

# Evolution and Application of Reinforcement Learning (RL) in Large Language Models (LLMs)

Networked Intelligence for Comprehensive Efficiency (NICE) Lab  
College of Information Science and Electronic Engineering  
Zhejiang University  
<http://nice.rongpeng.info/>



Aug. 12, 2025

# Content



## 1 Development Outline

- Overview: From Alignment to Advanced Reasoning
- Emergence: Alignment Tuning
- Development: Enhanced Reasoning

## 2 Evolutionary History of the System Framework

- RLHF\_v0
- RLHF\_v1
- RLHF\_v2
- RLHF\_v3
- System Framework

## 3 Technical Classifications

- Reward Design
- Critic Design
- Challenges in Reward Modeling

# Development Outline



# Overview

## From Alignment to Advanced Reasoning

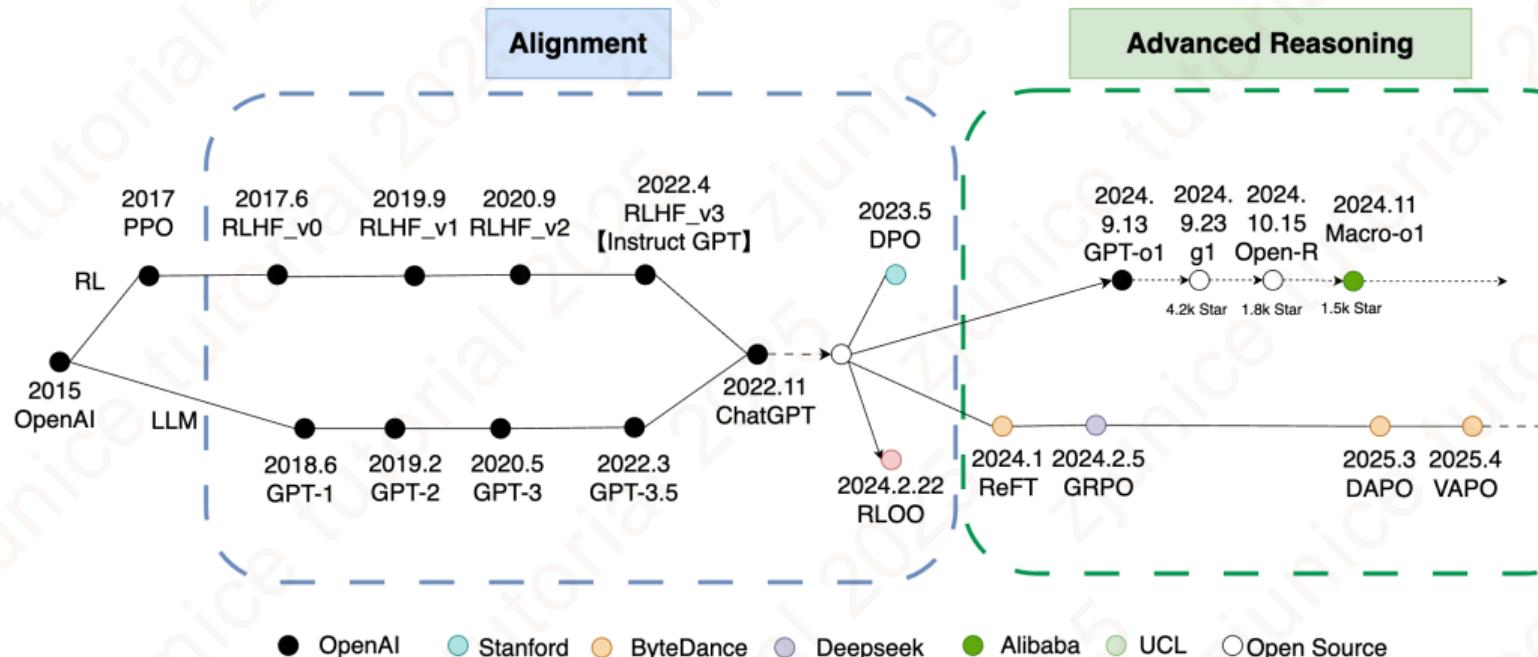


Figure: Roadmap of RL in LLM.

# Emergence: Alignment Tuning



■ **Goal:** Align the behavior of LLMs with **human preferences**.

e.g., *helpfulness, honesty, and harmlessness*.

■ **Motivation:** Reduce the reliance on costly human annotations.

■ **Typical Approaches:**

1 **With Reward Model (RM):**

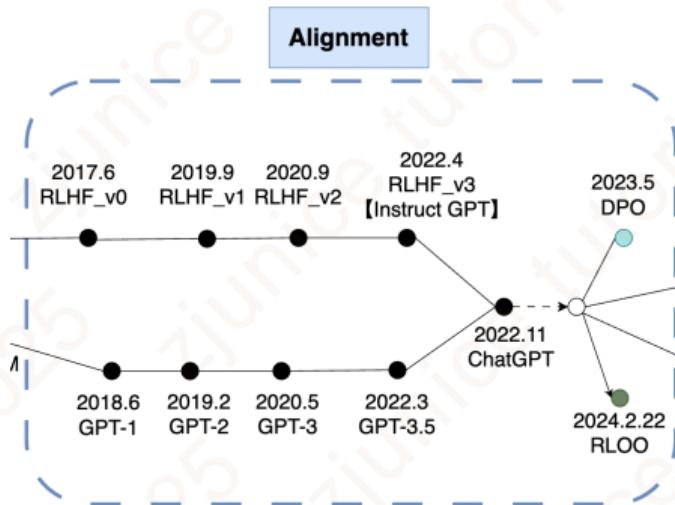
Train a separate RM to approximate human preferences; then fine-tune the LLM using PPO-based RLHF.

e.g., *InstructGPT*.

2 **Without RM:**

Directly learn from preference pairs via implicit objective optimization.

e.g., *DPO*.



# Development: Enhanced Reasoning



- **Goal:** Directly optimize final **task performance**.  
e.g., *accuracy, usefulness, reasoning quality*.
- **Motivation:** Overcome the limitations of Supervised Fine-Tuning (SFT), which is constrained by fixed reasoning paths and data quality.

