

# Reinforcement Learning in the Era Of LLMs

Nov. 2023

Hao Sun

University of Cambridge

[hs789@cam.ac.uk](mailto:hs789@cam.ac.uk)

---



# Content

- Essential Concepts in RL
  - Alignment: From SFT to RLHF
  - Prompting as an inverse-Alignment
-

# Some 'Terms' You May Have Heard About...

- Markov Decision Processes
  - States, Observations, Transitions (Dynamics), Actions, Rewards, Discount Factors...
  - Model-Based / Model Free
  - Value-Based / Policy-Based / Actor-Critic
  - On-Policy / Off-Policy
  - Online / Offline
  - Discrete Control / Continuous Control
-

# But those are not necessary...

- ~~Markov Decision Processes~~
  - ~~States, Observations, Transitions (Dynamics), Actions, Rewards, Discount Factors...~~
  - ~~Model-Based / Model Free~~
  - ~~Value-Based / Policy-Based / Actor-Critic~~
  - ~~On-Policy / Off-Policy~~
  - ~~Online / Offline~~
  - ~~Discrete Control / Continuous Control~~
-

# Essential Ideas

- What is RL?

*An **agent** learns from **trial and error**, to maximize a **cumulative reward**.*

- **agent**: can be human or neural networks. It has a policy (clinical guidelines, public policies)
  - **trial and error**: Online or Offline? Real or (data-drive) Simulator?
  - **cumulative reward**: the *Reward Hypothesis*
-

# Reward Can be Sparse...

- “Win the game”, but how?
- Imitation Learning (IL) Can Help.
- 1. Behavior Clone
  - Pros: **Simple, Does not need further interactions with the environment**
  - Cons: **Compounding error, multi-modality modeling**
- 2. GAIL: Generative Adversarial Imitation Learning [Ho et al. 2016]
  - Pros: **Solves the above cons.**
  - Cons: **Needs more interactions (the dynamics is accessible, but reward unknown)**

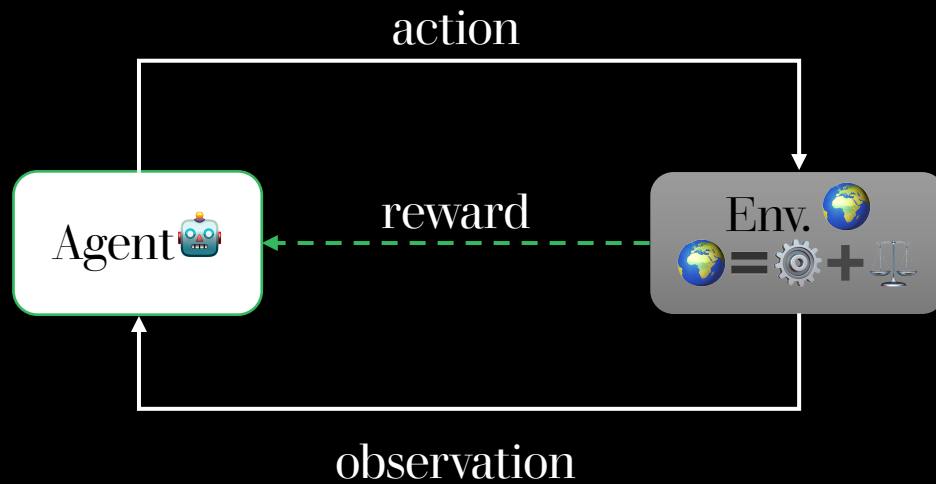
# Difference between Inverse RL and IL

- Inverse-RL  $\approx$  Imitation Learning, with an emphasis on explicit reward modeling

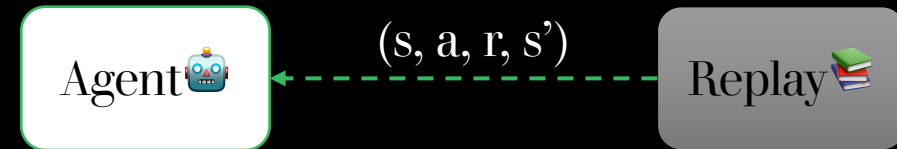
*Learning from **logged trial and error**, to find out **what cumulative reward** is being optimized.*

- **logs** can be either expert decisions or non-expert decisions. Extrapolation.
  - **trial and error** are offline data
  - **(the estimated) reward** can be used as an evaluator of trajectories/policies
-

