

Alignment of Language Models – Reward Maximization (Part -2)

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821

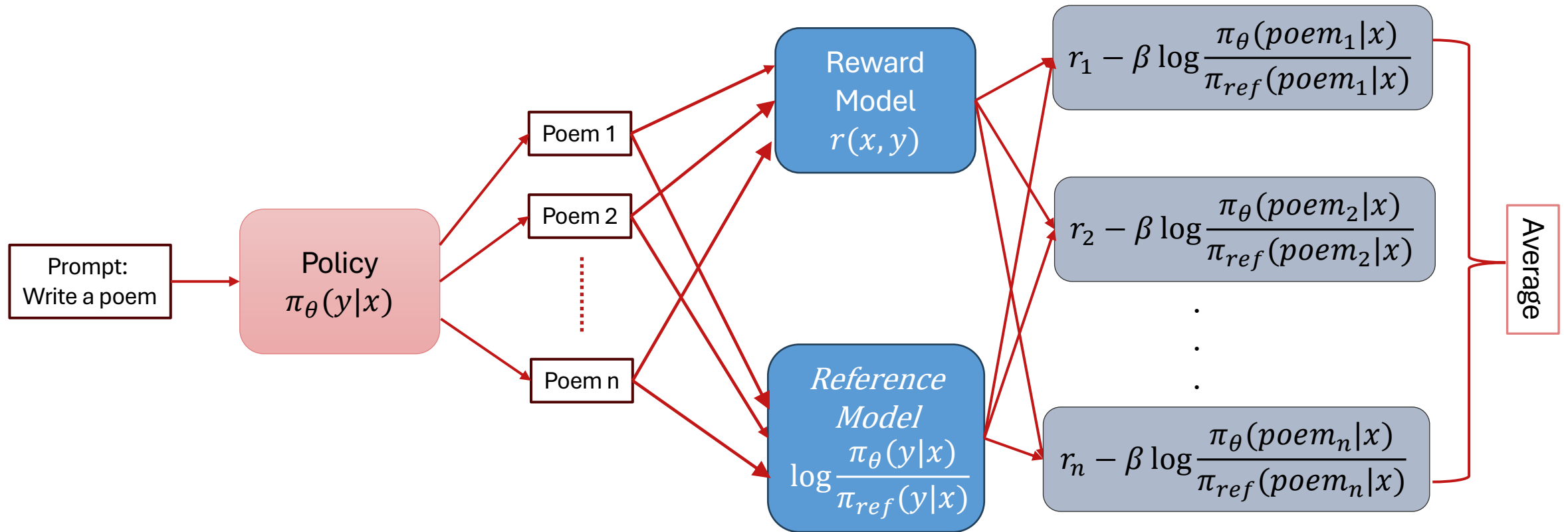


Gaurav Pandey
Research Scientist, IBM Research

Regularized reward maximization

- Maximize the reward
- Minimize the KL divergence
- Add a scaling factor β & combine

The regularized reward maximization objective



Regularized reward

$$E_{\pi_{\theta}(y|x)} \left[r(x, y) - \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{ref}(y|x)} \right] \equiv E_{\pi_{\theta}(y|x)} r_s(x, y)$$

$$\text{where } r_s(x, y) = r(x, y) - \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{ref}(y|x)}$$

- $r_s(x, y)$ is the regularized reward
- Maximizing the regularized reward ensures
 - High reward outputs as decided by the reward model
 - Outputs that have reasonable probability under the reference model



How to maximize – The REINFORCE algorithm?

- Compute the gradient of the objective.
- Train using Adam/Adagrad optimization algorithms

$$\nabla_{\theta} E_{\pi_{\theta}(y|x)} r_s(x, y)$$



Computing the derivative

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

- 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

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

- How to transform the derivative to this desired form?



