

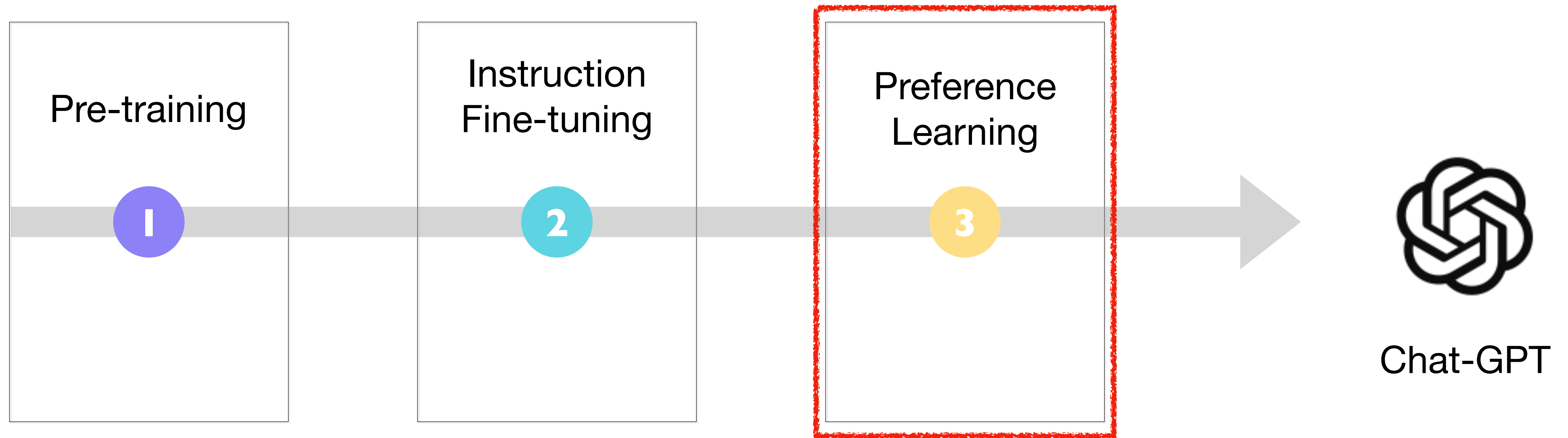


# Iterative Preference Learning for Large Language Model Post Training

Wei Xiong

University of Illinois Urbana-Champaign

# LLM training pipeline



# Outline

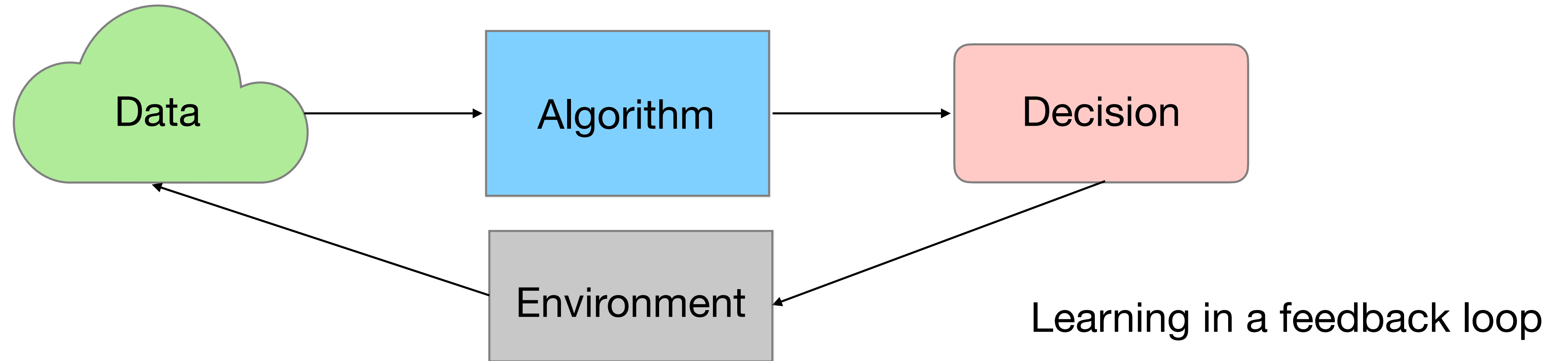
- Motivation: Preference Learning as **Sequential Decision Making**
- Introduction to Reinforcement Learning from Human Feedback (RLHF)
- Main Results: Online Iterative RLHF Framework
- Practical and Open-source Codebook: RLHFlow

# Supervised learning vs decision making



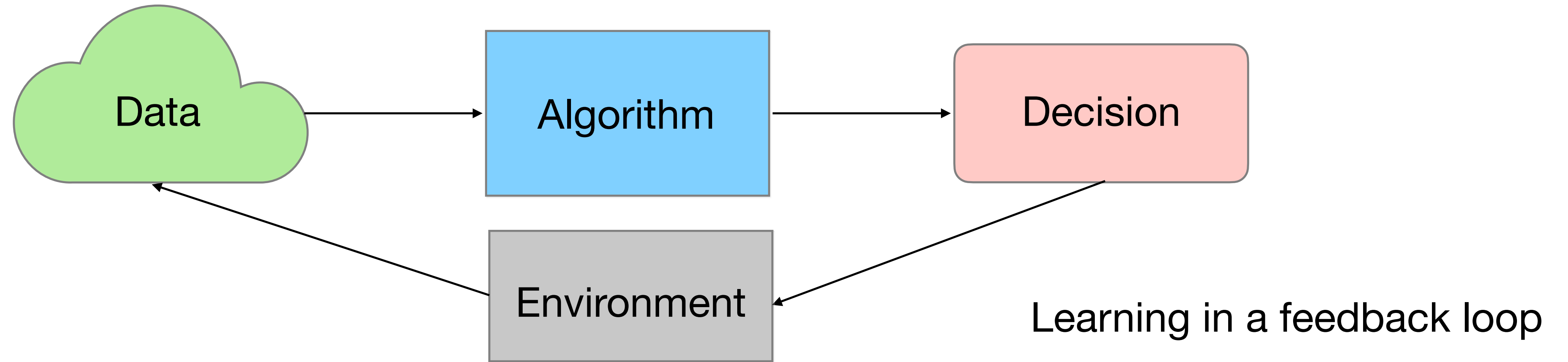
- **Supervised learning** predicts patterns from **passively** observed data
  - Image classification and speech recognition

# Supervised learning vs decision making

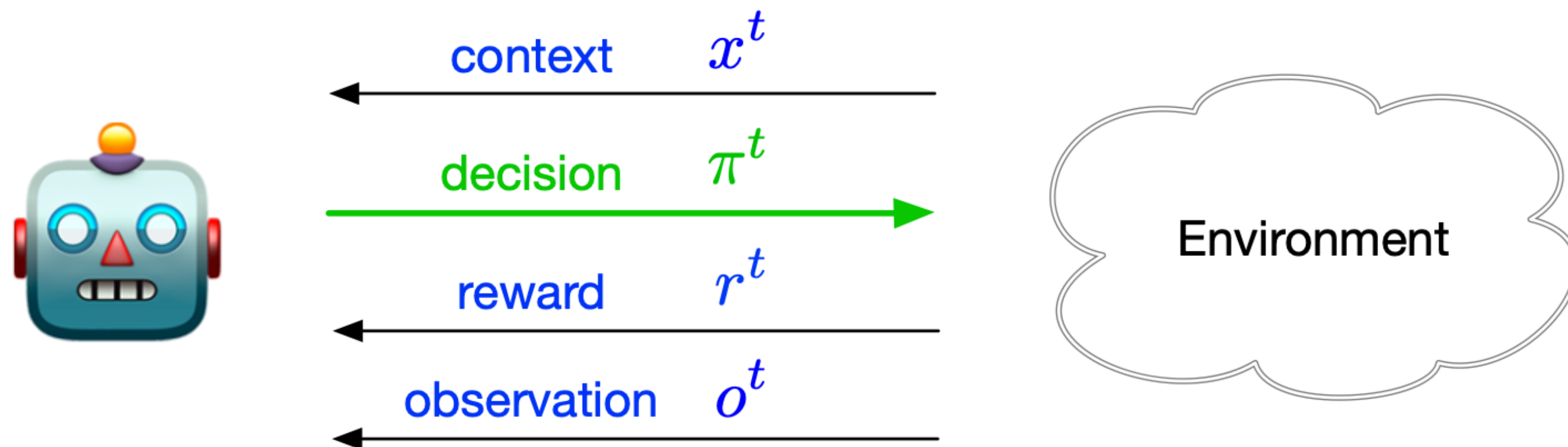


- **Supervised learning** predicts patterns from **passively** observed data
  - Image classification and speech recognition
- **Decision making** actively gathers information by **sequential interactions** with the environment
  - Recommendation system, robotics and game playing

# Preference learning as decision making

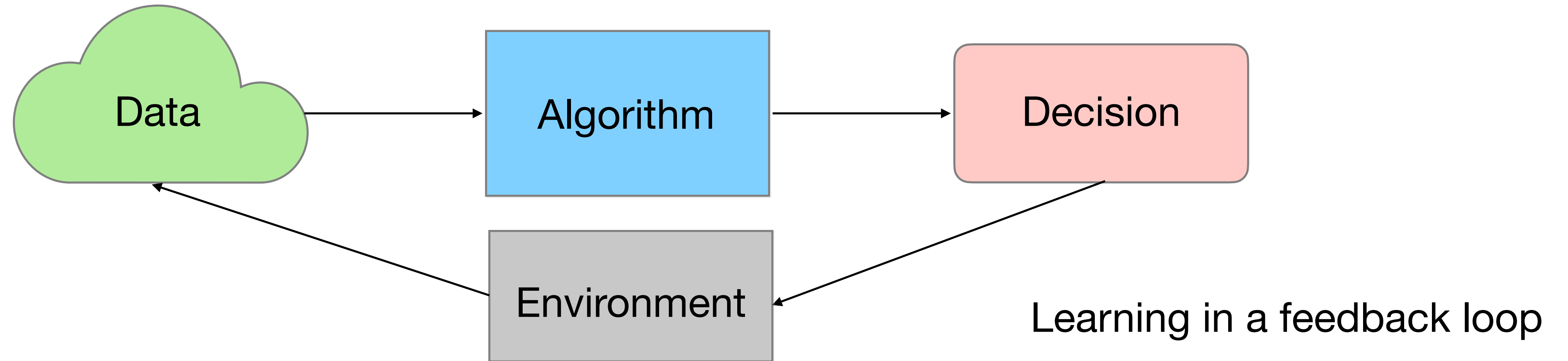


## Decision making for T steps



- **Context:** prompt
- **Decision:** (distribution) of response
- **Reward:** human preference (defined later)
- **Observation:** potential external message