## ■ Typical Approaches:

### ■ Without Training:

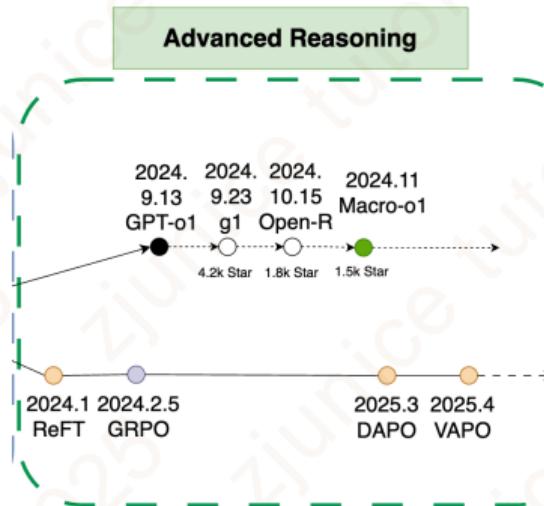
Reasoning path search methods.

e.g., *Monte Carlo Tree Search (MCTS)* +  
*Chain-of-Thought (CoT)/Chain-of-Action(CoA)*,  
*Prompt Engineering*.

### ■ With Training:

Fine-tuning via RL with rewards.

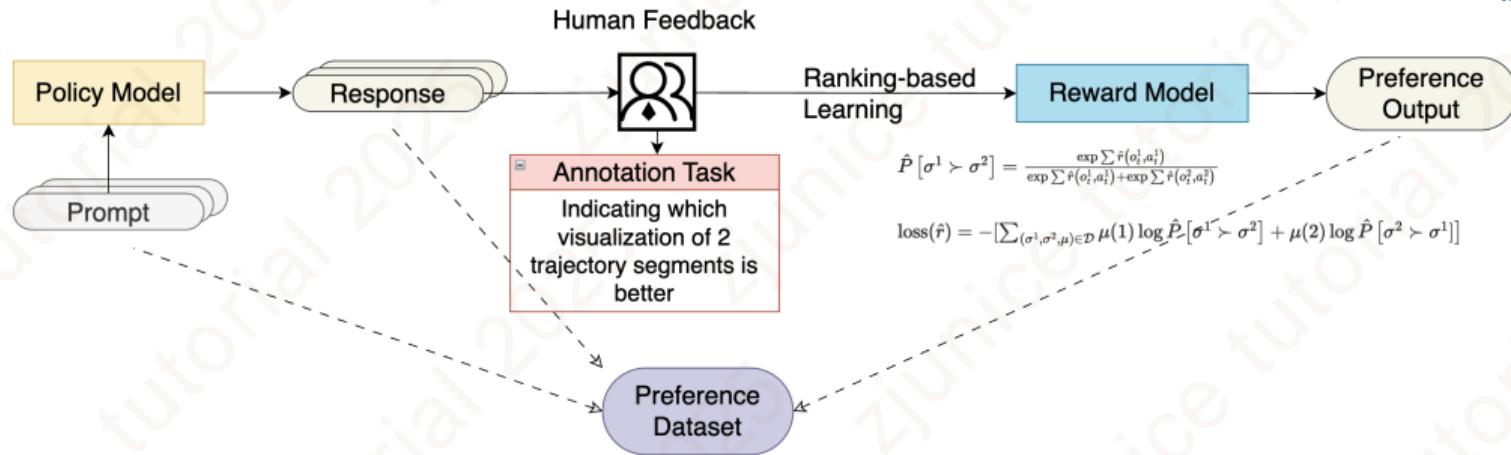
e.g., *ReFT, GRPO*.



# Evolutionary History of the System Framework



# System Framework 1: RLHF\_v0<sup>1</sup>

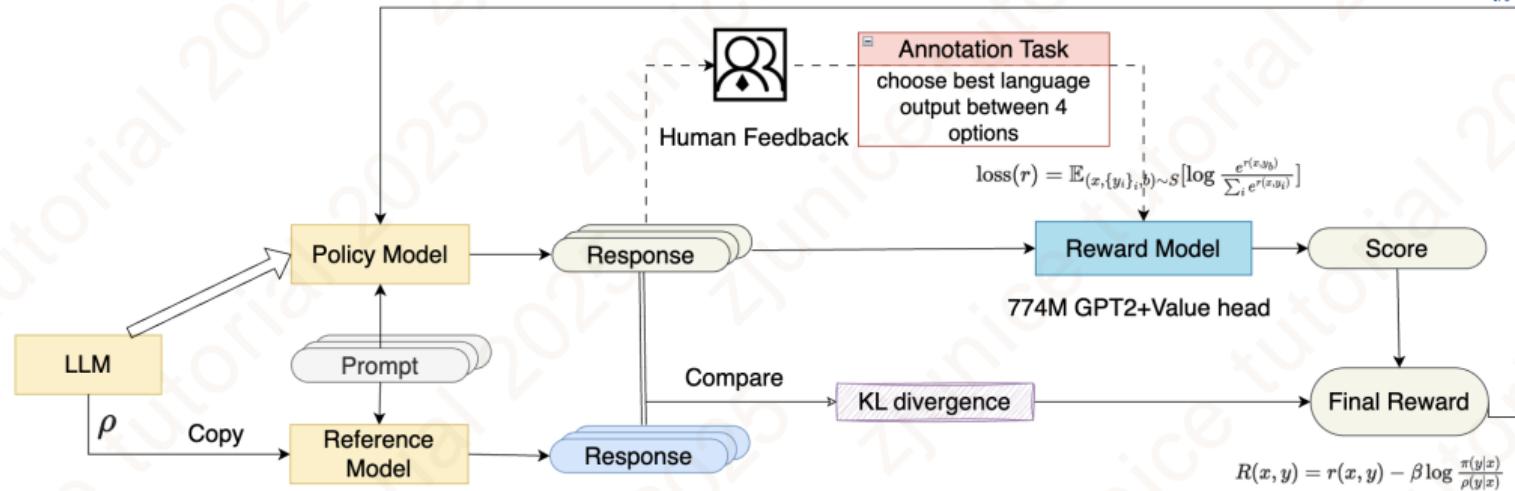


- **Task:** traditional RL benchmarks (e.g., Atari games).
- **Core Contribution:**
  - Preference-Based **Reward Modeling**;
  - Policy Optimization with Learned Rewards.

