

# Alignment of Language Models – Reward Maximization

Advances in Large Language Models

ELL8299 · AIL861



Gaurav Pandey  
Senior Research Scientist, IBM Research

# Reward Maximization

- Best-of-N policy
- Reward Maximization Objective Formulation
- Vanilla Policy Gradient
  - The REINFORCE gradient estimator
  - The problem of high variance - Baseline subtraction
  - Final algorithm
- Towards Actor/Critic methods
  - The problem of credit assignment
  - Q-functions, Value functions and advantage
  - Estimating the value function and advantage
  - Actor-Critic Family of algorithms



# Reward Maximization

- Best-of-N policy
- Reward Maximization Objective Formulation
- Vanilla Policy Gradient
  - The REINFORCE gradient estimator
  - The problem of high variance - Baseline subtraction
  - Final algorithm
- Towards Actor/Critic methods
  - The problem of credit assignment
  - Q-functions, Value functions and advantage
  - Estimating the value function and advantage
  - Final algorithm



# Best-of-N policy

Given

- Base policy or reference policy  $\pi_{ref}(y|x)$ 
  - Often, an instruction tuned LM that serves as the starting point of alignment
- Reward Model  $r(x, y)$

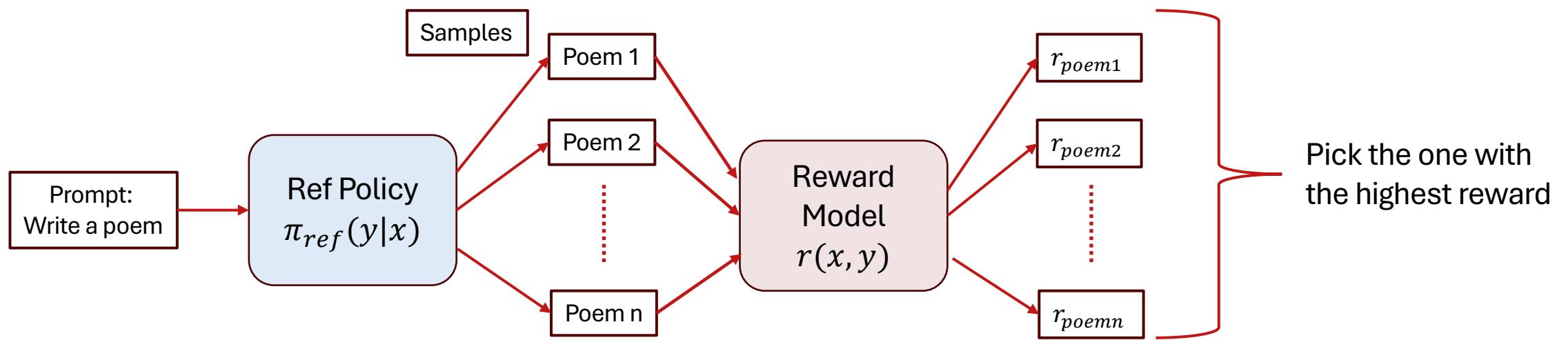
Aim – To generate outputs with high reward

**Solution**

- Sample multiple outputs from the policy  $\pi_{ref}$
- Score each output using the reward model
- Return the output with the highest reward



# Best-of-N policy



## Challenge:

- Too expensive during inference
- Can't get better iteratively
- Reward Model may not be available during inference
  - Verifiable rewards



# Reward Maximization

- Best-of-N policy
- Reward Maximization Objective Formulation
- Vanilla Policy Gradient
  - The REINFORCE gradient estimator
  - The problem of high variance - Baseline subtraction
  - Final algorithm
- Towards Actor/Critic methods
  - The problem of credit assignment
  - Q-functions, Value functions and advantage
  - Estimating the value function and advantage
  - Final algorithm



# The reward-maximization objective

Given

- Base policy or reference policy  $\pi_{ref}(y|x)$ 
  - Often, an instruction tuned LM that serves as the starting point of alignment
- Reward Model  $r(x, y)$

Aim

- To find a policy  $\pi_{\theta^*}(y|x)$ 
  - That generated outputs with high reward.
  - That stay close to the reference policy.



# Why care about closeness to $\pi_{ref}$ ?

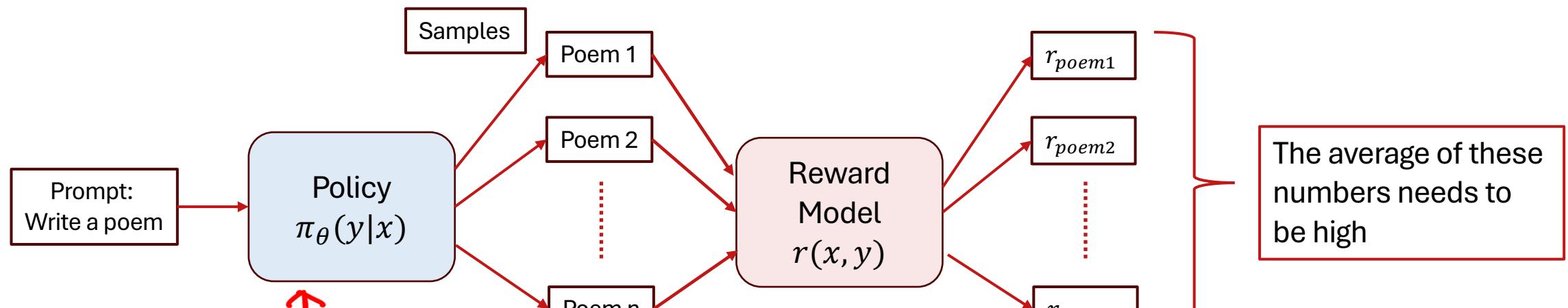
Reward Models are not perfect.

