$

- 1.5 and 2 are heuristically chosen, model are not sensitive to them
- This version perform worse than the clipped surrogate objective

PPO with Clipped Objective

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] \quad r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \quad r(\theta_{\text{old}}) = 1$$

Fluctuation happens when r changes too much -> limit r within a range?



$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

PPO with Clipped Objective

Algorithm 5 PPO with Clipped Objective

Input: initial policy parameters θ_0 , clipping threshold ϵ
for $k = 0, 1, 2, \dots$ **do**

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \cdots + \cdots + (\gamma\lambda)^{T-t+1}\delta_{T-1},$$

where $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$

Collect set of partial trajectories \mathcal{D}_k on policy $\pi_k = \pi(\theta_k)$

Estimate advantages $\hat{A}_t^{\pi_k}$ using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

by taking K steps of minibatch SGD (via Adam), where

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[\sum_{t=0}^T \left[\min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

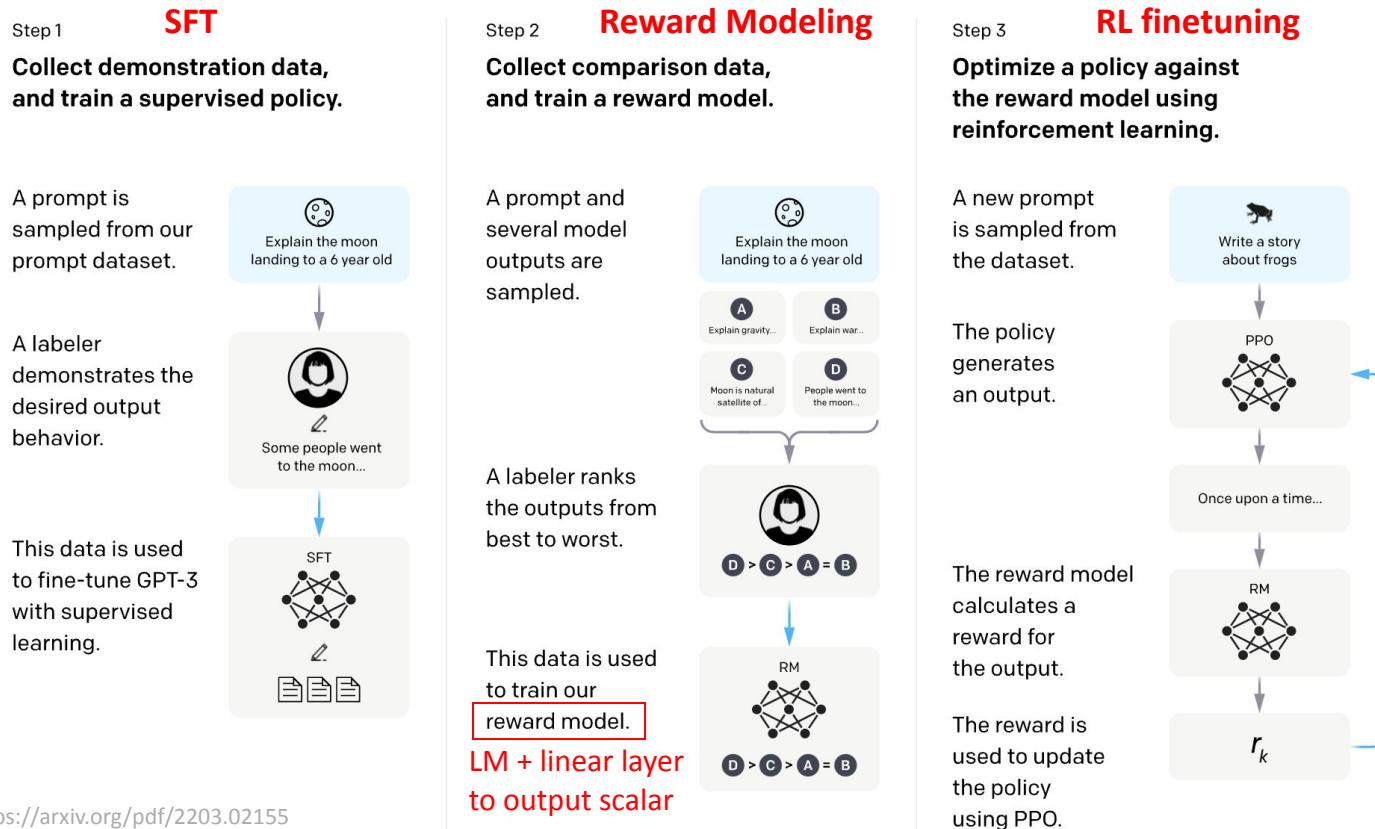
end for

algorithm	avg. normalized score
No clipping or penalty	-0.39
Clipping, $\epsilon = 0.1$	0.76
Clipping, $\epsilon = 0.2$	0.82
Clipping, $\epsilon = 0.3$	0.70
Adaptive KL $d_{\text{targ}} = 0.003$	0.68
Adaptive KL $d_{\text{targ}} = 0.01$	0.74
Adaptive KL $d_{\text{targ}} = 0.03$	0.71
Fixed KL, $\beta = 0.3$	0.62
Fixed KL, $\beta = 1.$	0.71
Fixed KL, $\beta = 3.$	0.72
Fixed KL, $\beta = 10.$	0.69

Table 1: Results from continuous control benchmark. Average normalized scores (over 21 runs of the algorithm, on 7 environments) for each algorithm / hyperparameter setting . β was initialized at 1.

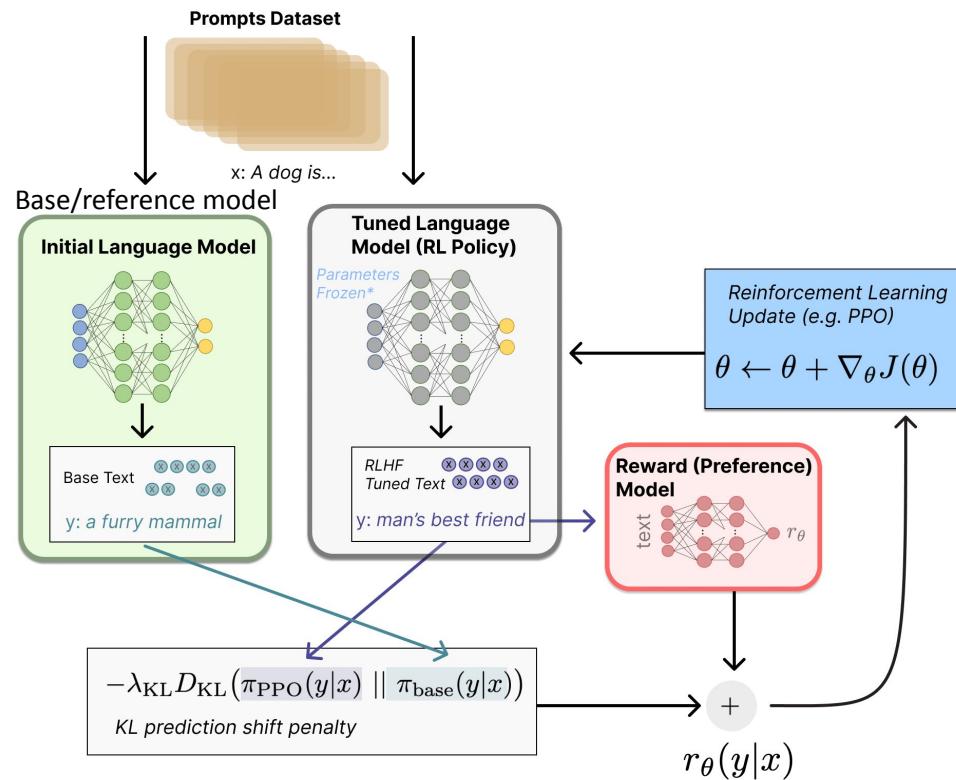
Reinforcement Learning with Human Feedback