# Exploration-exploitation trade-off



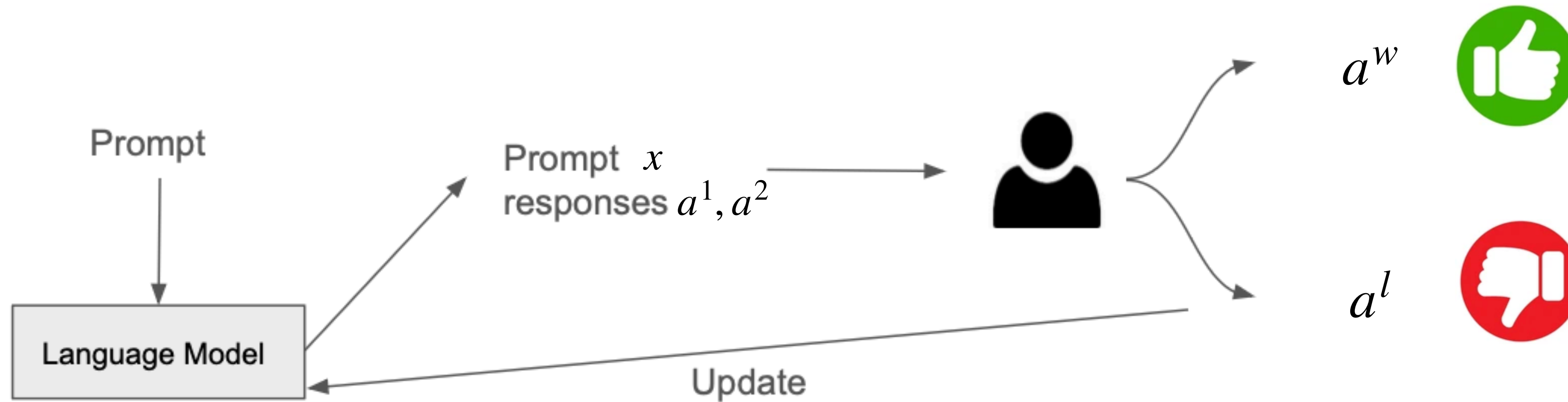
- Trade-off between **exploration** and **exploitation** in online sequential decision making:
  - want to focus on good decisions based on the history and avoid bad decisions to maximize rewards
  - need to try new decisions to learn the environment

**Main research problem: can we design principled preference learning algorithms under this online sequential decision making framework?**

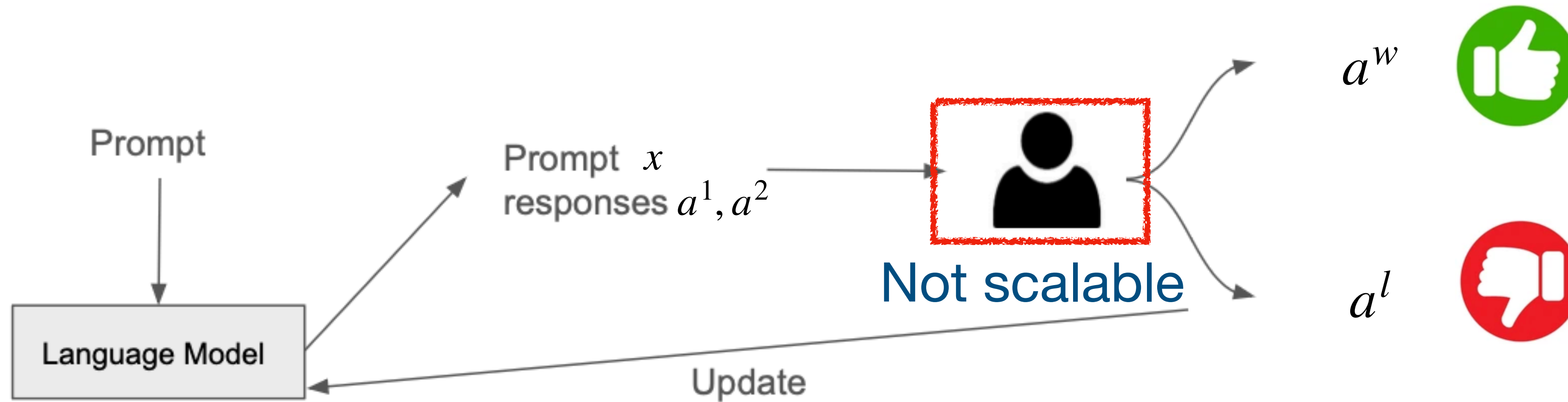
# Reinforcement Learning from Human Feedback



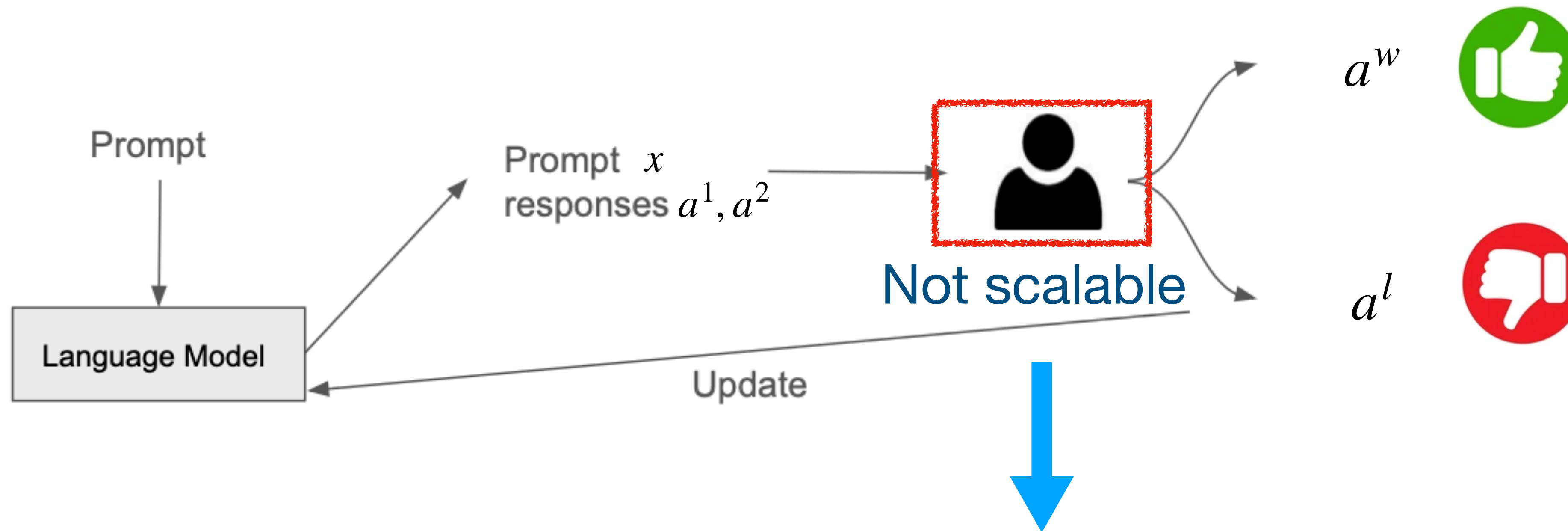
# Reinforcement learning from human feedback



# Reinforcement learning from human feedback



# Reinforcement learning from human feedback



Bradley-Terry (BT) model :

$$\mathcal{P}_{BT}^*(a^1 \succ a^2 \mid x, a^1, a^2) = \frac{e^{r^*(x, a^1)}}{e^{r^*(x, a^1)} + e^{r^*(x, a^2)}}$$

**Scalable:** we can query the reward as many times as we want

# Reinforcement learning from human feedback

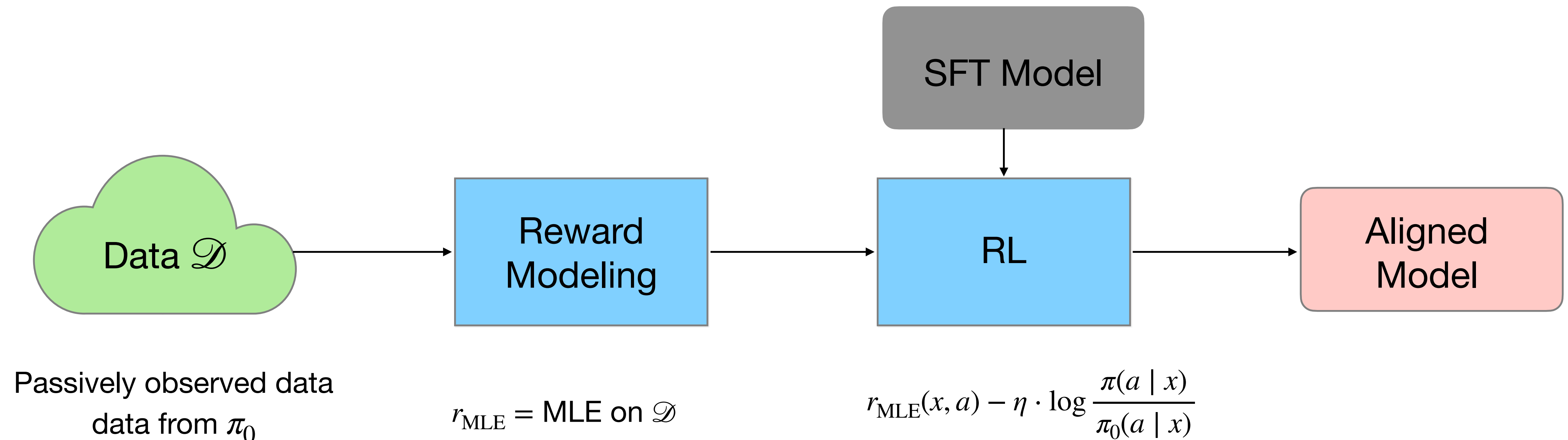
- Learning objective

$$\max_{\pi} J(\pi) = \max_{\pi} \mathbb{E}_{x \sim d_0} \left[ \underbrace{\mathbb{E}_{a \sim \pi(\cdot | x)} [r^*(x, a)]}_{\text{Optimize Reward}} - \underbrace{\eta \text{KL}(\pi(\cdot | x), \pi_0(\cdot | x))}_{\text{Stay Close to SFT Model } \pi_0} \right].$$

# Reinforcement learning from human feedback

- Learning objective

$$\max_{\pi} J(\pi) = \max_{\pi} \mathbb{E}_{x \sim d_0} \left[ \underbrace{\mathbb{E}_{a \sim \pi(\cdot | x)} [r^*(x, a)]}_{\text{Optimize Reward}} - \underbrace{\eta \text{KL}(\pi(\cdot | x), \pi_0(\cdot | x))}_{\text{Stay Close to SFT Model } \pi_0} \right].$$



