

# RL from Human Feedback

**Sham Kakade and Kianté Brantley**

**CS 2824: Foundations of Reinforcement Learning**

# Recap

We have covered a few RL algorithms REINFORCE, PPO;

# Recap

We have covered a few RL algorithms REINFORCE, PPO;

They all rely on a key and strong assumption: reward function/signal is given

# Motivation

Modern chatbots are pre-trained via next-token prediction on web data, followed by fine-tuning using human feedback (post-training)

# Outline

1. LLM as a policy

2. Learning reward functions from preference data

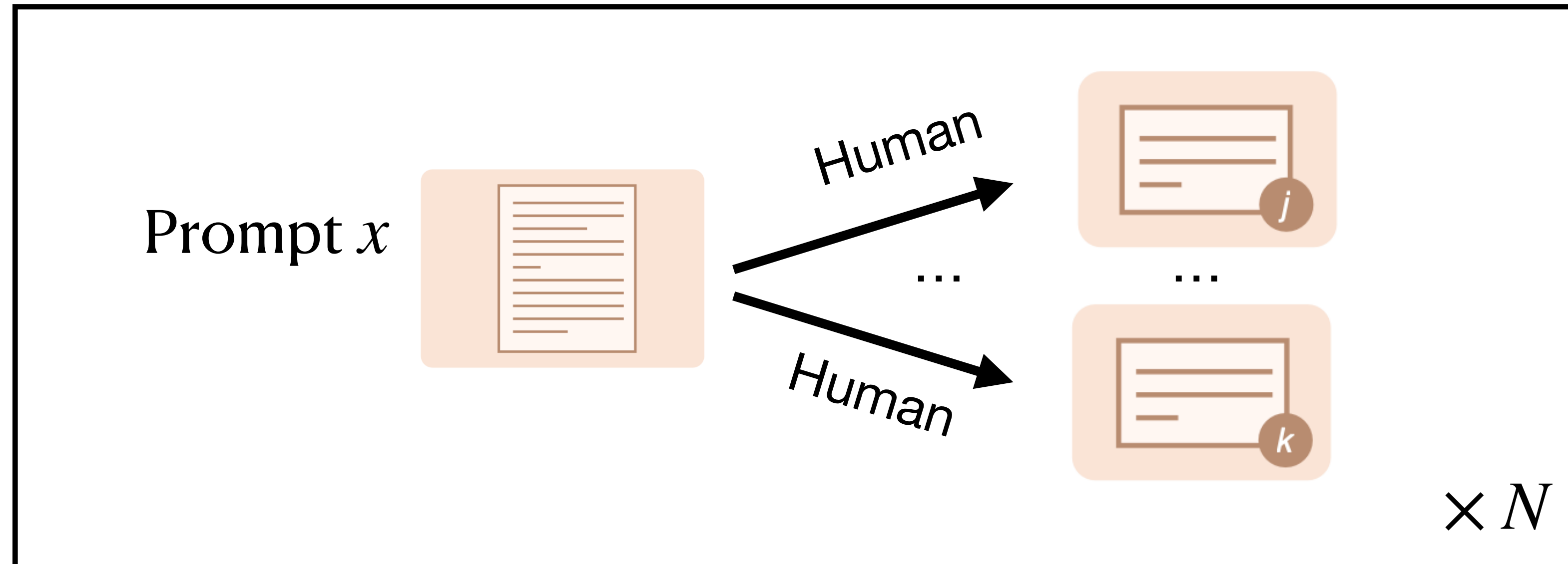
3. KL-regularized RL

4. DPO

5. REBEL

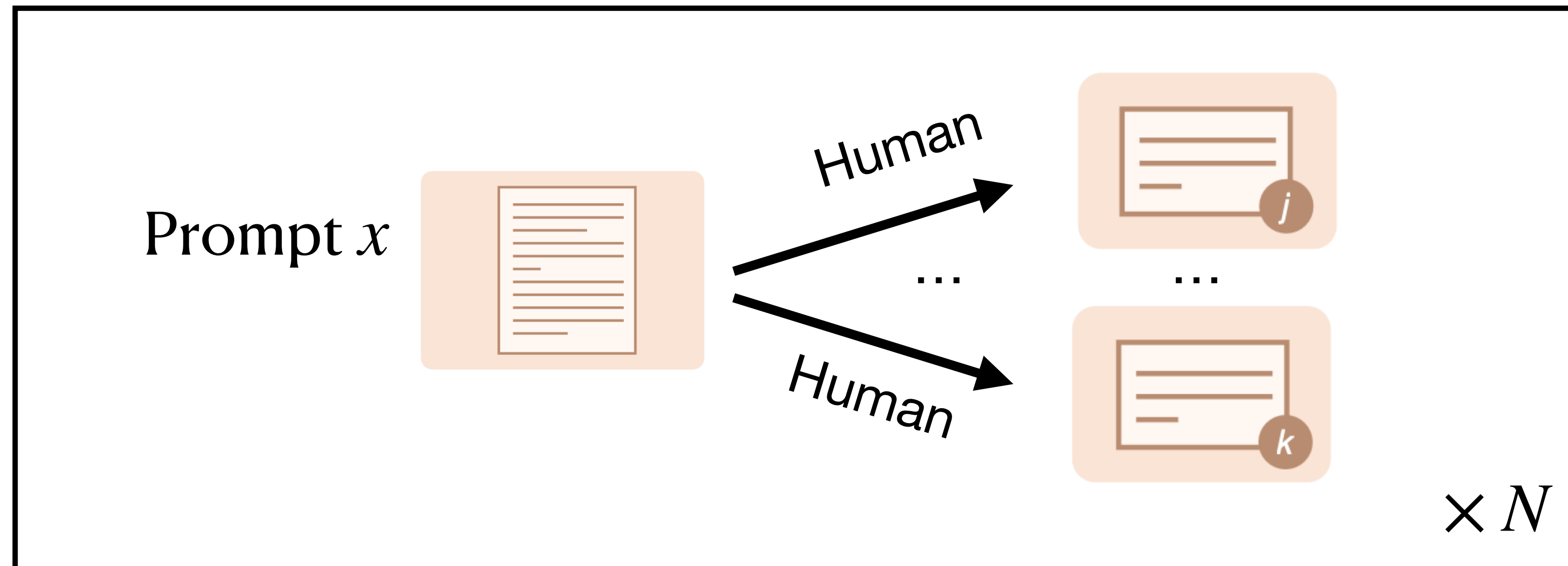
# The post-training Pipeline: Supervised Fine-tuning (SFT)