RLHF workflow



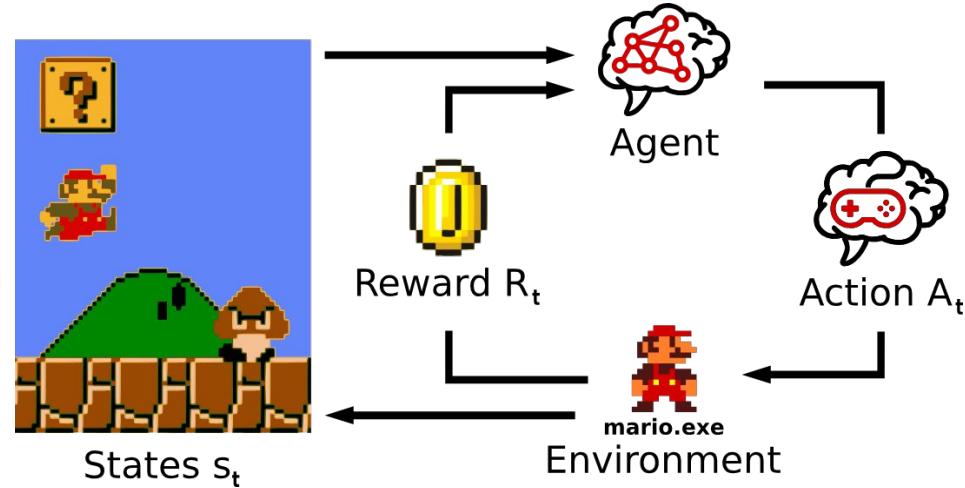
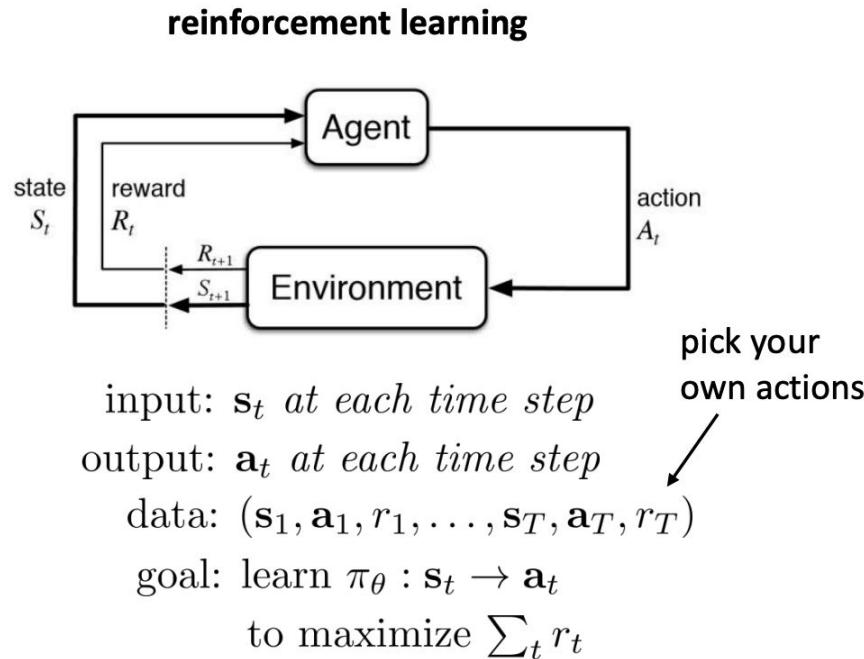
PPO in RLHF

- *Base model*: frozen, used to compute KL divergence
- *Policy model (actor model)*: the LM we want to finetune
- *Reward model*: generate an outcome reward for rollout



Direct Preference Optimization

Recap: Reinforcement Learning



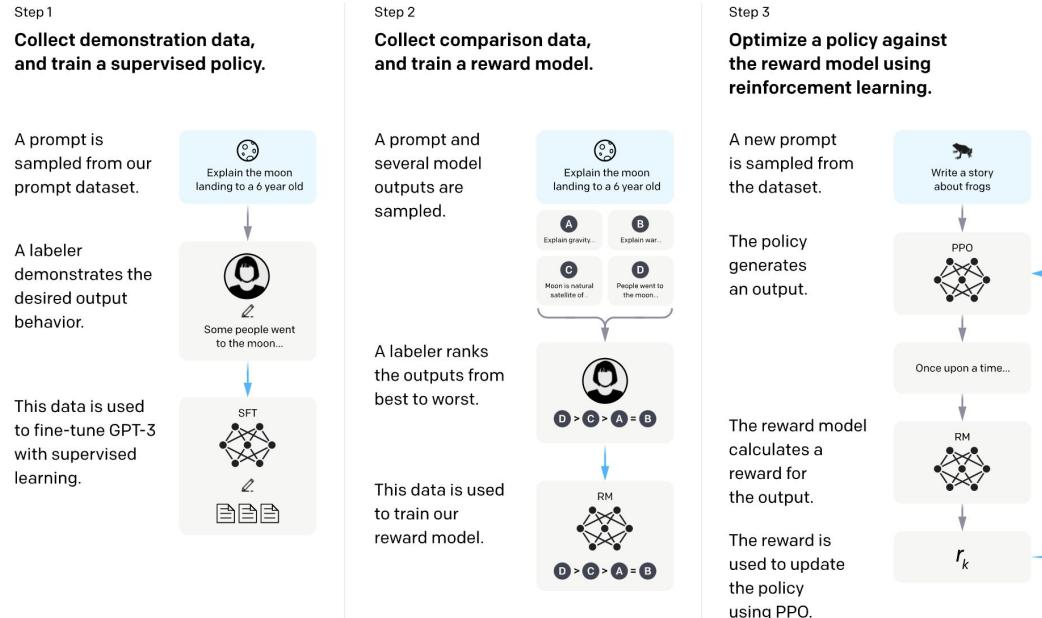
Recap: PPO and RLHF

PPO Objective functions:

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)] \right]$$

RLHF workflow:



Motivation

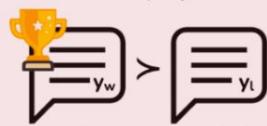
Building and training a reward model can be complicated

How to simplify the process of
RLHF?