- They have been trained to score only selected natural language outputs.
- The policy can hack the reward model – generate outputs with high reward but meaningless
- An input can have multiple correct outputs (Write a poem?)
  - Reward maximization can collapse the probability to 1 outputs
  - Staying close to  $\pi_{ref}$  can preserve diversity.



# Formulating the objective – Reward Maximization

- What does it mean for a policy to have high reward?



↑  
Train  $\theta$  such that

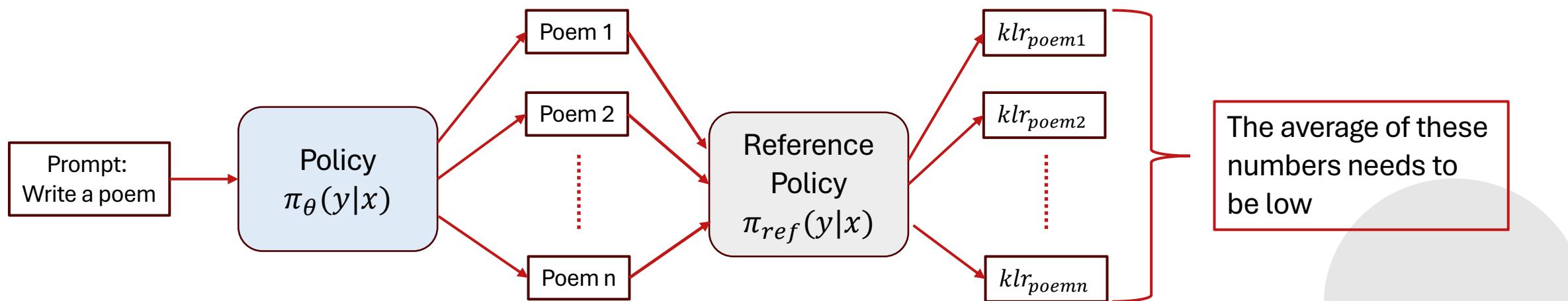
$$\mathbb{E}_{y \sim \pi_\theta(y|x)} r(x,y) \uparrow \approx \frac{1}{n} \sum_{i=1}^n r(x,y_i) \\ y_1, \dots, y_n \sim \pi_\theta(y|x)$$



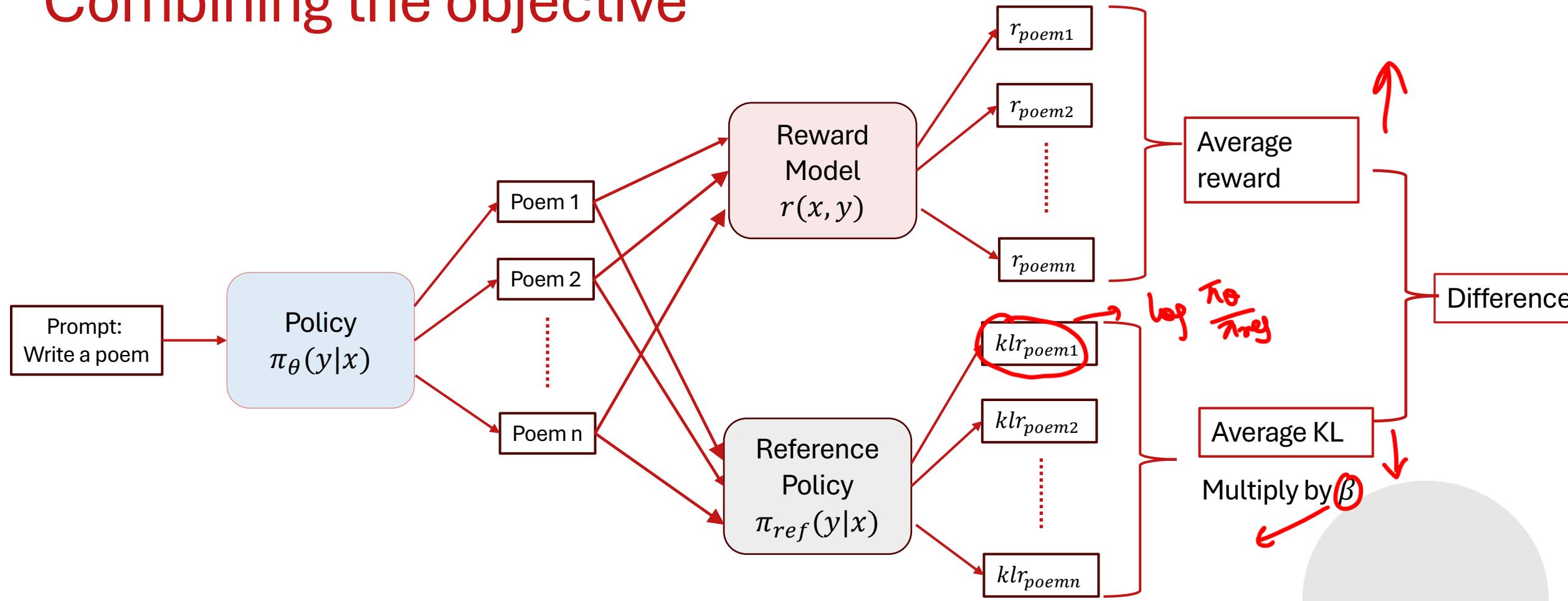
# Formulating the objective – closeness to $\pi_{ref}$

- How do we capture closeness to  $\pi_{ref}$ ?

$$KL(\pi_\theta(y|x) \parallel \pi_{ref}(y|x)) = \mathbb{E}_{y \sim \pi_\theta(y|x)} \left[ \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} \right]$$