# Graph: RL and Offline-RL



RL

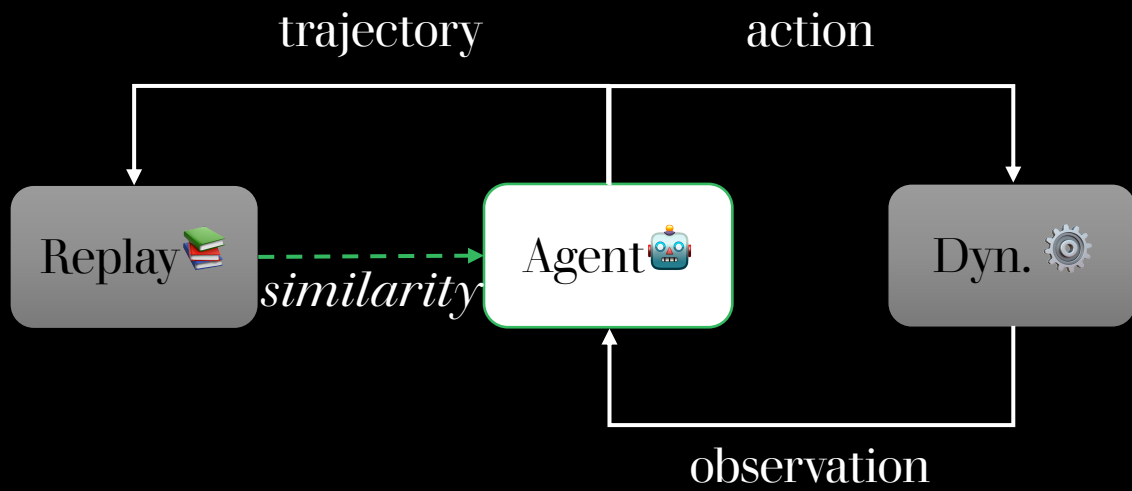


Offline-RL

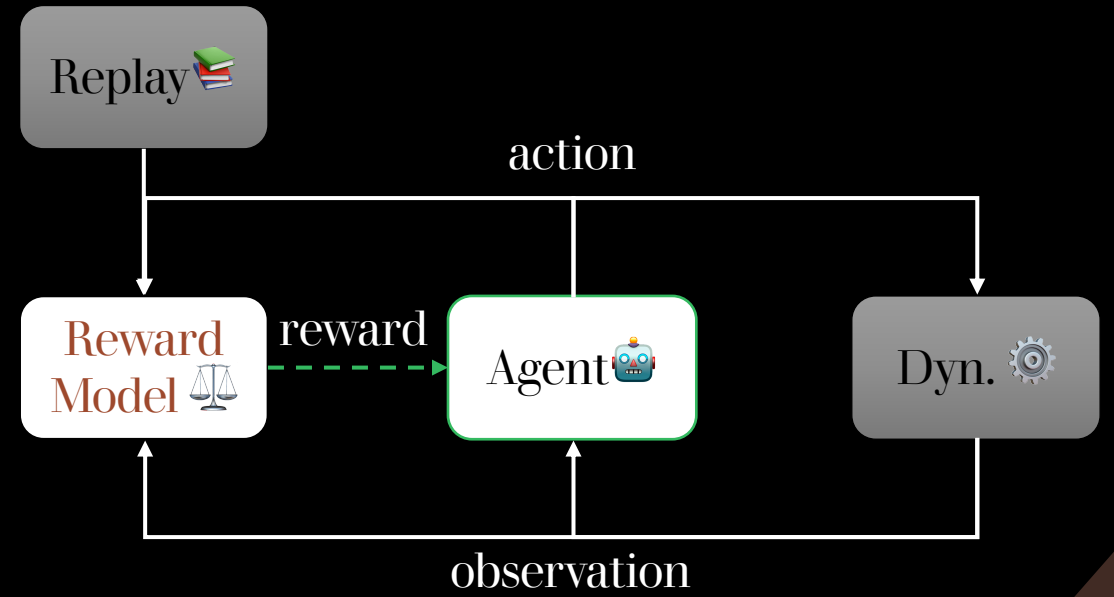
Env. = Dyn. + Rew.



# Graph: IL and IRL

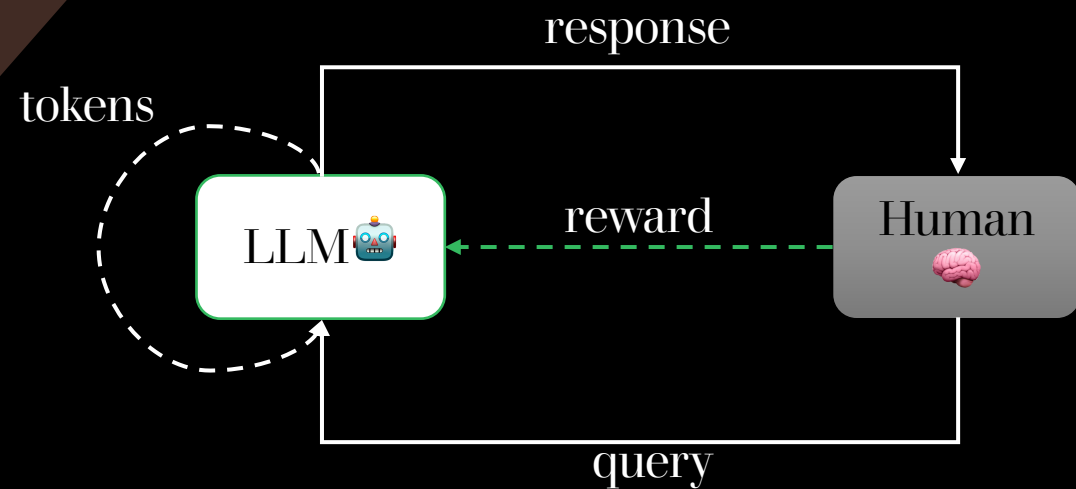


IL



IRL: Model Reward Explicitly

# LLM Alignment with Human Feedback



Alignment as RL

- Why RL?

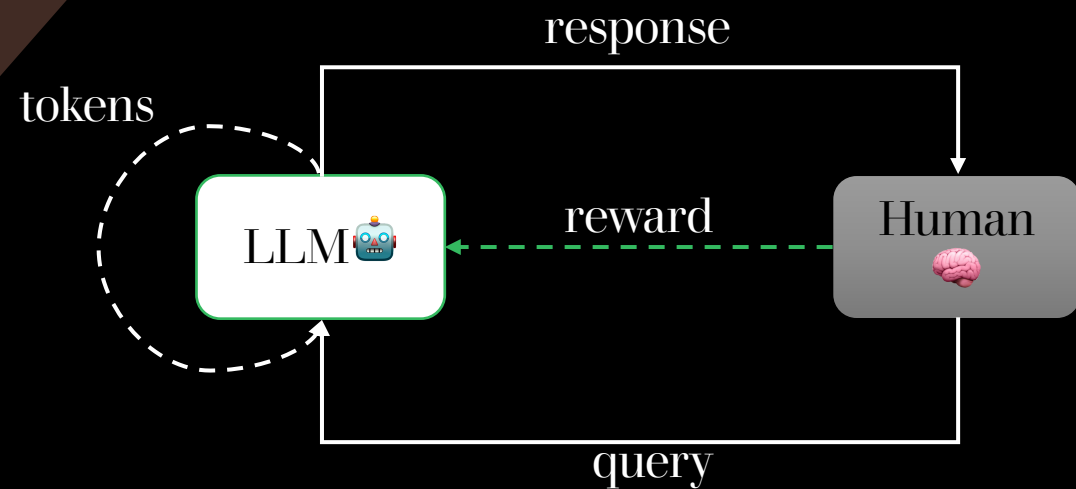
*We can not define the desired objective as a metric function.*

*We can not do back-prop through a black-box 🧠*

- Why not?