Collect instruction-response data



# The post-training Pipeline: Supervised Fine-tuning (SFT)

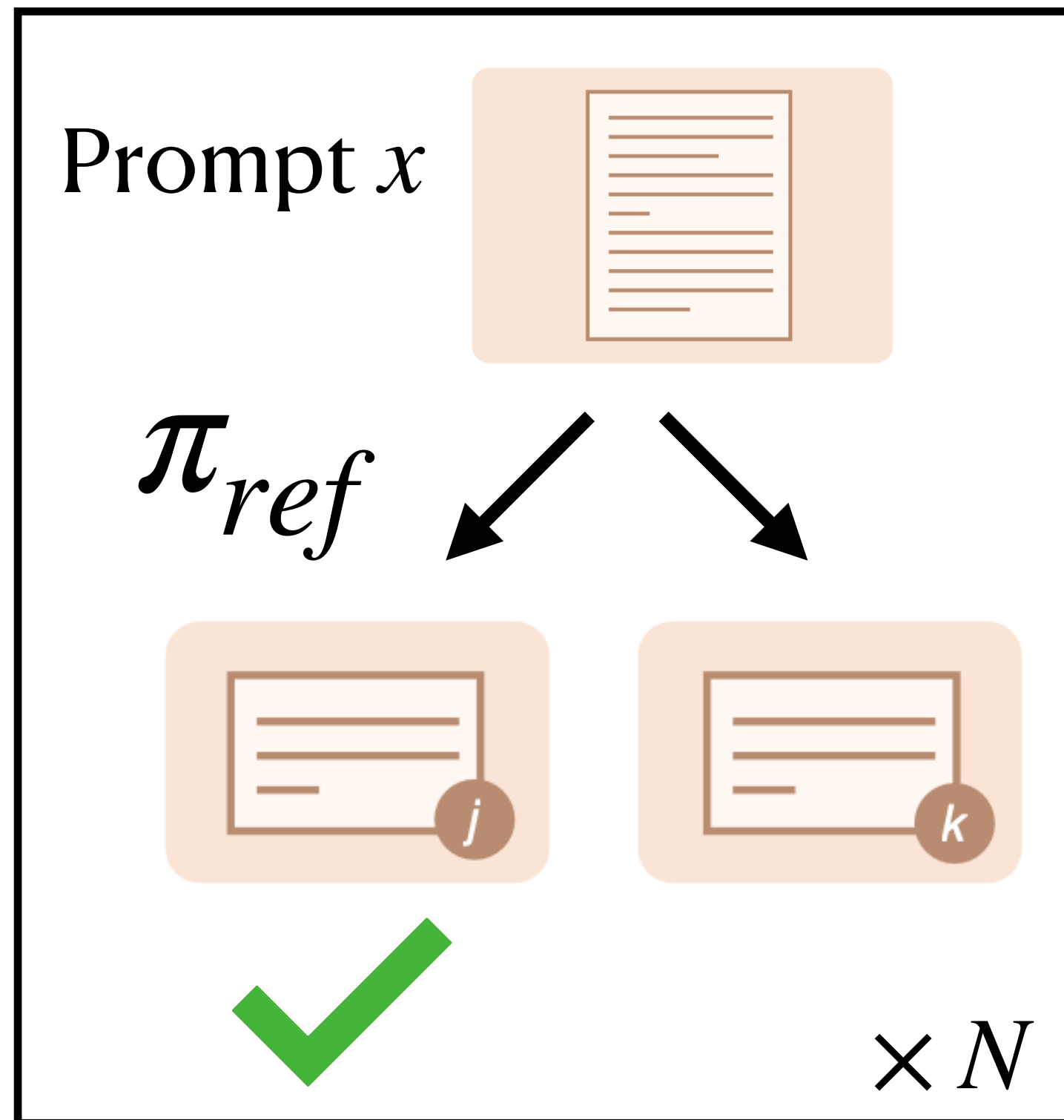
Collect instruction-response data



SFT: given prompts, train LLM to predict tokens in human responses

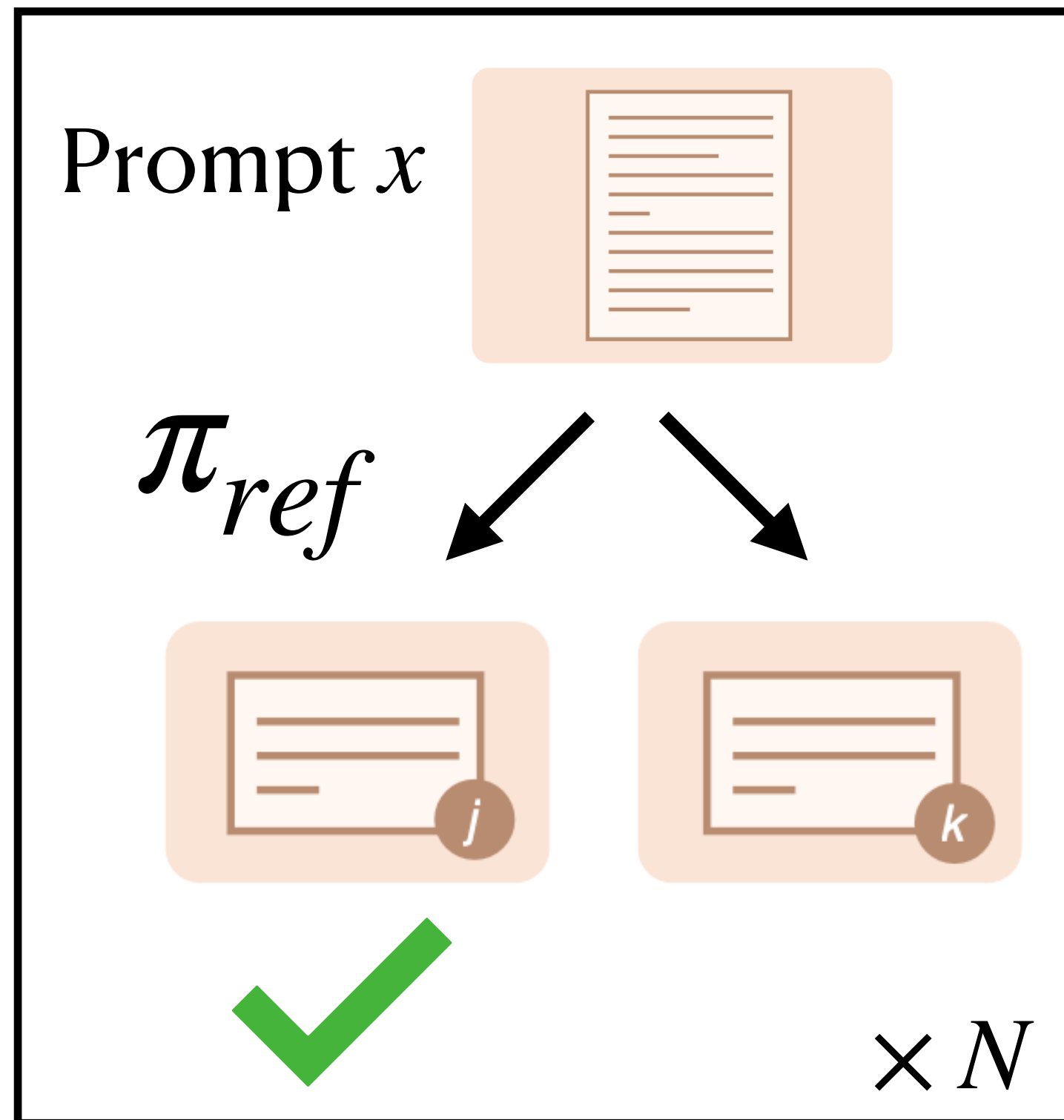
# The post-training Pipeline: RLHF

## 1. Collect preference dataset



# The post-training Pipeline: RLHF

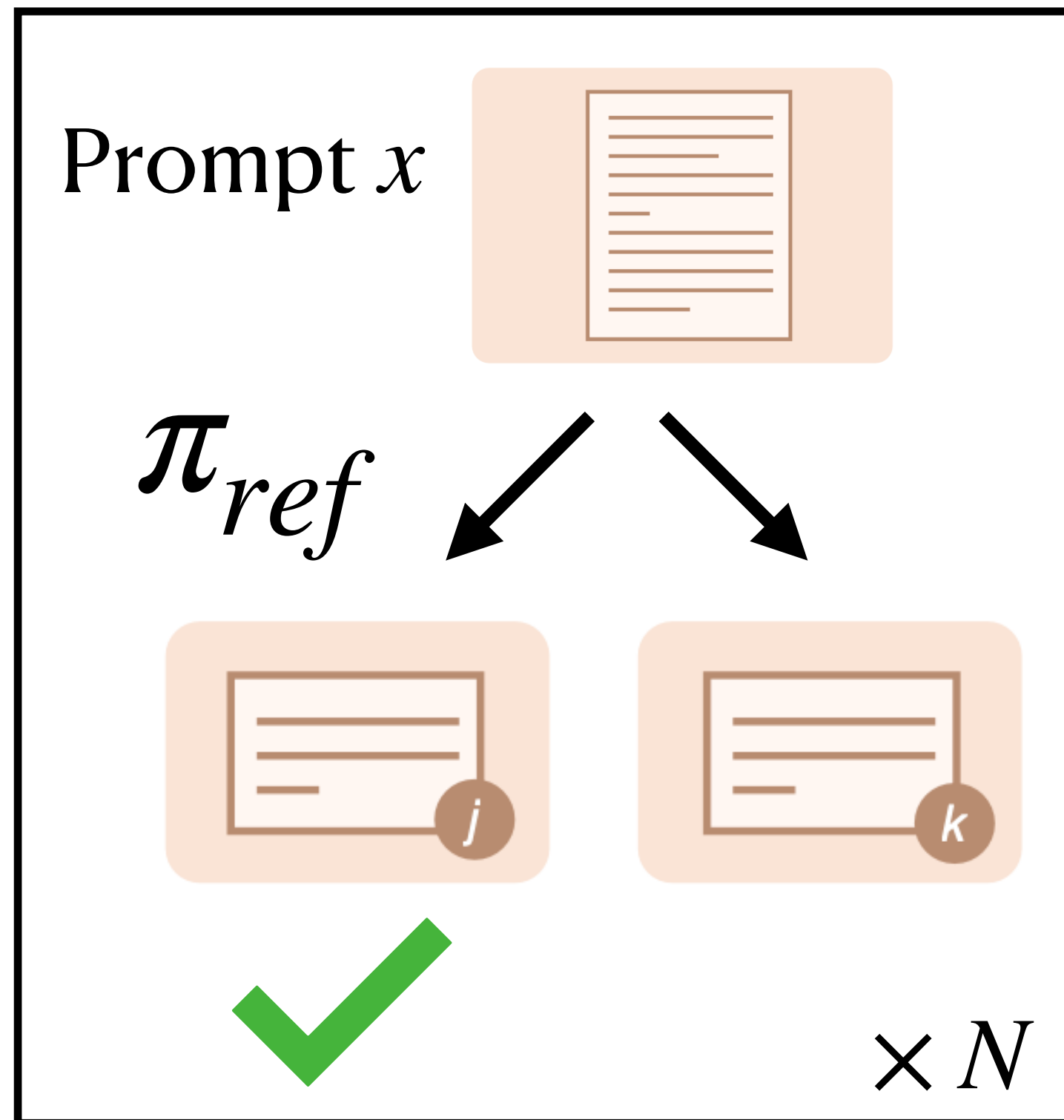
## 1. Collect preference dataset



$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

# The post-training Pipeline: RLHF

## 1. Collect preference dataset

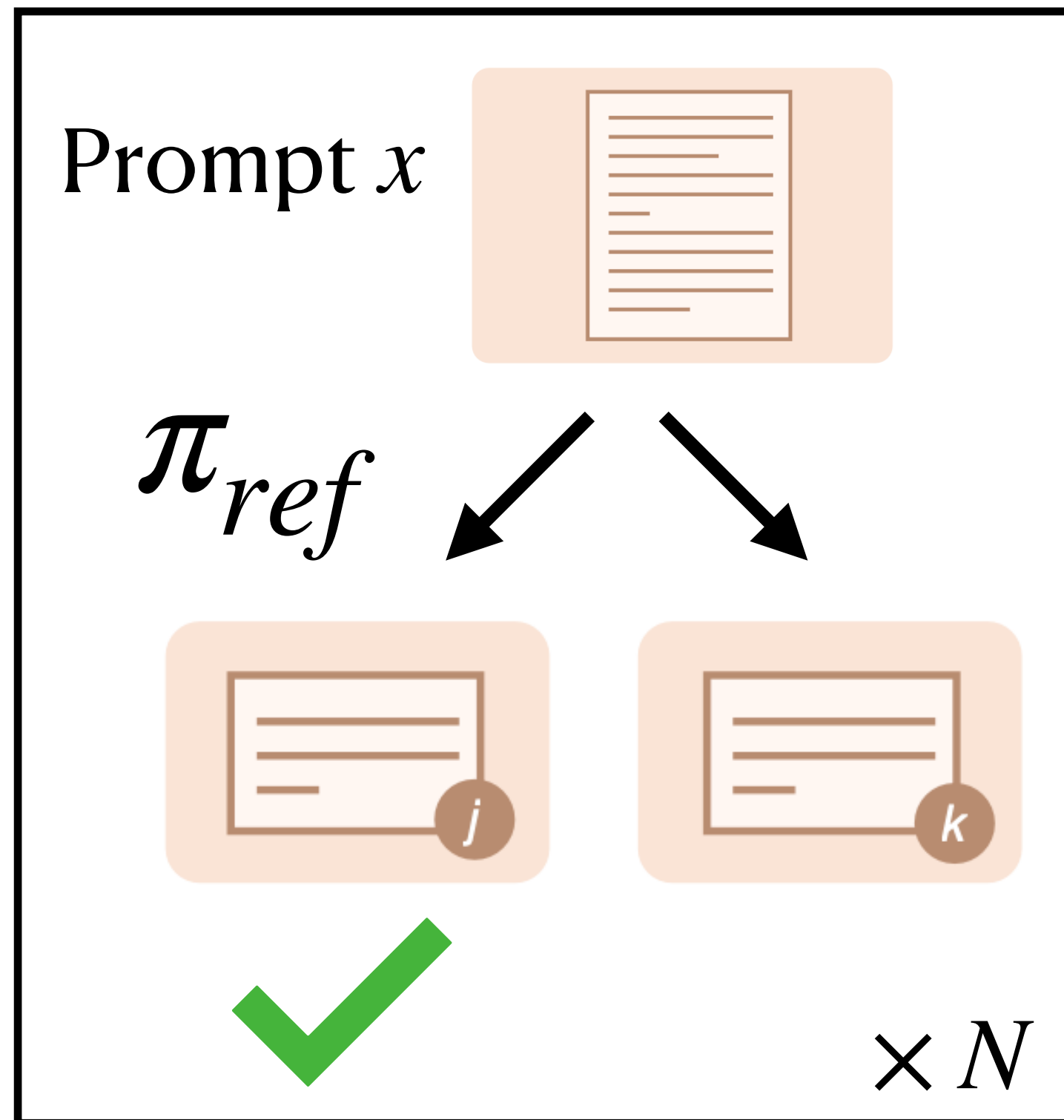