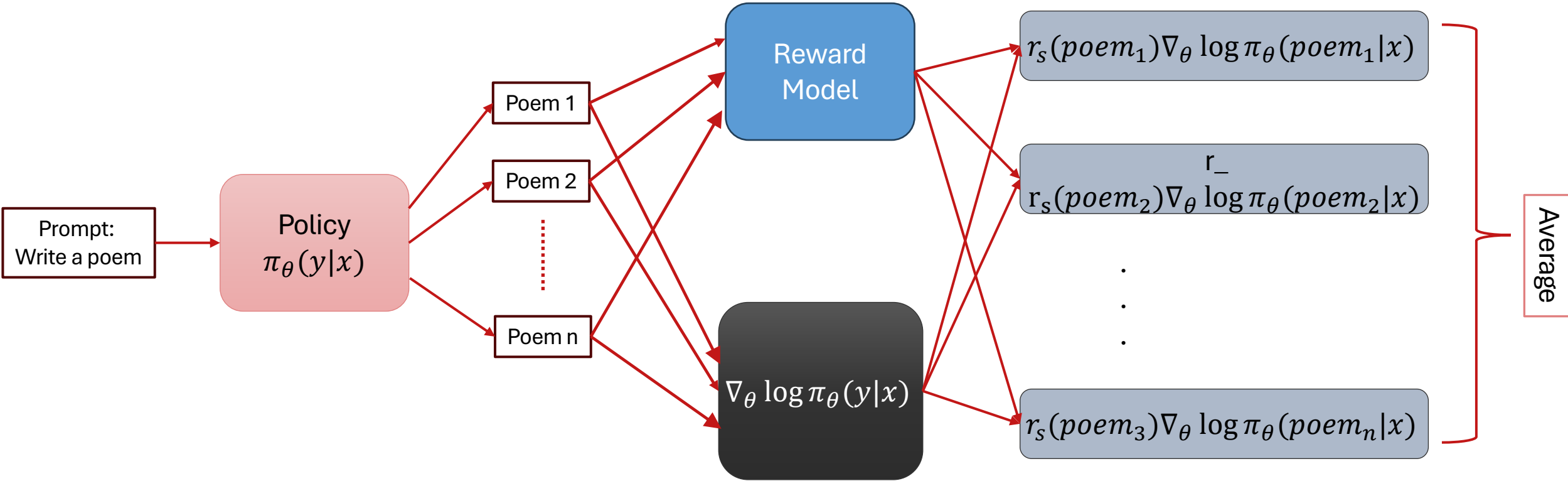
The log-derivative trick

$$\nabla_{\theta} \log \pi_{\theta}(y|x) =$$

Replacing it in the derivative, we get



Monte Carlo approximation



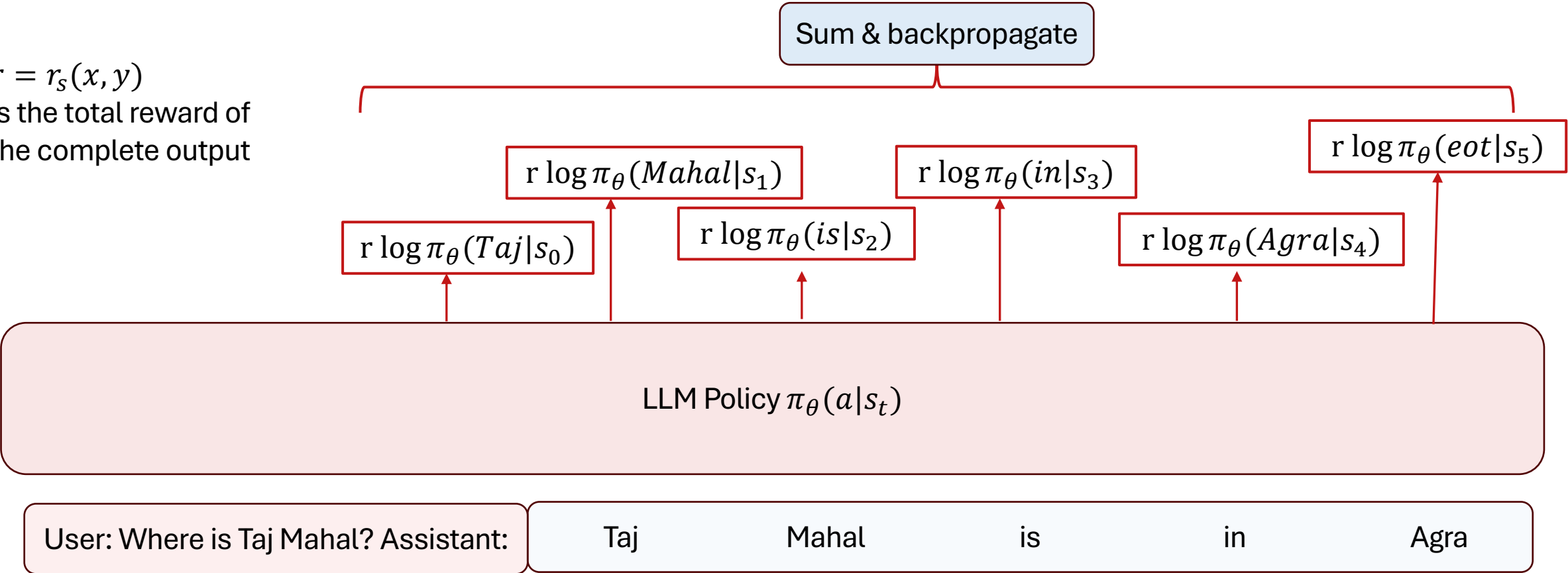
Expanding the gradient

- Let $y = (a_1, \dots, a_m)$ be the tokens of y .
- $r_s(x, y) \nabla_{\theta} \log \pi_{\theta}(y|x) =$



Implementing REINFORCE

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



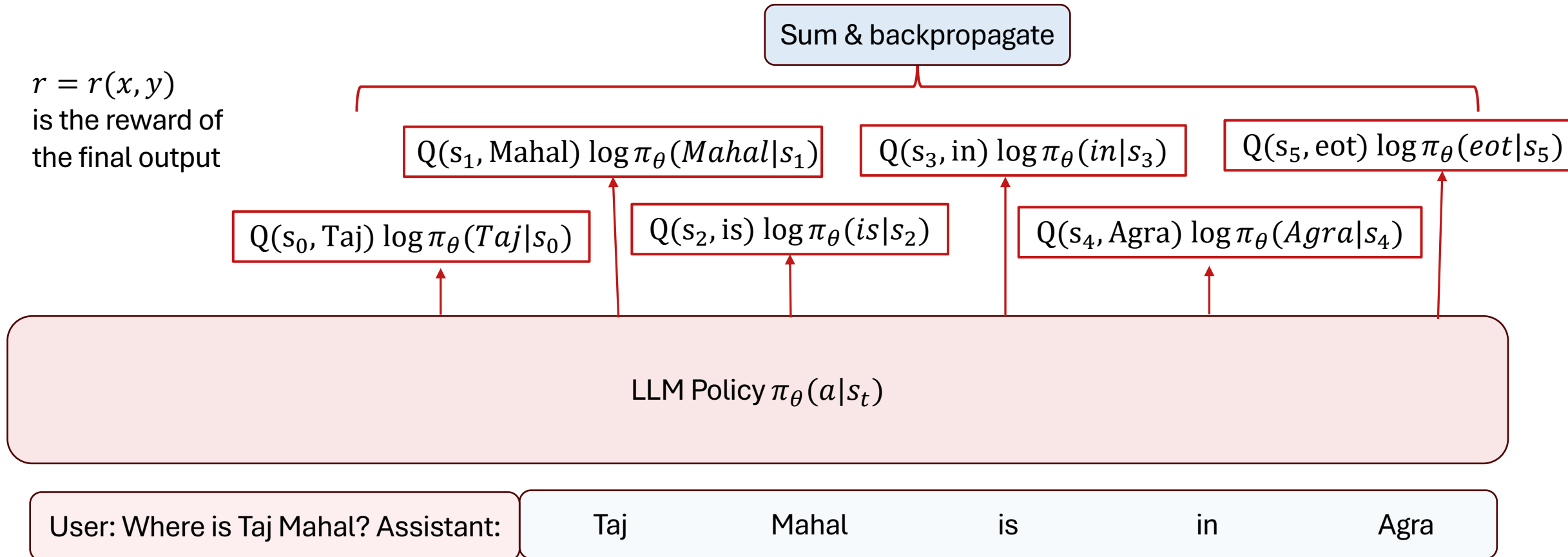
Problems with REINFORCE

- The reward at token “Taj” depends on the tokens generated in the future
- If the model had generated “Taj Mahal is in Paris”
 - The reward would be negative
 - The probability of generating Taj would be decreased
- If the model had generated “Taj Mahal is in Agra”
 - The reward would be positive
 - The probability of generating Taj would be increased
- This variance in the reward leads to unstable training.
- To reduce variance – take the average reward over all likely sequences (under the policy) that generate “Taj” for the first token.
- This is called the Q – *function*



REINFORCE 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.



Q-function & Value function

- The Q-function for a state-action pair is the average discounted cumulative reward received at the state after taking taking the specified action.
- The discount factor λ ensures that immediate rewards get higher weight.
- The Value function of a state is the average discounted cumulative reward received after reaching the state.