*Too expensive.*

# LLM Alignment with Human Feedback



Alignment as RL

- Why RL?

*We can not define the desired objective as a metric function.*

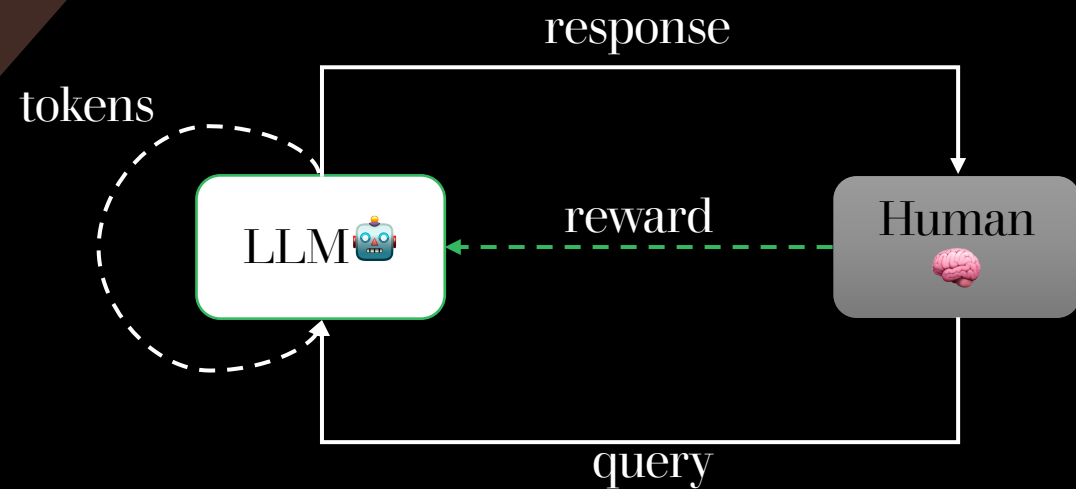
*We can not do back-prop through a black-box 🧠*

- Why not?

*Too expensive.*

*OpenAI: unless you have enough volunteers (users)...*

# LLM Alignment with GPT-4 Feedback?



Alignment as RL

- Why RL?

*We can not define the desired objective as a metric function.*

*We can not do back-prop through a black-box 🧠*

- Why not?

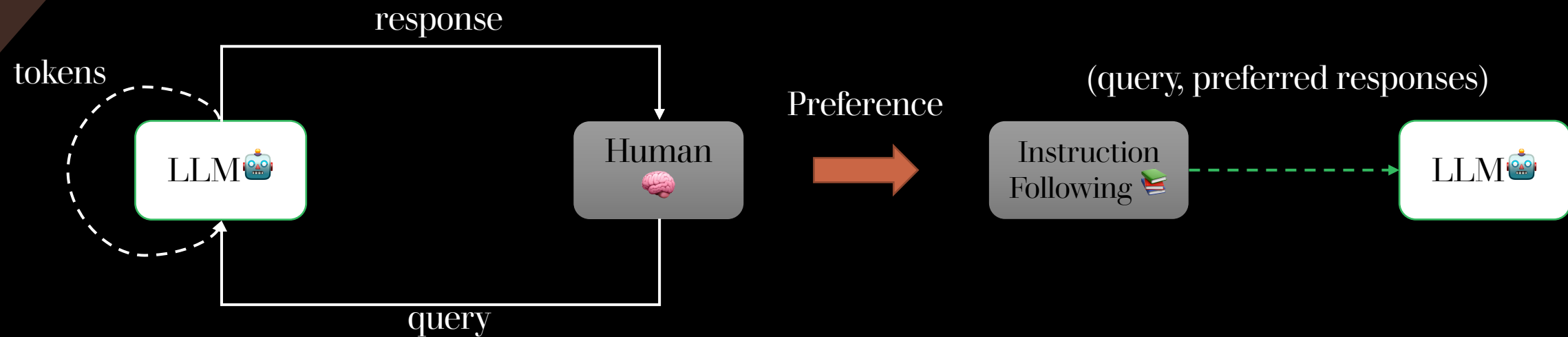
*Too expensive.*

*OpenAI: unless you have enough volunteers (users)...*

GPT-4?



