# Reinforcement Learning Enhanced LLMs: A Survey

**Shuhe Wang**♠, **Shengyu Zhang**♣, **Jie Zhang**★, **Runyi Hu**▲, **Xiaoya Li**♦
**Tianwei Zhang**▲, **Jiwei Li**♣, **Fei Wu**♣, **Guoyin Wang, Eduard Hovy**♠

## Abstract

This paper surveys research in the rapidly growing field of enhancing large language models (LLMs) with reinforcement learning (RL), a technique that enables LLMs to improve their performance by receiving feedback in the form of rewards based on the quality of their outputs, allowing them to generate more accurate, coherent, and contextually appropriate responses. In this work, we make a systematic review of the most up-to-date state of knowledge on RL-enhanced LLMs, attempting to consolidate and analyze the rapidly growing research in this field, helping researchers understand the current challenges and advancements. Specifically, we (1) detail the basics of RL; (2) introduce popular RL-enhanced LLMs; (3) review researches on two widely-used reward model-based RL techniques: Reinforcement Learning from Human Feedback (RLHF) and Reinforcement Learning from AI Feedback (RLAIF); and (4) explore Direct Preference Optimization (DPO), a set of methods that bypass the reward model to directly use human preference data for aligning LLM outputs with human expectations. We will also point out current challenges and deficiencies of existing methods and suggest some avenues for further improvements.

## 1 Introduction

Large language models (Jiang et al., 2023; OpenAI, 2023; Dubey et al., 2024) are sophisticated language models pre-trained on extensive text data, allowing them to produce coherent and fluent responses to diverse inputs. However, the interaction capabilities of these pre-trained LLMs can be inconsistent, sometimes leading to responses that,

while technically correct, may be harmful, biased, misleading, or irrelevant to users' needs. Therefore, it is crucial to align the outputs of pre-trained LLMs with human preferences before they can be effectively applied to various natural language tasks (Wang et al., 2023b; Wan et al., 2023; Sun et al., 2023c,b; Giray, 2023; Zhang, 2023; Long, 2023; Sun, 2023; Gao et al., 2023; Paranjape et al., 2023; Sun et al., 2023a; Diao et al., 2023; Wang et al., 2023a; Zhang et al., 2023b; Sun et al., 2023d; Liu et al., 2024d; Yao et al., 2024; Liu et al., 2024c; Lee et al., 2024; Kambhampati, 2024; Wang et al., 2024c).

Previously, a widely adopted approach for aligning the outputs of pre-trained LLMs with human preferences has been supervised fine-tuning (SFT) (Hu et al., 2021; Mishra et al., 2021; Wang et al., 2022; Du et al., 2022; Dettmers et al., 2023; Taori et al., 2023; Zhang et al., 2023a; Chiang et al., 2023; Xu et al., 2023; Peng et al., 2023; Mukherjee et al., 2023; Li et al., 2023; Ding et al., 2023; Luo et al., 2023; Wang et al., 2024d; Zhou et al., 2024). This method further trains LLMs on (Instruction, Answer) pairs, where "Instruction" represents the human prompt given to the model, and "Answer" is the target output that follows the instruction. SFT helps guide LLMs to produce responses that adhere to specific characteristics or domain knowledge, making it possible for humans to interactive with LLMs. Despite its effectiveness, SFT has limitations: during training, the model is constrained to learn specific answers we provide, with metrics like perplexity (PPL) penalizing synonym use. On one hand, this can hinder the LLM's ability to generalize, as tasks like writing and summarization have multiple valid phrasings. On the other hand, it may cause poor performance in aligning with human preferences, as no direct human feedback is incorporated into the training process.

To alleviate the above issues, reinforcement learning (RL) is adopted in aligning the outputs

---

♠The University of Melbourne, ♣Zhejiang University, ★CFAR and IHPC, A*STAR, Singapore, ▲Nanyang Technological University, ♦University of Washington
Email: shuhewang@student.unimelb.edu.au

Project page of this work can be found at: https://github.com/ShuheWang1998/Reinforcement-Learning-Enhanced-LLMs-A-Survey

**\* The latest update was on Dec. 3, 2024 (Version 1).**

of LLMs with human preferences, which can be decomposed into three steps: (1) First, before fine-tuning, a reward model (or reward function) is trained to approximate human preferences and score different LLM outputs; (2) Then, during each fine-tuning iteration, given a single instruction, the LLM generates multiple responses, each of which is scored by the trained reward model; (3) Finally, policy optimization, an RL optimization technique, updates the LLM's weights to improve predictions based on these preference scores. Fine-tuing LLMs with RL tackles the aforementioned issues simultaneously. In one line, rather than being restricted to learning a specific answer, RL adjusts the LLM based on various preference scores, rewarding any valid, well-phrased responses. In the other line, the reward model is designed to approximate human preferences, enabling direct training on human preferences and fostering the LLM's capacity for impressive creativity.

In this paper, we organize the most up-to-date state of knowledge on reinforcement learning (RL) in large language models (LLMs), attempting to consolidate and analyze the rapidly growing research in this field, helping researchers understand the current landscape, challenges, and advancements. Specifically,

- Section 2 presents the basics of reinforcement learning (RL) along with key terminologies, and outlines how the RL pipeline is adapted for LLMs.

- Section 3 introduces popular and powerful LLMs enhanced by reinforcement learning.