$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

## 2. Learn a reward model $\hat{r}$ using the data from step 1

# The post-training Pipeline: RLHF

1. Collect preference dataset



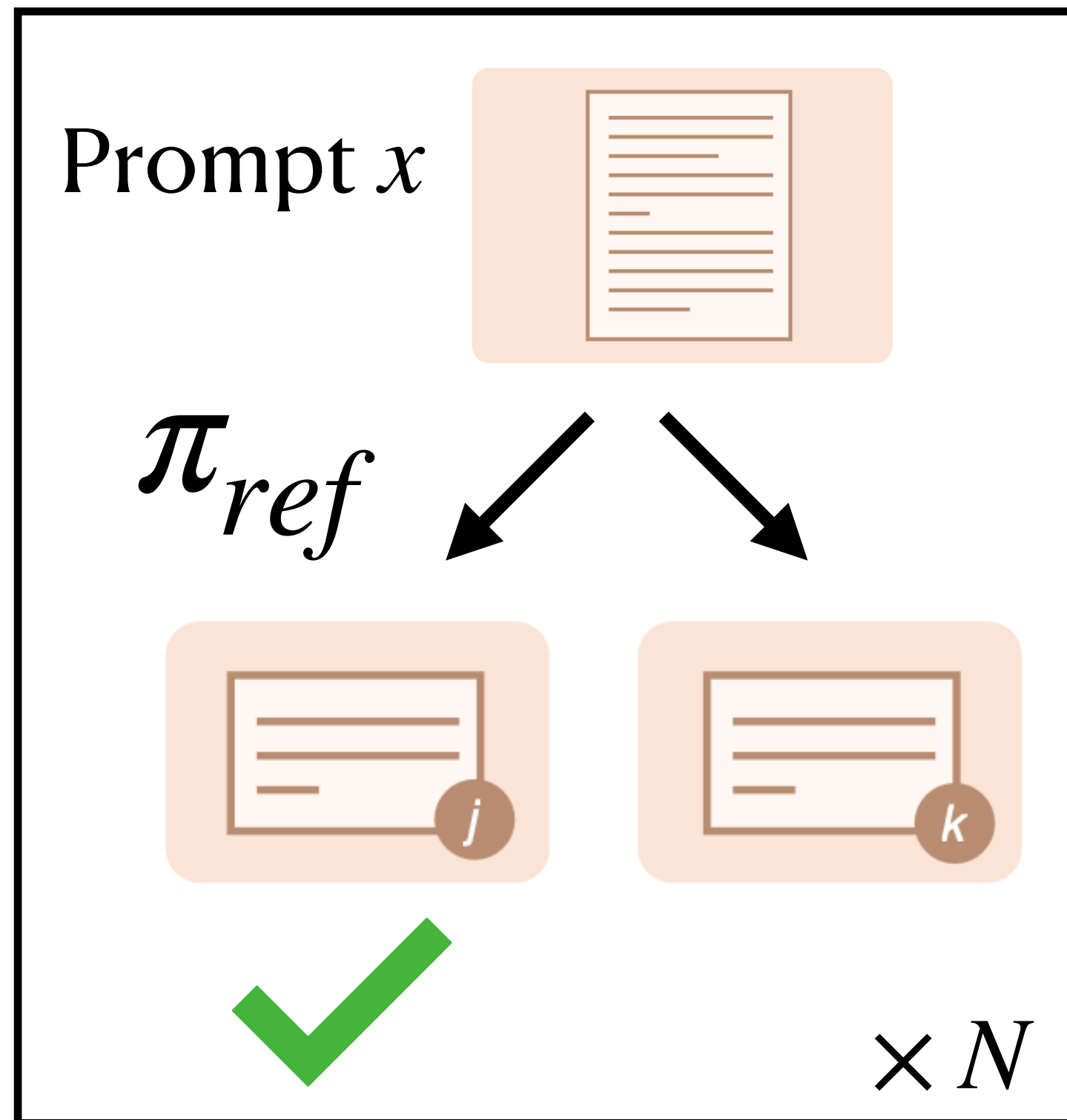
$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

2. Learn a reward model  $\hat{r}$  using the data from step 1

3. train policy via RL (e.g., PPO)

# The post-training Pipeline: RLHF

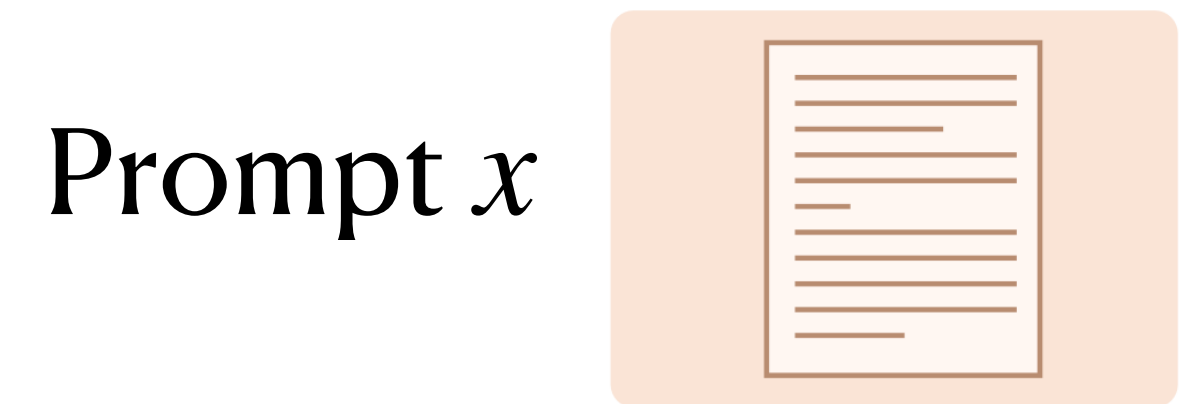
## 1. Collect preference dataset



$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

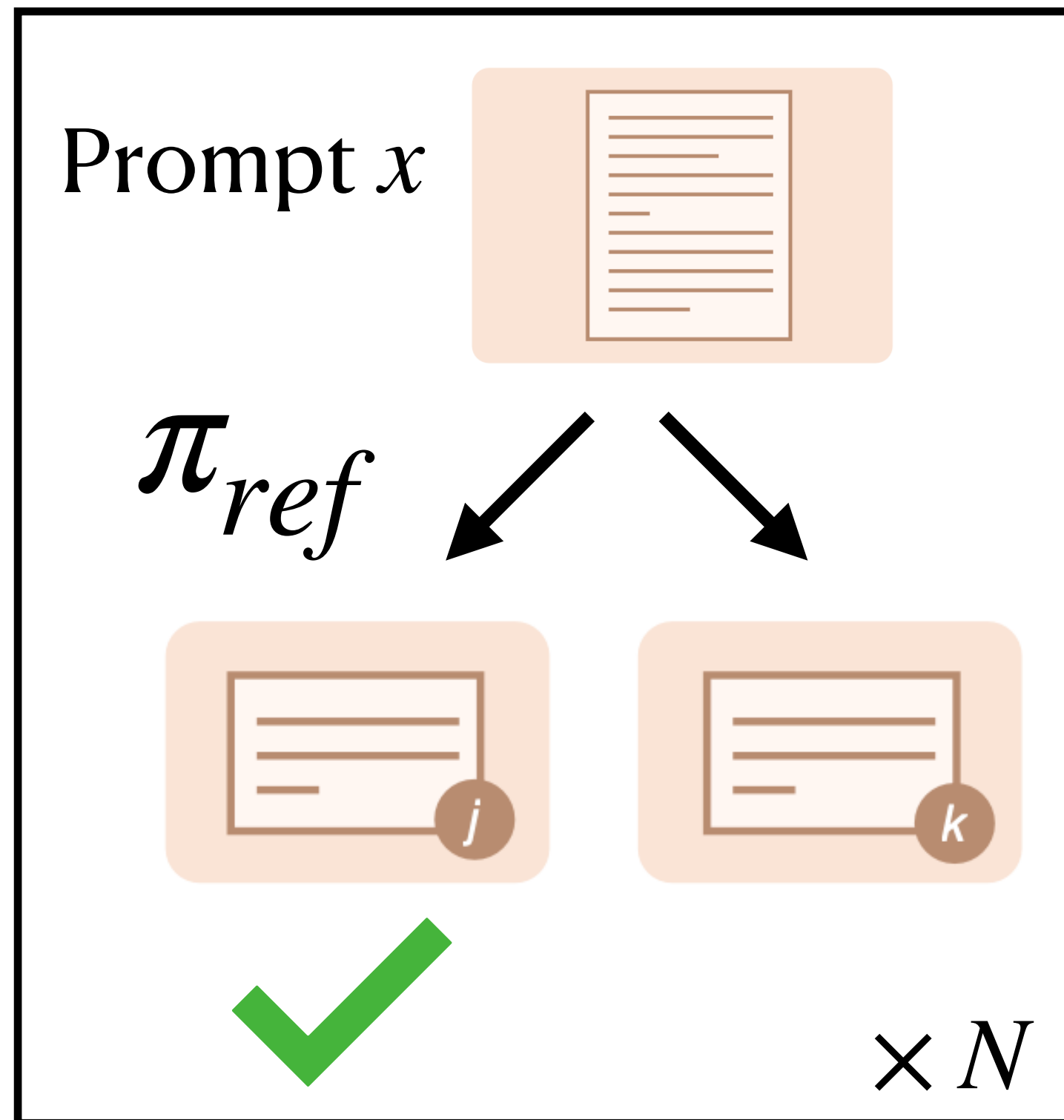
## 2. Learn a reward model $\hat{r}$ using the data from step 1