<sup>1</sup>P. F. Christiano et al., "Deep reinforcement learning from human preferences," 2017.



# System Framework 2: RLHF\_v1<sup>2</sup>

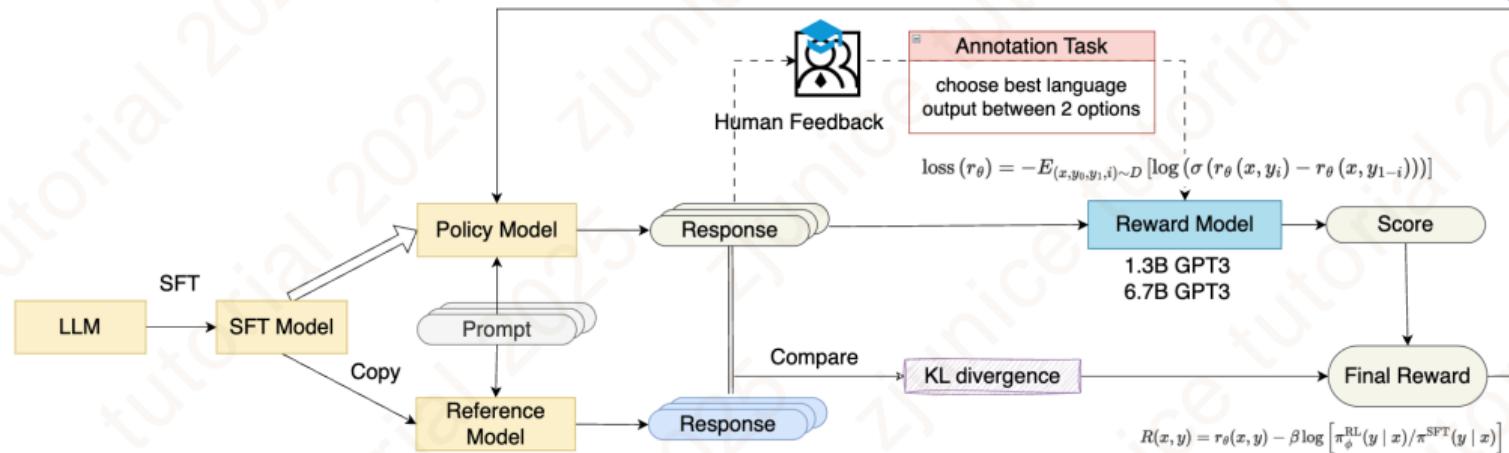


- **Task:** First RLHF for LLMs.
- **Core Contribution:**
  - Novel **Reward Model Initialization** and **KL Control**;
  - Online Data Collection to Mitigate Distribution Shift.

<sup>2</sup>D. M. Ziegler et al., "Fine-tuning language models from human preferences," 2019.



# System Framework 3: RLHF\_v2<sup>3</sup>



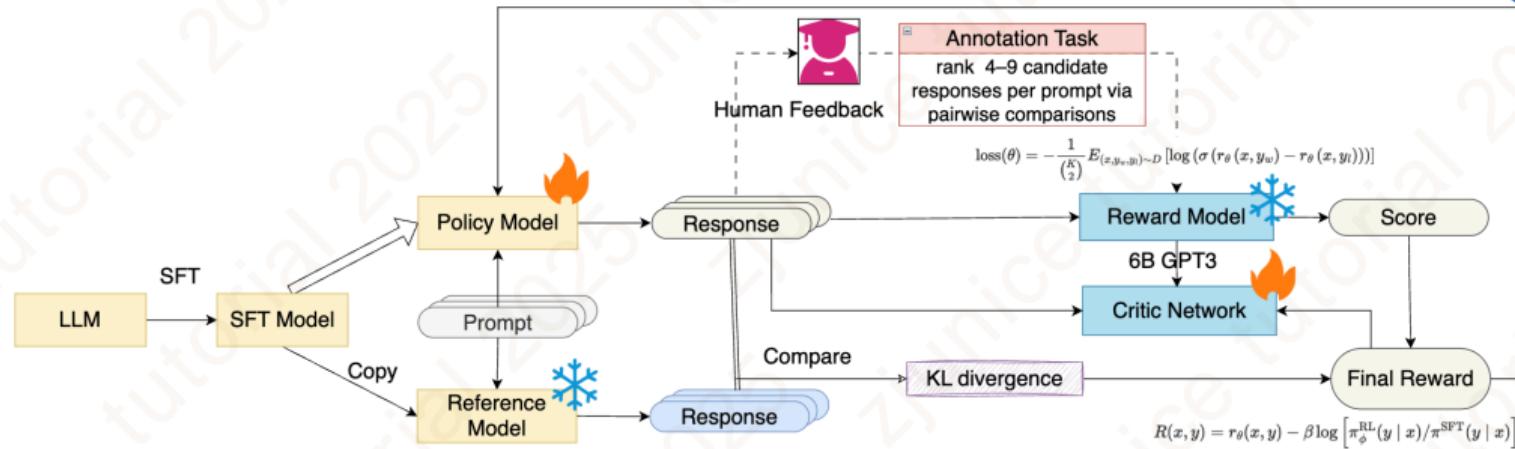
## Core Contribution:

- Expert Annotators & **Pairwise Preference Labeling**;
- Improved Alignment through **SFT + RLHF Pipeline**;
- Reward Model Initialization from SFT Policy.

<sup>3</sup>N. Stiennon et al., "Learning to summarize with human feedback," 2020.



# System Framework 4: RLHF\_v3 (InstructGPT)<sup>4</sup>



## Core Contribution:

- Higher-Quality Expert Labelers;
- Multi-Response Generation with **Full Pairwise Ranking** ;
- Scaling Reward Model with GPT-3.