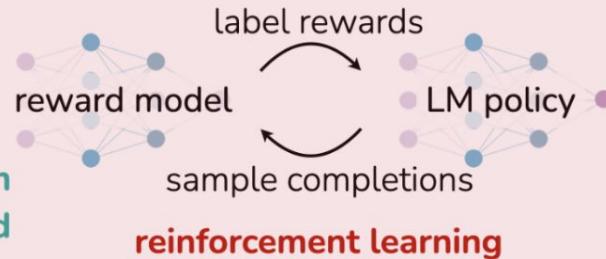
Method - Main idea

Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about
the history of jazz"



maximum likelihood



Direct Preference Optimization (DPO)

x: "write me a poem about
the history of jazz"

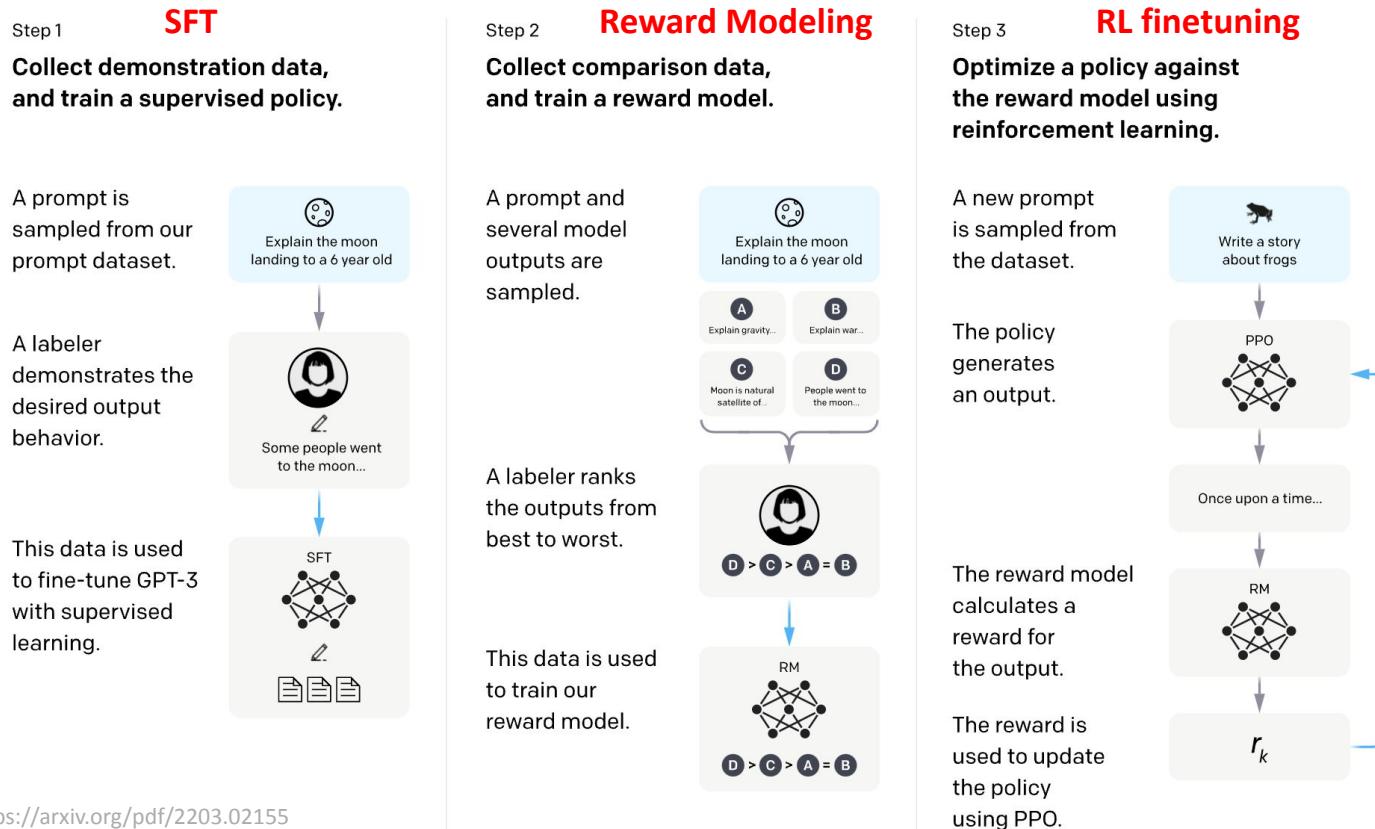


maximum likelihood



Directly learning from preference data
No need for explicitly fitting a reward model!