# Combining the objective



# Regularized reward maximization

- Maximize the reward

$$\mathbb{E}_{y \sim \pi_\theta(y|x)} r(x,y)$$

- Minimize the KL divergence

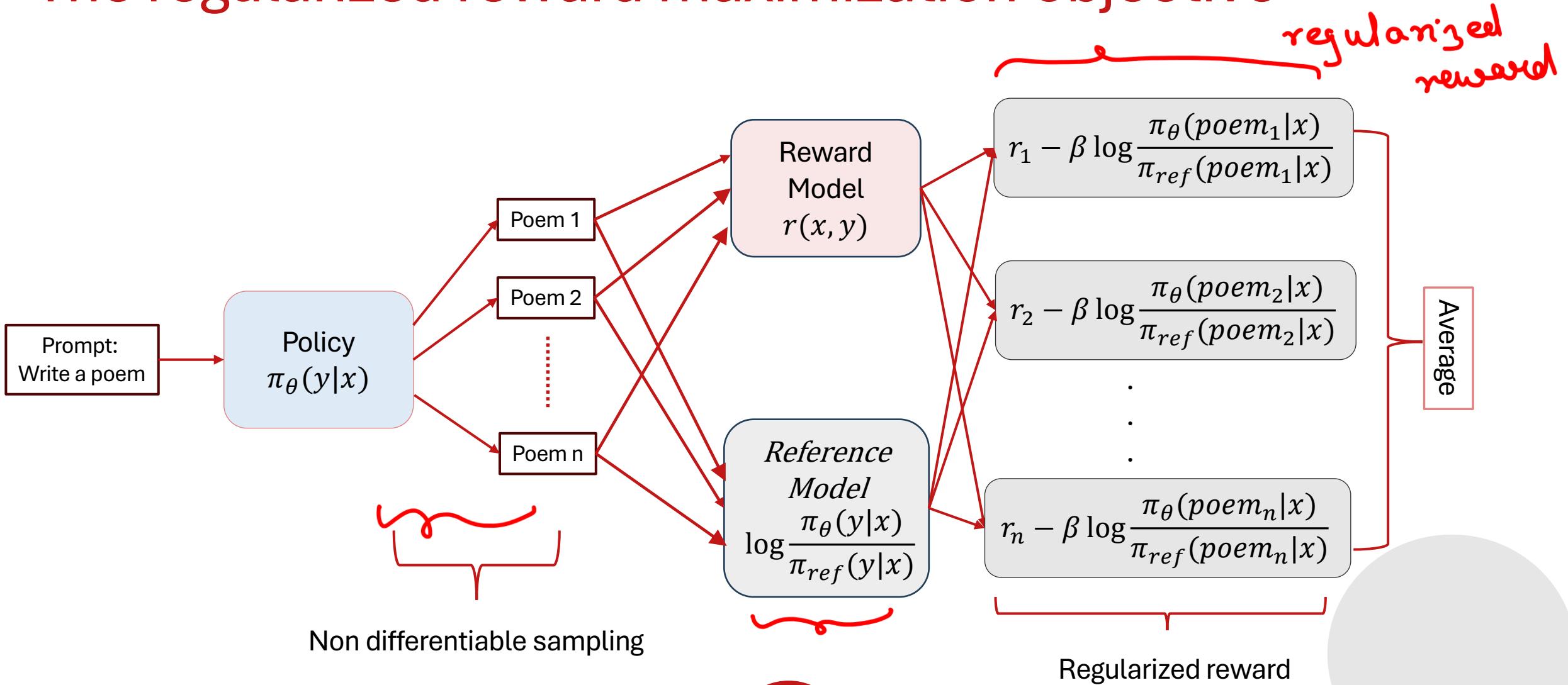
$$\mathbb{E}_{y \sim \pi_\theta(y|x)} \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$$

- Add a scaling factor  $\beta$  & combine

$$\mathbb{E}_{y \sim \pi_\theta(y|x)} \left[ r(x,y) - \beta \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \right]$$



# The regularized reward maximization objective



# Reward Maximization

- Best-of-N policy
- Reward Maximization Objective Formulation
- Towards Vanilla Policy Gradient
  - The REINFORCE gradient estimator
  - The problem of high variance - Baseline subtraction
  - Final algorithm
- Towards Actor/Critic methods
  - The problem of credit assignment
  - Q-functions, Value functions and advantage
  - Estimating the value function and advantage
  - Final algorithm



# How to maximize the objective?

- Derive the gradient of the objective.
- Estimate the gradient - the REINFORCE gradient estimator.
- Train using Adam/Adagrad optimization algorithms

$$\begin{aligned}\nabla_{\theta} E_{\pi_{\theta}}(y|x) r(x, y) &= \nabla_{\theta} \sum_{y \in Y} \pi_{\theta}(y|x) r(x, y) \\ &= \sum_{y \in Y} \nabla_{\theta} \pi_{\theta}(y|x) r(x, y).\end{aligned}$$

How to estimate this  
using samples ?



# Computing the gradient

$$\sum_{y \in Y} \nabla_{\theta} \pi_{\theta}(y|x) r(x, y)$$

*Not an expectation*

- Exact computation of the gradient is intractable
  - Output space is too large
- Can we approximate it using samples?
- To be able to do that, we need an expression of the form

$$E_{\pi_{\theta}}(y|x)[...] = \sum_{y \in Y} \pi_{\theta}(y|x) [...]$$

- How to transform the derivative to this desired form?