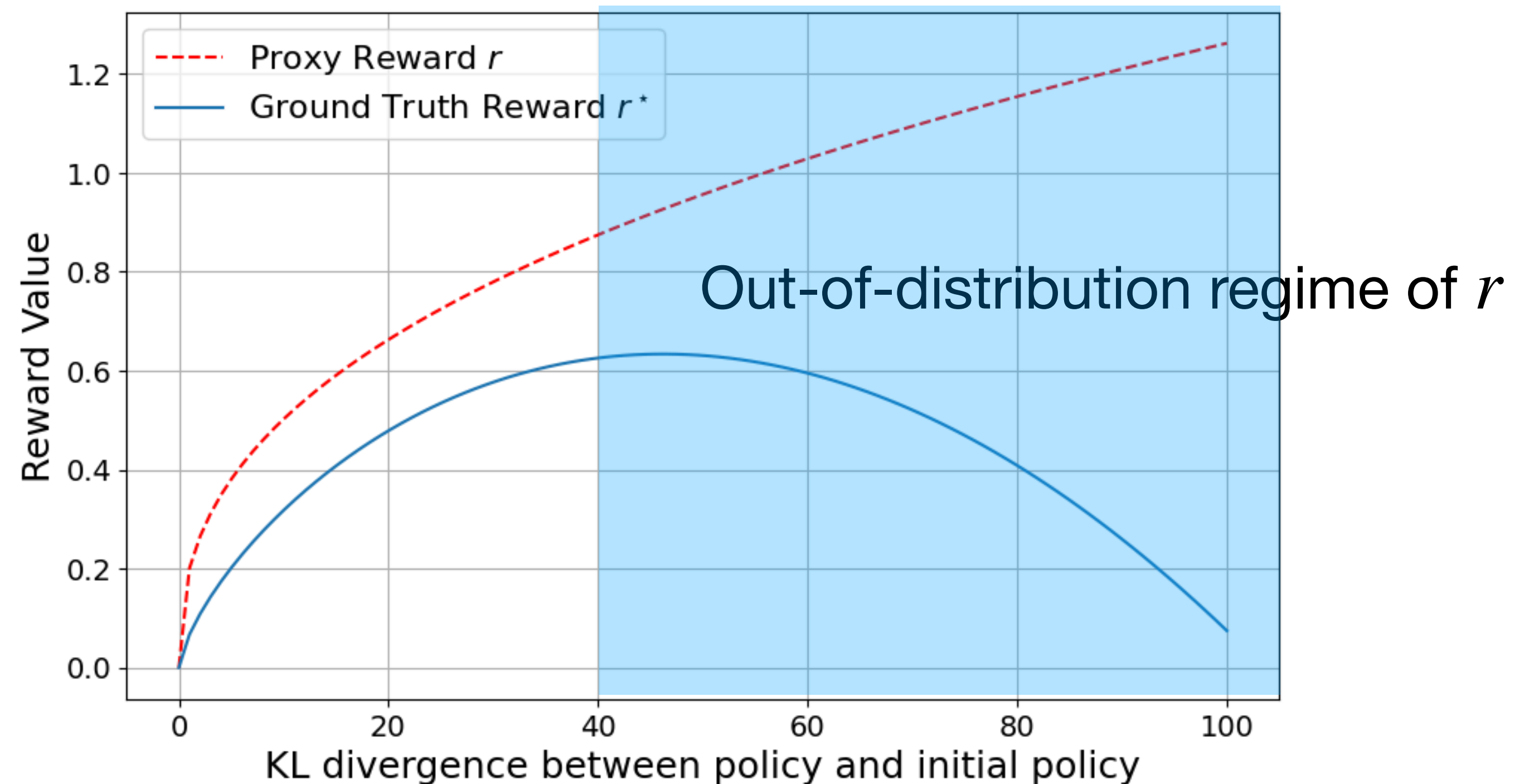
# Proxy reward over-optimization

A high proxy reward does not necessarily lead to a better performance.



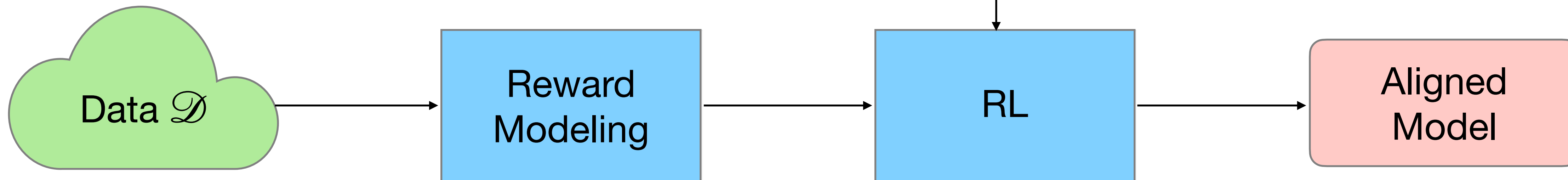
$$\mathcal{P}_{BT}^*(\cdot | x, a^1, a^2), r^* \longrightarrow r \text{ trained from } \mathcal{D}$$

The outputs of LLMs easily fall into **OOD** regime of proxy reward



# Online Iterative RLHF

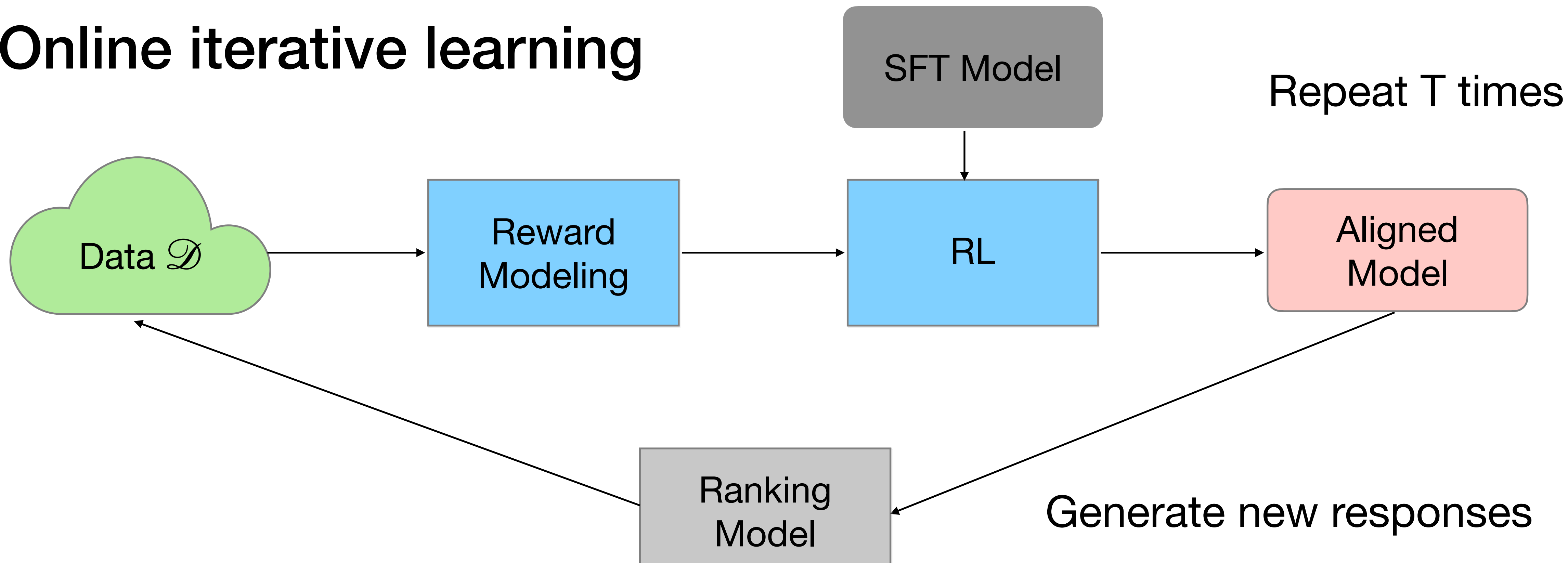
# Offline learning



Passively observed data  
data from  $\pi_0$



# Online iterative learning



Intuition: the new responses and their labels mitigate the OOD issue of proxy reward

# Online iterative RLHF with exploration

- For  $t = 1, 2, 3, \dots$       Divide the learning into T batches
- The main agent **exploits** the historical information:  $\pi_t^1 = \pi_{r_{t,\text{MLE}}}$  based on  $\mathcal{D}_{1:t-1}$

$$\pi_t^1 = \max_{\pi} \mathbb{E}_{x \sim d_0} \left[ \mathbb{E}_{a \sim \pi(\cdot | x)} [r_{t,\text{MLE}}(x, a)] - \eta \text{KL}(\pi(\cdot | x), \pi_0(\cdot | x)) \right].$$

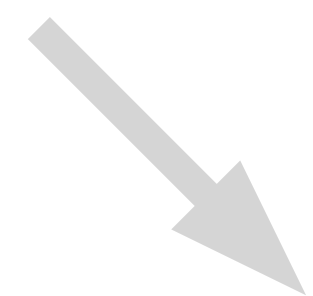
# Online iterative RLHF with exploration

- For  $t = 1, 2, 3, \dots$       Divide the learning into  $T$  batches
- The main agent **exploits** the historical information:  $\pi_t^1 = \pi_{r_{t,\text{MLE}}}$  based on  $\mathcal{D}_{1:t-1}$

$$\pi_t^1 = \max_{\pi} \mathbb{E}_{x \sim d_0} \left[ \mathbb{E}_{a \sim \pi(\cdot|x)} [r_{t,\text{MLE}}(x, a)] - \eta \text{KL}(\pi(\cdot|x), \pi_0(\cdot|x)) \right].$$

- The enhancer **explores** the environment by maximizing the uncertainty relative to  $\pi_t^1$

- $\pi_t^2 = \arg \max_{\pi' \in \Pi} \Gamma_t(\pi_t^1, \pi')$



Uncertainty estimator

- Collect  $m$  new samples  $x_{t,j}, a_{t,j}^1, a_{t,j}^2, y_{t,j} \sim (d_0, \pi_t^1, \pi_t^2, \mathcal{P}_{BT}^*)$  as  $\mathcal{D}_t$

# Uncertainty estimator

---

**Definition:** uncertainty estimator in linear case

**Suppose that**  $r = \langle \theta, \phi(x, a) \rangle : \theta, \phi(x, a) \in \mathbb{R}^d$ . For any two policies  $\pi_t^1, \pi_t^2$ , we define the information gain as

$$\Gamma_t(\pi_t^1, \pi_t^2) = C_{\dagger} \underbrace{\| \mathbb{E}_{\pi_t^1} \phi(x, a_t^1) - \mathbb{E}_{\pi_t^2} \phi(x, a_t^2) \|}_{\text{feature difference}} \|_{\Sigma_t^{-1}}$$

which is the projection of the **new feature difference** to **historical** feature covariance matrix.