<sup>4</sup>L. Ouyang et al., "Training language models to follow instructions with human feedback.", 2022.



# System Framework of RL in LLM

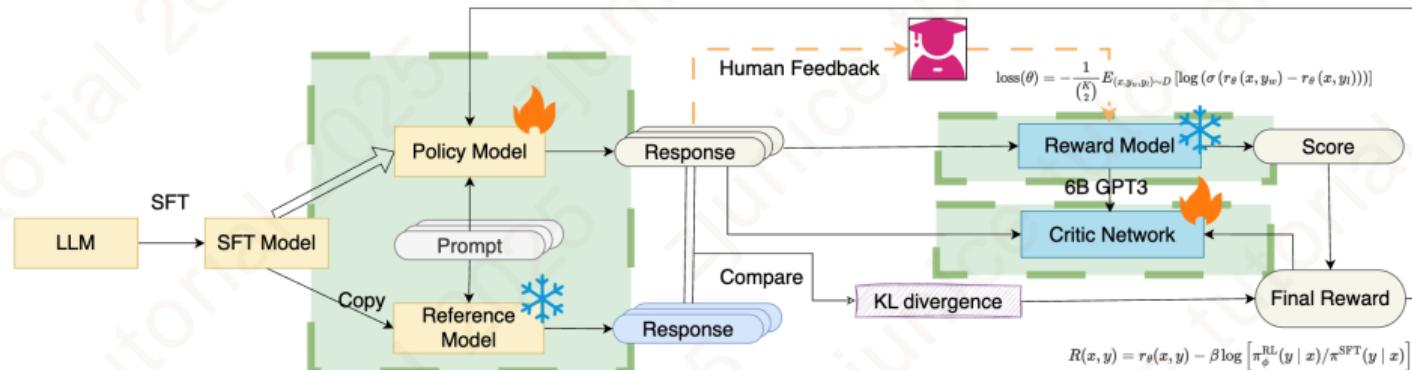


Figure: Complete workflow of RLHF.

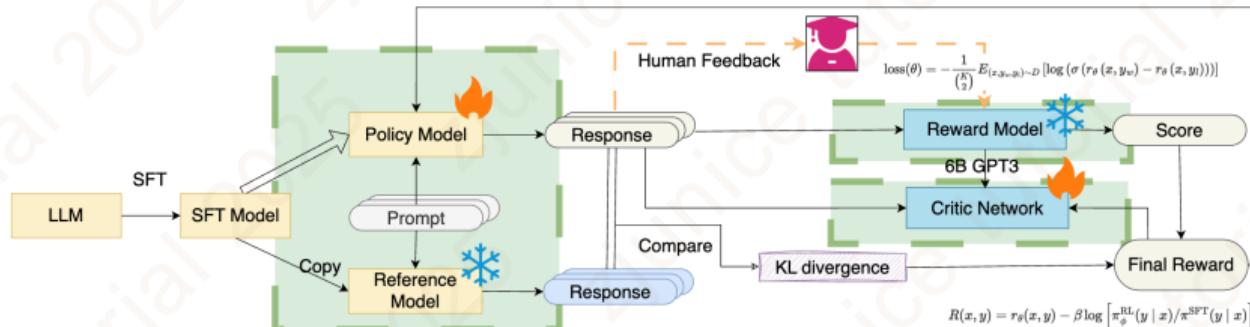
## Components

- Actor
- Critic (✗)
- Reward (✗)

# Technical Classifications



# Technical Classifications 1: Reward



■ **Implicit Rewards:** Learn from comparisons. e.g., [DPO](#)

■ **Explicit Rewards:**

■ **By source:** [source illustration](#)

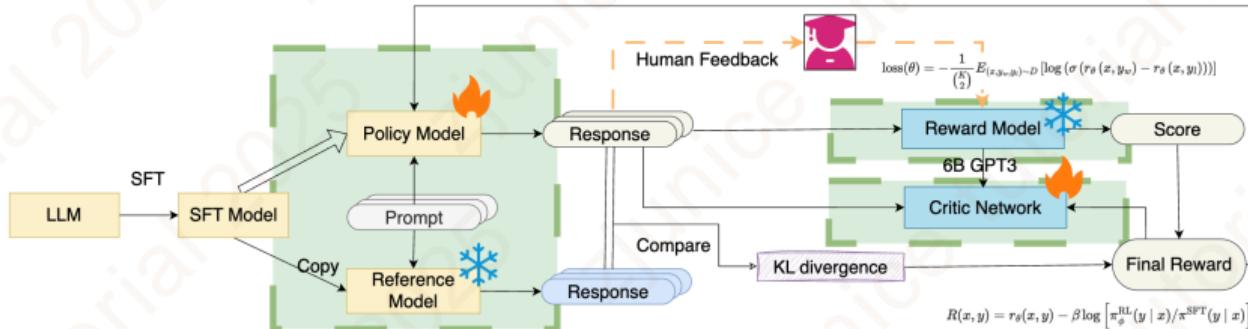
- 1 **External:** RM is a separate small LLM.
- 2 **Internal:** shares parameters with actor, with a reward head.

■ **By temporal granularity:** [PRM vs. ORM illustration](#)

- 1 **Outcome Reward Model (ORM):** Outcome-level. e.g., [GRPO](#)
- 2 **Process Reward Model (PRM):** Sequence- or token-level. e.g., [PRIME](#)



# Technical Classifications2: Critic Design



## With Critic:

Estimates expected value for current policy.