- Section 4 outlines the process of reinforcement learning from human feedback (RLHF), a training method that integrates reinforcement learning with human feedback to align LLMs with human values, preferences, and expectations.

- Section 5 reviews research on reinforcement learning from AI feedback (RLAIF), which presents a promising alternative or complement to RLHF by utilizing AI systems to provide feedback on the outputs of the LLM being trained, offering advantages in scalability, consistency, and cost-effectiveness.

- Section 6 provides an analysis of the challenges associated with RLHF and RLAIF.

- Section 7 discusses research on direct preference optimization (DPO), a series of methods that bypasses the reward model and directly utilizes human preference data to align LLM outputs with human expectations.

- Section 8 summarizes the current challenges and discusses opportunities for further improvement..

## 2 Basics: Reinforcement Learning for LLMs

In this section, we first detail the basics of reinforcement learning (RL) along with key terminologies, and then outline how the RL pipeline is adapted for LLMs.

### 2.1 Basics of Reinforcement Learning

Reinforcement Learning (RL) is a key approach in machine learning, focusing on how an agent engages with its environment to maximize cumulative rewards. Unlike supervised learning, which depends on labeled data, and unsupervised learning, which uncovers patterns in unlabeled data, RL emphasizes learning through direct feedback via trial and error. Below, we sequentially describe basic definitions and general pipeline of RL.

#### 2.1.1 Basic Definitions

Here, we use the training example in Figure 1 to illustrate the full process of RL. In this example, our goal is to train a robot to move from the bottom-left corner of a square to the top-right corner. Additionally, each grid cell has a reward score, and we aim to maximize the robot's total score. Before delving into the training process, we first introduce some relevant terms:

- **Agent:** An agent is the entity we train to make correct decisions. In this example, our goal is to train the robot to make movement decisions, so the robot is the agent.

- **Environment:** The environment is the external system that the agent interacts with. For our example, as the trained robot (agent) moves within the grid, the grid serves as the environment.

- **State:** The state represents the agent's position at each time $t$. For instance, at the beginning, at time $t_0$, the robot (agent) starts at the bottom-left corner, so the state at time $t_0$
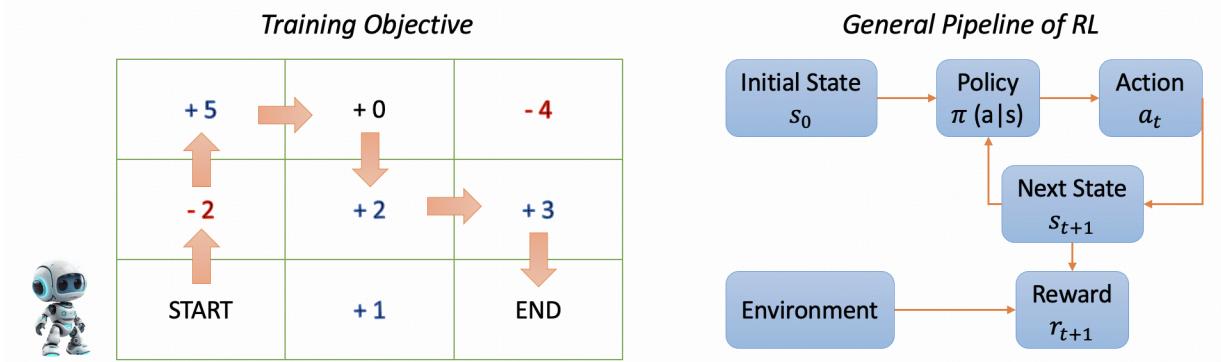
Figure 1: An example of the full process of RL. **Training Objective:** The goal is to train a robot to navigate from the bottom-left corner of a square to the top-right corner. Each grid cell is assigned a reward score, and the objective is to maximize the robot's overall score. **General Pipeline of RL:** The agent begins in an initial state $s_0$, and at each time step $t$, it selects an action $a_t$ based on its current state $s_t$. In response, the environment transitions to a new state $s_{t+1}$, and the agent receives a reward $r_t$.

is the bottom-left corner, represented by the coordinates $(0, 0)$.

- **Action(s):** Actions represent the possible choices available to the agent within the environment at each time $t$. For example, at the start, at time $t_0$, the robot (agent) can choose to move right or up, making these two actions available to the agent at $t_0$.

- **Reward(s):** Rewards are the signals or feedback provided by the environment to the agent based on the action it takes at each time $t$. For instance, at time $t_0$, the robot (agent) would receive a reward of +5 points for moving right, or a penalty of -1 point for moving up.

- **Policy:** A policy is a set of decision-making strategies that helps the agent choose an action at each time $t$. In practice, at time $t_0$, the policy represents a probability distribution that directs the robot (agent) to move right or up in order to maximize its cumulative rewards.

### 2.1.2 General Pipeline of RL

We have defined key terminologies used in RL, and in this section, we will continue to detail the general pipeline of RL.

As illustrated in Figure 1, the general reinforcement learning (RL) pipeline can be represented as a Markov Decision Process (MDP). Formally, the agent begins in an initial state $s_0$, and at each time step $t$, it selects an action $a_t$ based on its current state $s_t$. In response, the environment transitions to a new state $s_{t+1}$, and the agent receives a reward $r_t$. This cycle continues, with the agent's objective being to maximize the cumulative rewards it accumulates over time.