From Q-function to Advantage function

- For text generation using language models

$$s_{t+1} = (s_t, a_t)$$

- That is, once you have generated the next token, the next state is determined completely.

- Hence, the Q-function for a state-action pair can be written as

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma V(s_{t+1})$$

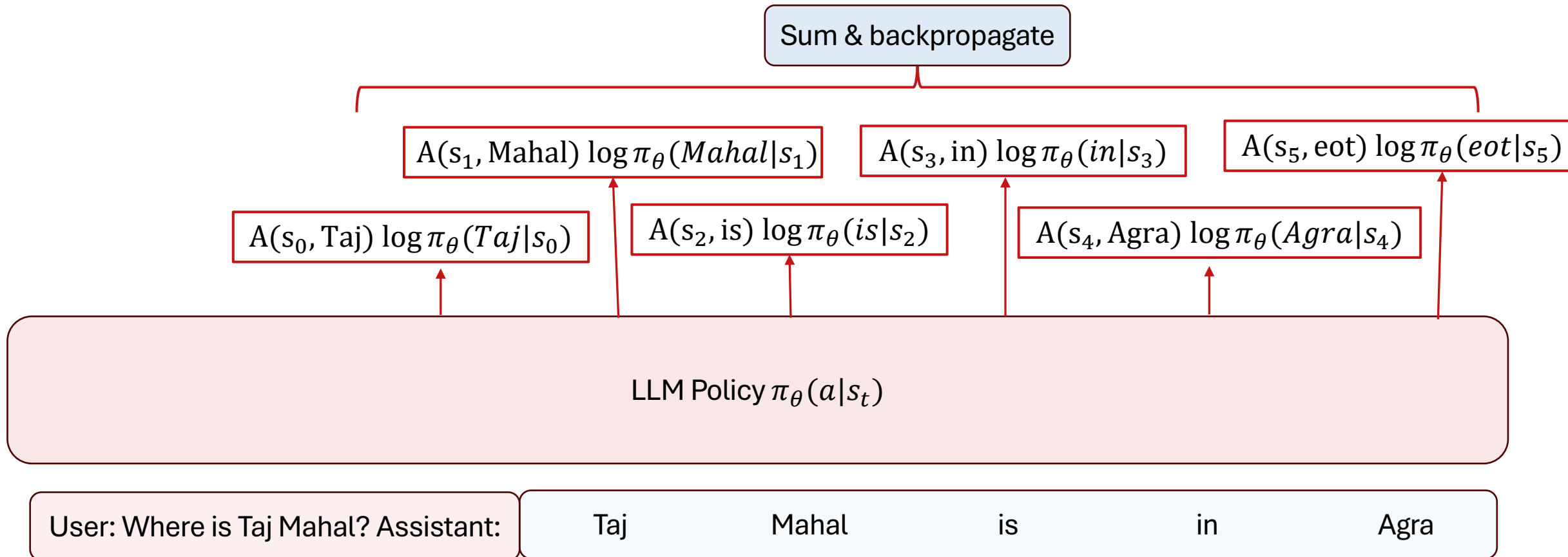
- To further reduce variance, the advantage function $A(s_t, a_t)$ is used instead of Q-function

$$\begin{aligned} A(s_t, a_t) &= Q(s_t, a_t) - V(s_t) \\ &= r(s_t, a_t) + \gamma V(s_{t+1}) - V(s_t) \end{aligned}$$