- 1 Initialized from RM. e.g., *InstructGPT*
- 2 Initialized from policy and perform Value-Pretraining. e.g., *VAPO*

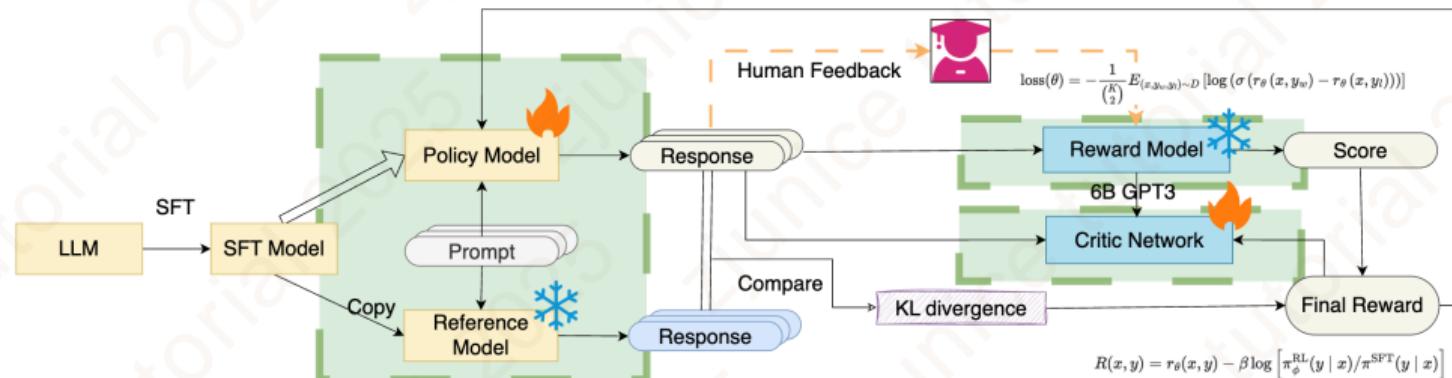
## Without Critic:

Relative ranking used as signal (no value function).

Key idea: Avoids bias from value estimation. e.g., *RLOO*, *DAPO*, *GRPO*



# Challenges in Reward Modeling



- **Distribution Shift/Overoptimization:** Reward model becomes misaligned.
- **Reward Hacking:** Model exploits reward shortcuts.
- **Reward Sparsity:** Good samples are rare.
- **Cost:** Human preference labels are expensive.
- **Non-additivity:** Mixed (PRM and ORM) reward scales cause instability.

# Thank you

Networked Intelligence for Comprehensive Efficiency (NICE) Lab  
College of Information Science and Electronic Engineering

Zhejiang University  
<https://nice.rongpeng.info>

## Appendix

# Source of Reward Model

## LLM + Value Head (Linear Layer)

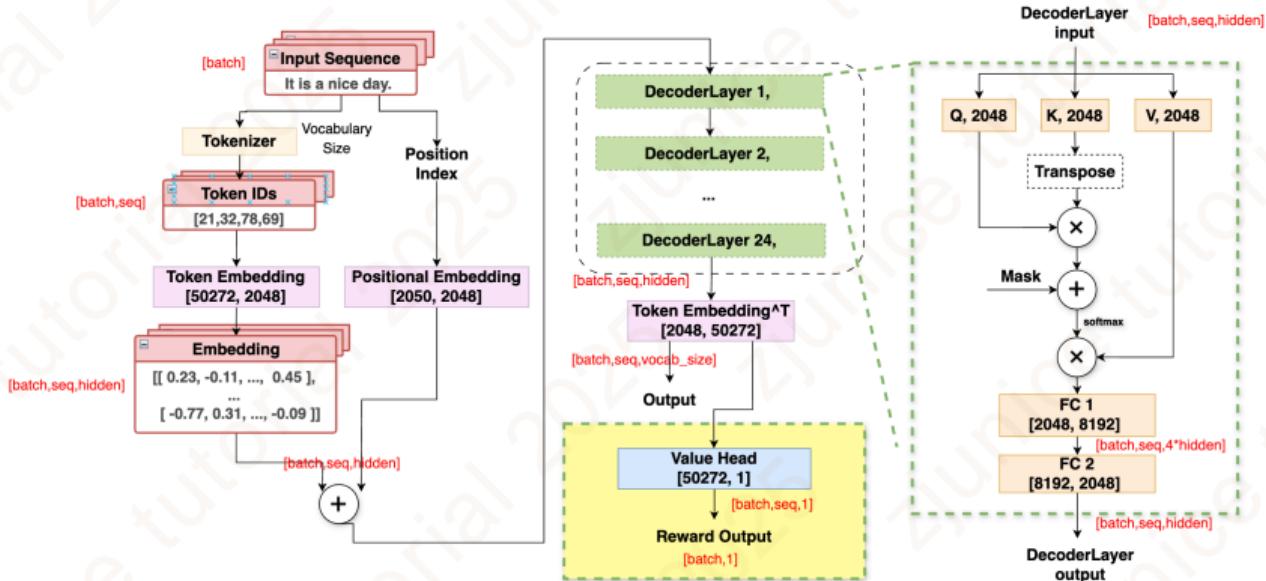


Figure: external vs. internal source of RM.

[Back Link to reward](#)

# Outcome Reward Model (ORM) vs. Process Reward Model (PRM)



**Problem:** Let  $p(x)$  be a monic polynomial of degree 4. Three of the roots of  $p(x)$  are 1, 2, and 3. Find  $p(0) + p(4)$ .

**Golden Answer:** 24

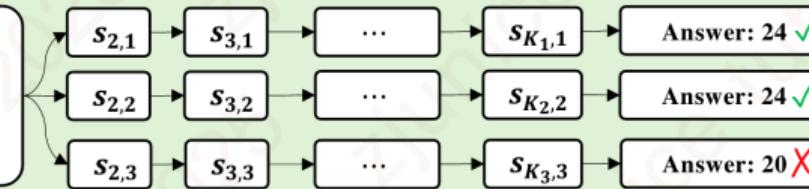
**Solution:**  $S = s_1, s_2, s_3, \dots, s_K$

**Answer:** 20 X

**(a) Outcome Annotation:**  $y_S = 0$

**Problem:** ....

**s<sub>1</sub>:** Since three of the roots of  $p(x)$  are 1, 2, and 3, we can write :  $p(x) = (x - 1)(x - 2)(x - 3)(x - r)$ .