Mapping into the specific example in Figure 1, at the initial time $t_0$, the robot starts at the bottom-left corner, denoted by the position (state) $s_0$. As time progresses, at each time step $t$, the robot chooses an action $a_t$ (either moving up or moving right). This action causes the robot to transition from its current position $s_t$ to a new position $s_{t+1}$, while earning a reward $t_t$. This cycle of movement and reward collection continues until the robot reaches the desired position (state) at the top-right corner, achieving the goal of maximum cumulative rewards.

### 2.2 RL for LLMs

We have outlined the general framework of RL above; now we will delve into the process of fine-tuning LLMs using RL. This approach aims to align LLMs with desired behaviors, enhance their performance, and ensure that their outputs are both effective and dependable.

In reinforcement learning (RL), there are six key components:agent, environment, state, action, reward, and policy. To apply RL for fine-tuning large language models (LLMs), the first step is to map these components to the LLM framework. LLMs are highly proficient at next-token prediction, where they take a sequence of tokens as input and predict the next token based on the given context. From an RL perspective, we can view the LLM itself as the policy. The current textual sequence represents the state, and based on this state, the LLM generates an action—the next token. This action updates the state, creating a new state that
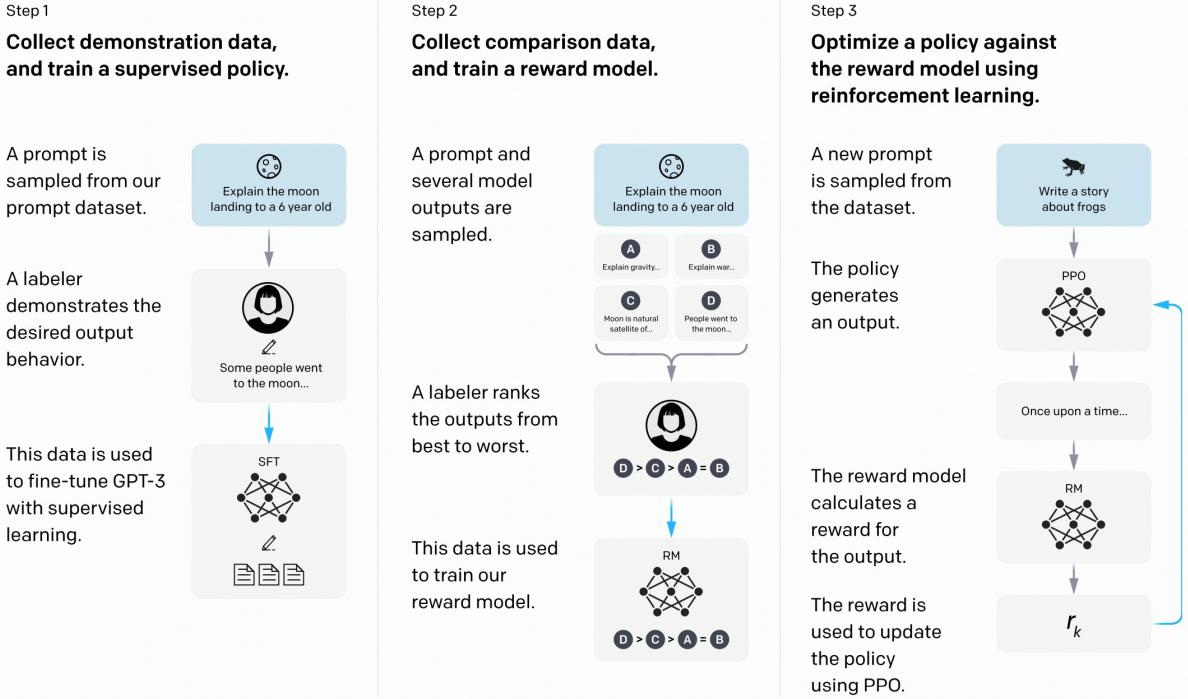
Figure 2: The framework of RL for LLMs proposed by Ouyang et al. (2022).

incorporates the newly added token. After generating a complete textual sequence, a reward is determined by assessing the quality of the LLM's output using a pre-trained reward model.

Figure 2 illustrates the specific RL framework for LLMs as proposed by (Ouyang et al., 2022). Ouyang et al. (2022) starts with an instruction-tuned model trained through supervised learning, enabling it to generate structured responses to human instructions. Then, Ouyang et al. (2022) applies the following two steps:

**Step 1: Collect comparison data, and train a reward model.** Ouyang et al. (2022) collects a dataset of comparisons between outputs of the instruction-tuned model, where labelers indicate which output they prefer for a given input. Then, the collected dataset is used to train a reward model (RM) to predict the human-preferred output.

**Step 2: Optimize a policy against the reward model using PPO.** Ouyang et al. (2022) leverages the output of the RM as a scalar reward, and fine-tunes the instruction-tuned model to optimize this reward using the PPO algorithm (Schulman et al., 2017).