---

$$\Sigma_t = \lambda C_{\dagger}^2 I + \sum_{s=1}^{t-1} \mathbb{E}_{x \sim d_0, a^1 \sim \pi_s^1, a^2 \sim \pi_s^2} (\phi(x, a^1) - \phi(x, a^2))^\top (\phi(x, a^1) - \phi(x, a^2))$$

# Theoretical result

---

**Theorem:** Guarantee for the online iterative preference learning

If we run the online iterative RLHF with batch size  $m = O(d/\epsilon^2)$  for  $T = \tilde{\Omega}(d)$  times, with probability at least  $1 - \delta$ , we can find a  $t_0 \in [T]$  such that

$$J(\pi^\star) - J(\pi_{t_0}^1) + \eta \text{KL}(\pi^\star, \pi_{t_0}^1) \leq \epsilon$$

where  $J(\pi) = \mathbb{E}_{d_0, \pi}[r^\star(x, a) - \eta \text{KL}(\pi, \pi_0)]$ .

---

# Theoretical result

---

**Theorem:** Guarantee for the online iterative preference learning

If we run the online iterative RLHF with batch size  $m = O(d/\epsilon^2)$  for  $T = \tilde{\Omega}(d)$  times, with probability at least  $1 - \delta$ , we can find a  $t_0 \in [T]$  such that

$$J(\pi^\star) - J(\pi_{t_0}^1) + \eta \text{KL}(\pi^\star, \pi_{t_0}^1) \leq \epsilon$$

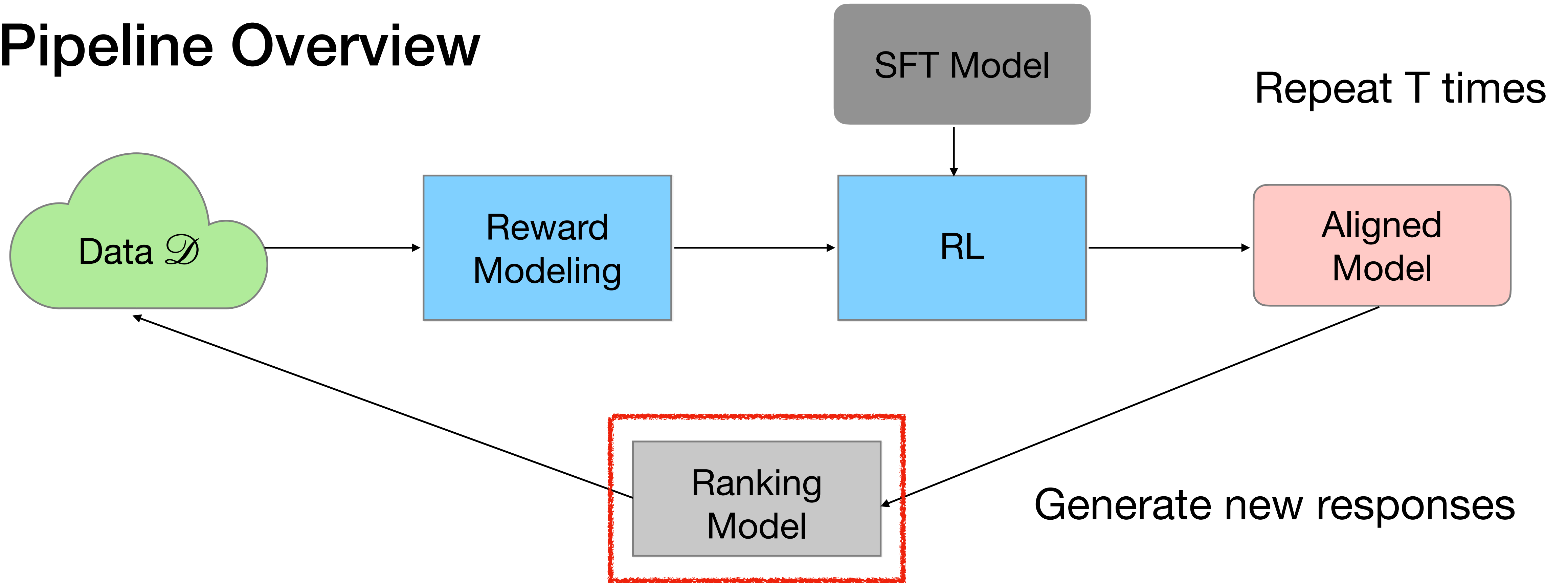
where  $J(\pi) = \mathbb{E}_{d_0, \pi}[r^\star(x, a) - \eta \text{KL}(\pi, \pi_0)]$ .

---

- + The algorithm is provably efficient
- - Iterative human feedback is expensive to collect for open-source project
- - It is not clear how to construct the uncertainty estimator for general neural network

# **RLHFlow: Open-source Online Iterative RLHF**

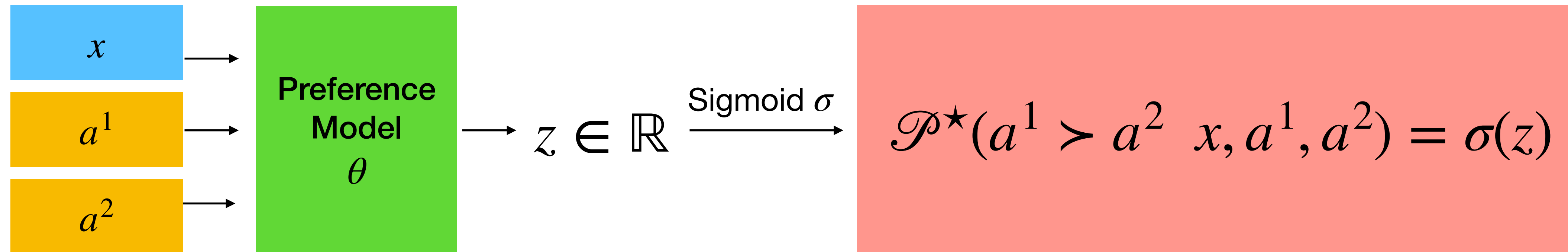
# Pipeline Overview



- A mixture of different types of ranking models on open-source data
- Heuristic rule: length penalty, final result checking for MATH/Coding...



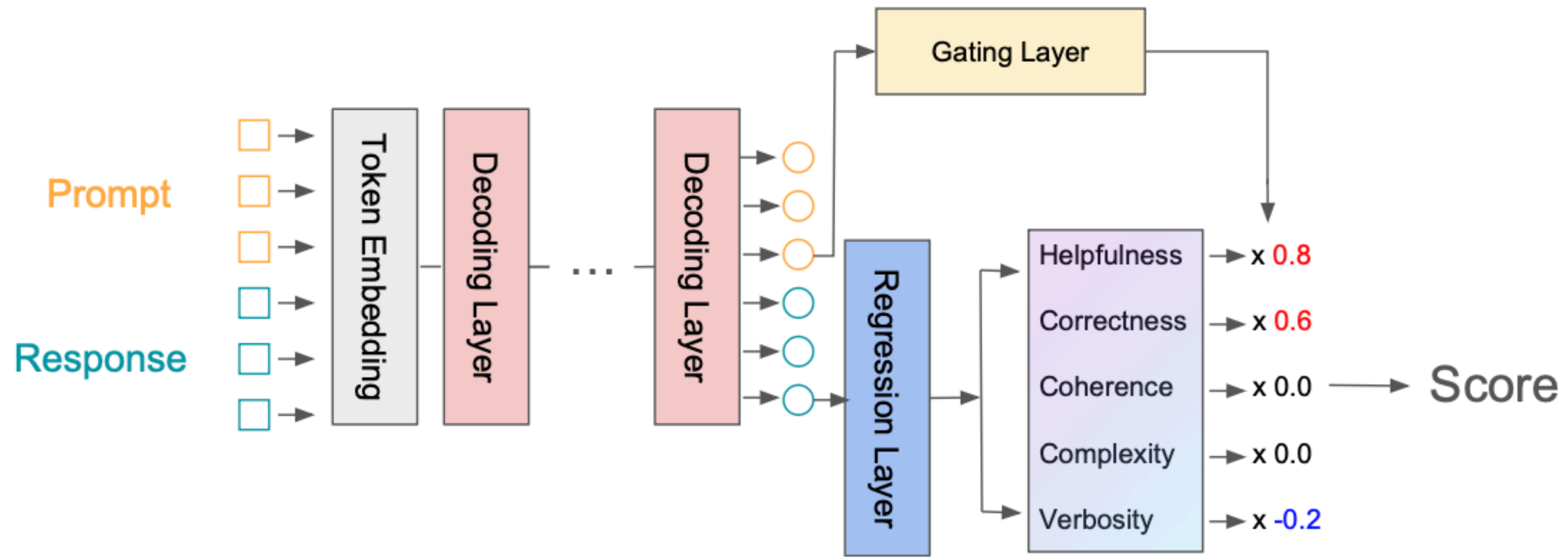
# Next-token prediction as pairwise preference model



instruction = [CONTEXT]  $\{x\}$  [RESPONSE A]  $\{a^1\}$  [RESPONSE B]  $\{a^2\}$

$$\mathbb{P}(a^1 \succ a^2 \mid x, a^1, a^2) = \mathbb{P}(A \mid \text{instruction})$$