# LLM Alignment with Human Feedback Logs



Preference Data Generation

Alignment as Offline-RL

# Alignment as Offline-RL: How to Learn?



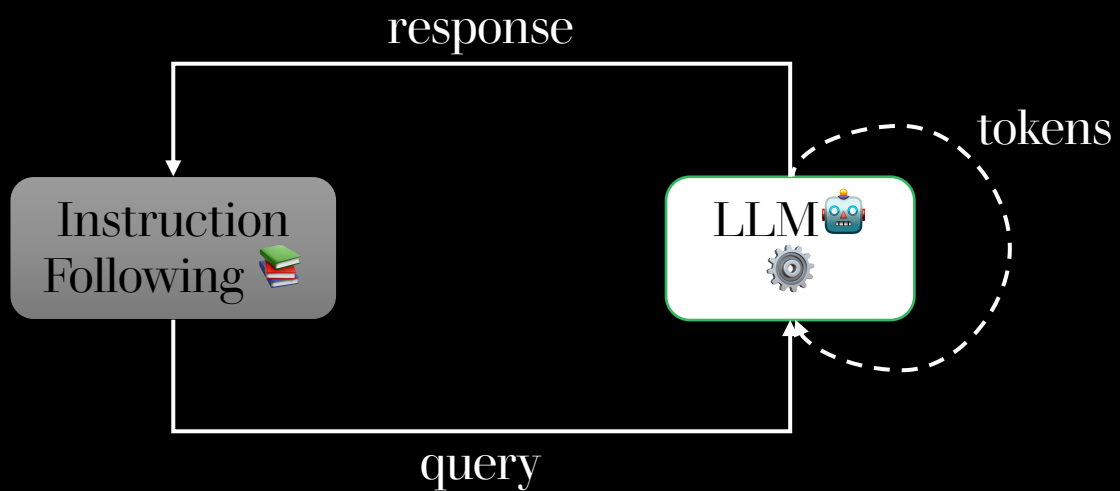
- We can always use behavior clone.  
*BC = SFT, supervised-fine-tuning*
- It is simple, stable, efficient.
- But language modeling is not 1-step decision.  
*Compounding errors*

Alignment as Offline-RL

---

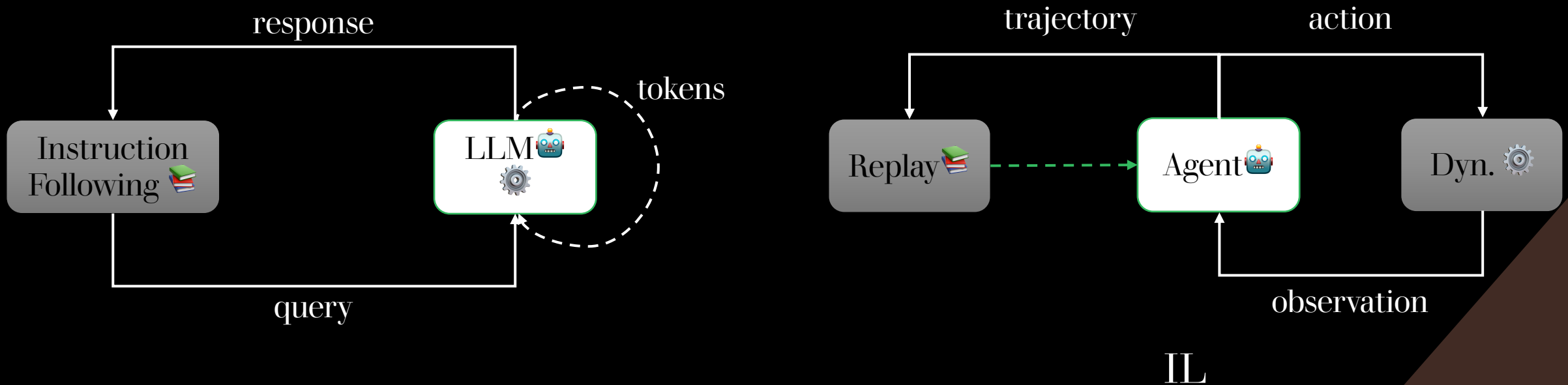
# What Makes LLM Alignment Special?

- The transition dynamics is deterministic and known!



# What Makes LLM Alignment Special?

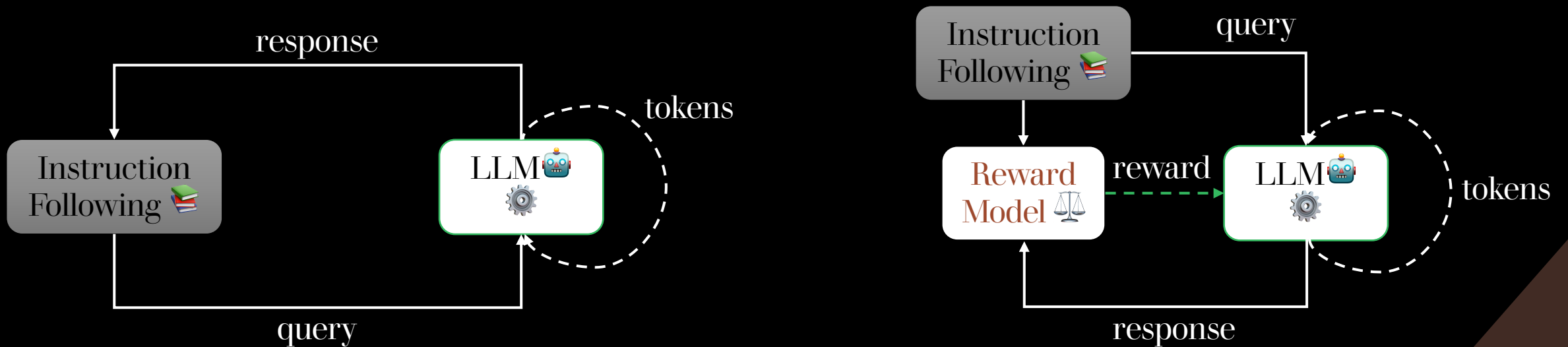
- The transition dynamics is deterministic and known!
- Recall the framework of Imitation Learning.





# RLHF: Solving Offline-RL via **Online** Inverse RL

- Inverse RL: learn the reward model, then optimize the policy.



# SFT vs RLHF\*: from the RL Perspective

- LLM alignment with logged human feedback (preference) can be interpreted as
  1. *Offline-RL --- solve it with behavior clone --- SFT*
  2. *Imitation Learning --- solve it with IRL --- RLHF*
- We can always do both:  
*RLHF uses SFT as a warm-start (OpenAI Alignment Paper)*
- Potential Alternatives? Someone would try GAIL...