# The log-derivative trick

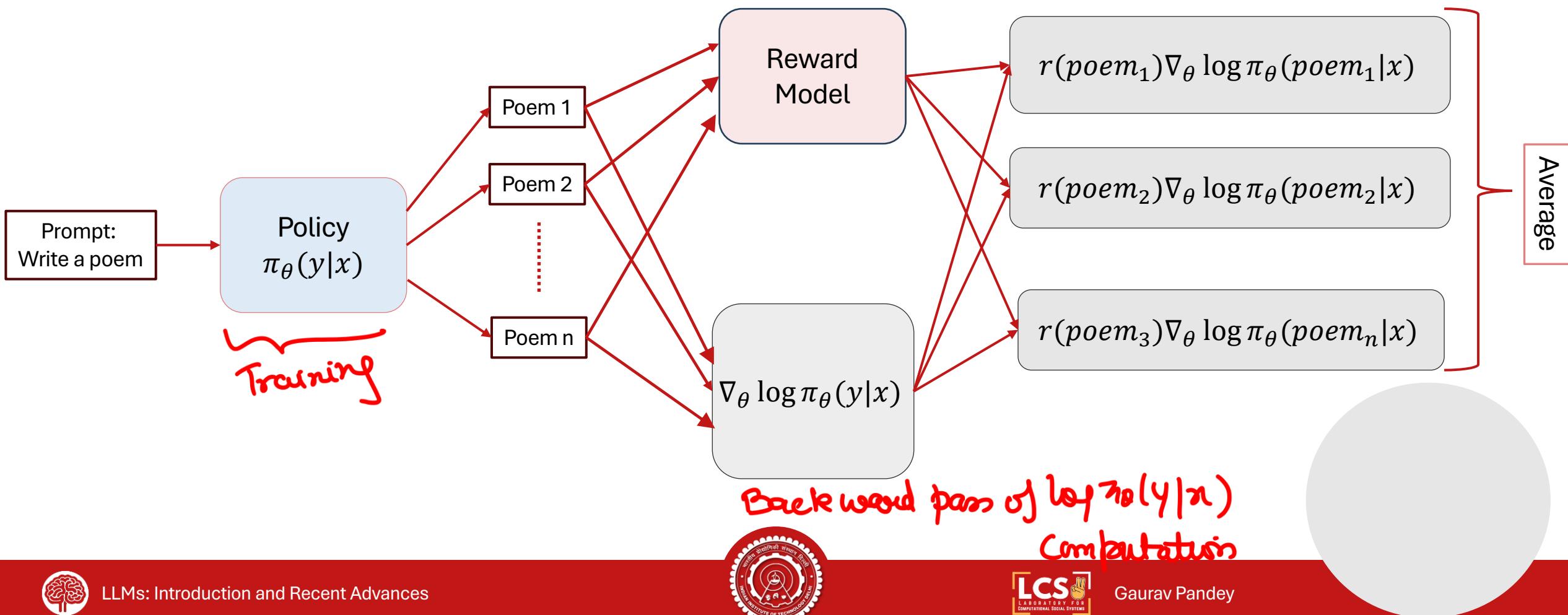
$$\nabla_{\theta} \log \pi_{\theta}(y|x) = \frac{1}{\pi_{\theta}(y|x)} \nabla_{\theta} \underbrace{\pi_{\theta}(y|x)}_{\pi_{\theta}(y|x)} \Leftrightarrow \nabla_{\theta} \pi_{\theta}(y|x) = \underbrace{\pi_{\theta}(y|x)}_{\pi_{\theta}(y|x)} \nabla_{\theta} \log \pi_{\theta}(y|x)$$

- Replacing it in the gradient, we get

$$\begin{aligned} \sum_{y \in Y} \nabla_{\theta} \pi_{\theta}(y|x) r(x,y) &= \sum_{y \in Y} \pi_{\theta}(y|x) \left[ \nabla_{\theta} \log \pi_{\theta}(y|x) r(x,y) \right] \\ &= \mathbb{E}_{y \sim \pi_{\theta}(y|x)} \left[ r(x,y) \nabla_{\theta} \log \pi_{\theta}(y|x) \right] \end{aligned}$$



# Estimating the gradient – REINFORCE



# Estimating the gradient – REINFORCE

$$\nabla_{\theta} E_{\pi_{\theta}}(y|x) r(x, y) \approx \frac{1}{n} \sum_{i=1}^n r(x, y_i) \nabla_{\theta} \log \pi_{\theta}(y_i|x)$$

Where  $y_1, \dots, y_n$  are samples from the policy  $\underline{\pi_{\theta}(y|x)}$ .

- The estimator is very noisy – different samples give different gradient estimates.



# Expanding the gradient

- Let  $y = (a_1, \dots, a_m)$  be the tokens of  $y$ .

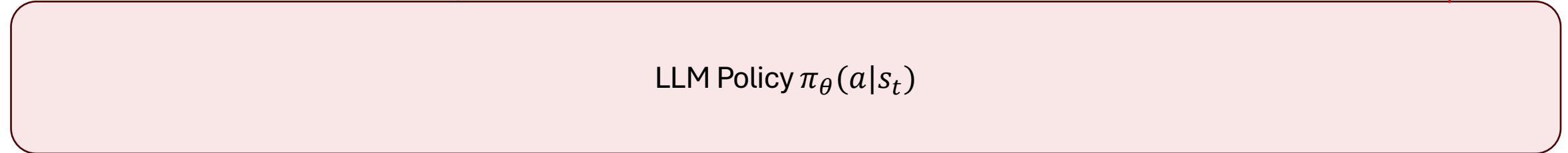
$$\begin{aligned} r(x, y) \nabla_{\theta} \log \pi_{\theta}(y|x) &= \underbrace{r(x, y)}_{\gamma(x, y)} \nabla_{\theta} \left[ \sum_{j=1}^m \log \pi_{\theta}(a_j | a_1, \dots, a_{j-1}, x) \right] \\ &= r(x, y) \sum_{j=1}^m \nabla_{\theta} \log \pi_{\theta}(a_j | a_1, \dots, a_{j-1}, x) \\ &= \sum_{j=1}^m \underbrace{r(x, y)}_{\gamma(x, y)} \nabla_{\theta} \log \pi_{\theta}(a_j | \underbrace{a_1, \dots, a_{j-1}}_{s_j}, x) \end{aligned}$$