## 3. train policy via RL (e.g., PPO)



# The post-training Pipeline: RLHF

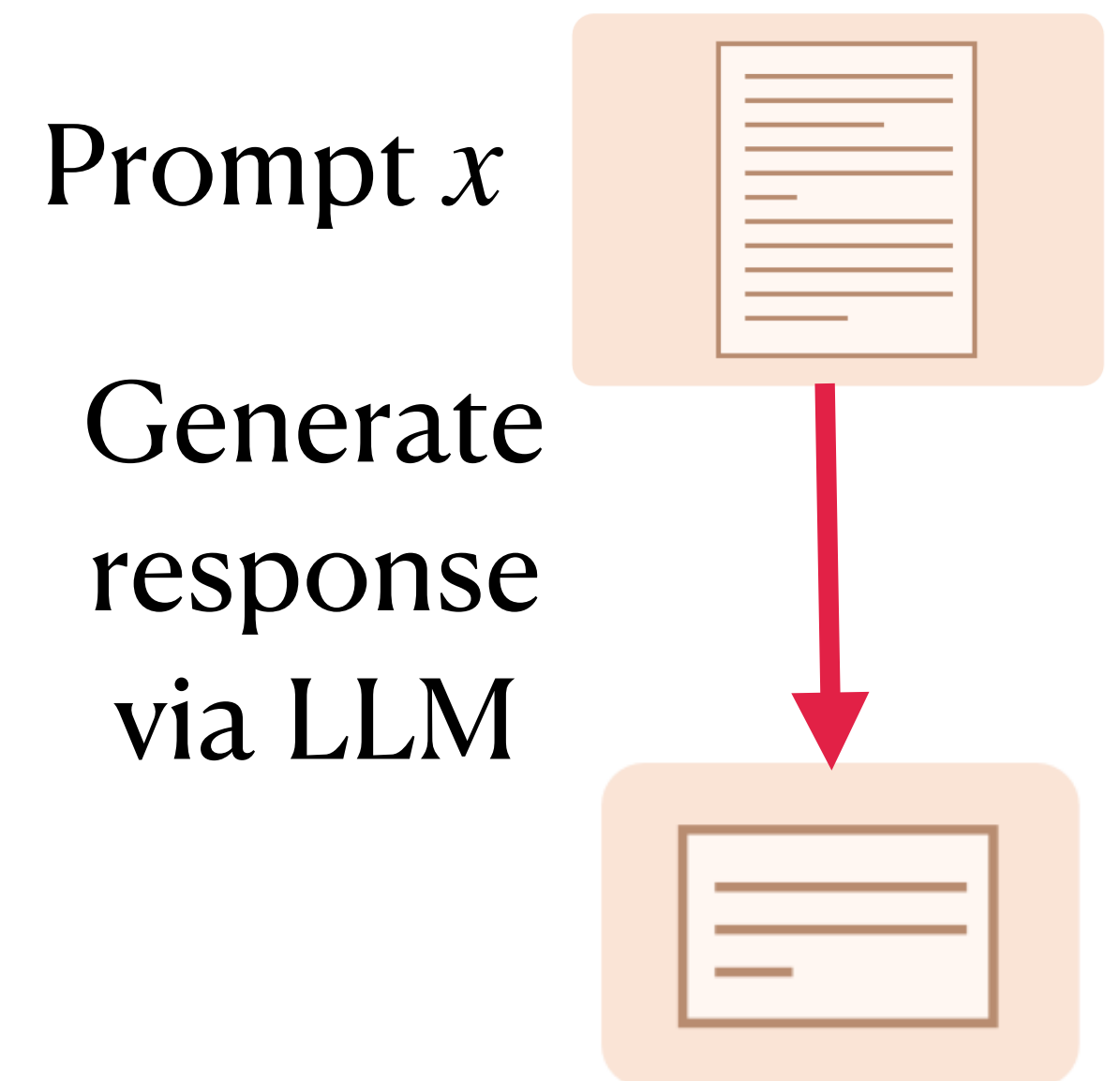
## 1. Collect preference dataset



$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

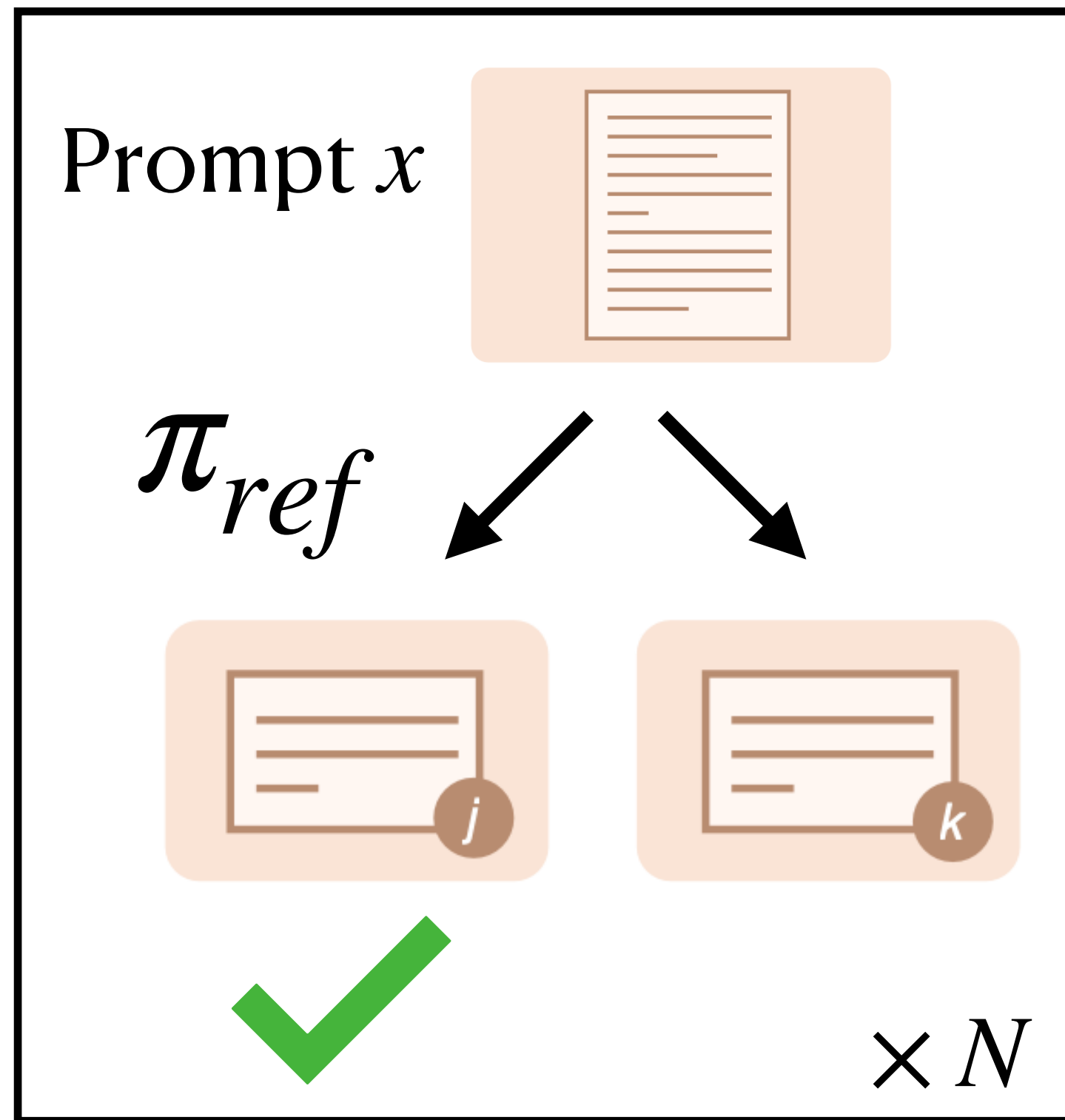
## 2. Learn a reward model $\hat{r}$ using the data from step 1