---

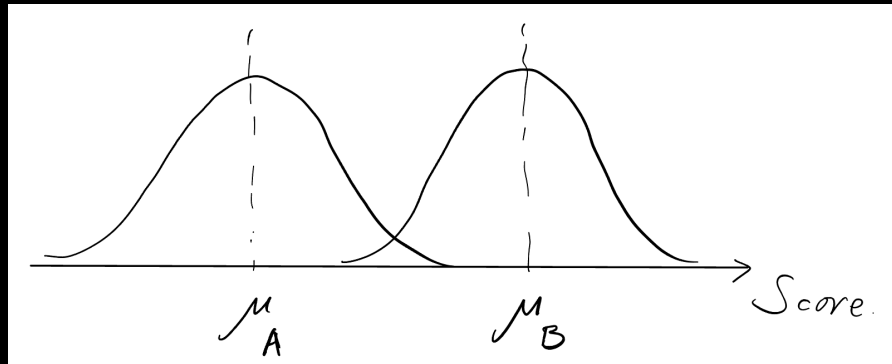
\* OpenAI's SFT is based on a separated high-quality response written by human.

# RLHF

- Underlying assumptions:
    1. *Learning a reward model is statistically **easier** than directly learning aligned LLMs.*
    2. *There are some higher-level metrics that **can not** be captured by token-level distances.*
  - Two steps
    1. *Reward Learning (Response Evaluation)*
    2. *LLM Optimization (Response Optimization)*
-

# Step 1: RM --- The Bradley-Terry Model

- Ranking is better than scoring, because the latter is noisier.
- Bradley-Terry Model is used to turn preferences into scores.



$$P(x_A > x_B | x_A \sim N(\mu_A, \sigma_A^2), x_B \sim N(\mu_B, \sigma_B^2))$$

$$P(x_A > x_B) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{\mu_A - \mu_B}{\sqrt{2(\sigma_A^2 + \sigma_B^2)}}\right)$$

# The Bradley-Terry Model

- In RLHF (and also many MOBA games), people use a slightly different function form

$$S_i = \exp[\beta_i]$$

$$P(i \succ j) = \frac{1}{1+S_j/S_i} = \frac{1}{1+e^{-(\beta_i-\beta_j)}} = \sigma(\beta_i - \beta_j) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{\beta_i - \beta_j}{2}\right)$$

- Equivalently: [DPO paper, Equation (1)]

$$p^*(y_1 \succ y_2 | x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))}$$

- Practical Optimization: Binary Classification

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

---

# Is B-T Model a Good Choice?

- In paired games, BT-model is used to attribute scores to different players

*evaluating player ability*

*fairness of game (trade off between waiting time)*

*number of players  $\ll$  number of games*

*applied in an online manner --- error correction!*

*reward is comparable: the same game*

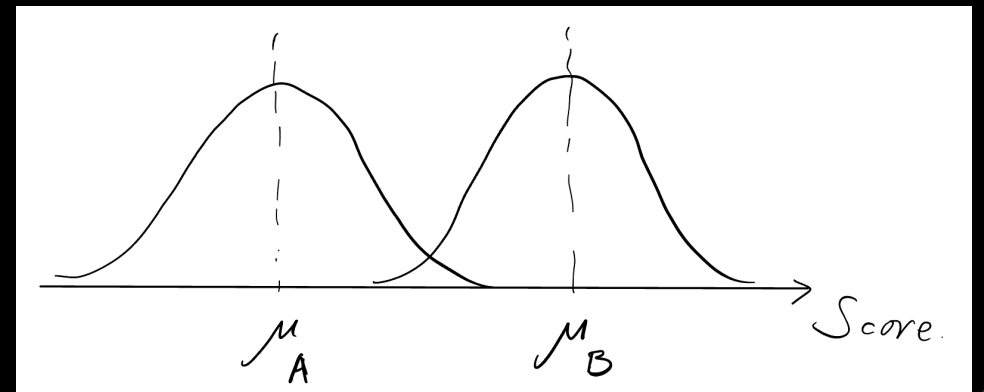
- In RLHF, labor annotation is noisy, and can be biased

*every labor + query- (paired) response = a game*

*number of players (query)  $\sim$  number of games*

*offline manner --- no error correction*

*is the reward value really comparable? E.g., Some queries can be toxic.*



# Reward Model Overoptimization

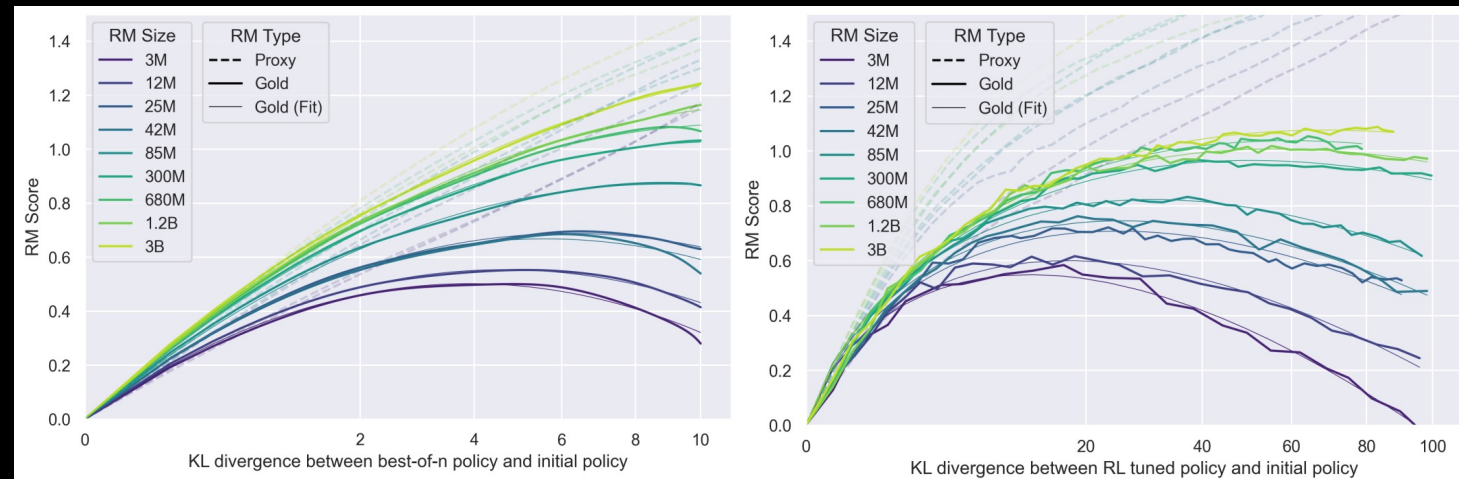
- Size of RM?