# Implementing the REINFORCE gradient estimate

$r = r(x, y)$   
is the total reward of  
the complete output

*reward model*  
 $-ve \beta KL$



# The Vanilla Policy Gradient Algorithm

- Repeat until convergence
  - Sample a batch of prompts  $B$
  - For each prompt  $x$ , sample one or more outputs -  $y_1, \dots, y_n \sim \underline{\pi_\theta(y|x)}$
  - For each output  $y$ 
    - Compute the total reward of the output  $r(x, y)$  ✓
    - Compute the log-probability of the output  $\log \pi_\theta(y|x)$
  - To get the loss, multiply the reward with the negative log-probability of the output
    - For each output of the prompt
    - For each prompt in the batchand average.
  - Backpropagate the loss to compute the gradient.
  - Use Adam or any other optimizer to update the weights.



# Motivating the baseline

- Let  $X$  and  $Y$  be random variable defined as below

$$P(X = x) = \begin{cases} \frac{1}{10}, & 95 < x \leq 105 \\ 0, & \text{otherwise.} \end{cases} \quad P(Y = y) = \begin{cases} \frac{1}{2}, & x \in \{-1,1\} \\ 0, & \text{otherwise} \end{cases}$$

- What do samples of  $XY$  look like?

95, -96, 103, -104 *(very +ve or very -ve)*

- The variance is too high.

- What do samples of  $(X - 100)Y$  look like?

$\tilde{E}[X]$  -5, 4, 3, -4 *(variance is low)*

- Centering  $X$  can reduce the variance.



# Baseline Subtraction to reduce variance

- The REINFORCE gradient estimate is an average of the product of 2 quantities:
  - The reward of the response  $r(x, y)$
  - The gradient of the log-prob of output under the policy  $\nabla_{\theta} \log \pi_{\theta}(y|x)$
- By centering the reward, we can reduce the variance of the estimator.
- Define  $\bar{r}(x) = \frac{1}{n} \sum_{i=1}^n r(x, y_i)$
- The centered REINFORCE gradient estimate is given by

$$\frac{1}{n} \sum_{i=1}^n \underbrace{(r(x, y_i) - \bar{r}(x))}_{\hat{r}(x, y_i)} \nabla_{\theta} \log \pi_{\theta}(y_i|x)$$



# The Vanilla Policy Gradient Algorithm (with baseline)

- Repeat until convergence
  - Sample a batch of prompts  $B$
  - For each prompt  $x$ , sample one or more outputs -  $y_1, \dots, y_n \sim \pi_\theta(y|x)$
  - For each output  $y$ 
    - Compute the total reward of the output  $r(x, y)$
    - Compute the log-probability of the output  $\log \pi_\theta(y|x)$
  - Compute the average reward (per prompt or per batch)
    - Subtract it from the reward of each output to get baselined reward
    - $\hat{r}(x, y) = r(x, y) - \bar{r}(x, y)$
  - To get the loss, multiply and average the **baselined** reward with the negative log-probability of the output
  - Backpropagate the loss to compute the gradient.
  - Use Adam or any other optimizer to update the weights.



# Reward Maximization

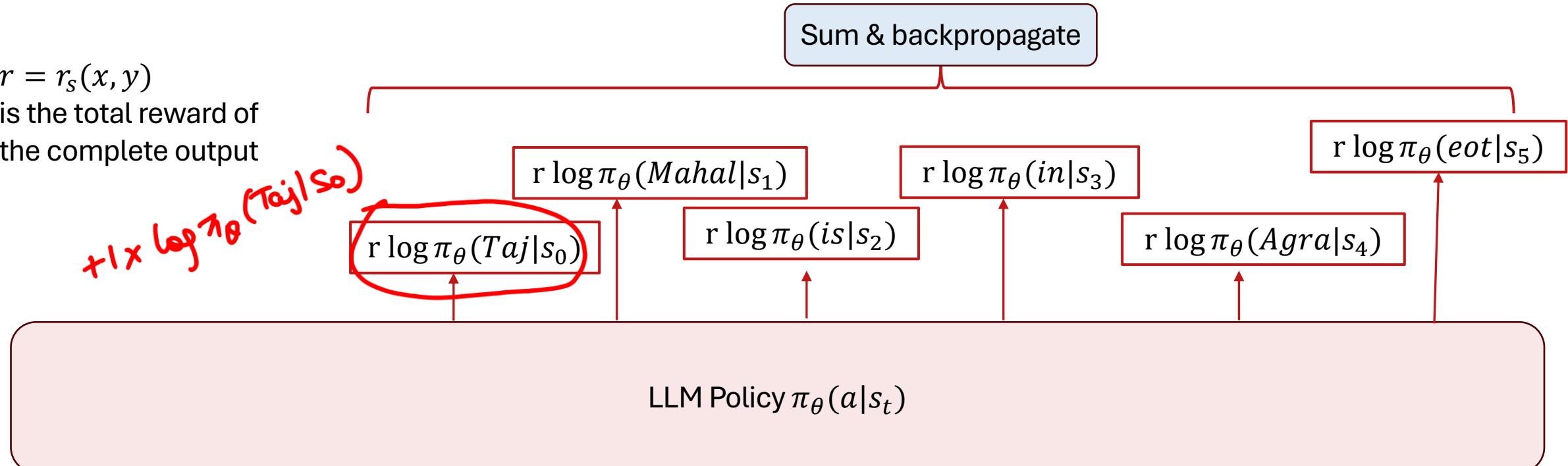
- Best-of-N policy
- Reward Maximization Objective Formulation
- Vanilla Policy Gradient
  - The REINFORCE gradient estimator
  - Baseline subtraction
  - Final algorithm
- Towards Actor/Critic methods
  - The problem of credit assignment
  - Q-functions, Value functions and advantage
  - Estimating the value function and advantage
  - Actor-Critic Family of algorithms