**(b): Process Annotation:**  $y_{s_1}^{SE} = \frac{2}{3}$  ;  $y_{s_1}^{HE} = 1$

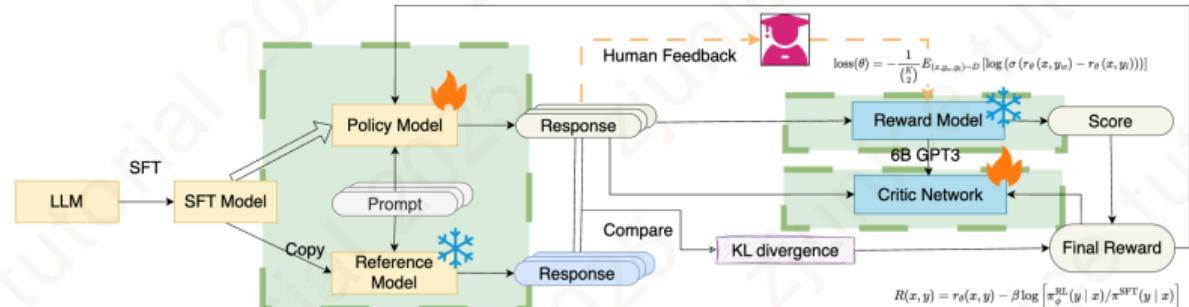
$s_i$ : the  $i$ -th step of the solution  $S$ .     $s_{ij}$ : the  $i$ -th step of the  $j$ -th finalized solution.

**Figure:** Example of different RM in Math-Shepherd<sup>5</sup>.

[Back Link to reward](#)

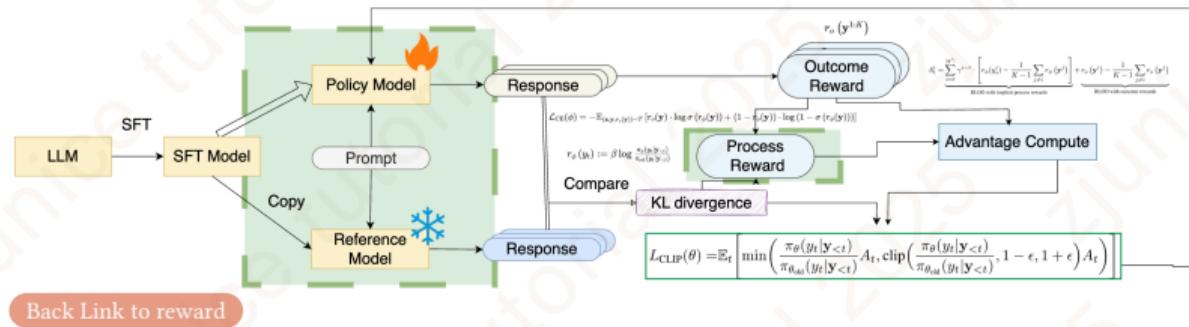
<sup>5</sup>P. Wang et al., "Math-shepherd: Verify and reinforce llms step-by-step without human annotations," in Proc. ACL, 2024.

# PRIME: Process Reinforcement through IMplicit rEwards<sup>6</sup>



## Component

- Actor
- Critic
- Reward

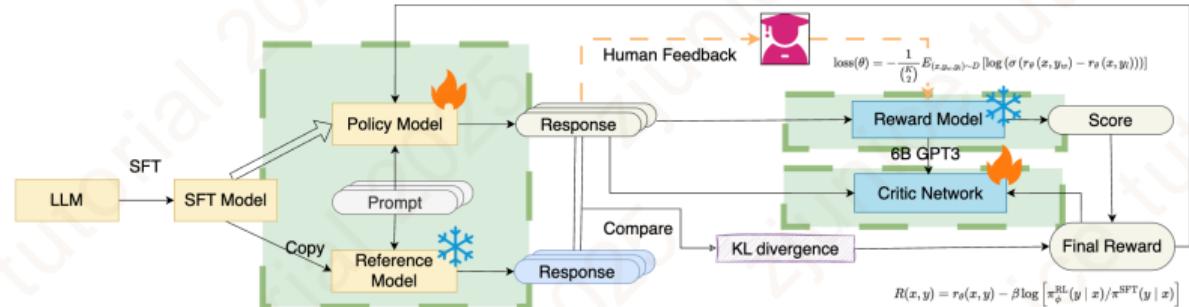


## PRIME

- PRM  
Supervision from ORM.

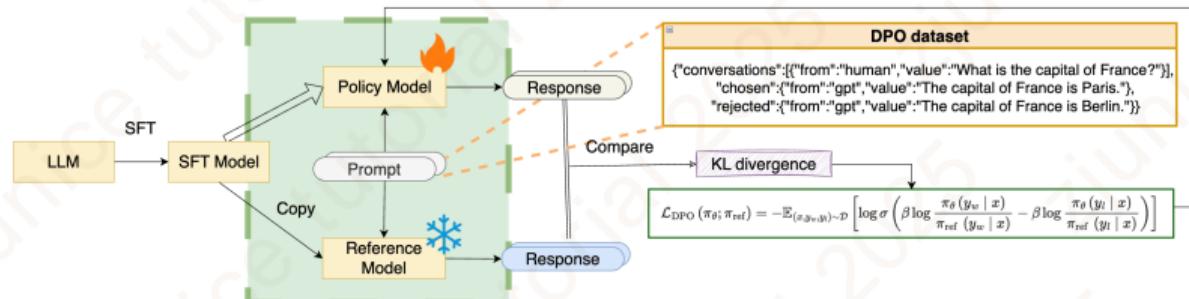
<sup>6</sup>G. Cui et al., "Process reinforcement through implicit rewards," *arXiv preprint arXiv:2502.01456*, 2025 .