## 3 Popular LLMs Enhanced by RL

Recent popular LLMs with strong capabilities almost all leverage reinforcement learning (RL) to further enhance their performance during the post-training process. The RL methods adopted by these models can be typically divided into two main lines: 1. Traditional RL approaches, such as Reinforcement Learning from Human Feedback (RLHF) and Reinforcement Learning from AI Feedback (RLAIF). These methods require training a reward model and involve a complex and often unstable process, using algorithms like Proximal Policy Optimization (PPO) (Schulman et al., 2017) to optimize the policy model. Models like InstructGPT (Ouyang et al., 2022), GPT-4 (OpenAI, 2023), and Claude 3 (Anthropic, 2024) follow this approach. 2. Simplified approaches, such as Direct Preference Optimization (DPO) (Rafailov et al., 2024) and Reward-aware Preference Optimization (RPO) (Adler et al., 2024). These methods discard the reward model, offering a stable, performant, and computationally efficient solution. Models like Llama 3 (Dubey et al., 2024), Qwen 2 (Yang et al., 2024a), and Nemotron-4 340B (Adler et al., 2024) follow this approach. In this section, we provide a detailed description of each model, starting with a brief overview of these RL enhanced LLMs and followed by an explanation of how RL is applied in their post-training process. An overview of these RL Enhanced LLMs is shown in Tab 1.

| RL Enhanced LLMs | Organization | # Params | RL Methods |
|---|---|---|---|
| Instruct-GPT (Ouyang et al., 2022) | OpenAI | 1.3B, 6B, 175B | RLHF, PPO |
| GPT-4 (OpenAI, 2023) | OpenAI | - | RLHF, PPO, RBRM |
| Gemini (Team et al., 2023) | Google | - | RLHF |
| InternLM2 (Cai et al., 2024) | 上海人工智能实验室 | 1.8B, 7B, 20B | RLHF, PPO |
| Claude 3 (Anthropic, 2024) | ANTHROP\C | - | RLAIF |
| Reka (Team et al., 2024c) | Reka | 7B, 21B | RLHF, PPO |
| Zephyr (HuggingFaceH4, 2024) | Argilla | 141B-A39B | ORPO |
| Phi-3 (Abdin et al., 2024) | Microsoft | 3.8B, 7B, 14B | DPO |
| DeepSeek-V2 (Liu et al., 2024a) | deepseek | 236B-A21B | GRPO |
| ChatGLM (GLM et al., 2024) | ZHIPU·AI | 6B, 9B | ChatGLM-RLHF |
| Nemotron-4 340B (Adler et al., 2024) | NVIDIA | 340B | DPO, RPO |
| Llama 3 (Dubey et al., 2024) | Meta | 8B, 70B, 405B | DPO |
| Qwen2 (Yang et al., 2024a) | Alibaba | (0.5-72)B, 57B-A14B | DPO |
| Gemma2 (Team et al., 2024b) | Google | 2B, 9B, 27B | RLHF |
| Starling-7B (Zhu et al., 2024) | Berkeley UNIVERSITY OF CALIFORNIA | 7B | RLAIF, PPO |
| Athene-70B (Nexusflow, 2024) | Nexusflow | 70B | RLHF |
| Hermes 3 (Teknium et al., 2024) | NOUS RESEARCH | 8B, 70B, 405B | DPO |
| o1 (OpenAI, 2024b) | OpenAI | - | RL through CoT |

Table 1: An overview of RL Enhanced LLMs. The format '141B-A39B' refers to MoE models with 141B total and 39B active parameters.

## 3.1 InstructGPT

InstructGPT (Ouyang et al., 2022) is a series of language models fine-tuned from GPT-3 (Brown et al., 2020) by OpenAI, using human feedback to better align with human intent. The series includes models in three sizes: 1.3 B, 6 B, and 175 B parameters. The model is first fine-tuned using supervised learning with prompts collected from the OpenAI API or written by labelers and corresponding labeler demonstrations, then further refined using reinforcement learning from human feedback (RLHF). Human evaluations reveal that InstructGPT outputs are preferred over GPT-3. Notably, the 1.3B parameter InstructGPT model is favored over the 175B GPT-3, despite having 100 times fewer parameters. Additionally, InstructGPT demonstrates improved truthfulness and reduced toxic outputs, with minimal performance trade-offs on public NLP datasets.

Before applying reinforcement learning (RL), the authors train a 6B reward model (RM) initial-ized from the supervised fine-tuned (SFT) model, with the final unembedding layer removed. This RM is trained using comparison data ranked by labelers. During the RL phase, they fine-tune the SFT model to optimize the scalar reward output from the RM using the PPO algorithm (Schulman et al., 2017). To address performance regressions on public NLP datasets, they experiment with mixing pre-training gradients with PPO gradients, resulting in models known as PPO-ptx.

## 3.2 GPT-4

GPT-4 (OpenAI, 2023), developed by OpenAI, is a large multimodal model that can process both image and text inputs to produce text outputs. It excels at understanding and generating natural language, particularly in complex and nuanced scenarios. Evaluations show that GPT-4 performs exceptionally well on a range of human-designed exams, often surpassing the majority of human test takers. Additionally, it outperforms earlier large language models and most state-of-the-art systems, which

frequently rely on benchmark-specific training or hand-engineered solutions.

GPT-4 leverages RLHF methods, as outlined in InstructGPT (Ouyang et al., 2022) which we have describe in Sec 3.1, in the post-training alignment stage. To steer the models more effectively towards appropriate refusals at a finer level, the authors further use a zero-shot GPT-4 classifier as the rule-based reward model (RBRM). This RBRM provides an additional reward signal to the GPT-4 policy model during PPO fine-tuning on a subset of training prompts. The RBRM takes a prompt (optional), the policy model's output, and a human-written rubric (e.g., a set of rules in multiple-choice style) as input, then classifies the output according to the rubric. Through this approach, GPT-4 is rewarded for refusing harmful content and for appropriately responding to known-safe prompts.