RLHF workflow



Reward modeling phase in RLHF

Feedback comes as **preferences over model samples**:

$$\mathcal{D} = \{x^i, y_w^i, y_l^i\}$$

↑
Prompt
↑
Preferred response
↑
Dispreferred response

Bradley-Terry Model connects rewards to preferences:

$$p(y_w \succ y_l \mid x) = \sigma(r(x, \overset{\swarrow}{y_w}) - r(\overset{\searrow}{x}, y_l))$$

Reward assigned to preferred and dispreferred responses

Train the reward model by **minimizing negative log likelihood**:

$$\mathcal{L}_R(\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$$

Method - Recap on RLHF

RL fine-tuning phase: Learn a policy that **maximize reward** while **stay close** to the reference policy

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(y | x) || \pi_{\text{ref}}(y | x)]$$



PPO objective used in RLHF paper

Method - Recap on RLHF

RL fine-tuning phase: Learn a policy that **maximize reward** while **stay close** to the reference policy

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(y | x) || \pi_{\text{ref}}(y | x)]$$

Reward from
reward model Policy learnt
from PPO Policy right
after SFT

Method - Recap on RLHF

RL fine-tuning phase: Learn a policy that **maximize reward** while **stay close** to the reference policy

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(y | x) || \pi_{\text{ref}}(y | x)]$$

Maximizes the rewards

Prevents the learnt policy from changing too much from the SFT model

Direct Preference Optimization

RLHF Objective

(get **high reward**, stay close to reference model)

Closed-form Optimal Policy

(write **optimal policy** as function of **reward function**; from prior work)

$$\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [r(x, y)] - \beta \mathbb{D}_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))$$

any reward function

$$\pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$

$$\sum_y \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$

intractable!

Direct Preference Optimization

RLHF Objective

(get **high reward**, stay close to reference model)

Closed-form Optimal Policy

(write **optimal policy** as function of **reward function**; from prior work)

Rearrange

(write **any reward function** as function of **optimal policy**)

$$\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [r(x, y)] - \beta \mathbb{D}_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))$$

any reward function

$$\pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} r(x, y) \right)$$



$$r(x, y) = \underbrace{\beta \log \frac{\pi^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x)}_{\text{some parameterization of a reward function}}$$

Direct Preference Optimization

A loss function on
reward functions



A transformation
between reward
functions and policies



A loss function
on policies

Derived from the Bradley-Terry model of human preferences:

$$\mathcal{L}_R(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$$

$$r_{\pi_\theta}(x, y) = \beta \log \frac{\pi_\theta(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x)$$

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

Takeaways

- DPO removes the RL training loop from RLHF
- DPO is simpler and computationally cheaper than PPO
- DPO is optimizing the same objective as RLHF

Is RLHF really RL?