## 3. train policy via RL (e.g., PPO)



# The post-training Pipeline: RLHF

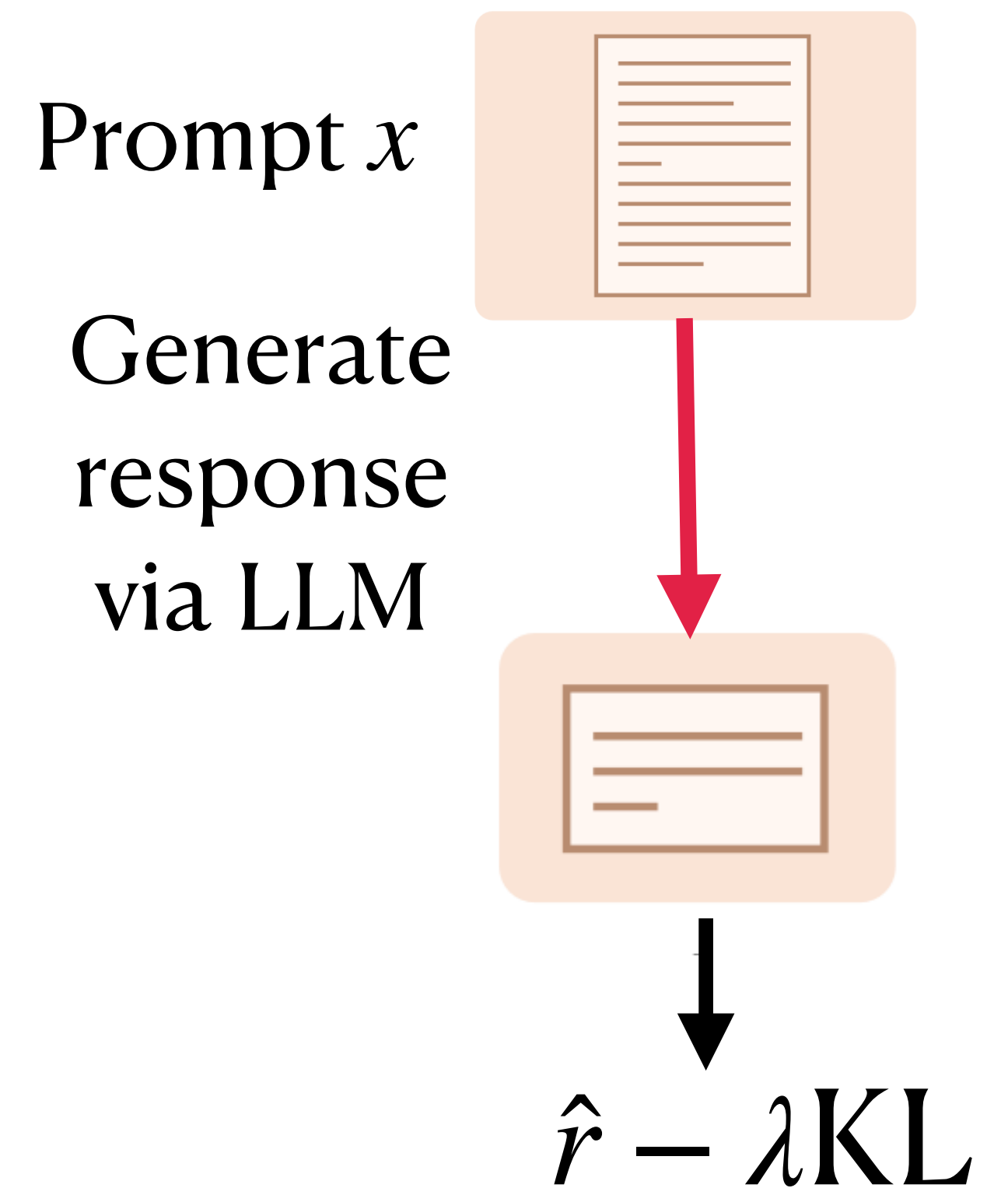
## 1. Collect preference dataset



$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

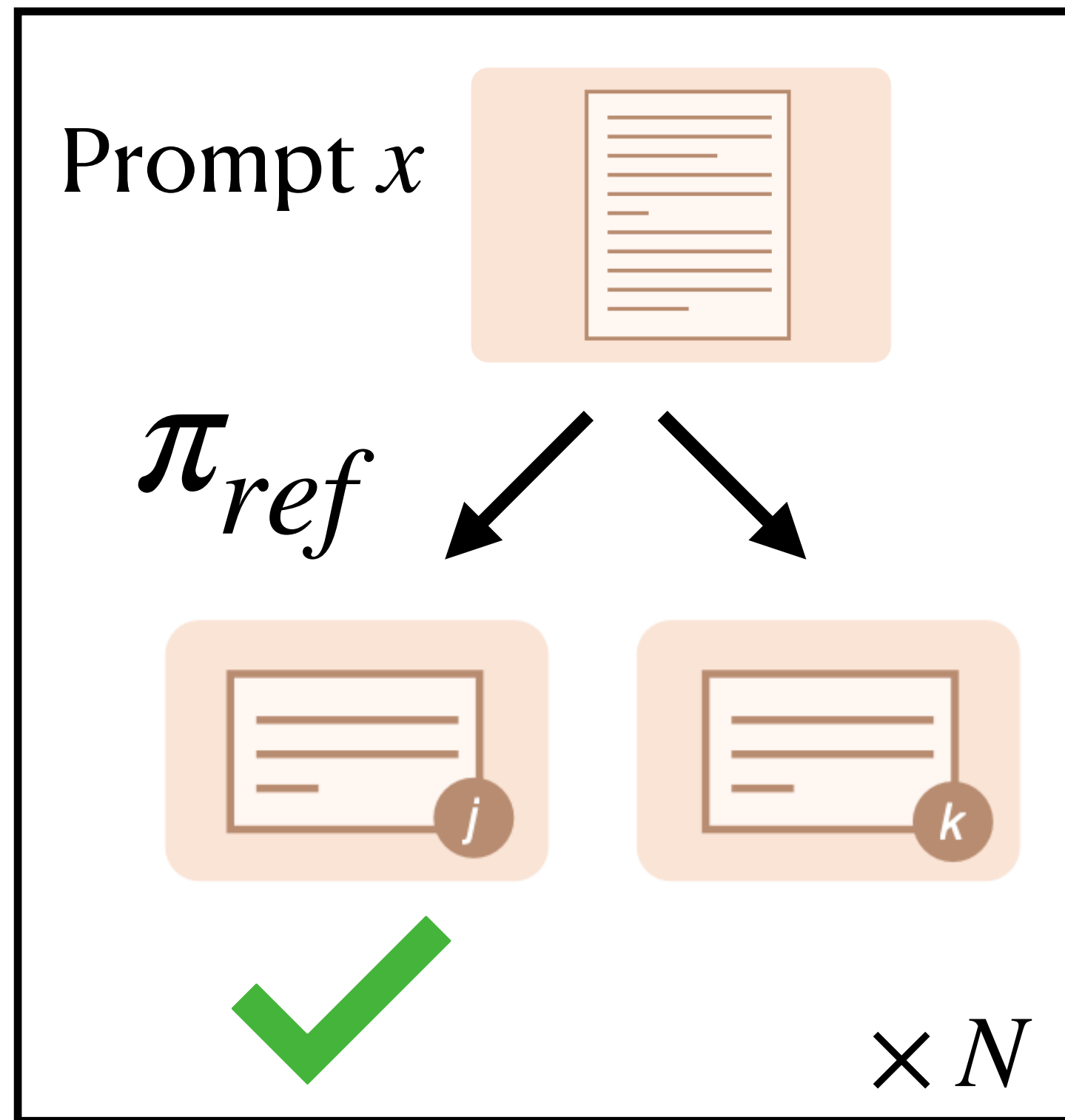
## 2. Learn a reward model $\hat{r}$ using the data from step 1

## 3. train policy via RL (e.g., PPO)



# The post-training Pipeline: RLHF

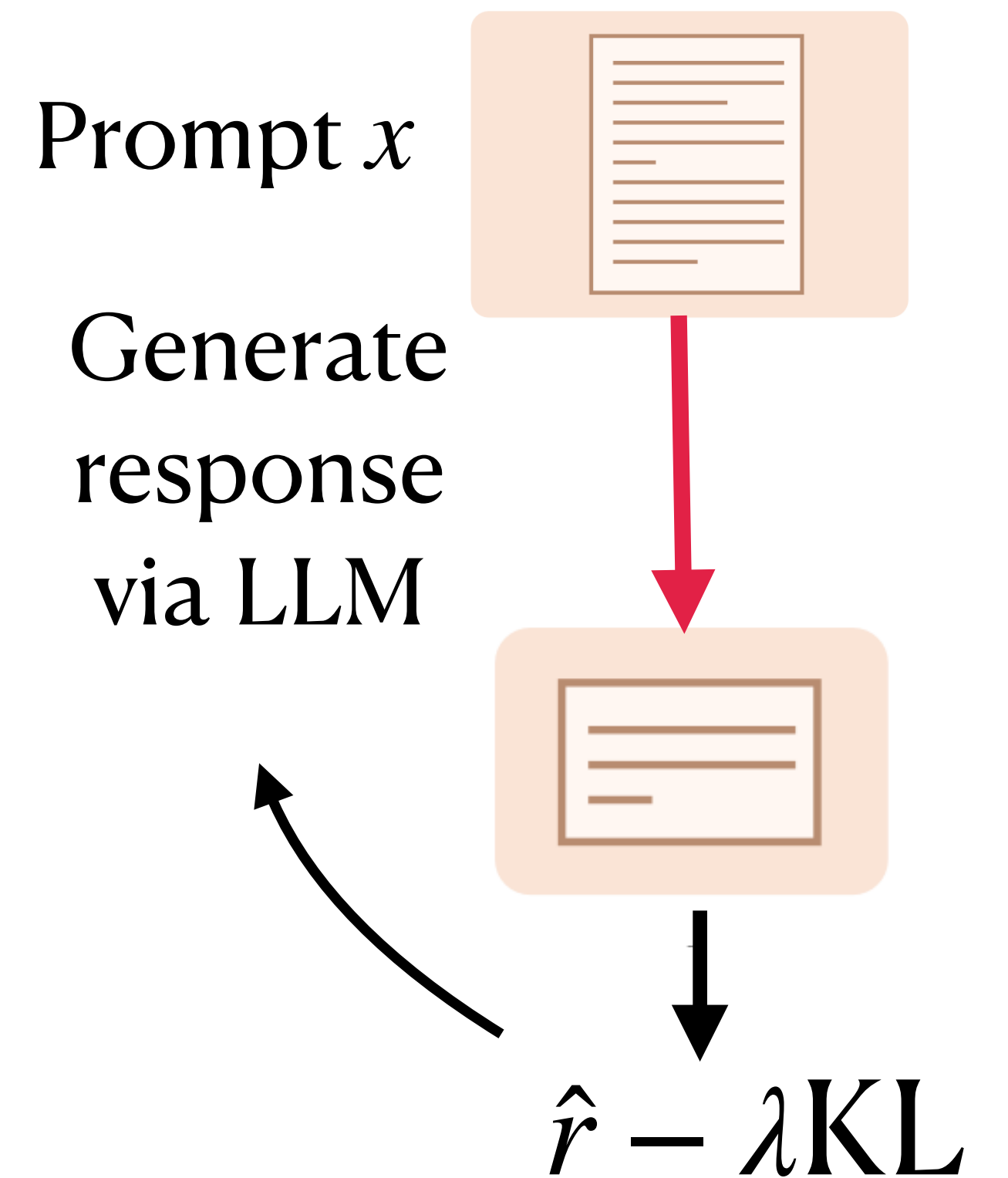
## 1. Collect preference dataset



$$\mathcal{D}_{off} = \{x, \tau, \tau', z\}$$

## 2. Learn a reward model $\hat{r}$ using the data from step 1

## 3. train policy via RL (e.g., PPO)



# What's the benefit of RLHF over SFT?

**Evaluation is often easier than generation**

# What's the benefit of RLHF over SFT?

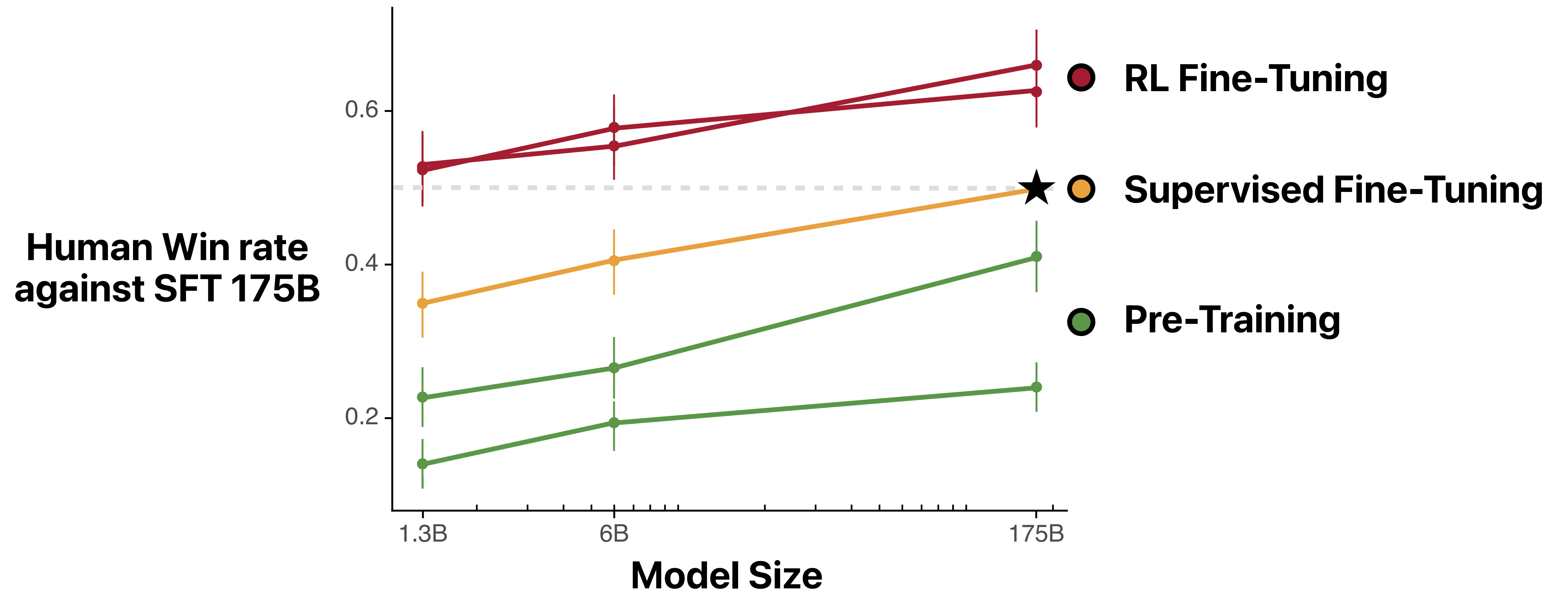
**Evaluation is often easier than generation**

Given a high quality reward, RLHF can often make model outperform humans:

# What's the benefit of RLHF over SFT?

Evaluation is often easier than generation

Given a high quality reward, RLHF can often make model outperform humans:



# The MDP formulation of text generation

**Initial state**  $s_0$ : prompt  $x$

**Action**: token  $y$ ; action space: all possible tokens

**State**: prompt + generated tokens, e.g.,  $s_h = (x, y_0, y_1, \dots, y_{h-1})$

**Transition**: concatenation, i.e., given  $s_h$  and  $y_h$ ,  $s_{h+1} = (s_h, y_h)$

**Terminate**: either hits the maximum content length or hits the special EOS token

# Outline

1. LLM as a policy

2. Learning reward functions from preference data

3. KL-regularized RL

4. DPO

5. REBEL

# Learning reward from human data

Reward design can be challenging in RL

# Bradley-Terry Model

Assume there is a ground truth reward  $r^\star(x, \tau)$  (i.e., high reward means response is good)

# Bradley-Terry Model

Assume there is a ground truth reward  $r^*(x, \tau)$  (i.e., high reward means response is good)

The BT model assumes that humans generate labels based on the following probabilistic model:

$$P(\tau \text{ is preferred over } \tau' \text{ given } x) = \frac{1}{1 + \exp\left(-\left(\frac{r^*(x, \tau) - r^*(x, \tau')}{\Delta(\tau, \tau')}\right)\right)}$$