# Multi-head reward model with MoE aggregation



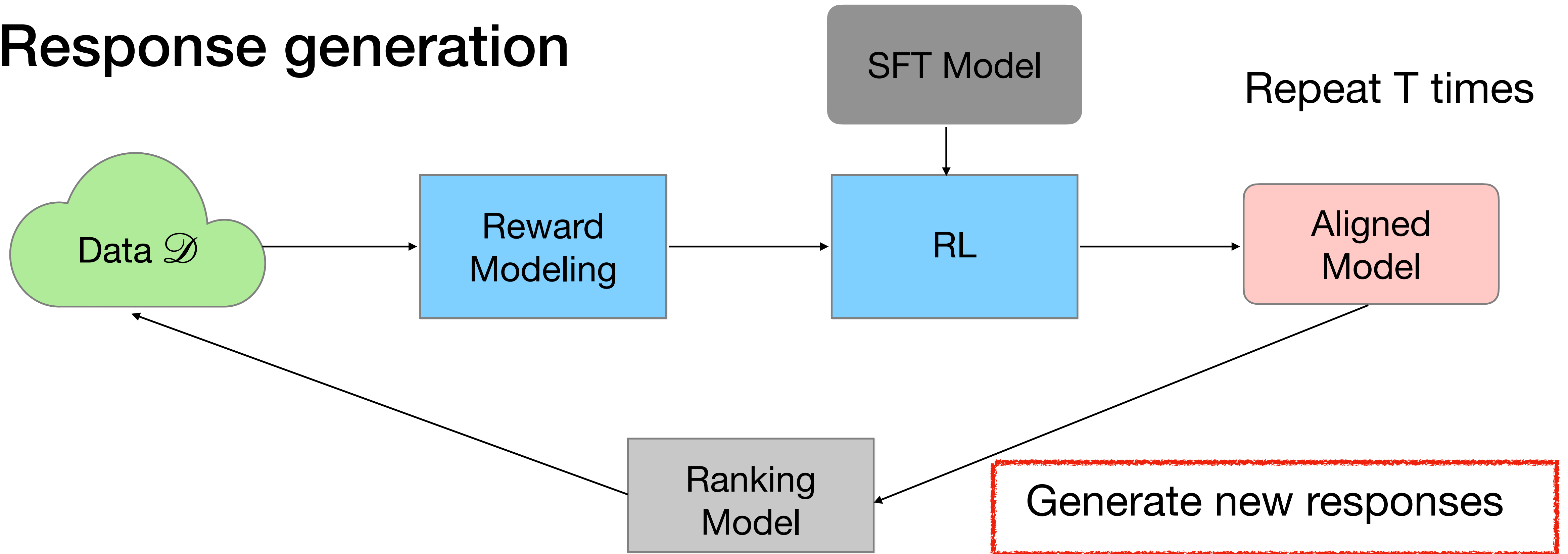
- Multi-head reward modeling from different criteria
- MoE aggregation with the coefficient determined by the embedding of the prompt

# Reward modeling: reward benchmark results

▲	Model	▲	Model Type	▲	Score	▲	Chat	▲	Chat Hard	▲	Safety	▲	Reasoning	▲
1	<a href="#">nvidia/Nemotron-4-340B-Reward</a> *		Custom Classifier		92.2		95.8		87.1		92.2		93.6	
2	<a href="#">RLHFlow/ArmoRM-Llama3-8B-v0.1</a>		Custom Classifier		90.8		96.9		76.8		92.2		97.3	
3	<a href="#">Cohere May 2024</a> *		Custom Classifier		89.5		96.4		71.3		92.7		97.7	
4	<a href="#">nvidia/Llama3-70B-SteerLM-RM</a> *		Custom Classifier		89.0		91.3		80.3		93.7		90.6	
5	<a href="#">facebook/Self-taught-Llama-3-70B</a> *		Generative		88.7		96.9		84.0		91.5		82.5	
6	<a href="#">google/gemini-1.5-pro-0514</a> *		Generative		88.1		92.3		80.6		87.5		92.0	
7	<a href="#">google/flame-1.0-24B-july-2024</a> *		Generative		88.1		92.2		75.7		90.7		93.8	
8	<a href="#">RLHFlow/pair-preference-model-LLaMA3-8B</a>		Custom Classifier		87.1		98.3		65.8		89.7		94.7	
9	<a href="#">Cohere March 2024</a> *		Custom Classifier		87.1		94.7		65.1		90.3		98.2	
10	<a href="#">openai/gpt-4o-2024-08-06</a>		Generative		86.7		96.1		76.1		88.1		86.6	
11	<a href="#">openai/gpt-4-0125-preview</a>		Generative		85.9		95.3		74.3		87.2		86.9	
12	<a href="#">openai/gpt-4-turbo-2024-04-09</a>		Generative		85.1		95.3		75.4		87.1		82.7	
13	<a href="#">openai/gpt-4o-2024-05-13</a>		Generative		84.7		96.6		70.4		86.7		84.9	

The models serve as the ranking models for 30+ follow-up preference learning research projects.

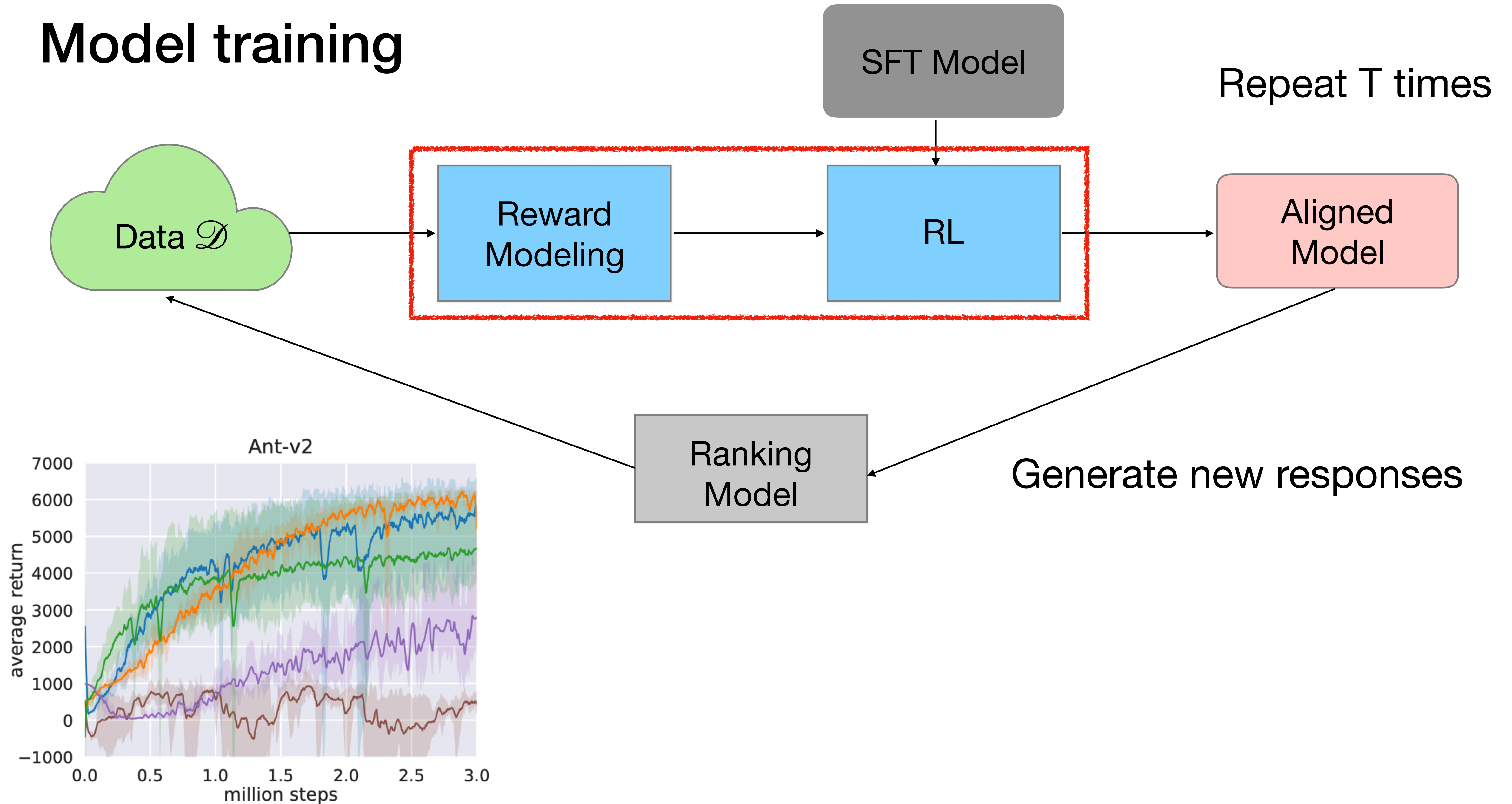
# Response generation



- Heuristic strategies to maximize sample diversity
  - Sample n responses and use the **best** one and the **worst** one to construct a pair
  - Tuning sampling parameter like the temperature
  - Collect data by different checkpoints

$$\Gamma_t(\pi_t^1, \pi_t^2) = C_{\dagger} \underbrace{\| \mathbb{E}_{\pi_t^1} \phi(x, a_t^1) - \mathbb{E}_{\pi_t^2} \phi(x, a_t^2) \|_{\Sigma_t^{-1}}}_{\text{feature difference}}$$

# Model training



# Direct preference optimization (DPO)

- Gibbs distribution

$$\pi_r(\cdot | x) = \max_{\pi} \left[ \mathbb{E}_{a \sim \pi(\cdot | x)} [r(x, a)] - \eta \text{KL}(\pi(\cdot | x), \pi_0(\cdot | x)) \right] = \frac{1}{Z(x)} \cdot \pi_0(\cdot | x) \cdot \exp\left(\frac{1}{\eta} r(x, \cdot)\right)$$

- Re-parameterize reward by policy:

$$r(x, a) = \underbrace{\eta \log \frac{\pi_r(a | x)}{\pi_0(a | x)}}_{\text{Implicit reward}} + \cancel{\eta \log Z(x)}$$

$$Z(x) = \sum_{a \in \mathcal{A}} \pi_0(a | x) \cdot \exp\left(\frac{1}{\eta} r(x, a)\right)$$