# DPO: Direct Preference Optimization<sup>7</sup>



## Component

- Actor
- Critic
- Reward



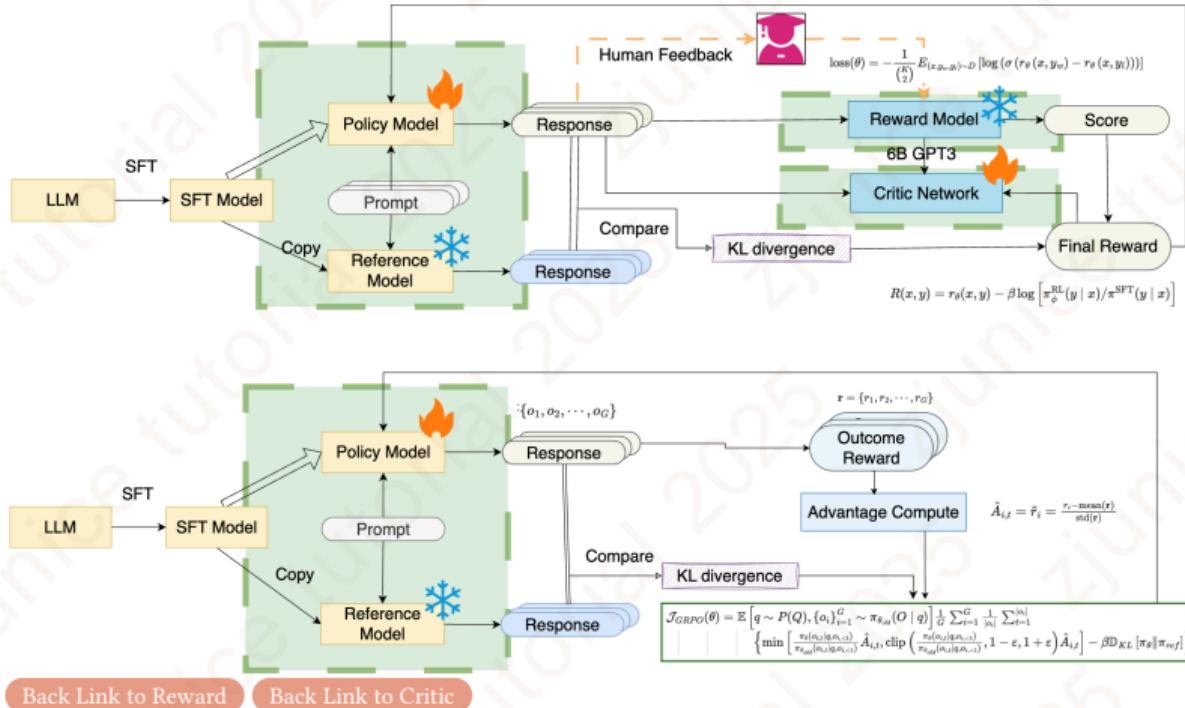
## DPO

- pairwise comparison

[Back Link to Reward](#)

<sup>7</sup>R. Rafailov et al., “Direct preference optimization: Your language model is secretly a reward model,” 2023.

# GRPO: Group Relative Policy Optimization<sup>8,9</sup>



## Component

- Actor
- Critic
- Reward

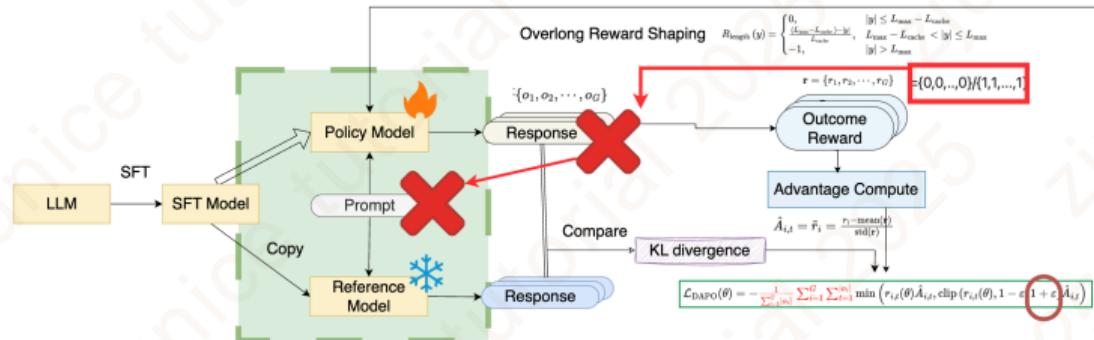
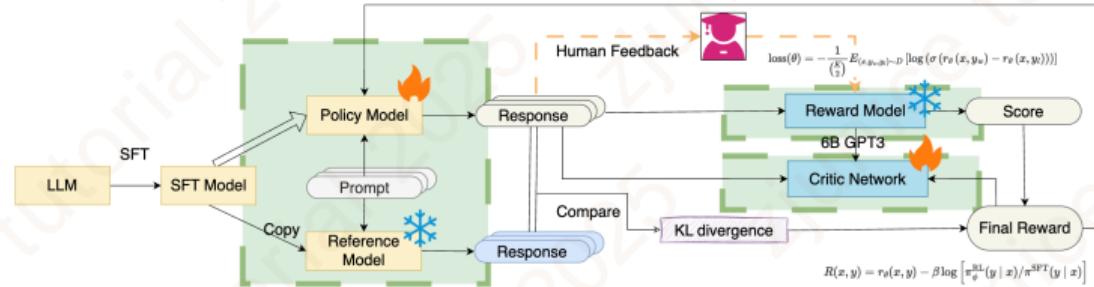
## GRPO

- Average group advantage.

<sup>8</sup>Z. Shao et al., “Deepseekmath: Pushing the limits of mathematical reasoning in open language models,” *arXiv preprint arXiv:2402.03300*, 2024 .

<sup>9</sup>D. Guo et al., “Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,” *arXiv preprint arXiv:2501.12948*, 2025 .