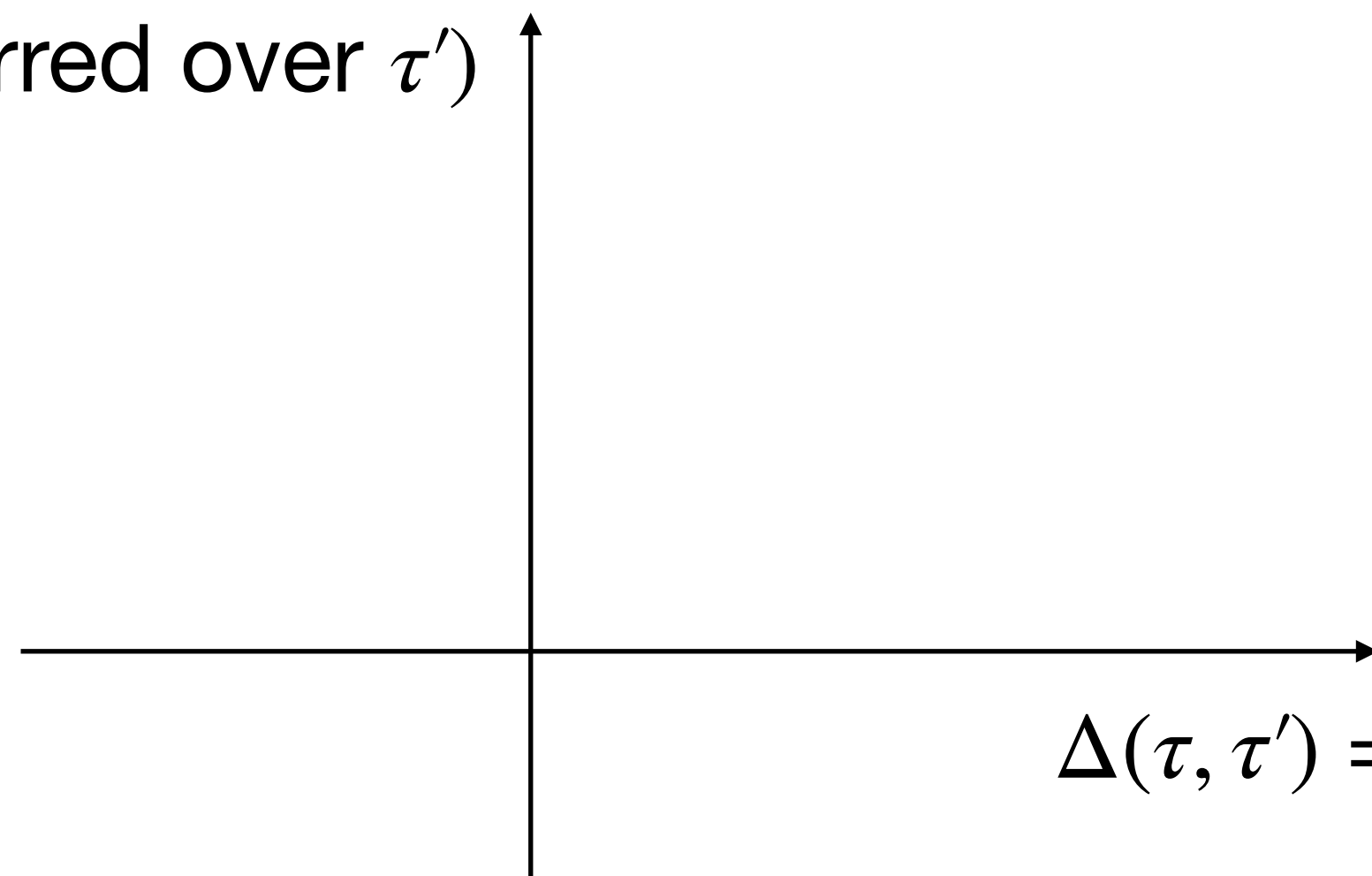
# Bradley-Terry Model

Assume there is a ground truth reward  $r^*(x, \tau)$  (i.e., high reward means response is good)

The BT model assumes that humans generate labels based on the following probabilistic model:

$$P(\tau \text{ is preferred over } \tau' \text{ given } x) = \frac{1}{1 + \exp\left(-\left(\frac{r^*(x, \tau) - r^*(x, \tau')}{\Delta(\tau, \tau')}\right)\right)}$$

$P(\tau \text{ preferred over } \tau')$



$$\Delta(\tau, \tau') = r^*(x, \tau) - r^*(x, \tau')$$

# Learning reward based on the Bradley-Terry assumption

Given a preference dataset  $\mathcal{D} = \{x, \tau, \tau', z\}$ , where label  $z \in \{1, -1\}$  is generated via BT on  $r^\star$   
(1 indicates  $\tau$  is preferred over  $\tau'$ ; -1 otherwise)

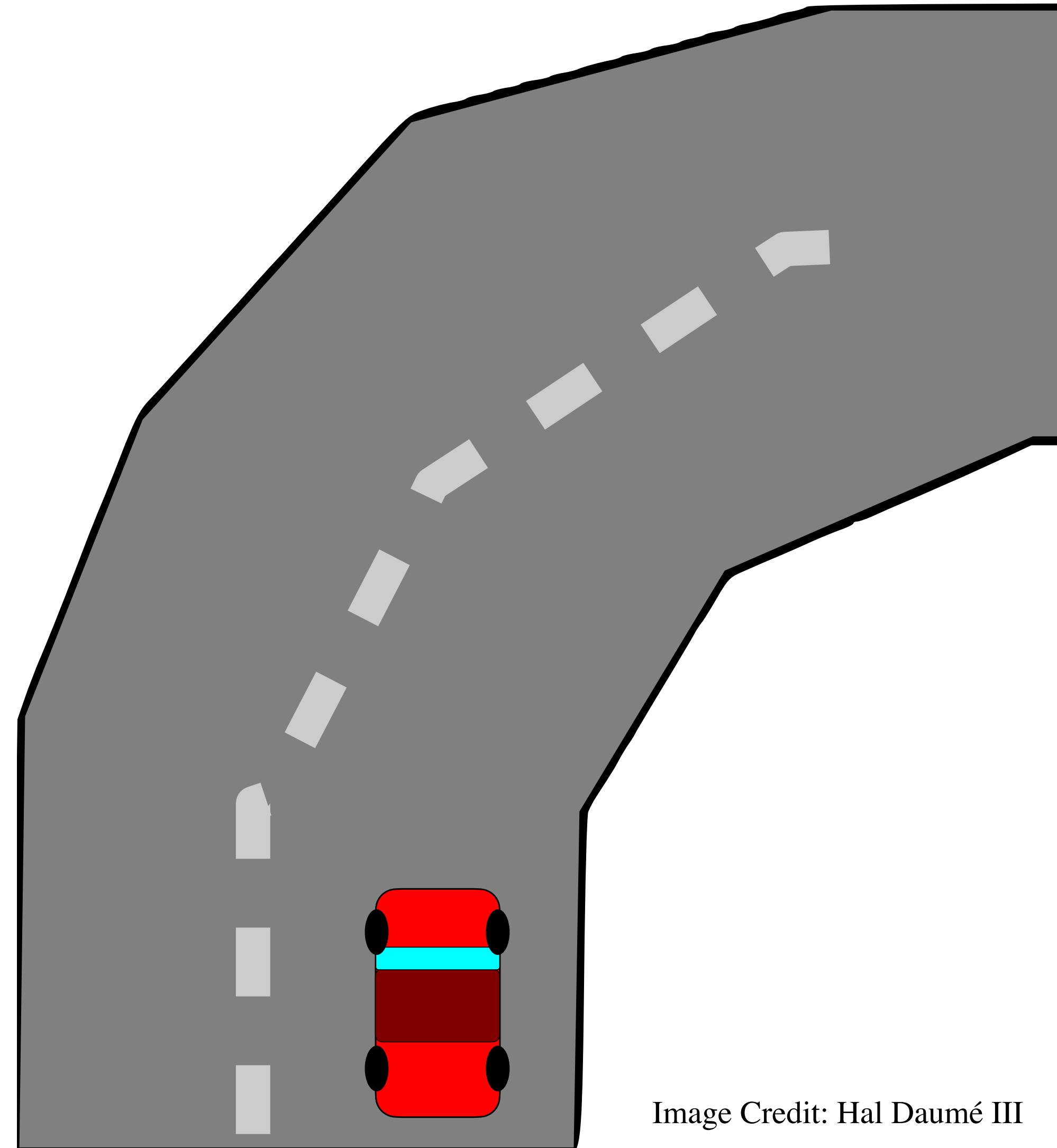
Learn a reward model via MLE:

$$\max_r \sum_{x, \tau, \tau', z} \ln \frac{1}{1 + \exp(-z \cdot (r(x, \tau) - r(x, \tau'))))}$$

# Outline

1. LLM as a policy
2. Learning reward functions from preference data
3. KL-regularized RL
4. DPO
5. REBEL

**RL is very good at reward hacking**



## To avoid reward hacking

We form the following KL regularized RL objective

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

## To avoid reward hacking

We form the following KL regularized RL objective

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

$\beta$  : controls the strength of KL-reg;

## To avoid reward hacking

We form the following KL regularized RL objective

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

$\beta$  : controls the strength of KL-reg;

“stay close” to the SFT policy  $\pi_{ref}$ .

# To avoid reward hacking

We form the following KL regularized RL objective

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

$\beta$  : controls the strength of KL-reg;

“stay close” to the SFT policy  $\pi_{ref}$ .

Q: Why this can help avoid reward hacking?

# How to optimize the KL-reg RL objective

A simple heuristic is to add KL to reward

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$

# How to optimize the KL-reg RL objective

A simple heuristic is to add KL to reward

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$
$$= \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \underbrace{\left( \hat{r}(x, \tau) - \beta \ln \frac{\pi(\tau | x)}{\pi_{ref}(\tau | x)} \right)}_{:= r_{new}(x, \tau)} \right]$$

# How to optimize the KL-reg RL objective

A simple heuristic is to add KL to reward

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$

$$= \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \underbrace{\left( \hat{r}(x, \tau) - \beta \ln \frac{\pi(\tau | x)}{\pi_{ref}(\tau | x)} \right)}_{:= r_{new}(x, \tau)} \right]$$

Run PG (reinforce or PPO) w/  
 $r_{new}(x, \tau)$  as the reward signal

# How to optimize the KL-reg RL objective

A simple heuristic is to add KL to reward

$$J(\pi_\theta) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$

$$= \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \underbrace{\left( \hat{r}(x, \tau) - \beta \ln \frac{\pi(\tau | x)}{\pi_{ref}(\tau | x)} \right)}_{:= r_{new}(x, \tau)} \right]$$

Run PG (reinforce or PPO) w/  
 $r_{new}(x, \tau)$  as the reward signal

***Remark: it works, but it  
is not the exact gradient***

# Outline

1. LLM as a policy
2. Learning reward functions from preference data
3. KL-regularized RL
4. DPO
5. REBEL

# When models are large...

RM + PPO can be hard to optimize...

At least need to maintain 4 big models in GPU RAM (RM,  $\pi$ , V,  $\pi_{ref}$ ...)

# KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

What's the  $\arg \max_{\pi} J(\pi)$  ?

---

## KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

What's the  $\arg \max_{\pi} J(\pi)$  ?

---

Consider on a  $(x, \tau)$  pair, what is  $\partial J(\pi) / \partial \pi(\tau | x)$  ?

## KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$

What's the  $\arg \max_{\pi} J(\pi)$  ?