- LM with different sizes are used:

*OpenAI: 6B RM for 175B LM*

*DeepMind: 75B RM for 75B LM*

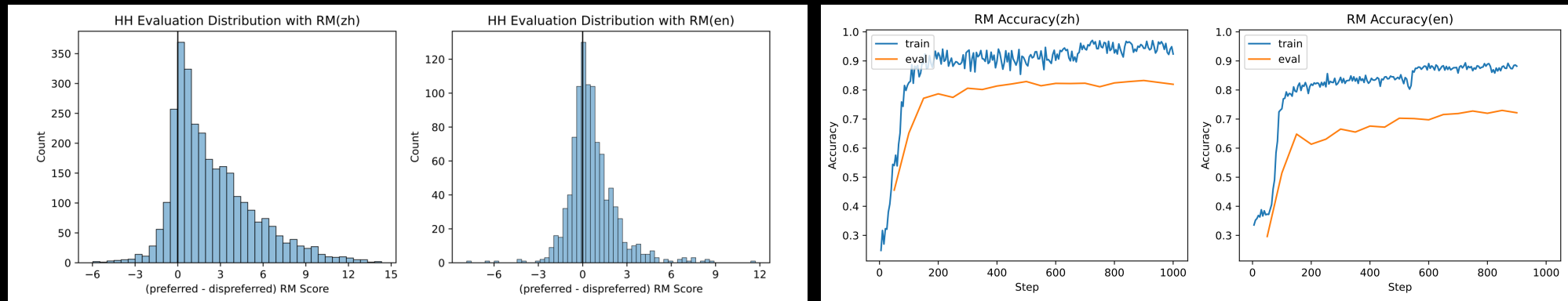
*Larger Model is Better*



- Core Idea: The RM should be able to **understand** responses.

# Some Evidence: RM Quality & Data Quality

- Model: LLaMA-7B / OpenChineseLLaMA
- Dataset (en): HH-RLHF 118k helpful + 42k harmless as train, 7.5k as val., 1k as test.
- Dataset (zh): annotated 31k helpful + 8k harmless, 30k train, 6k val., 3k as test.
- Results:





# Step 2. Learning with RM

- RL Algorithms

- PPO: ([Secrets of RLHF in Large Language Models](#))
- Best-of-N: Empirically better than PPO (but is not parameterized)

*Challenges:*

- multiple LLMs required. e.g., reference model, actor, critic, reward model.
- not stable, hard to train, sensitive to hyper-params & seeds.
- episodic return is a very sparse reward. (trajectory-level return)
- no dynamic programming structure

- Evolution Strategies can be a scalable alternative [Salimans et al. 2017]

- RAFT
- RRHF

*Sample a batch, and select the best using RM*

# Challenges of RLHF

- RM:
    - *Reward Overoptimization --- use larger model (other efficient choices?)*
    - *Noisy and Offline dataset --- use clean dataset (how to clean-up existing ones?)*
  - Policy Learning:
    - *Sparse reward --- dense reward or better credit assignment (e.g., hindsight)*
    - *Is it possible to learn a dense reward model? --- maybe use a new dataset*
    - *Conservative update: do not need to change the LLM too much (e.g., BoN KL-div)*
-

# DPO: Implicit Imitation Learning

- Soft-Q-Learning: Not using arg-max, but soft-max. (exponential sum over q-values)
- Motivation: better exploration (max-ent RL)

$$\pi_{\text{MaxEnt}}^* = \arg \max_{\pi} \sum_t \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_{\pi}} [r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_t))],$$

$$\eta(\pi) = \mathbb{E}_{a \sim \pi, r} [r] - \tau D_{\text{KL}}[\pi \parallel \bar{\pi}]$$

- A similar objective in RLHF (Eqn.3 in DPO):

**RL Fine-Tuning Phase:** During the RL phase, we use the learned reward function to provide feedback to the language model. In particular, we formulate the following optimization problem

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta}(y | x) \parallel \pi_{\text{ref}}(y | x)] \quad (3)$$

# DPO: Implicit Imitation Learning

- The optimal policy has a closed-form (expressed as function of reference policy and reward)

$$\eta(\pi) = \mathbb{E}_{a \sim \pi, r} [r] - \tau D_{\text{KL}} [\pi \parallel \bar{\pi}]$$

$$\pi_{\bar{r}}^{\mathcal{B}}(a) = \bar{\pi}(a) \exp(\bar{r}(a)/\tau) / \underbrace{\mathbb{E}_{a' \sim \bar{\pi}} [\exp(\bar{r}(a')/\tau)]}_{\text{normalizing constant}}.$$

- DPO put this result in preference-based learning, and cancel-out the normalizing constant

$$r(x, y) = \beta \log \frac{\pi_r(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x).$$

$$p^*(y_1 \succ y_2 | x) = \frac{1}{1 + \exp \left( \beta \log \frac{\pi^*(y_2|x)}{\pi_{\text{ref}}(y_2|x)} - \beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)} \right)}$$