### 3.3 Gemini

Gemini (Team et al., 2023) represents a family of advanced multimodal models developed by Google, distinguished by their impressive capabilities. The initial version, Gemini 1.0, comes in three sizes—Ultra, Pro, and Nano—ranging from large to small in terms of performance. Each size is tailored to address specific computational constraints and application needs. Notably, Gemini Ultra, the most powerful variant, achieves state-of-the-art results in 30 out of 32 benchmarks and is the first model to attain human expert-level performance on MMLU (Hendrycks et al., 2020), while setting new records across all 20 multimodal benchmarks.

Gemini implements a post-training process that utilizes an optimized feedback loop, collecting human-AI interactions to drive continuous improvement in key performance areas. During the post-training's RLHF phase, an iterative approach is adopted wherein reinforcement learning (RL) incrementally enhances the reward model (RM). Concurrently, the RM undergoes continuous refinement through systematic evaluation and data collection. This dynamic interplay promotes ongoing advancement in both RL and RM, leading to progressively improved performance over time.

### 3.4 InternLM2

InternLM2 (Cai et al., 2024) is an open-source series of large language models developed by Shanghai AI Laboratory, available in three sizes: 1.8B, 7B, and 20B. The model demonstrates superior performance across six dimensions and 30 bench-marks, including long-context modeling and open-ended subjective evaluations, thanks to innovative pre-training and optimization techniques.

To further enhance alignment, InternLM2 employs a novel strategy called Conditional Online Reinforcement Learning from Human Feedback (COOL RLHF) with the use of PPO. This approach addresses two key challenges. The first is preference conflict, where it is difficult to satisfy two preferences, such as helpfulness and harmlessness, simultaneously. The second challenge is reward hacking, which becomes more problematic as the model's scale increases and its policy becomes more powerful. COOL RLHF introduces a Conditional Reward mechanism that reconciles diverse preferences by allowing a single reward model to dynamically adjust its focus based on specific conditional prompts, effectively integrating multiple preferences. Additionally, COOL RLHF incorporates a multi-round Online RLHF strategy with two distinct pathways: a Fast Path for immediate, targeted improvements and a Slow Path for long-term, comprehensive refinement of the reward model. This approach enables the model to quickly adapt to new human feedback while reducing the risk of reward hacking.

### 3.5 Claude 3

Claude 3 (Anthropic, 2024) is a family of large multimodal models developed by Anthropic, which demonstrates strong performance across benchmark evaluations. It comprises three models with varying abilities and speeds: the largest, Claude 3 Opus; the mid-sized, Claude 3 Sonnet; and the smallest, Claude 3 Haiku. The Claude 3 models show strong benchmark performance, setting new standards in reasoning, math, and coding. Claude 3 Opus achieves state-of-the-art results on evaluations such as GPQA (Rein et al., 2023), MMLU (Hendrycks et al., 2020), and MMMU (Yue et al., 2024). Claude 3 Haiku matches or surpasses Claude 2 in most text tasks, while Sonnet and Opus perform significantly better.

The authors use a technique called Constitutional AI (Bai et al., 2022) to align Claude 3 with human values during reinforcement learning (RL). In the RL stage, Constitutional AI follows a process similar to RLHF, but instead of human preferences for harmlessness, it uses AI feedback, known as RLAIF. Specifically, it distills language model interpretations of a set of rules and principles into a hybrid human/AI preference model (PM), using hu-

man labels for helpfulness and AI labels for harmlessness. Afterwards, they fine-tune the supervised learning model using RL with this PM, resulting in a policy trained by RLAIF.

## 3.6 Zephyr 141B-A39B

Zephyr 141B-A39B (HuggingFaceH4, 2024) is the newest addition to the Zephyr (Tunstall et al., 2023) series of language models, developed through a collaboration between Argilla, KAIST, and Hugging Face. This model is a Mixture of Experts (MoE) with a total of 141 billion parameters, 39 billion of which are active, fine-tuned from Mixtral-8x22B-v0.1 (Mistral AI, 2024).

Zephyr 141B-A39B employs a novel alignment algorithm known as Odds Ratio Preference Optimization (ORPO) (Hong et al., 2024). ORPO is a straightforward, unified alignment approach that discourages the model from adopting undesired generation styles during supervised fine-tuning. Notably, ORPO does not require an SFT warm-up phase, a reward model, or a reference model, making it highly resource-efficient. The method works by adding an odds ratio-based penalty to the standard SFT negative log-likelihood loss, enabling the model to distinguish between preferred and non-preferred response styles.

## 3.7 DeepSeek-V2

DeepSeek-V2 (Liu et al., 2024a), developed by DeepSeek-AI, is a powerful Mixture-of-Experts (MoE) language model designed for economical training and efficient inference. It features innovative architectures such as Multi-head Latent Attention (MLA) and DeepSeekMoE. With 236 billion total parameters, of which 21 billion are activated per token, it supports a context length of up to 128K tokens. The model is pre-trained on a high-quality, multi-source corpus of 8.1 trillion tokens. Evaluations show that DeepSeek-V2, along with its chat versions, maintains top-tier performance among open-source models, despite having only 21 billion activated parameters.