---

Consider on a  $(x, \tau)$  pair, what is  $\partial J(\pi) / \partial \pi(\tau | x)$  ?

$$\frac{\partial J(\pi)}{\partial \pi(\tau | x)} = \hat{r}(x, \tau) - \beta \left( \ln \pi(\tau | x) - \ln \pi_{ref}(\tau | x) + 1 \right)$$

# KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \mid \pi_{ref}(\cdot | x) \right) \right]$$

What's the  $\arg \max_{\pi} J(\pi)$  ?

---

Consider on a  $(x, \tau)$  pair, what is  $\partial J(\pi) / \partial \pi(\tau | x)$  ?

$$\frac{\partial J(\pi)}{\partial \pi(\tau | x)} = \hat{r}(x, \tau) - \beta \left( \ln \pi(\tau | x) - \ln \pi_{ref}(\tau | x) + 1 \right)$$

$$\pi(\tau | x) \propto \pi_{ref}(\tau | x) \exp \left( \hat{r}(x, \tau) / \beta \right)$$

## KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot | x)} \hat{r}(x, \tau) - \beta \text{KL} \left( \pi(\cdot | x) \middle| \pi_{ref}(\cdot | x) \right) \right]$$

What's the  $\arg \max_{\pi} J(\pi)$  ?

---

Consider on a  $(x, \tau)$  pair, what is  $\partial J(\pi) / \partial \pi(\tau | x)$  ?

$$\frac{\partial J(\pi)}{\partial \pi(\tau | x)} = \hat{r}(x, \tau) - \beta \left( \ln \pi(\tau | x) - \ln \pi_{ref}(\tau | x) + 1 \right)$$

$$\pi(\tau | x) \propto \pi_{ref}(\tau | x) \exp(\hat{r}(x, \tau) / \beta)$$

$$\pi(\tau | x) = \pi_{ref}(\tau | x) \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right) / Z(x), \text{ where } Z(x) = \mathbb{E}_{\tau \sim \pi_{ref}(\cdot | x)} \exp(\hat{r}(x, \tau) / \beta)$$

# Can we parameterize RM using policies?

In sum, the optimal policy of the KL-reg RL objective is:

$$\hat{\pi}(\tau | x) = \frac{\pi_{ref}(\tau | x) \cdot \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right)}{Z(x)}$$

---

# Can we parameterize RM using policies?

In sum, the optimal policy of the KL-reg RL objective is:

$$\hat{\pi}(\tau | x) = \frac{\pi_{ref}(\tau | x) \cdot \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right)}{Z(x)}$$

---

$$\ln \hat{\pi}(\tau | x) = \ln \pi_{ref}(\tau | x) - \ln Z(x) + \frac{\hat{r}(x, \tau)}{\beta}$$

# Can we parameterize RM using policies?

In sum, the optimal policy of the KL-reg RL objective is:

$$\hat{\pi}(\tau | x) = \frac{\pi_{ref}(\tau | x) \cdot \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right)}{Z(x)}$$

---

$$\ln \hat{\pi}(\tau | x) = \ln \pi_{ref}(\tau | x) - \ln Z(x) + \frac{\hat{r}(x, \tau)}{\beta}$$

$$\hat{r}(x, \tau) = \beta \left( \ln \frac{\hat{\pi}(\tau | x)}{\pi_{ref}(\tau | x)} + \ln Z(x) \right)$$

# Can we parameterize RM using policies?

In sum, the optimal policy of the KL-reg RL objective is:

$$\hat{\pi}(\tau | x) = \frac{\pi_{ref}(\tau | x) \cdot \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right)}{Z(x)}$$

---

$$\ln \hat{\pi}(\tau | x) = \ln \pi_{ref}(\tau | x) - \ln Z(x) + \frac{\hat{r}(x, \tau)}{\beta}$$

$$\hat{r}(x, \tau) = \beta \left( \ln \frac{\hat{\pi}(\tau | x)}{\pi_{ref}(\tau | x)} + \ln Z(x) \right) \quad \text{Not done yet, this } Z(x) \text{ technically contains } \hat{r}!$$

# Can we parameterize RM using policies?

In sum, the optimal policy of the KL-reg RL objective is:

$$\hat{\pi}(\tau | x) = \frac{\pi_{ref}(\tau | x) \cdot \exp\left(\frac{\hat{r}(x, \tau)}{\beta}\right)}{Z(x)}$$

---

$$\ln \hat{\pi}(\tau | x) = \ln \pi_{ref}(\tau | x) - \ln Z(x) + \frac{\hat{r}(x, \tau)}{\beta}$$

$$\hat{r}(x, \tau) = \beta \left( \ln \frac{\hat{\pi}(\tau | x)}{\pi_{ref}(\tau | x)} + \ln Z(x) \right)$$

**Not done yet, this  $Z(x)$  technically contains  $\hat{r}$ !**  
**But  $\ln Z(x)$  is a shift that is independent of  $\tau$ ...**

# DPO

1. Take any policy  $\pi_\theta$ , we can use it to model the reward difference:

$$r_\theta(\tau | x) - r_\theta(\tau' | x) := \beta \left( \ln \frac{\pi_\theta(\tau | x)}{\pi_{ref}(\tau | x)} - \ln \frac{\pi_\theta(\tau' | x)}{\pi_{ref}(\tau' | x)} \right)$$

# DPO

1. Take any policy  $\pi_\theta$ , we can use it to model the reward difference:

$$r_\theta(\tau | x) - r_\theta(\tau' | x) := \beta \left( \ln \frac{\pi_\theta(\tau | x)}{\pi_{ref}(\tau | x)} - \ln \frac{\pi_\theta(\tau' | x)}{\pi_{ref}(\tau' | x)} \right)$$

2. Now plug this into the MLE loss we had for learning the reward difference:

# DPO

1. Take any policy  $\pi_\theta$ , we can use it to model the reward difference:

$$r_\theta(\tau | x) - r_\theta(\tau' | x) := \beta \left( \ln \frac{\pi_\theta(\tau | x)}{\pi_{ref}(\tau | x)} - \ln \frac{\pi_\theta(\tau' | x)}{\pi_{ref}(\tau' | x)} \right)$$

2. Now plug this into the MLE loss we had for learning the reward difference:

$$\arg \max_{\theta} \sum_{x, \tau, \tau', z} \ln \frac{1}{1 + \exp \left( -z \cdot (r_\theta(x, \tau) - r_\theta(x, \tau')) \right)}$$

# DPO

DPO optimizes policy  $\pi_\theta$  directly using the following loss:

$$\arg \max_{\theta} \sum_{x, \tau, \tau', z} \ln \frac{1}{1 + \exp \left( -z \cdot \beta \left( \ln \frac{\pi_\theta(\tau | x)}{\pi_{ref}(\tau | x)} - \ln \frac{\pi_\theta(\tau' | x)}{\pi_{ref}(\tau' | x)} \right) \right)}$$

# Outline

1. LLM as a policy
2. Learning reward functions from preference data
3. KL-regularized RL
4. DPO
5. REBEL

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau | x) = \pi_t(\tau | x) \exp(r(x, \tau) / \beta) / Z(x)$$

---

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau | x) = \pi_t(\tau | x) \exp(r(x, \tau) / \beta) / Z(x)$$

---

1. Take log on both sides and rearrange terms, we get

$$r(x, \tau) = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} + \ln Z(x) \right)$$

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau | x) = \pi_t(\tau | x) \exp(r(x, \tau) / \beta) / Z(x)$$

---

1. Take log on both sides and rearrange terms, we get

$$r(x, \tau) = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} + \ln Z(x) \right)$$

2. Instead of modeling reward, we model reward difference to cancel  $Z(x)$ :

$$r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

This indicates  $\pi_{t+1}$  is **the minimizer** of the following least square regression problem:

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

This indicates  $\pi_{t+1}$  is **the minimizer** of the following least square regression problem:

$$\mathbb{E}_{x, \tau, \tau'} \left( \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right) - (r(x, \tau) - r(x, \tau')) \right)^2$$

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

This indicates  $\pi_{t+1}$  is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$  should be the minimizer  
regardless of what the  
distribution is;

$$\mathbb{E}_{x, \tau, \tau'} \left( \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right) - (r(x, \tau) - r(x, \tau')) \right)^2$$

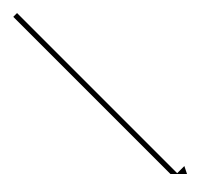
# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

This indicates  $\pi_{t+1}$  is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$  should be the minimizer  
regardless of what the  
distribution is;


$$\mathbb{E}_{x, \tau, \tau'} \left( \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right) - (r(x, \tau) - r(x, \tau')) \right)^2$$

In practice, we often use

$$x, \tau \sim \pi_t(\cdot | x), \tau' \sim \pi_t(\cdot | x)$$

# Reparameterization

We obtained the following relationship between  $r$  and  $\pi_{t+1}$  &  $\pi_t$ :

$$\forall(x, \tau, \tau') : r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right)$$

This indicates  $\pi_{t+1}$  is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$  should be the minimizer  
regardless of what the  
distribution is;

$$\mathbb{E}_{x, \tau, \tau'} \left( \beta \left( \ln \frac{\pi_{t+1}(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi_{t+1}(\tau' | x)}{\pi_t(\tau' | x)} \right) - \frac{(r(x, \tau) - r(x, \tau'))}{\text{Relative reward}} \right)^2$$

In practice, we often use

$$x, \tau \sim \pi_t(\cdot | x), \tau' \sim \pi_t(\cdot | x)$$

Relative reward

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given  $\pi_t$ , we compute  $\pi_{t+1}$  via least square regression:

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given  $\pi_t$ , we compute  $\pi_{t+1}$  via least square regression:

$$\pi_{t+1} = \arg \min_{\pi} \mathbb{E}_{x, (\tau, \tau') \sim \pi_t(\cdot|x)} \left( \beta \left( \ln \frac{\pi(\tau|x)}{\pi_t(\tau|x)} - \ln \frac{\pi(\tau'|x)}{\pi_t(\tau'|x)} \right) - (r(x, \tau) - r(x, \tau')) \right)^2$$

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given  $\pi_t$ , we compute  $\pi_{t+1}$  via least square regression:

$$\pi_{t+1} = \arg \min_{\pi} \mathbb{E}_{x, (\tau, \tau') \sim \pi_t(\cdot | x)} \left( \beta \left( \ln \frac{\pi(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi(\tau' | x)}{\pi_t(\tau' | x)} \right) - (r(x, \tau) - r(x, \tau')) \right)^2$$

sample  $\tau, \tau'$  from the latest policy  
 $\pi_t(\cdot | x)$ , independently;

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given  $\pi_t$ , we compute  $\pi_{t+1}$  via least square regression:

$$\pi_{t+1} = \arg \min_{\pi} \mathbb{E}_{x, (\tau, \tau') \sim \pi_t(\cdot | x)} \left( \beta \left( \ln \frac{\pi(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi(\tau' | x)}{\pi_t(\tau' | x)} \right) - \frac{(r(x, \tau) - r(x, \tau'))}{\text{Relative reward}} \right)^2$$

sample  $\tau, \tau'$  from the latest policy  
 $\pi_t(\cdot | x)$ , independently;