- “End-to-end”—direct optimization / a smart idea / overoptimization? (B-o-128>DPO)

# Instruction Following by Prompting

- Prompt engineering is an effective approach in eliciting the abilities of LLMs
- In-context Learning/Fine-tuning

*Few-Shot Prompting + In-Context Learning*

*Zero-Shot Prompting*

- Examples:

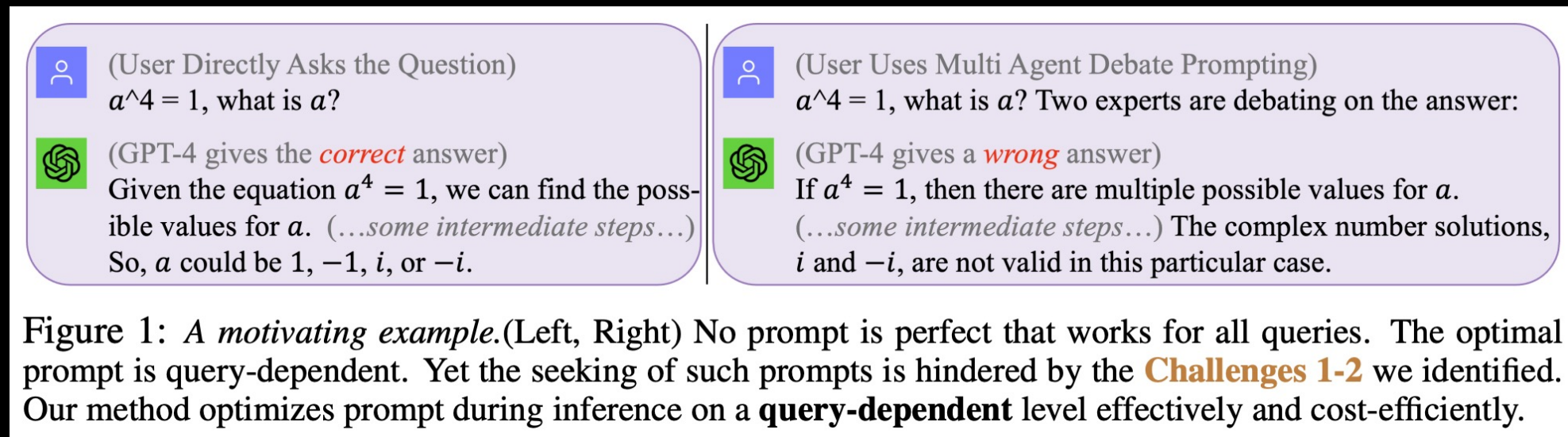
*CoT: let's think step by step...*

*OPRO: take a deep breath ...*

- How to design? Previous approaches: learning from **trial and error**.
-

# Prompt-OIRL

- Offline Prompt Evaluation and Optimization with IRL [Sun et al. 2025]



# Prompting Using RL?

- RL is Expensive.

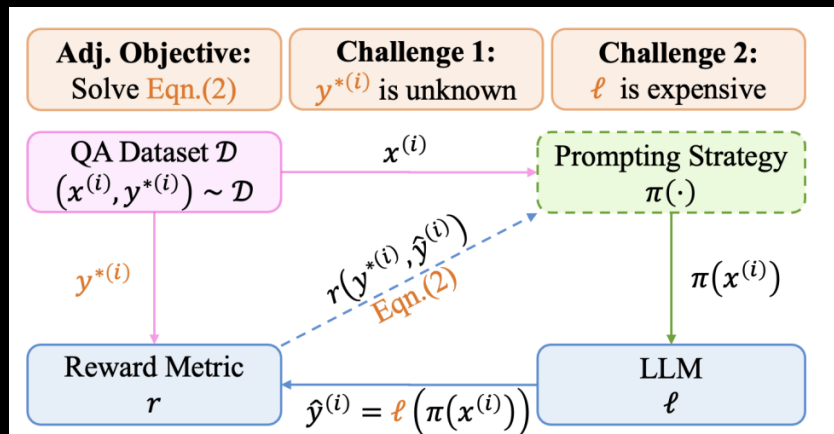


Figure 2: *The Adjusted Objective and Challenges in prompt optimization.* We use blue to denote fixed functions, pink for datasets, and green for functions to be optimized. Solid lines show the flow of outputs, and dashed lines denote the learning process.

How about using existing expert demonstrations?