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



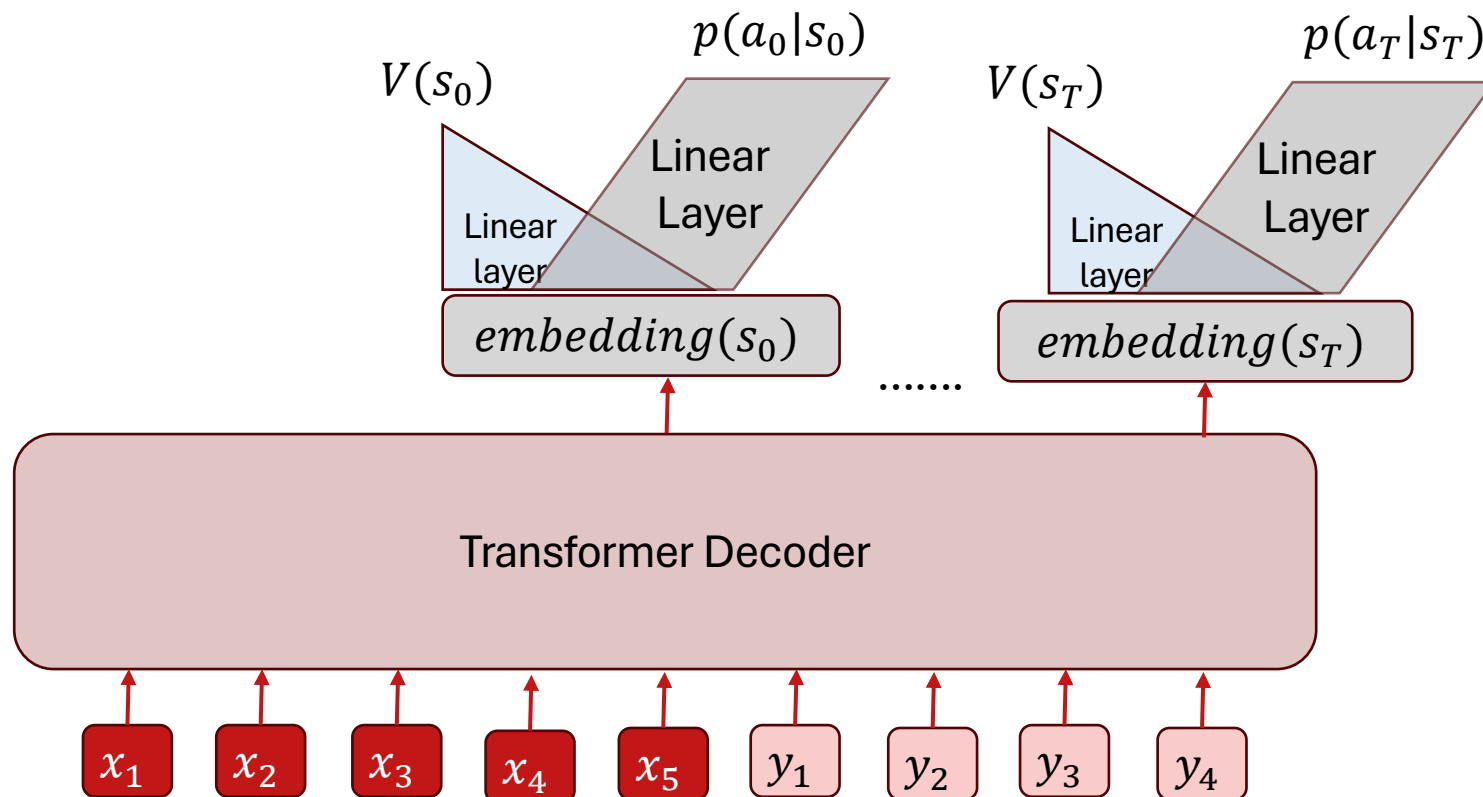
REINFORCE with advantage functions



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



Implementing the Value function



Learning the Value function

- Given an input x , sample $y = (a_0, \dots, a_T)$ from the policy $\pi_\theta(y|x)$
- Compute the cumulative discounted reward for each time-step

$$R_t =$$

- Minimize the mean-squared error



Vanilla Policy Gradient

- Repeat until convergence
 - Sample a batch of prompts B
 - For each prompt, sample one-or more outputs
 - For each output $y = (a_1, \dots, a_T)$
 - Compute the reward r_t at each token a_t
 - Compute cumulative discounted reward R_t for each token
 - Compute the value & advantage function A_t for each token
 - Apply few gradient updates using REINFORCE with the advantage values computed above
 - Apply few gradient updates to train the value function by minimizing the MSE.

Credit: <https://spinningup.openai.com/en/latest/algorithms/ppo.html>



Problems

- Sampling from the policy after every update can be challenging.
- Solution: Sample from an older fixed policy instead