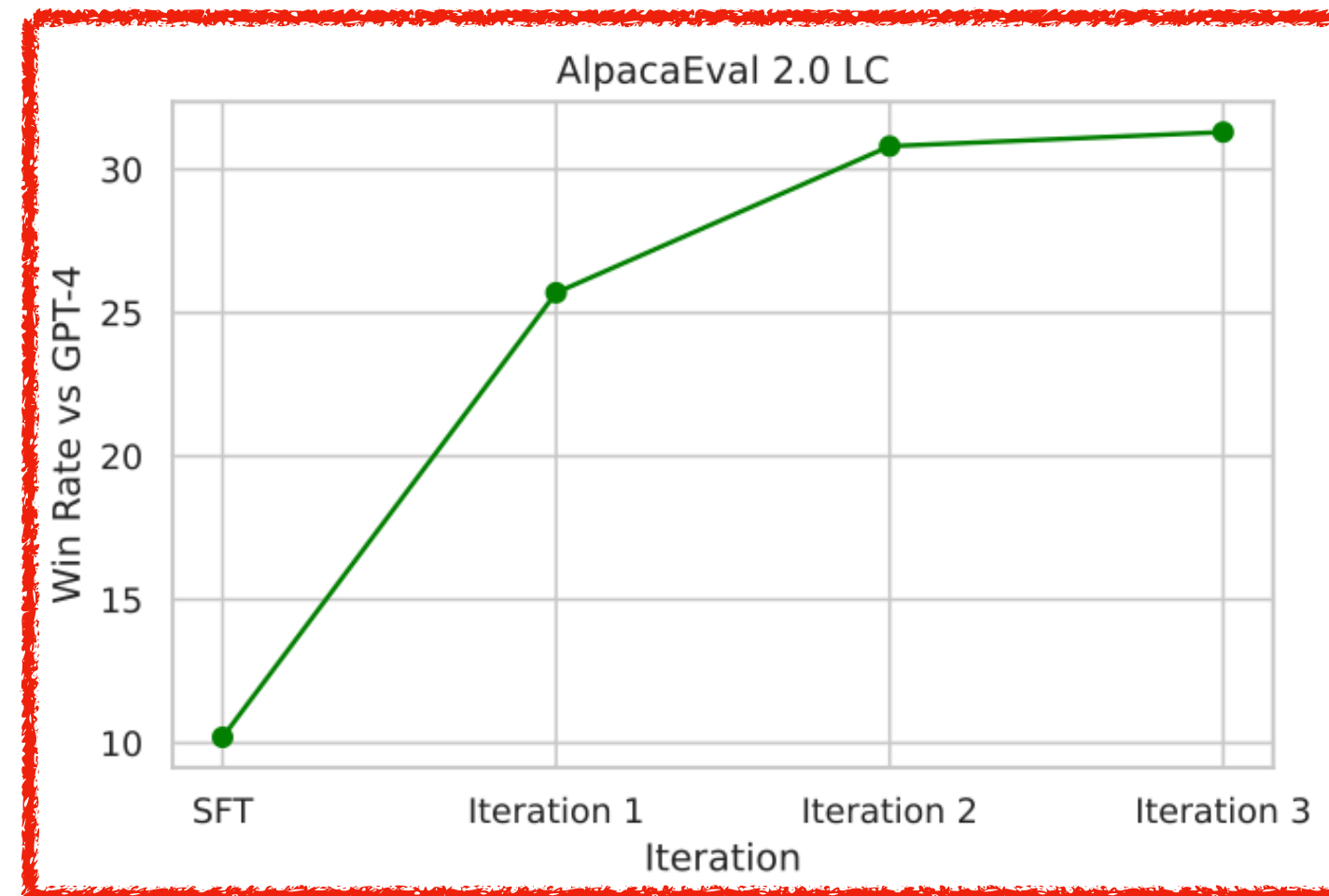
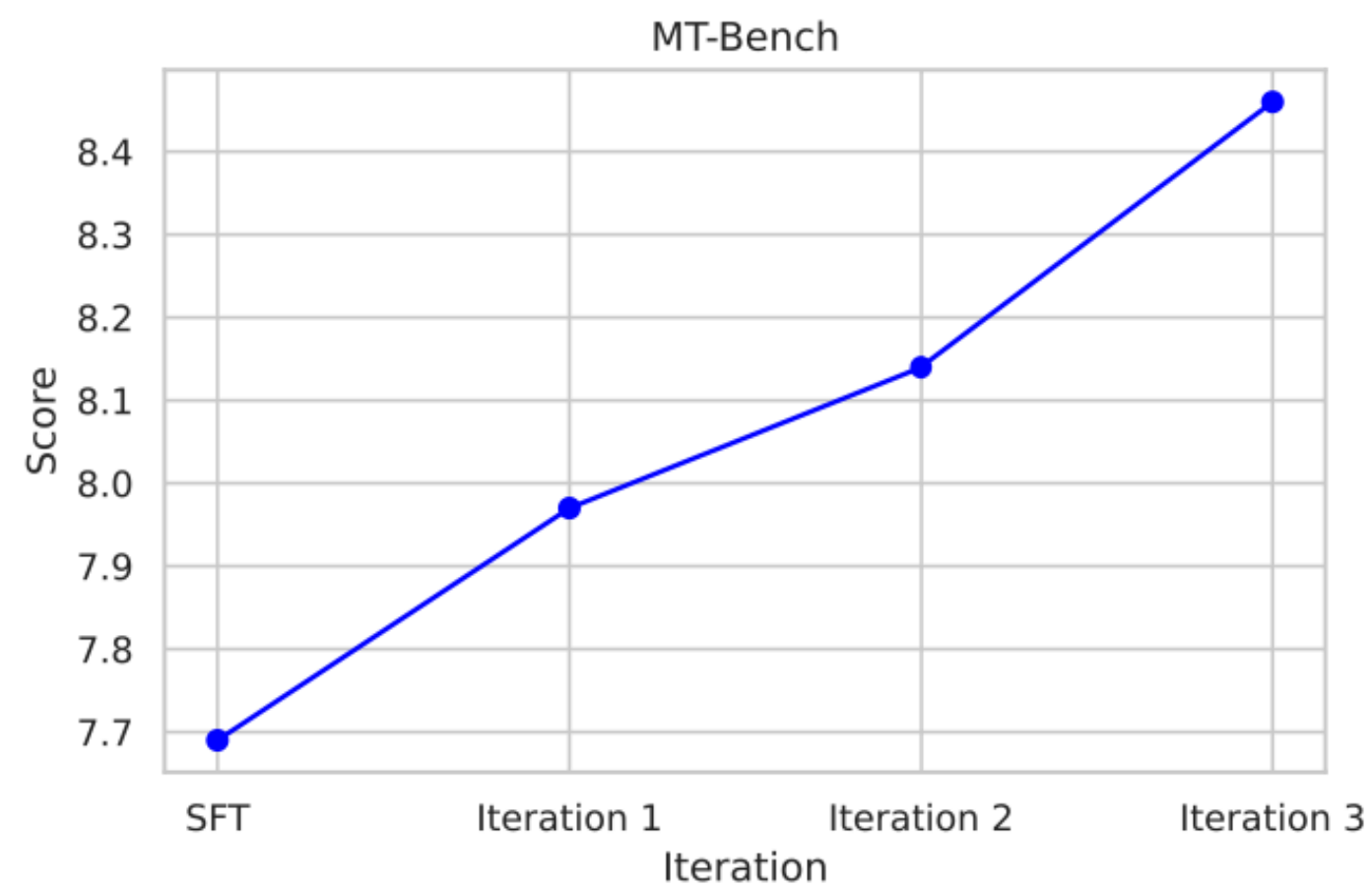
- MLE in reward space -> policy optimization:

$$\ell_{\text{reward}}(r_{\theta}) = \sum_{(x, a^w, a^l) \in \mathcal{D}} \log \left( \sigma \left( r_{\theta}(x, a^w) - r_{\theta}(x, a^l) \right) \right)$$

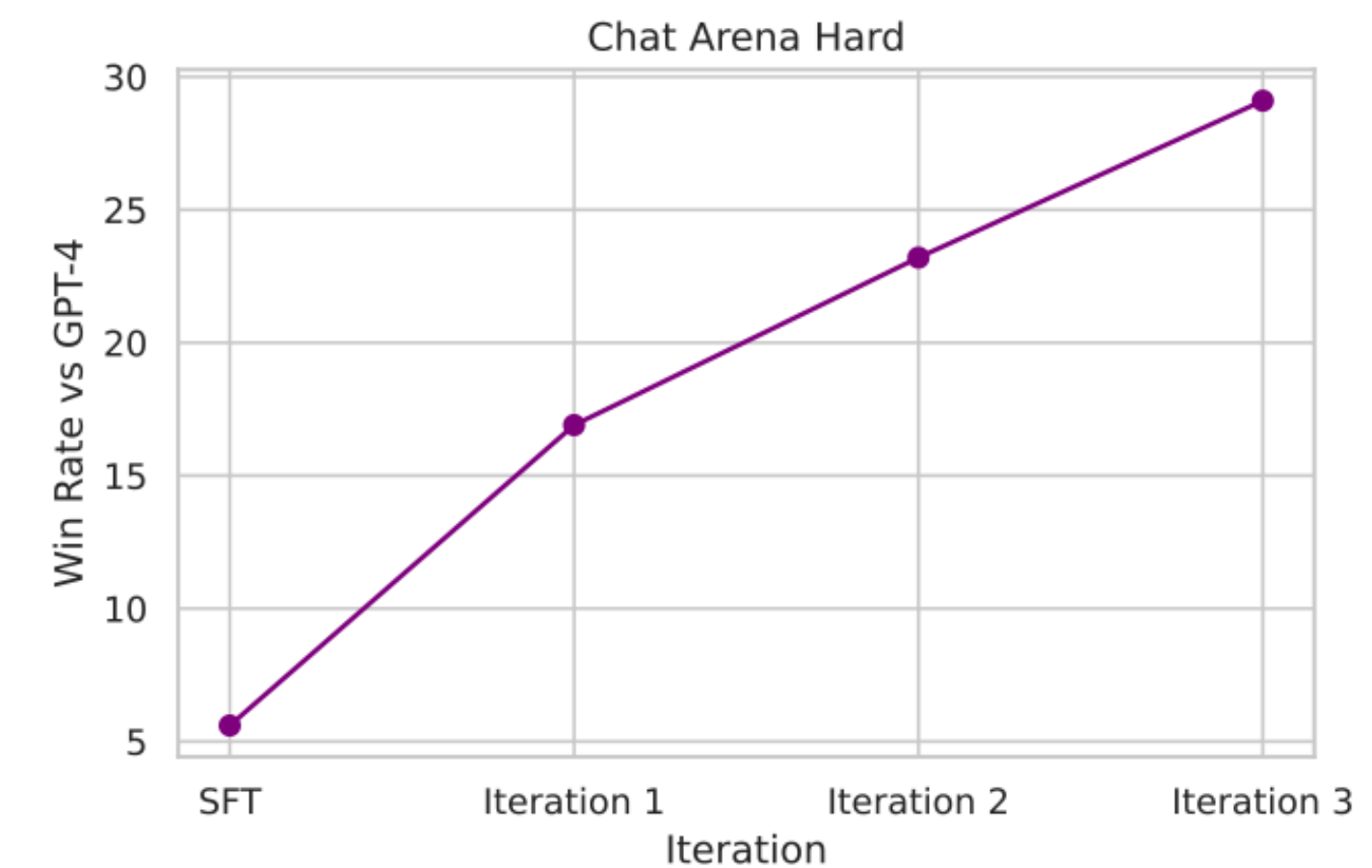
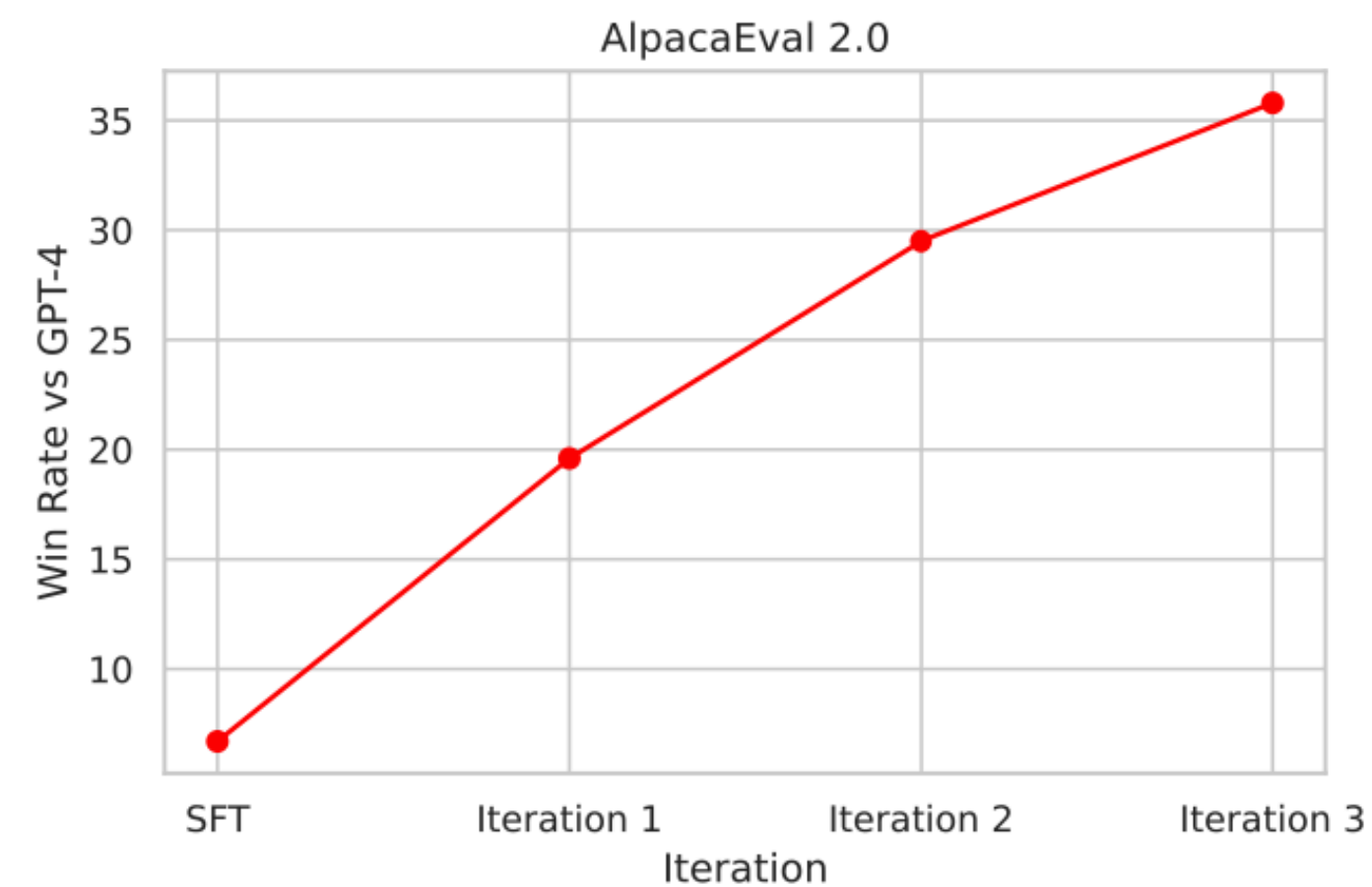