DeepSeek-V2 is optimized using Group Relative Policy Optimization (GRPO) (Shao et al., 2024) during the RL phase to reduce training costs. Unlike traditional RL methods that use a critic model of similar size to the policy model, which increases training expenses, GRPO foregoes the critic model and estimates the baseline from scores computed on a group of outputs for the same question. Additionally, a two-stage RL training strategy is employed:

the first stage focuses on reasoning alignment, and the second on human preference alignment, as the authors find these stages exhibit distinct characteristics.

## 3.8 ChatGLM

ChatGLM (GLM et al., 2024), developed by Zhipu AI, represents an evolving series of large language models. The latest version in this series is GLM-4, which includes variants such as GLM-4, GLM-4-Air, and GLM-4-9B. These models are pre-trained on a dataset of over 10 trillion tokens, predominantly in Chinese and English, and are subsequently post-trained through a combination of supervised fine-tuning (SFT) and RLHF to achieve advanced alignment quality. Evaluation results indicate that GLM-4 rivals or even surpasses GPT-4 (OpenAI, 2023) on general benchmarks like MMLU, and demonstrates superior performance in Chinese-specific alignments as measured by Align-Bench (Liu et al., 2023b).

The reinforcement learning phase involves the ChatGLM-RLHF (Hou et al., 2024) pipeline, which enhances alignment with human preferences. This pipeline comprises three primary components: gathering human preference data, training a reward model, and optimizing policy models. To support large-scale training, ChatGLM-RLHF includes methods to reduce reward variance for stable training, leverages model parallelism with fused gradient descent, and applies regularization constraints to prevent catastrophic forgetting in large language models. Experimental results confirm that ChatGLM-RLHF yields substantial improvements in alignment-focused tasks compared to the supervised fine-tuned version of ChatGLM.

## 3.9 Nemotron-4 340B

Nemotron-4 340B (Adler et al., 2024) is a family of models released by NVIDIA, consisting of Nemotron-4-340B-Base, Nemotron-4-340B-Instruct, and Nemotron-4-340B-Reward. The Nemotron-4-340B-Base model is trained on 9 trillion tokens from a high-quality dataset. In the alignment process to develop Nemotron-4-340B-Instruct, over 98% of the data used is synthetically generated by the model. Evaluations demonstrate that these models perform competitively with open-access models across a broad range of evaluation benchmarks.

During the preference fine-tuning phase, both DPO (Rafailov et al., 2024) and a new alignment

algorithm, Reward-aware Preference Optimization (RPO), are employed to improve the model through multiple iterations. RPO addresses a limitation in DPO, where the quality difference between selected and rejected responses is not considered, leading to overfitting and the forgetting of valuable responses. RPO uses an implicit reward from the policy network to approximate this gap, enabling the model to better learn from and retain superior feedback.

### 3.10 Llama 3

Llama 3 (Dubey et al., 2024), developed by Meta, is a collection of open-source foundational language models available in sizes of 8 billion, 70 billion, and 405 billion parameters. It is trained on a significantly larger corpus consisting of approximately 15 trillion multilingual tokens, a notable increase compared to the 1.8 trillion tokens used for Llama 2 (Touvron et al., 2023). Extensive empirical evaluations demonstrate that Llama 3 achieves performance comparable to leading models, such as GPT-4 (OpenAI, 2023), across a diverse range of tasks.

The post-training process for aligning Llama 3 with human feedback involves six rounds of iterative refinement. Each round includes supervised fine-tuning (SFT) followed by DPO, with the final model being an average of the outputs from all rounds. For each round, a reward model (RM) is trained on newly collected preference annotation data, targeting a wide range of capabilities built upon the pre-trained checkpoint. After SFT, DPO is applied to further optimize the SFT models, using recent preference data batches obtained from the best-performing models of previous rounds. To enhance the stability of DPO training, two key adjustments are implemented: masking out formatting tokens in the DPO loss and introducing regularization via an NLL (negative log-likelihood) loss.

### 3.11 Qwen2

Qwen2 (Yang et al., 2024a), developed by Alibaba, is a series of large language models ranging from 0.5 billion to 72 billion parameters in dense configurations, as well as a Mixture-of-Experts variant with 57 billion parameters, of which 14 billion are activated per token. It is pre-trained on a high-quality, large-scale dataset containing over 7 trillion tokens, covering a wide array of domains and languages. Extensive evaluations show that Qwen2 outperforms most prior open-weight models, including its predecessor Qwen1.5, and delivers competitive results across a range of benchmarks, including language understanding, generation, multilingual proficiency, coding, mathematics, and reasoning.

The preference fine-tuning process for Qwen2 consists of two main stages: offline and online learning. In the offline stage, Qwen2 is optimized using DPO, which aims to maximize the likelihood difference between two responses to the same prompt, based on a pre-compiled preference dataset. In the online stage, the model improves continuously in real-time by utilizing preference pairs selected by the reward model from multiple responses generated by the current policy model. Additionally, the Online Merging Optimizer (Lu et al., 2024) is employed to minimize alignment costs.

### 3.12 Gemma 2

Gemma 2 (Team et al., 2024b), developed by Google, is the latest addition to the Gemma family of lightweight, state-of-the-art open models, with sizes ranging from 2 billion to 27 billion parameters. The model incorporates several well-established modifications to the Transformer architecture, including interleaving local-global attentions (Beltagy et al., 2020) and group-query attention (Ainslie et al., 2023). Experiments demonstrate that these models deliver the best performance for their size and even provide competitive alternatives to models 2-3 times larger.