# The REINFORCE gradient estimate

$$r = r_s(x, y)$$

is the total reward of  
the complete output



# The problem of credit assignment

- The reward of generating a token (like “Taj”) depends on how the future sequence unfolds.
  - Example: “Taj Mahal is in Paris” → negative reward.
  - Example: “Taj Mahal is in Agra” → positive reward.
- This makes credit assignment hard: should “Taj” be encouraged or discouraged?
- The **Q-function** solves this by taking the expected reward of generating “Taj” and then continuing according to the policy.
- In other words, the Q-function averages over all possible continuations that follow “Taj,” giving a stable estimate for learning.

The cumulative discounted reward at time  $t$  is also known as reward-to-go or return at time  $t$

↓  
using the policy.



# The Q-function

- Given a state  $s$  and action  $a$ , the Q function for a policy  $\pi$  is defined as below:

$$Q(s, a) = \mathbb{E}_{a_{t+1}, a_{t+2}, \dots} \left[ r(s, a) + \gamma \mathbb{E}(s_{t+1}, a_{t+1}) + \gamma^2 \mathbb{E}(s_{t+2}, a_{t+2}) \dots \right]$$

- $s$ : state (context before “Taj”)
- $a$ : action (choose token “Taj”)
- $r_t$ : reward at step  $t$
- Expectation is over all future actions sampled from policy  $\pi$

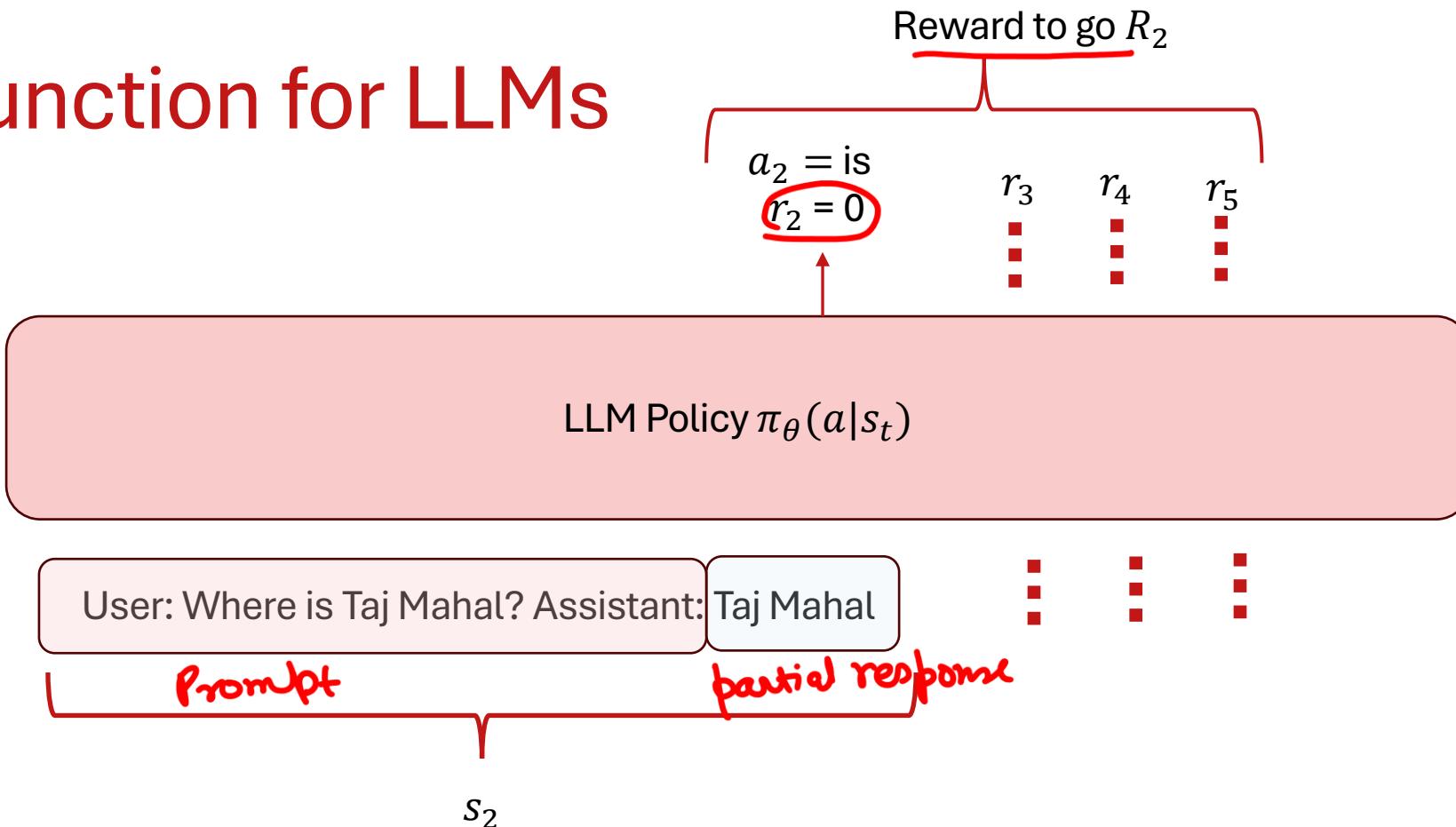


# Why discounting?

- Rewards that you get immediately after an action should get higher weight.
- What happens with discounting?
  - If  $\gamma = 1$ : final reward is propagated **equally to all tokens**.
  - If  $\gamma < 1$ : earlier tokens get **less credit**, later tokens get more.
- Response generation works well with  $\gamma = 1$ .
- Discounting is useful mainly when there are **intermediate rewards** (games, robotics)