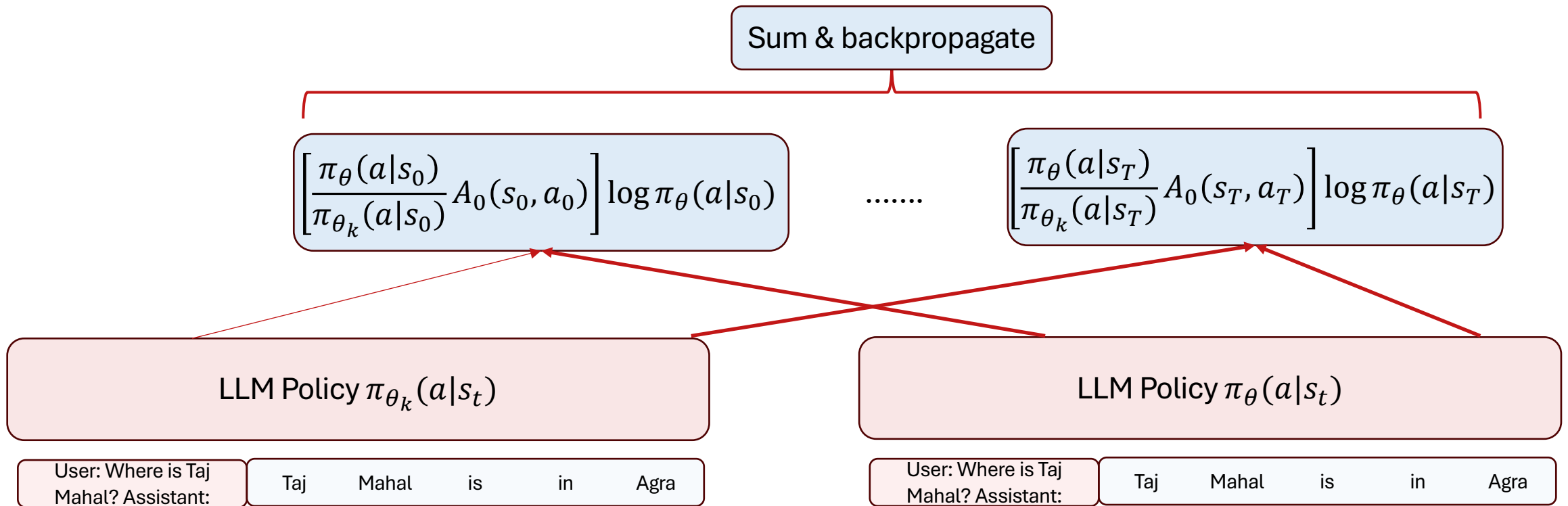


REINFORCE with importance weights



REINFORCE with importance weights

The term in the square brackets is kept constant during backpropagation. In Pytorch, this means using `.detach()` function



Proximal Policy Optimization

- Keeping the batch of prompts & outputs fixed, how much can we update the policy?
- If we update too much, the importance weights can change drastically.
- PPO-CLIP

$$(1 - \epsilon) \leq \frac{\pi_{\theta}(a_t | s_t)}{\pi_k(a_t | s_t)} \leq (1 + \epsilon)$$

- This ensures that the no matter how many updates are done to π_{θ} , it stays close to π_{θ_t}



PPO-CLIP

$$(1 - \epsilon) \leq \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_t}(a_t|s_t)} \leq (1 + \epsilon)$$

To achieve above, maximize the following

- When advantage is positive

$$\max_{\theta} \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)}, (1 + \epsilon) \right) A_t(s_t, a_t)$$

- When advantage is negative

$$\max_{\theta} \max \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)}, (1 - \epsilon) \right) A_t(s_t, a_t)$$

Credit: <https://spinningup.openai.com/en/latest/algorithms/ppo.html>



PPO-CLIP with +ve advantage

$$\max_{\theta} \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_t}(a_t|s_t)}, (1 + \epsilon) \right) A_t(s_t, a_t)$$



PPO-CLIP with -ve advantage

$$\max_{\theta} \max \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_t}(a_t|s_t)}, (1 - \epsilon) \right) A_t(s_t, a_t)$$



The PPO-CLIP algorithm

- For $k = 1$ to K
 - Sample a batch of prompts B
 - For each prompt, sample one-or more outputs from $\pi_{\theta_k}(y|x)$
 - For each output $y = (a_1, \dots, a_T)$
 - Compute the reward r_t at each token a_t
 - Compute cumulative discounted reward R_t for each token
 - Compute the value & advantage function A_t for each token
 - Apply few gradient updates using REINFORCE PPO-CLIP with the advantage values computed above
 - Apply few gradient updates to train the value function by minimizing the MSE.

Credit: <https://spinningup.openai.com/en/latest/algorithms/ppo.html>



Things to remember

- The log-derivative trick should be used to compute gradient in REINFORCE
- The log-probability of the tokens should be weighed by the advantage function to reduce variance
- Importance weights should be used to allow sampling from a fixed policy
- The importance weights should be clipped to prevent large gradient updates.