$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}) = - \sum_{(x, a^w, a^l) \in \mathcal{D}} \log \sigma \left( \eta \left[ \log \frac{\pi_{\theta}(a^w | x)}{\pi_0(a^w | x)} - \log \frac{\pi_{\theta}(a^l | x)}{\pi_0(a^l | x)} \right] \right).$$

# Main result: state-of-the-art chat model

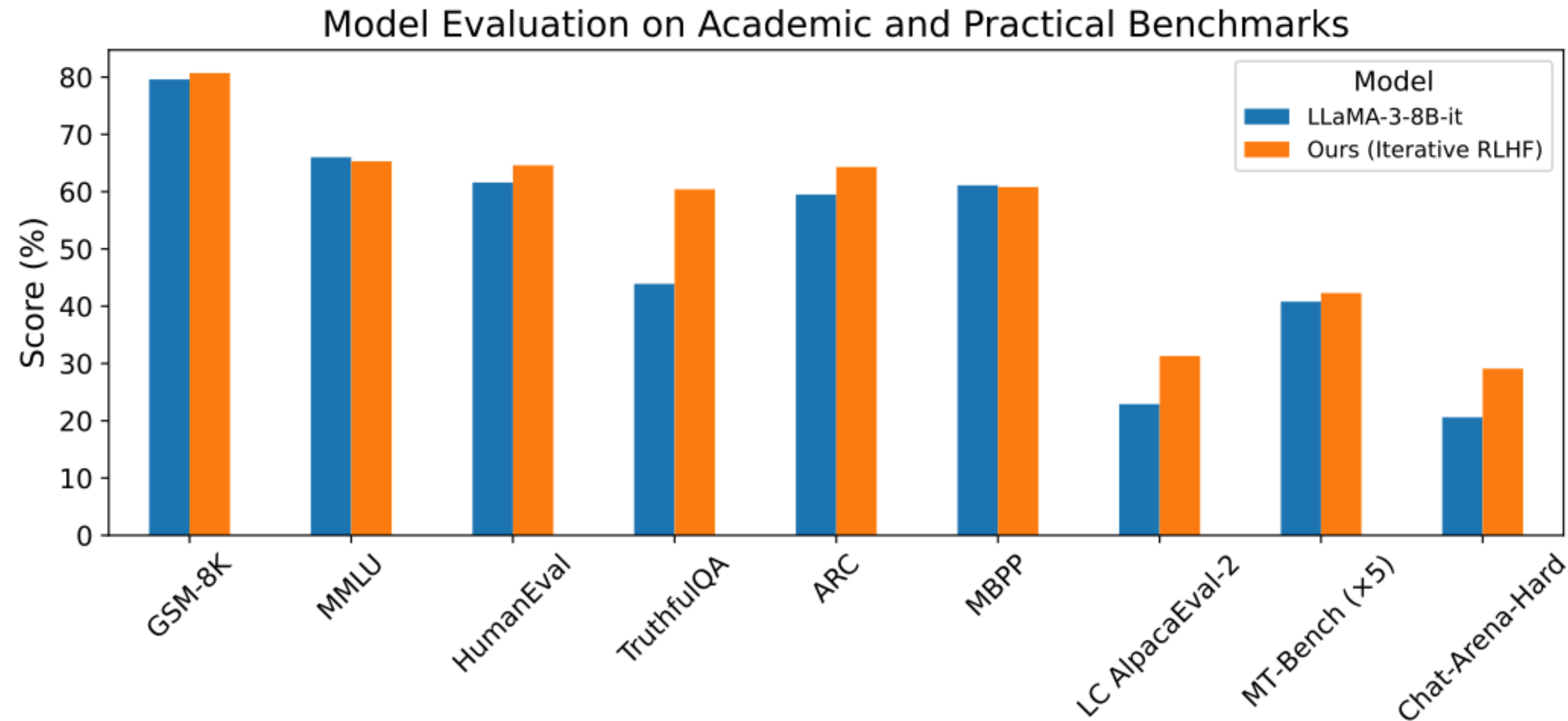


Base model: LLaMA3-8B-SFT



Evaluation results of models at different iterations on instruction-following benchmarks.

# Main result: state-of-the-art chat model



Evaluation results on standard academic and instruction-following benchmarks.



# Multi-turn tool-integrated reasoning

*Prompt:*

**User:** Let

$$f(x) = \begin{cases} ax + 3, & \text{if } x > 2, \\ x - 5 & \text{if } -2 \leq x \leq 2, \\ 2x - b & \text{if } x < -2. \end{cases}$$

Find  $a + b$  if the piecewise function is continuous (which means that its graph can be drawn without lifting your pencil from the paper) ( $s_1 = x$ )

**Assistant:** The function must be continuous at 2 and  $-2$ . That means  $x - 5 = ax + 3$  at  $x = 2$  and  $x - 5 = 2x - b$  at  $x = -2$ . Let's now solve for  $x$  in each case.

# Multi-turn tool-integrated reasoning

*Prompt:*

**User:** Let

$$f(x) = \begin{cases} ax + 3, & \text{if } x > 2, \\ x - 5 & \text{if } -2 \leq x \leq 2, \\ 2x - b & \text{if } x < -2. \end{cases}$$

Find  $a + b$  if the piecewise function is continuous (which means that its graph can be drawn without lifting your pencil from the paper) ( $s_1 = x$ )

**Assistant:** The function must be continuous at 2 and  $-2$ . That means  $x - 5 = ax + 3$  at  $x = 2$  and  $x - 5 = 2x - b$  at  $x = -2$ . Let's now solve for  $x$  in each case.

```
1 from sympy import symbols, Eq, solve
2 ## define the variables
3 x, a, b = symbols('x a b')
4
5
6 ## the function must be continuous at 2
7 eq1 = Eq(x-5, a*x + 3)
8 eq1_sol = solve(eq1, x)
9
10 ## the function must be continuous at -2
11 eq2 = Eq(x-5, 2*x - b)
12 eq2_sol = solve(eq2, x)
13
14 ## print the solutions
15 print(eq1_sol)
16 print(eq2_sol)
```

( $a_1$ )

# Multi-turn tool-integrated reasoning

*Prompt:*

**User:** Let

$$f(x) = \begin{cases} ax + 3, & \text{if } x > 2, \\ x - 5 & \text{if } -2 \leq x \leq 2, \\ 2x - b & \text{if } x < -2. \end{cases}$$

Find  $a + b$  if the piecewise function is continuous (which means that its graph can be drawn without lifting your pencil from the paper) ( $s_1 = x$ )

**Assistant:** The function must be continuous at 2 and  $-2$ . That means  $x - 5 = ax + 3$  at  $x = 2$  and  $x - 5 = 2x - b$  at  $x = -2$ . Let's now solve for  $x$  in each case.

```
1 from sympy import symbols, Eq, solve
2 ## define the variables
3 x, a, b = symbols('x a b')
4
5
6 ## the function must be continuous at 2
7 eq1 = Eq(x-5, a*x + 3)
8 eq1_sol = solve(eq1, x)
9
10 ## the function must be continuous at -2
11 eq2 = Eq(x-5, 2*x - b)
12 eq2_sol = solve(eq2, x)
13
14 ## print the solutions
15 print(eq1_sol)
16 print(eq2_sol)
```