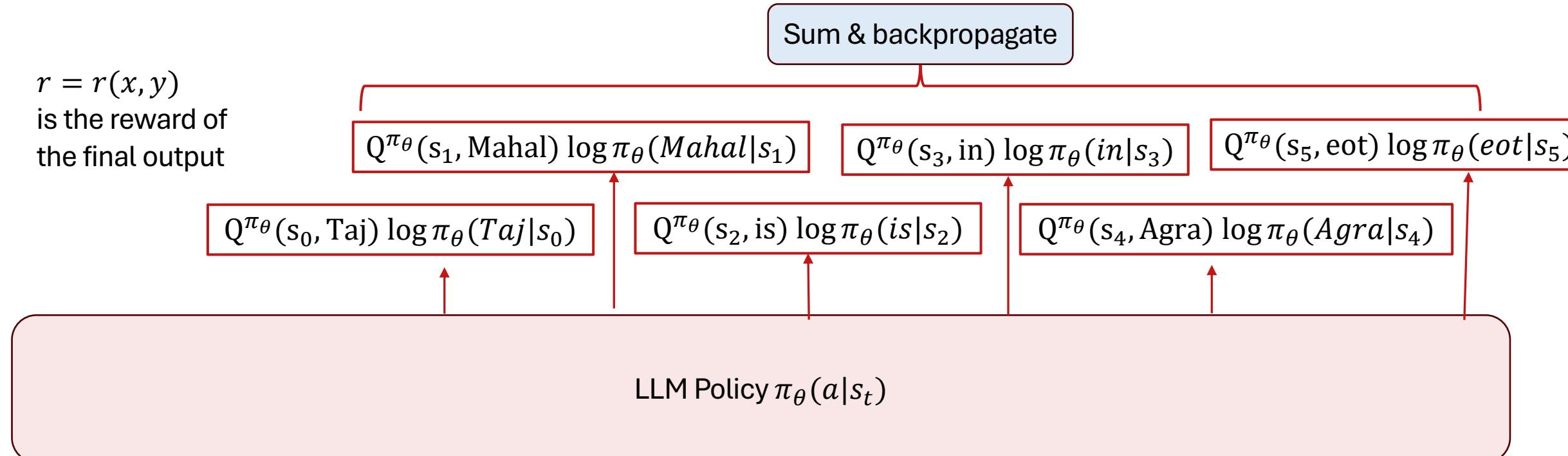
# The Q-function for LLMs



# REINFORCE gradient estimation with Q functions

$$r = r(x, y)$$

is the reward of  
the final output



Doesn't matter what gets generated in the future. The “reward” at token “Taj” is fixed.



# Value functions – A baseline for the Q functions

- Mean-centering the Q-function can further reduce variance in the estimator.
- The expected Q-function at a state  $s$  under policy  $\pi$  is referred to as the value function  $V^\pi(s)$

$$V^\pi(s) = E_{\pi(a|s)}Q(s, a) \quad \checkmark$$

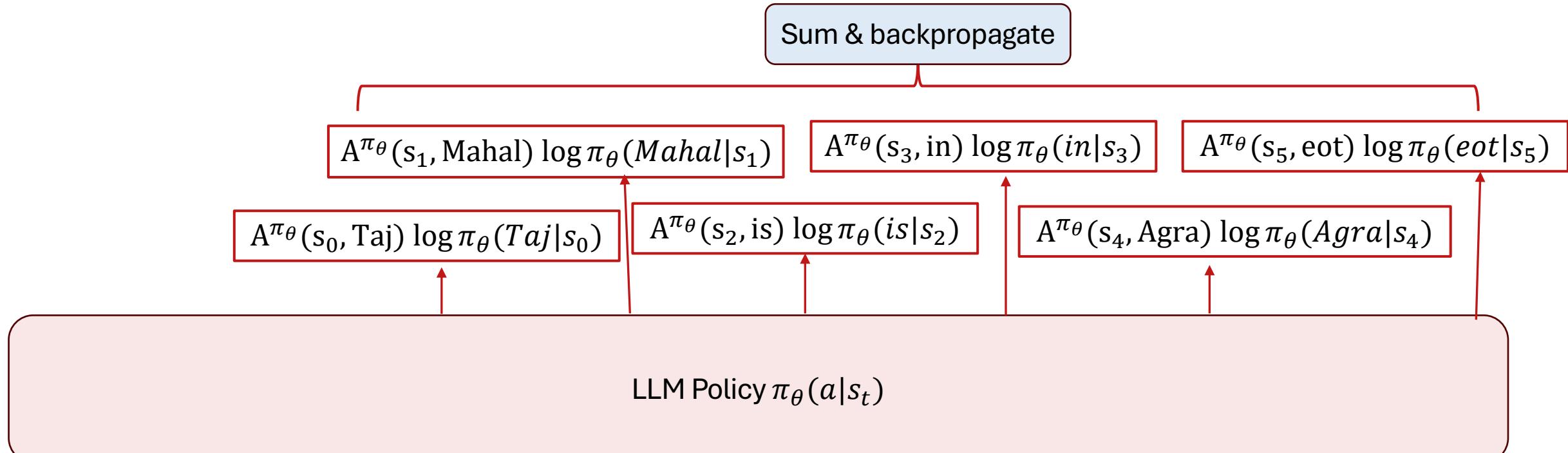
- The baseline subtracted Q is referred to as advantage function.

$$A^\pi(s, a) = \underline{Q^\pi(s, a) - V^\pi(s)} \rightarrow \text{Advantage function}$$

- Intuitively, advantage function captures contribution of the current action over an average action at the same state.