Similar to Gemma 1.1 (Team et al., 2024a), during the post-training RLHF phase, the authors use a high-capacity model as an automatic rater to tune hyperparameters and mitigate reward hacking (Amodei et al., 2016; Skalse et al., 2022). However, unlike Gemma 1.1, they employ a reward model that is an order of magnitude larger than the policy model. This reward model is specifically designed to focus on conversational capabilities, with an emphasis on multi-turn interactions.

### 3.13 Starling-7B

Starling-7B (Zhu et al., 2024) is a strong 7-billion-parameter chat model developed by UC Berkeley, focused on alignment with human preferences for helpfulness and harmlessness. It is fine-tuned from Openchat-3.5 (Wang et al., 2024a) using RLAIF on a high-quality preference dataset called Nectar, which comprises 3.8 million pairwise comparisons

generated by prompting GPT-4 to rank responses. As a result, the model's score on MT-Bench improves from 7.81 to 8.09, its score on AlpacaEval increases from 88.51% to 91.99%, and its human evaluation ELO on Chatbot Arena (Chiang et al., 2024) rises from 1072 to 1087.

The authors introduce several improvements to the PPO algorithm during the RLAIF process to enhance training stability and robustness. First, they introduce a constant positive reward for length control to prevent excessive verbosity. This adjustment helps address the issue where a highly negative reward from the reward model during the early stages can cause the policy model to become overly verbose after only a few gradient updates. Second, they pretrain the critic model to reduce early performance drops due to a randomly initialized critic. Third, they conduct full parameter tuning on both the actor and critic models, as opposed to tuning only the top four layers, to maximize performance improvements during the reinforcement learning stage.

## 3.14 o1

OpenAI's o1 (OpenAI, 2024b) is a newly developed large language model optimized for complex reasoning, utilizing reinforcement learning for its training. Before producing responses, o1 engages in an extensive internal thought process, enabling it to excel across various reasoning tasks. The model significantly surpasses GPT-4o (OpenAI, 2024a) in many challenging tasks: ranks in the 89th percentile on Codeforces for competitive programming, places among the top 500 participants in the AIME for mathematics, and surpasses PhD-level accuracy in scientific benchmarks such as GPQA.

The training of o1 involves a large-scale reinforcement learning algorithm that emphasizes productive thinking through a detailed chain of thought (CoT) (Wei et al., 2023), implemented with high data efficiency. To preserve the model's unfiltered reasoning ability, no policy compliance or user preference training is applied to its internal thought processes, which also provides a unique opportunity to understand the model's raw thought process. This approach allows o1 to refine its strategies, correct errors, and deconstruct complex problems during training. Notably, the model's performance improves with increased training compute and with more extensive test-time computation.

## 3.15 Others

**Reka Core, Flash, and Edge:** Team et al. (2024c) are powerful multimodal language models developed from scratch by Reka. Reka Edge and Reka Flash are dense models with 7B and 21B parameters, respectively, outperforming many larger models and offering exceptional performance for their compute class. The flagship model, Reka Core, competes with leading models like GPT-4v, Gemini, and Claude 3 in both automated and blind human evaluations. During post-training, following supervised fine-tuning, Reka models undergo multiple rounds of RLHF using PPO to enhance alignment further.

**Phi-3:** Abdin et al. (2024) is a series of language models introduced by Microsoft, comprising phi-3-mini, phi-3-small, and phi-3-medium. Remarkably, the smallest model, phi-3-mini, is trained on 3.3 trillion tokens yet contains only 3.8 billion parameters, making it compact enough for deployment on a mobile device. Despite its relatively small size, phi-3-mini demonstrates performance comparable to larger models like Mixtral 8x7B and GPT-3.5, achieving 69% on MMLU and a score of 8.38 on MT-bench in both academic benchmarks and internal testing. During post-training, the authors employ DPO to guide phi-3 away from undesired behavior by treating those outputs as "rejected" responses.

**Athene-70B:** Nexusflow (2024) is a powerful chat model fine-tuned from Llama-3-70B (Dubey et al., 2024), developed by Nexusflow. It achieves an impressive Arena-Hard-Auto score of 77.8%, placing it close to leading proprietary models like GPT-4o (79.2%) and Claude-3.5-Sonnet (79.3%). This marks a significant leap from its predecessor, Llama-3-70B-Instruct, which scored 46.6%. This progress is attributed to Nexusflow's targeted post-training approach, which enhances the model's performance. Specifically, Nexusflow curates high-quality preference data based on internal benchmark evaluations covering instruction following, coding, creative writing, and multilingual tasks. This data is then used for targeted RLHF, resulting in substantial performance gains over Llama-3-70B-Instruct.