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given  $\pi_t$ , we compute  $\pi_{t+1}$  via least square regression:

$$\pi_{t+1} = \arg \min_{\pi} \mathbb{E}_{x, (\tau, \tau') \sim \pi_t(\cdot | x)} \left( \underbrace{\beta \left( \ln \frac{\pi(\tau | x)}{\pi_t(\tau | x)} - \ln \frac{\pi(\tau' | x)}{\pi_t(\tau' | x)} \right)}_{\text{Regressor}} - \underbrace{(r(x, \tau) - r(x, \tau'))}_{\text{Relative reward}} \right)^2$$

sample  $\tau, \tau'$  from the latest policy  
 $\pi_t(\cdot | x)$ , independently;

# REBEL Theoretical Analysis

**Assumption 1:** over  $T$  iterations, we have the following for some  $\epsilon$ :

$$\mathbb{E}_{x, y \sim \pi_t(\cdot|x), y' \sim \pi_t(\cdot|x)} \left( \frac{1}{\eta} \left( \ln \frac{\pi_{t+1}(y|x)}{\pi_t(y|x)} - \ln \frac{\pi_{t+1}(y'|x)}{\pi_t(y'|x)} \right) - (r(x, y) - r(x, y')) \right)^2 \leq \epsilon$$

# REBEL Theoretical Analysis

**Assumption 1:** over  $T$  iterations, we have the following for some  $\epsilon$ :

$$\mathbb{E}_{x, y \sim \pi_t(\cdot|x), y' \sim \pi_t(\cdot|x)} \left( \frac{1}{\eta} \left( \ln \frac{\pi_{t+1}(y|x)}{\pi_t(y|x)} - \ln \frac{\pi_{t+1}(y'|x)}{\pi_t(y'|x)} \right) - (r(x, y) - r(x, y')) \right)^2 \leq \epsilon$$

**Data Coverage:** Given a test policy  $\pi$ , we denote the concentrability coefficient as:

$$C_{\pi_t \rightarrow \pi} = \max_{x, y} \frac{\pi(y|x)}{\pi_t(y|x)}$$

$\pi_t$  covers  $\pi$  if  $C_{\pi_t \rightarrow \pi} < +\infty$

# REBEL Theorem

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

$$C_{\max} = \max_{\pi \in \{\pi_1, \dots, \pi_T\}} C_{\pi \rightarrow \pi^*}$$

# Theorem – Lemma 1

**Assume**  $\max_{x,y,t} |A_t(x, y)| \leq A \in \mathbb{R}^+$ , and  $\pi_0(\cdot | x)$  is uniform over  $\mathcal{Y}$ .

**Then,  $\eta = \sqrt{\ln(|\mathcal{Y}|)/(A^2 T)}$ , for the sequence of policies computed by REBEL, we have:**

$$\forall \pi, x : \sum_{t=0}^{T-1} \mathbb{E}_{y \sim \pi(\cdot | x)} A_t(x, y) \leq 2A \sqrt{\ln(|\mathcal{Y}|) T}.$$

# Theorem – Proof Lemma 1

**Start with  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
where  $Z_t(x)$  is the normalization constant**

# Theorem – Proof Lemma 1

**Start with  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
where  $Z_t(x)$  is the normalization constant**

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

# Theorem – Proof Lemma 1

**Start with  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
where  $Z_t(x)$  is the normalization constant**

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

$$\log \pi_{t+1}(y | x) = \log \pi_t(y | x) + \eta A_t(x, y) - \log Z_t(x),$$

## Theorem – Proof Lemma 1

**Start with**  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
**where**  $Z_t(x)$  **is the normalization constant**

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

$$\log \pi_{t+1}(y | x) = \log \pi_t(y | x) + \eta A_t(x, y) - \log Z_t(x),$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

# Theorem – Proof Lemma 1

**Start with**  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
**where**  $Z_t(x)$  **is the normalization constant**

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

$$\log \pi_{t+1}(y | x) = \log \pi_t(y | x) + \eta A_t(x, y) - \log Z_t(x),$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \log \pi(y | x) - \log \pi(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

# Theorem – Proof Lemma 1

**Start with**  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
**where**  $Z_t(x)$  **is the normalization constant**

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

$$\log \pi_{t+1}(y | x) = \log \pi_t(y | x) + \eta A_t(x, y) - \log Z_t(x),$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \log \pi(y | x) - \log \pi(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

# Theorem – Proof Lemma 1

Start with  $\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x)$ ,  
where  $Z_t(x)$  is the normalization constant

$$\pi_{t+1}(y | x) = \pi_t(y | x)\exp(\eta A_t(x, y))/Z_t(x),$$

$$\log \pi_{t+1}(y | x) = \log \pi_t(y | x) + \eta A_t(x, y) - \log Z_t(x),$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

$$E_{y \sim \pi} \left[ \log \pi_{t+1}(y | x) \right] = E_{y \sim \pi} \left[ \log \pi_t(y | x) \right] + E_{y \sim \pi} \left[ \log \pi(y | x) - \log \pi(y | x) \right] + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} \left[ \eta A_t(x, y) \right] - E_{y \sim \pi} \left[ \log Z_t(x) \right]$$

term

# Theorem – Proof Lemma 1

$$\log Z_t(x)$$

$$\log Z_t(x) = \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} \exp(\eta A_t(x, y))$$

# Theorem – Proof Lemma 1

$$\log Z_t(x)$$

$$\log Z_t(x) = \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} \exp(\eta A_t(x, y))$$

**Using  $\exp(x) \leq 1 + x + x^2$  for any  $x \leq 1$ :**

# Theorem – Proof Lemma 1

$$\log Z_t(x)$$

$$\log Z_t(x) = \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} \exp(\eta A_t(x, y))$$

**Using  $\exp(x) \leq 1 + x + x^2$  for any  $x \leq 1$ :**

$$\log Z_t(x) \leq \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} (1 + \eta A_t(x, y) + \eta^2 A_t(x, y)^2)$$

# Theorem – Proof Lemma 1

$$\log Z_t(x)$$

$$\log Z_t(x) = \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} \exp(\eta A_t(x, y))$$

**Using  $\exp(x) \leq 1 + x + x^2$  for any  $x \leq 1$ :**

$$\begin{aligned} \log Z_t(x) &\leq \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} (1 + \eta A_t(x, y) + \eta^2 A_t(x, y)^2) \\ &\leq \ln(1 + 0 + \eta^2 A^2) \end{aligned}$$

# Theorem – Proof Lemma 1

$$\log Z_t(x)$$

$$\log Z_t(x) = \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} \exp(\eta A_t(x, y))$$

**Using  $\exp(x) \leq 1 + x + x^2$  for any  $x \leq 1$ :**

$$\log Z_t(x) \leq \ln \mathbb{E}_{y \sim \pi_t(\cdot|x)} (1 + \eta A_t(x, y) + \eta^2 A_t(x, y)^2)$$

$$\leq \ln(1 + 0 + \eta^2 A^2)$$

$$\leq \eta^2 A^2$$

## Theorem – Proof Lemma 1

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} [\eta A_t(x, y)] - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] = D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) - E_{y \sim \pi} [\log Z_t(x)]$$

## Theorem – Proof Lemma 1

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} [\eta A_t(x, y)] - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] = D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \eta^2 A^2$$

## Theorem – Proof Lemma 1

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} [\eta A_t(x, y)] - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] = D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \eta^2 A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_{t+1}) - \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \frac{1}{\eta} \sum_{t=0}^{T-1} \eta^2 A^2$$

## Theorem – Proof Lemma 1

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} [\eta A_t(x, y)] - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] = D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \eta^2 A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_{t+1}) - \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \frac{1}{\eta} \sum_{t=0}^{T-1} \eta^2 A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_0) - D_{\text{KL}}(\pi || \pi_T) - T\eta A^2$$

## Theorem – Proof Lemma 1

$$-D_{\text{KL}}(\pi || \pi_{t+1}) = -D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} [\eta A_t(x, y)] - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] = D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) - E_{y \sim \pi} [\log Z_t(x)]$$

$$E_{y \sim \pi} [\eta A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_{t+1}) - D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \eta^2 A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_{t+1}) - \frac{1}{\eta} \sum_{t=0}^{T-1} D_{\text{KL}}(\pi || \pi_t) + E_{y \sim \pi} - \frac{1}{\eta} \sum_{t=0}^{T-1} \eta^2 A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_0) - D_{\text{KL}}(\pi || \pi_T) - T\eta A^2$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq D_{\text{KL}}(\pi || \pi_0) - T\eta A^2 \leq \frac{\ln |y|}{\eta} + T\eta A^2$$

# Theorem – Proof Lemma 1

$$\text{Set } \eta = \sqrt{\frac{\ln(|y|)}{A^2 T}}$$

$$\sum_{t=0}^{T-1} E_{y \sim \pi} [A_t(x, y)] \leq 2A \sqrt{\ln(|y|) T}$$

# Theorem – Proof

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

## Theorem – Proof

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

**Consider a comparator policy  $\pi^*$ . We start with the performance difference between  $\pi^*$  and the uniform mixture policy.**

# Theorem – Proof

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

**Consider a comparator policy  $\pi^*$ . We start with the performance difference between  $\pi^*$  and the uniform mixture policy.**

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y))$$

# Theorem – Proof

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

**Consider a comparator policy  $\pi^*$ . We start with the performance difference between  $\pi^*$  and the uniform mixture policy.**

$$\begin{aligned} & \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y)) \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A_t(x, y)) + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y) - A_t(x, y)) \end{aligned}$$

# Theorem – Proof

**Under assumption 1, after  $T$  many iterations, among the learned policies  $\pi_1, \dots, \pi_T$ , there must exist a policy  $\pi_t$  such that:**

$$\forall \pi^* : \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} r(x, y) - \mathbb{E}_{x, y \sim \pi_t(\cdot|x)} r(x, y) \leq O\left(\sqrt{\frac{1}{T}} + \sqrt{C_{\max} \epsilon}\right)$$

**Consider a comparator policy  $\pi^*$ . We start with the performance difference between  $\pi^*$  and the uniform mixture policy.**

$$\begin{aligned} & \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y)) \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A_t(x, y)) + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{x, y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y) - A_t(x, y)) \\ &\leq 2A \sqrt{\frac{\ln(|\mathcal{Y}|)}{T}} + \frac{1}{T} \sum_{t=0}^{T-1} \sqrt{\mathbb{E}_x \mathbb{E}_{y \sim \pi^*(\cdot|x)} (A^{\pi_t}(x, y) - A_t(x, y))^2}. \end{aligned}$$

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

Reward Model (RM) is learned from human feedback (i.e., pair-wise preference)

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

Reward Model (RM) is learned from human feedback (i.e., pair-wise preference)

RM learning is based on the Bradley-Terry model

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

Reward Model (RM) is learned from human feedback (i.e., pair-wise preference)

RM learning is based on the Bradley-Terry model

KL regularization is important to avoid hacking the learned RM

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

Reward Model (RM) is learned from human feedback (i.e., pair-wise preference)

RM learning is based on the Bradley-Terry model

KL regularization is important to avoid hacking the learned RM

DPO reparameterizes the reward difference via BT

# Summary

RLHF is a tool for post-training LLMs so that llms can understand and follow human instructions

Reward Model (RM) is learned from human feedback (i.e., pair-wise preference)

RM learning is based on the Bradley-Terry model

KL regularization is important to avoid hacking the learned RM

DPO reparameterizes the reward difference via BT

REBEL reparameterizes the reward difference without BT