( $a_1$ )

$$s_h = (x, a_1, o_1, a_2, o_2, \dots, a_{h-1}, o_{h-1})$$

**User:** “output  $[-8/(a - 1)] [b - 5]$ ” ( $o_1$ )

# Preference learning for multi-turn reasoning task

- Trajectory:  $x$ -prompt,  $a$ -action,  $o$ -external observation,  $s_h$ : history at the beginning of step  $h$

$$\tau = (x, \underbrace{a_1, o_1, \dots, o_{H-1}, a_H}_y).$$

- Trajectory preference

$$\mathcal{P}_{BT}^*(y^1 \succ y^2 \mid x, y^1, y^2) = \frac{e^{r^*(x, y^1)}}{e^{r^*(x, y^1)} + e^{r^*(x, y^2)}}$$

- Learning objective

$$\arg \max_{\pi} J(\pi; \mathcal{M}^*, \pi_0) = \mathbb{E}_{x \sim d_0} \mathbb{E}_{a_h \sim \pi_h(\cdot \mid s_h), o_h \sim \mathbb{P}_h(\cdot \mid s_h, a_h)} \left[ r^*(x, y) - \eta \sum_{h=1}^H \text{KL}(\pi_h(\cdot \mid s_h), \pi_{0,h}(\cdot \mid s_h)) \right].$$

# Multi-turn direct preference learning

- Re-parameterization trick to connect the **model** with the **policy**

$$r(s_H, a_H) = \underbrace{\eta \sum_{h=1}^H \log \frac{\pi_{\mathcal{M},h}(a_h | s_h)}{\pi_{0,h}(a_h | s_h)}}_{\text{term (A)}} + \underbrace{\frac{Z(x)}{\cancel{\quad}}}_{\text{term (B)}} + \underbrace{\sum_{h=1}^{H-1} \left[ V_{\mathcal{M},h+1}(s_{h+1}) - \mathbb{E}_{o_h \sim \mathbb{P}_h(\cdot | s_h, a_h)} V_{\mathcal{M},h+1}(s_{h+1}) \right]}_{\text{term (C)}}.$$

- Term (C) is not zero except for

$V_{\mathcal{M},h}$  : optimal V value function under  $\mathcal{M} = (r, \mathbb{P})$

- $H = 1$ : original DPO
- $o_h$  is deterministic given the history

Implementation: run DPO but mask out the external messages.

# Main result: improving reasoning ability

Base Model	Method	with Tool	GSM8K	MATH	AVG
Gemma-1.1-it-7B	SFT <sup>†</sup>	✓	77.5	46.1	61.8
Gemma-1.1-it-7B	RAFT	✓	79.2	47.3	63.3
Gemma-1.1-it-7B	Iterative Single-turn DPO	✓	81.7	48.9	65.3
Gemma-1.1-it-7B	Iterative M-DPO + fixed reference	✓	79.9	48.0	64.0
Gemma-1.1-it-7B	M-DPO Iteration 1	✓	81.5	49.1	65.3
Gemma-1.1-it-7B	M-DPO Iteration 2	✓	82.5	49.7	66.1
Gemma-1.1-it-7B	M-DPO Iteration 3	✓	83.9 ↑6.4	51.2 ↑5.1	67.6
CodeGemma-1.1-it-7B	SFT <sup>†</sup>	✓	77.3	46.4	61.9
CodeGemma-1.1-it-7B	RAFT	✓	78.8	48.4	63.6
CodeGemma-1.1-it-7B	Iterative Single-turn DPO	✓	79.1	48.9	64.0
CodeGemma-1.1-it-7B	Iterative M-DPO	✓	81.5 ↑4.2	50.1 ↑3.7	65.8
Mistral-7B-v0.3	SFT <sup>†</sup>	✓	77.8	42.7	60.3
Mistral-7B-v0.3	RAFT	✓	79.8	43.7	61.8
Mistral-7B-v0.3	Iterative Single-turn DPO	✓	79.8	45.1	62.5
Mistral-7B-v0.3	Iterative M-DPO	✓	82.3 ↑4.5	47.5 ↑4.8	64.9
Gemma-2-it-9B	SFT <sup>†</sup>	✓	84.1	51.0	67.6
Gemma-2-it-9B	RAFT	✓	84.2	52.6	68.4
Gemma-2-it-9B	Iterative Single-turn DPO	✓	85.2	53.1	69.2
Gemma-2-it-9B	Iterative Single-turn KTO	✓	85.4	52.9	69.2
Gemma-2-it-9B	Iterative M-DPO	✓	86.3 ↑2.2	54.5 ↑3.5	70.4

Prompt: training set MATH and GSM8K  
Reward: binary reward by checking the answer

# A practical and **open-source** codebook

**RM-Bradley-Terry** >  
We train the reward model as the maximum likelihood estimat...

- sfairXC/FsfairX-LLaMA3-RM-v0.1**  
Text Classification • Updated Apr 24 • 17k • 42
- hendrydong/preference\_700K**  
Viewer • Updated Apr 17 • 700k • 3.34k • 2
- weqweasdas/RM-Mistral-7B**  
Text Classification • Updated Mar 31 • 4.74k • 20
- weqweasdas/preference\_dataset\_mixture2\_an...**

**PM-pair** >  
This is a collection of materials for training pairwise preference...

- RLHFlow/pair-preference-dataset-mix1**  
Viewer • Updated May 6 • 548k • 71 • 1
- RLHFlow/pair-preference-model-LLaMA3-8B**  
Text Generation • Updated May 24 • 8.19k • 32
- RLHFlow/pair\_preference\_model\_dataset**  
Viewer • Updated Apr 20 • 699k • 285 • 1

**Online RLHF** >  
Datasets, code, and models for online RLHF (i.e., iterative DPO)

- RLHFlow/prompt-collection-v0.1**  
Viewer • Updated May 7 • 179k • 26 • 6
- RLHFlow/pair-preference-model-LLaMA3-8B**  
Text Generation • Updated May 24 • 8.19k • 32
- sfairXC/FsfairX-LLaMA3-RM-v0.1**  
Text Classification • Updated Apr 24 • 17k • 42
- RLHFlow/llama-sft** private

**RLHFlow Reward Models** >  
Reward models trained by RLHFlow codebase (https://github.c...

- RLHFlow/ArmoRM-Llama3-8B-v0.1**  
Text Classification • Updated Jul 18 • 64.8k • 115
- RLHFlow/pair-preference-model-LLaMA3-8B**  
Text Generation • Updated May 24 • 8.19k • 32
- sfairXC/FsfairX-LLaMA3-RM-v0.1**  
Text Classification • Updated Apr 24 • 17k • 42
- RLHF Workflow: From Reward Modeling to Online RLHF

**RLHF-Reward-Modeling** Public  
Recipes to train reward model for RLHF.  
Python 629 54

**Online-RLHF** Public  
A recipe for online RLHF.  
Python 374 42

Dataset

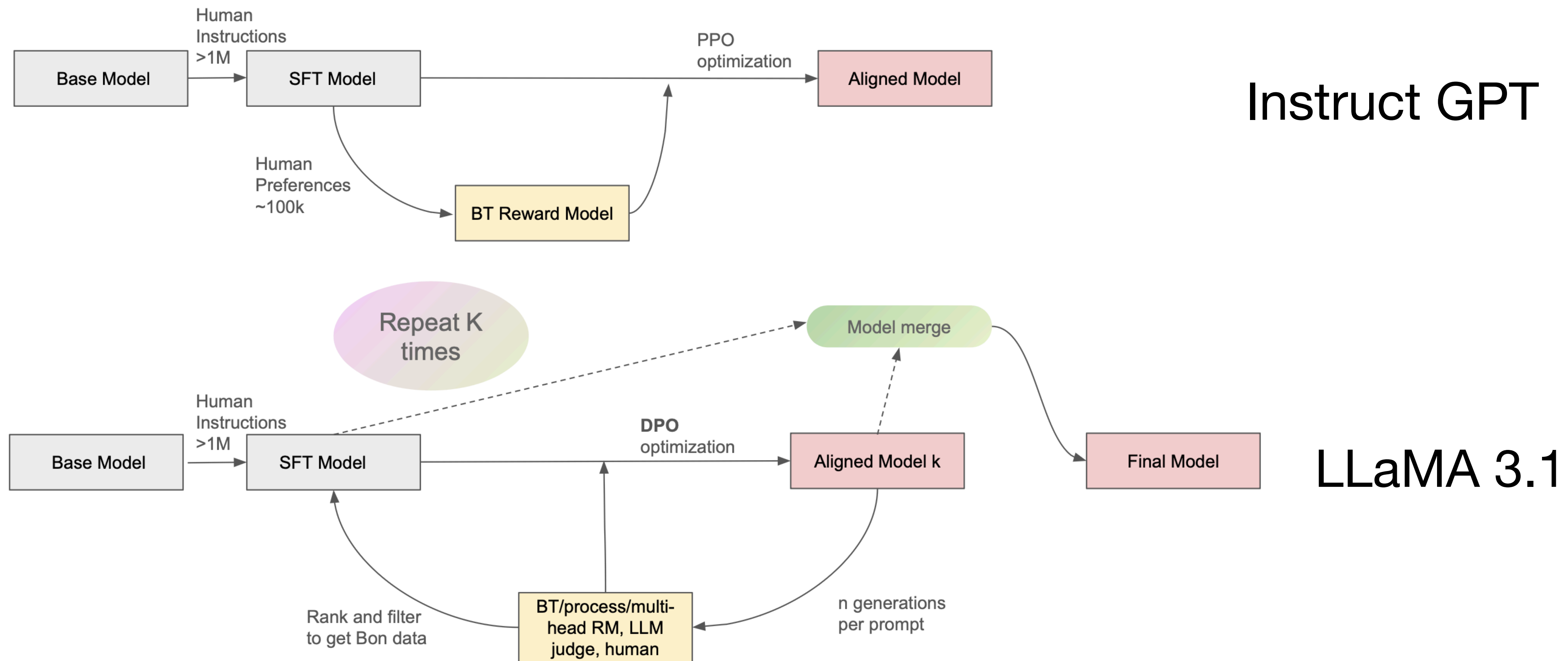
Training code

Hyper-parameters

Final models

# Takeaway

- RLHF benefits from **continuous online exploration** through interactions with the rater
- Online iterative direct preference learning is a robust recipe to make good chatbot





**Thanks for listening!**