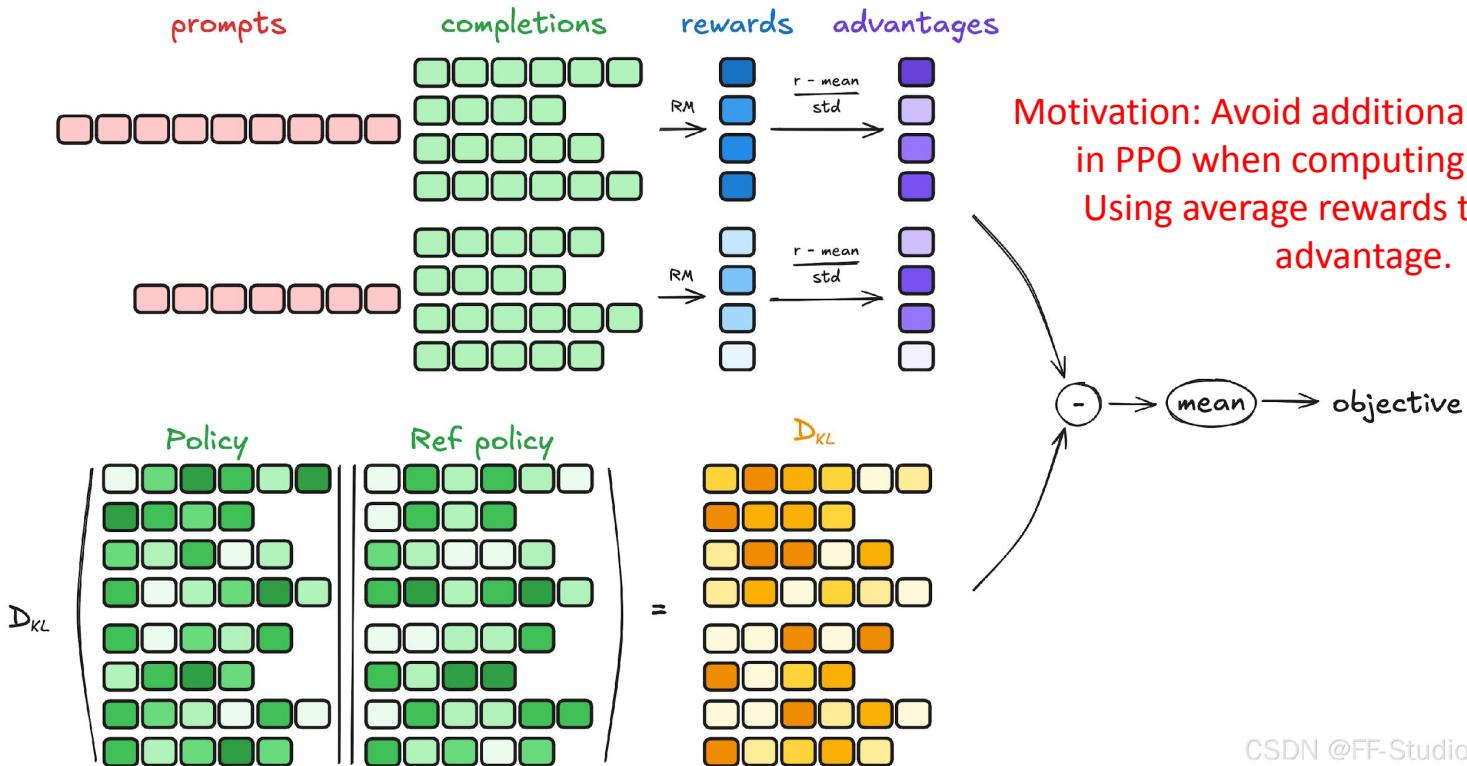
Contextual Bandit?

- Reward is only given after the sequence is completely generated
- LM itself can be regarded as a direct mapping from inputs to distributions over completions, rather than a sequential decision-making agent in the space of tokens