# DAPO: Decoupled Clip and Dynamic sAmpling Policy Optimization<sup>10</sup>



[Back Link to Critic](#)

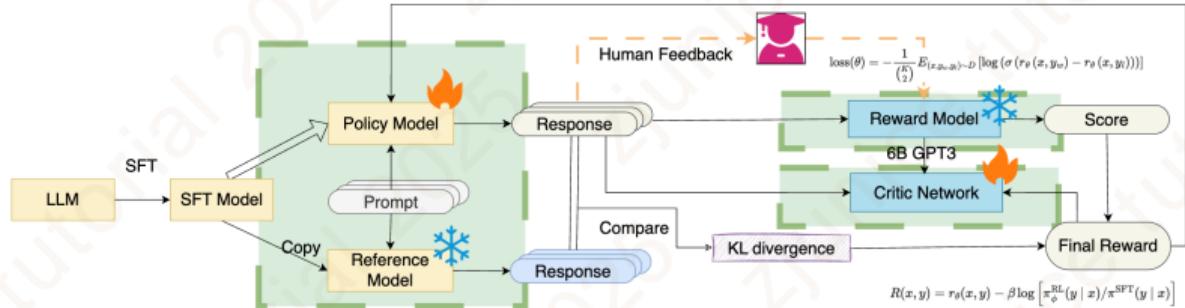
- Critic
- Reward

## DAPO

- Remove KL.
- Clip-Higher.
- Dynamic Sampling.
- Token-Level Policy Gradient Loss.
- Overlong Reward Shaping

<sup>10</sup>Q. Yu et al., "Dapo: An open-source llm reinforcement learning system at scale," *arXiv preprint arXiv:2503.14476*, 2025 .

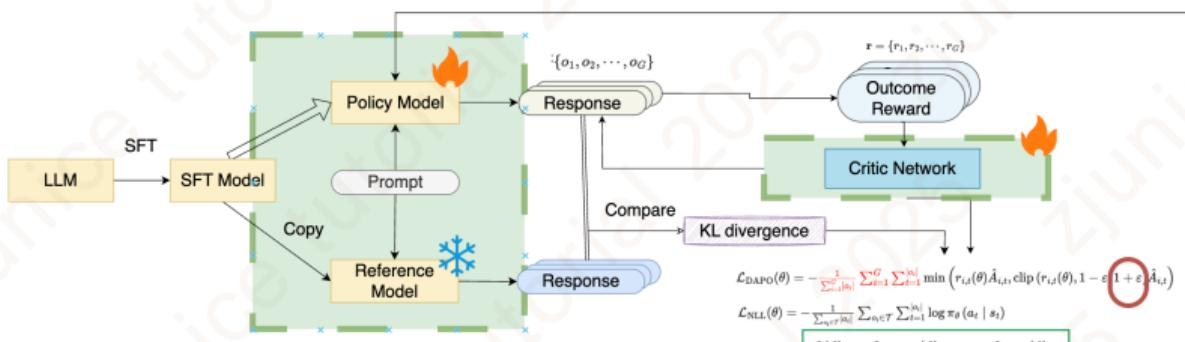
# VAPO: Value Augmented proximal Policy Optimization<sup>11</sup>



$$R(x, y) = r_\theta(x, y) - \beta \log \left[ \frac{\pi_\phi^\text{RL}(y | x)}{\pi_\phi^\text{SFT}(y | x)} \right]$$

## Component

- Actor
- Critic
- Reward



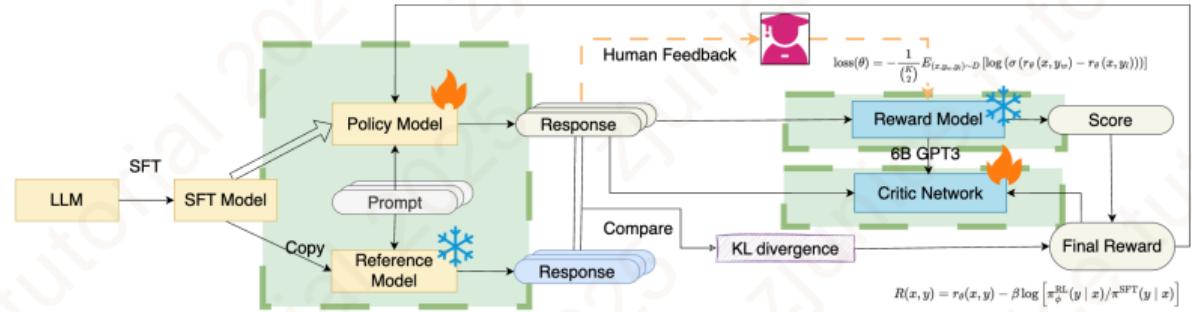
Back Link to Critic

## VAPO

- Value-Pretraining.
- Positive Example LM Loss.

<sup>11</sup>Y. Yue et al., "Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks," *arXiv preprint arXiv:2504.05118*, 2025.

# RLOO:REINFORCE Leave-One-Out<sup>12</sup>

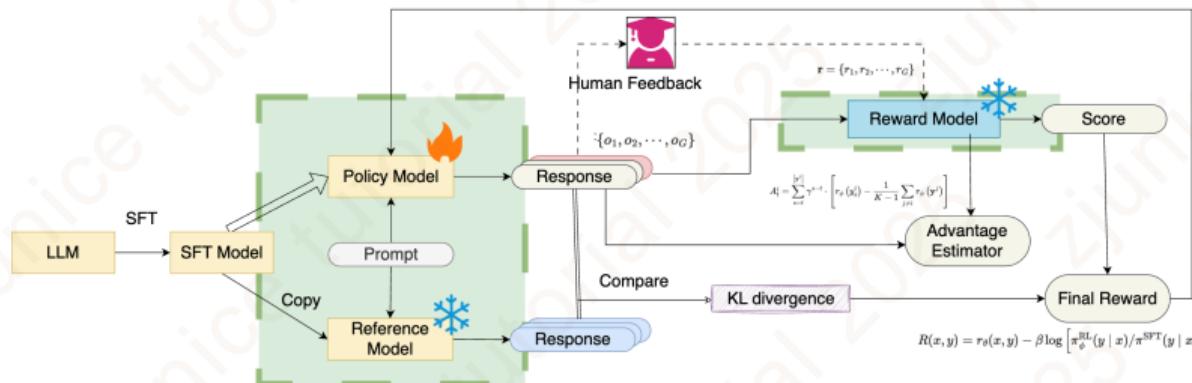


## Component

- Actor
- Critic
- Reward

## RLOO

- The average of the other samples in the group.



[Back Link to Critic](#)

<sup>12</sup>A. Ahmadian et al., "Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms," *arXiv preprint arXiv:2402.14740*, 2024.