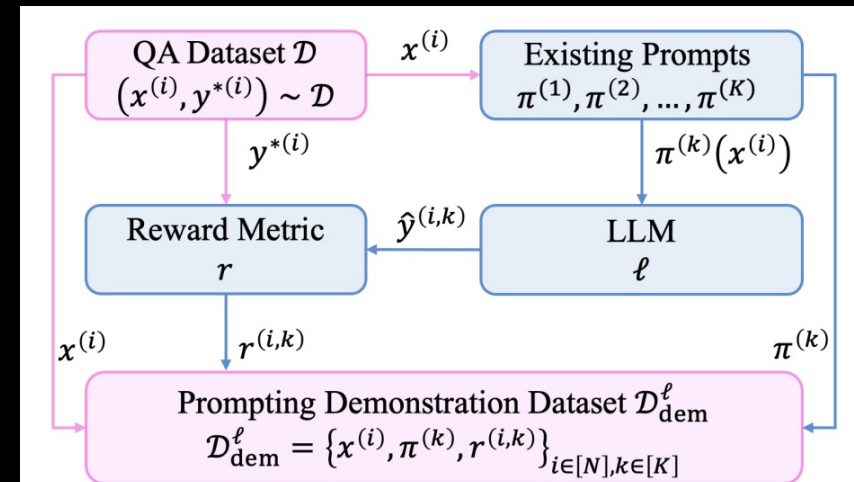
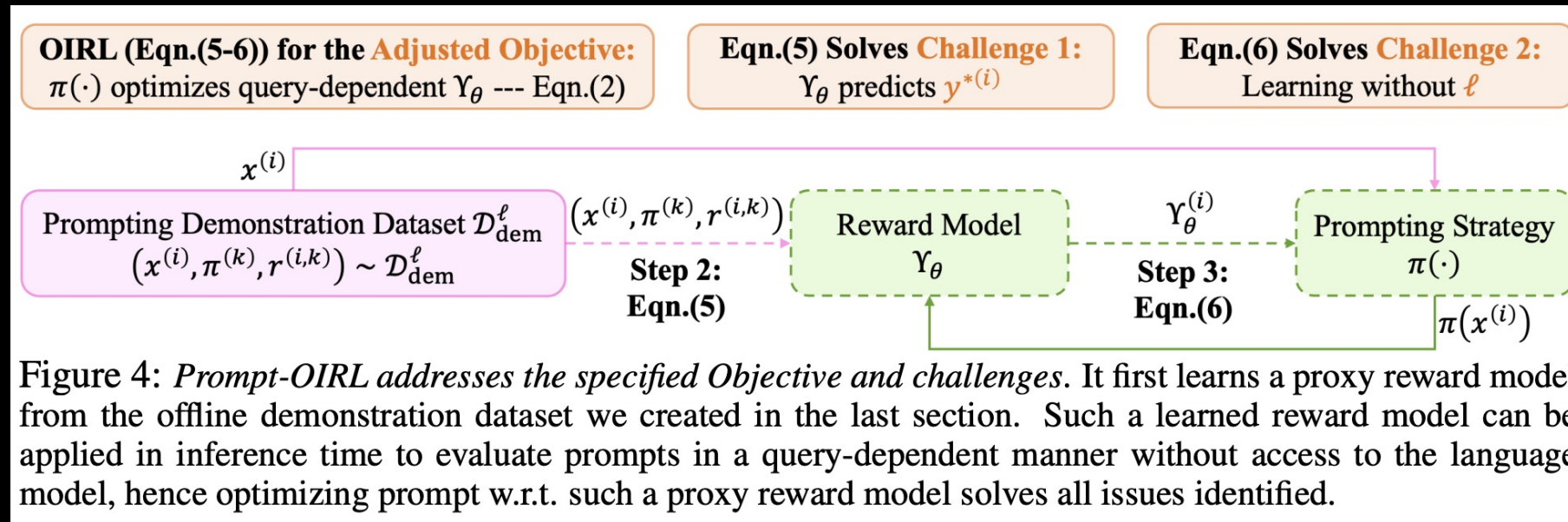


Figure 3: *The offline demonstration dataset is generated as a by-product of evaluating existing (query-agnostic) prompts.*

# Offline Inverse RL --- RLAIIF

- Reward model estimate the *preference* of LLMs



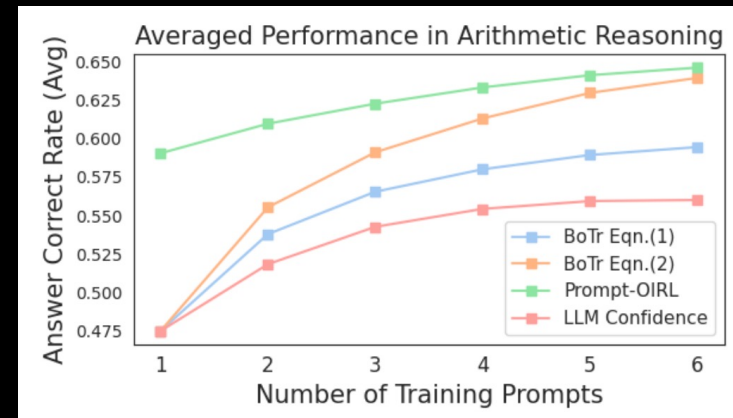


# Results

- Experiments on Arithmetic Reasoning Datasets (GSM8K, SVAMP, MAWPS)
- TakeAways:
  1. *Prompt-OIRL further improve the ability of LLMs in inference.*
  2. *It is extremely cheap to train and deploy Prompt-OIRL.*

Table 2: prompts used in offline training dataset collection.

No.	Effective Prompts Discovered by Experts and Algorithms	Explanation
1	“The answer is:”	direct prompting
2	“Let’s think step by step:”	zero-shot CoT (Kojima et al., 2022)
3	“Let’s work this out in a step by step way to be sure we have the right answer:”	APE discovered (Zhou et al., 2022b)
4	“First, decompose the question into several sub-questions that need to be solved, and then solve each question step by step:”	Least-to-most (Zhou et al., 2022a)
5	“Imagine three different experts are answering this question. All experts will write down 1 step of their thinking, and then share it with the group. Then all experts will go on to the next step, etc. If any expert realizes they’re wrong at any point then they leave.”	Tree-of-thought (Hulbert, 2023)
6	“3 experts are discussing the question, trying to solve it step by step, and make sure the result is correct:”	multi-agent debate (Liang et al., 2023)



# Summary

- RL is learning from trial and errors to maximize a cumulative reward.
  - Define reward function can be easy, but exploration of RL is hard.
  - With expert demonstrations, IL can improve learning efficiency.
  - Behavior Clone is the simplest IL, but it suffers from compounding errors.
  - IRL first learns a RM, and then use the learned RM to optimize policy.
  - SFT is behavior clone, RLHF is online IRL.
  - Assumptions under Reward Model Learning --- The Bradley-Terry Model
  - Overoptimization in Reward Model Learning
  - Given an RM, there are multiple approaches to optimize LLMs to align with human.
  - DPO, and Soft-Q-Learning
  - (Our Recent Work) Prompt-OIRL is able to effectively and efficiently perform offline prompt evaluation and optimization.
-