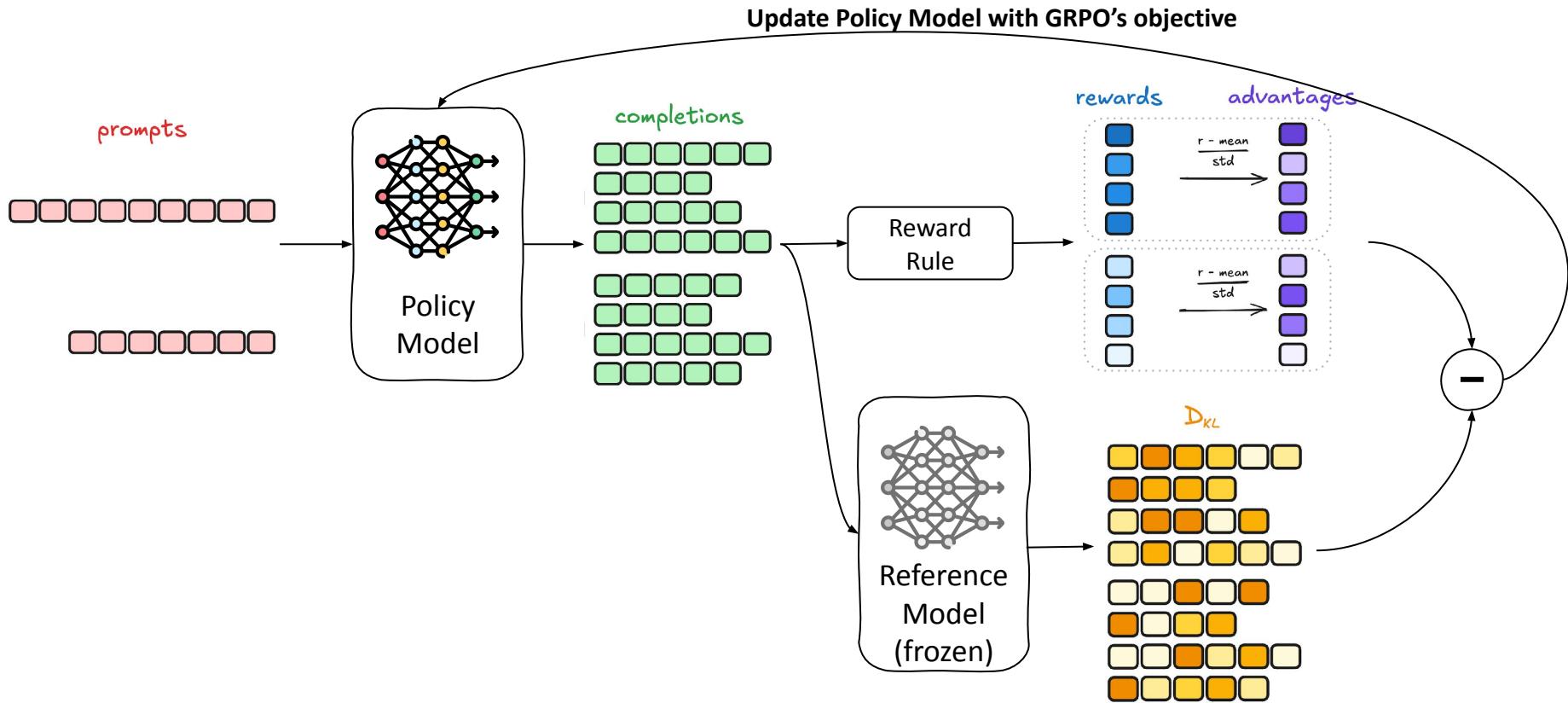
Offline RL / Supervised learning?

- The environment consists the policy(LLM) itself
- No new info from the environment for LLM except for the preference data

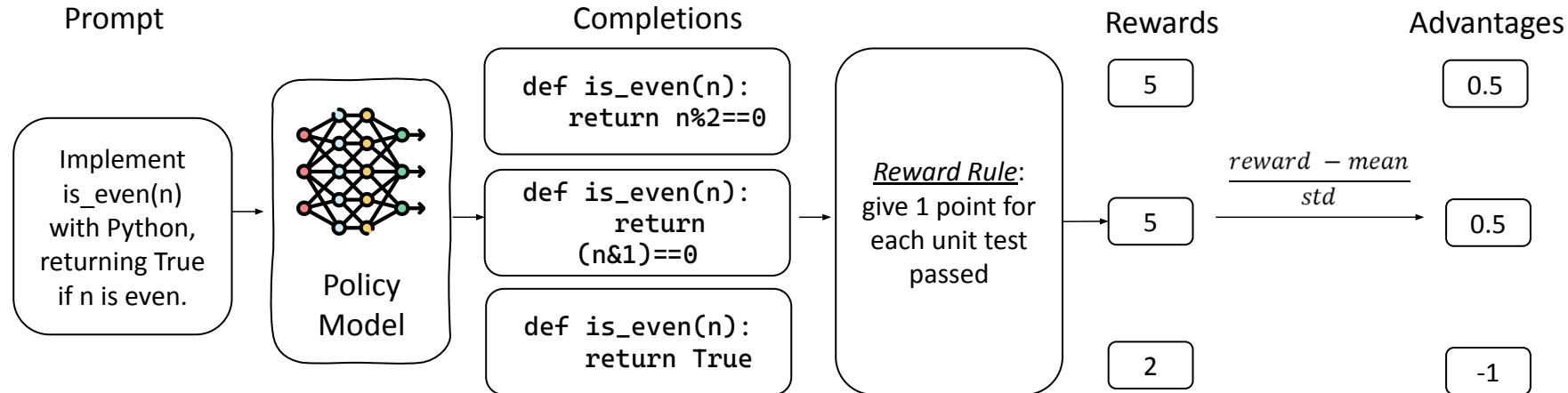
Group Relative Policy Optimization (GRPO)



GRPO Workflow



GRPO – An example



Dive a little bit into the RL infra side...