# REINFORCE with advantage functions



Doesn't matter what gets generated in the future. The “reward” at token “Taj” is fixed.



# Estimating Q and Value functions for an LLM

- Given: an input  $x$ , sample  $y = \underbrace{(a_0, \dots, a_T)}$  from the policy  $\pi_\theta(y|x)$
- The Q-function for a state-action pair is estimated using a single sample

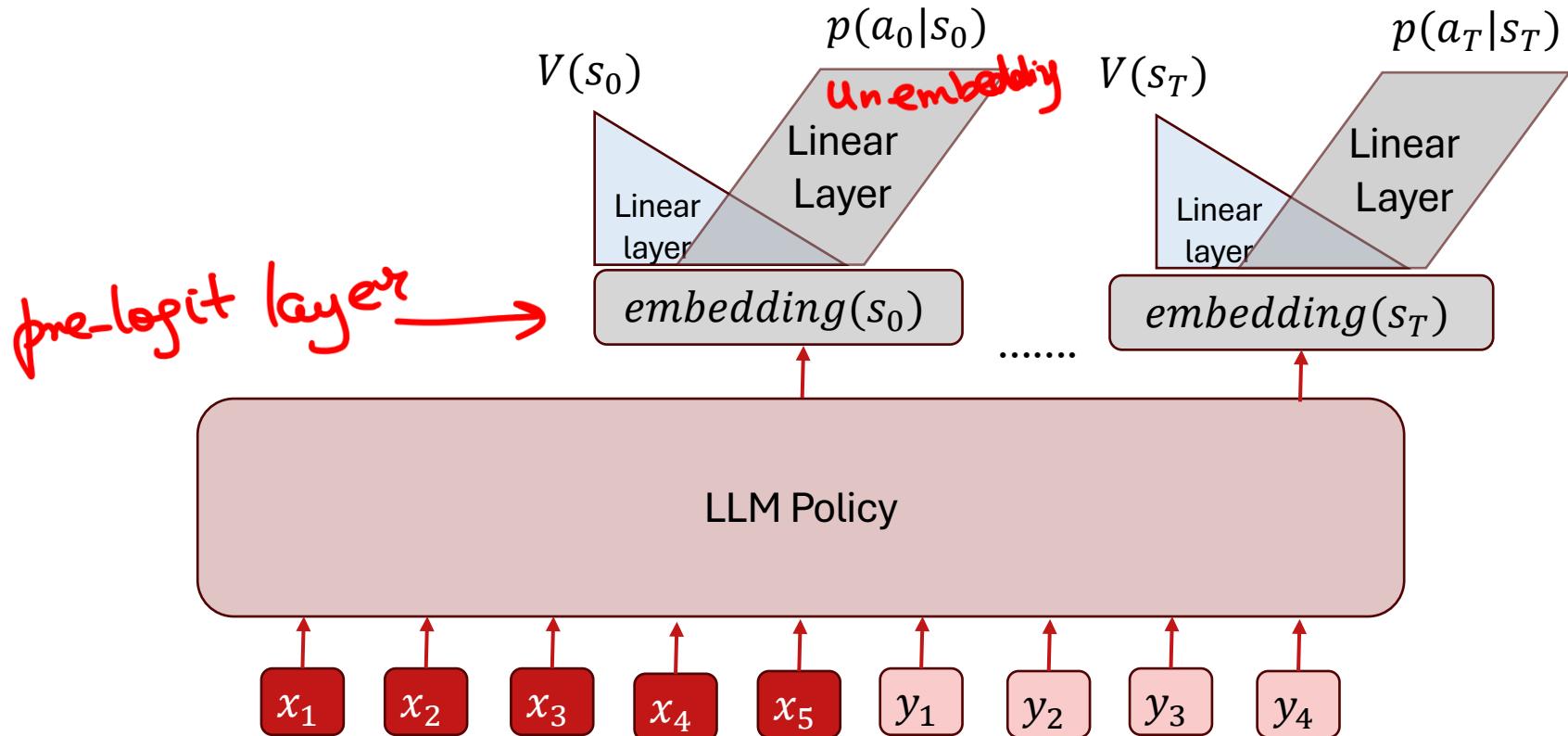
$$Q(s_t, a_t) = \underbrace{r(s_t, a_t) + r(s_{t+1}, a_{t+1}) \dots}_{\text{...}} + r(s_T, a_T)$$

- The value function is estimated by training a neural network
  - The cumulative discounted reward for each time-step is computed.
  - MSE between predicted value and the cumulative discounted reward is minimized.

$$\frac{1}{T} \sum_{t=0}^T (V_p(s_t) - [r(s_t, a_t) + r(s_{t+1}, a_{t+1}) \dots])^2.$$



# The Value function network $V_\phi$



# The Actor-Critic Family

- The actor is the policy network that chooses actions given states.
  - It is parameterized as  $\pi_\theta$
- The critic is the value function approximator  $V_\phi$
- The actor is trained using the following gradient estimate.

$$E[\nabla_\theta \log \pi_\theta(a_t | s_t)(R_t - V_\phi(s_t))] \quad \}$$

- The critic is trained to minimize the MSE

$$E[(R_t - V_\phi(s_t))^2] \quad \}$$

- This family of algorithms includes PPO.