AREAL: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning

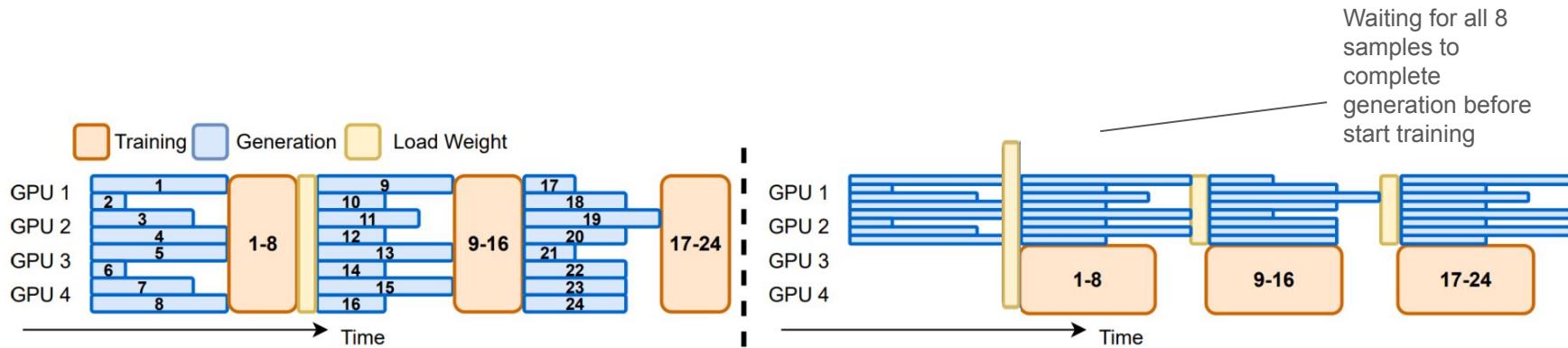
**Wei Fu¹², Jiaxuan Gao¹², Xujie Shen², Chen Zhu², Zhiyu Mei¹²,
Chuiyi He², Shusheng Xu¹², Guo Wei², Jun Mei², Jiashu Wang²³,
Tongkai Yang², Binhang Yuan³, Yi Wu¹²**

¹ IIIS, Tsinghua University, ² Ant Research, ³ HKUST
fuwth17@gmail.com, jxwuyi@gmail.com

Motivation to improve RL training system for reasoning LM

RL can be used to improve reasoning abilities in LLMs

- On-policy algorithms require rollouts to be generated by the latest policy
- Naive synchronous pipeline
 - alternates the generation phase and training phase strictly
 - all nodes wait for the slowest task in the batch to finish before training can start
 - chain-of-thought lengths vary greatly, creating a “longest task bottleneck”
 - leads to low GPU utilization as many GPUs sit idle waiting for others to finish generation.



AReAL's async generation

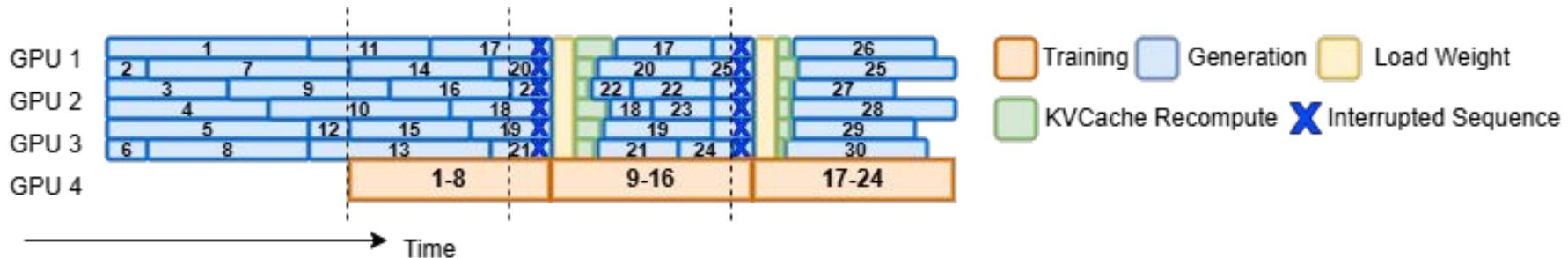


Figure 3: Illustration of generation management in AREAL. Vertical lines show the ready time for the next step training. Blue crosses show the interrupted requests when new parameters arrive.

Data staleness → staleness control during generation task submission

Inconsistent policy → decoupled PPO objective proposed in [1]

Two challenges to make it asynchronous

Data Staleness:

The policy that generated the rollout might be too old for the current training policy (off-policy issue)

Solution: don't start new rollout task until the current tasks+predefined staleness tolerance have been consumed(used in training)

Inconsistent policy:

Policies used for a single rollout may contains multiple different policies

Solution: decoupled PPO objective