**Hermes 3:** Teknium et al. (2024) is a series of neutrally-aligned generalist instruction and tool-use models with advanced reasoning and creative capabilities, developed by Nous Research. It is
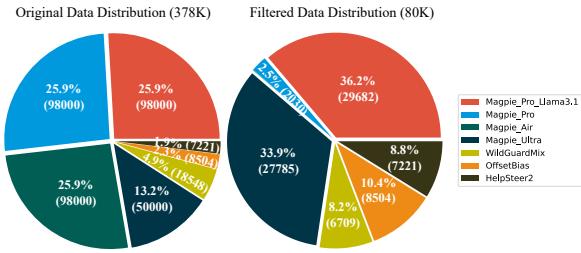
Figure 3: The composition of the Skywork-Reward. The figure is copied from Liu et al. (2024b).

finetuned from Llama 3.1 (Dubey et al., 2024) in 8B, 70B, and 405B variants and the largest model, Hermes 3 405B, sets the state-of-the-art performance among open-weight models across several public benchmarks. Hermes is trained on diverse synthetic reasoning tasks and creative applications such as role playing and writing. It is designed to precisely and neutrally follow system and instruction prompts, unlike many commercial models that may decline instructions for moral reasons. To further align Hermes, the authors leverage DPO and train a LoRA (Hu et al., 2021) adapter instead of fine-tuning the entire model, significantly reducing GPU memory usage for both the reference and trained models.

# 4 RLHF: Reinforcement Learning from Human Feedback

Reinforcement learning from human feedback (RLHF) is a training approach that combines reinforcement learning (RL) with human feedback to align LLMs with human values, preferences, and expectations. RLHF consists of two main components: (1) **Collecting Human Feedback to Train Reward Model**, where human evaluators provide feedback on the LLM's outputs by scoring or ranking responses based on factors such as quality and relevance. This feedback is then used to train a reward model that predicts the quality of the outputs and serves as the reward function in the RL process; and (2) **Preference Optimization Using Human Feedback**, where the trained reward model guides the optimization of the LLM's outputs to maximize predicted rewards, aligning the LLM's behavior with human preferences. Below, we will illustrate these two components via recent research studies.

## 4.1 Collecting Human Feedback to Train Reward Model

**Skywork-Reward (Liu et al., 2024b).** Skywork-Reward is a carefully designed dataset containing 80,000 high-quality preference pairs, curated through effective data selection and filtering strategies. As shown in Figure 3, the original dataset, with 378,000 preference pairs, is significantly refined into a compact, high-quality dataset of 80,000 pairs. Despite being significantly smaller than existing datasets, it achieves exceptional quality through rigorous cleaning, consistency checks, model-based scoring to filter out low-quality samples, and manual reviews. Covering a diverse range of tasks such as instruction following, code generation, and multilingual handling, Skywork-Reward serves as the foundation for models like Skywork-Reward-Gemma-27B, which excel on benchmarks[1]. By enabling language models to better understand human preferences, Skywork-Reward helps LLMs become more accurate and useful in real-world applications.

**TÜLU-V2-mix (Ivison et al., 2023).** TÜLU-V2-mix is designed to enhance instruction-following capabilities in large language models, offering a diverse dataset that improves the model's generalization and execution abilities across multi-domain tasks. It covers a wide range of tasks, including question answering, code generation, translation, and multi-turn conversations, with a strong emphasis on multilingual adaptability and handling complex real-world scenarios. Skywork-Reward, on the other hand, is designed to align models with human preferences using preference pairs, helping models learn to generate user-preferred responses, such as fluent and coherent text. While TÜLU-V2-mix excels in generalization across a wide range of tasks, Skywork-Reward specializes in optimizing user-centric outputs. Together, they address complementary goals for advancing language model capabilities.

## 4.2 Preference Optimization Using Human Feedback

Once the reward model is trained, it is used to guide the fine-tuning of the original LLM through reinforcement learning. The main objective is to improve the LLM's behavior based on the predicted rewards, making it more likely to generate outputs

---

[1]https://huggingface.co/spaces/allenai/reward-bench

that align with human preferences. Recent research (Ouyang et al., 2022; Yuan et al., 2023; Dong et al., 2024; Ahmadian et al., 2024) has shown that this process can be broken down into two key steps:

**(1) Rewarding:** In this step, the LLM generates multiple outputs in response to a given instruction. Each output is then passed through the trained reward model, which assigns a scalar score that approximates human preferences.

**(2) Policy Optimization:** In this step, the LLM is fine-tuned by adjusting its parameters to maximize the predicted reward, using the Proximal Policy Optimization (PPO) (Schulman et al., 2017) or Trust Region Policy Optimization (TRPO) (Schulman, 2015) algorithm.

These two steps—rewarding and policy optimization—can be iterated, meaning that the process of generating outputs, rewarding them with the trained reward model, and fine-tuning the LLM to maximize rewards can be repeated multiple times. With each iteration, the LLM's performance improves as it refines its behavior to better align with human preferences. This iterative cycle allows the LLM to continuously adapt and optimize its responses, ultimately leading to more effective and aligned outputs.

## 5 RLAIF: Reinforcement Learning from AI Feedback

Reinforcement learning from AI feedback (RLAIF) serves as a promising alternative or supplement to RLHF that leverages AI systems—often more powerful or specialized LLMs (e.g., GPT-4 (OpenAI, 2024a))—to provide feedback on the outputs of the LLM being trained. This approach provides benefits such as scalability, consistency, and cost efficiency while minimizing reliance on human evaluators. Below, we explore several methods for substituting human feedback with AI feedback in reinforcement learning, highlighting approaches: (1) Distilling AI Feedback to Train Reward Model, (2) Prompting LLMs As a Reward Function, and (3) Self-Rewarding.

### 5.1 Distilling AI Feedback to Train Reward Model

Beyond manually collected data, distilling datasets from pre-trained LLMs presents an efficient alternative. By leveraging the outputs of powerful LLMs like GPT-4, researchers can build a bridge between manual curation and autonomous evaluation.

**UltraFeedback (Cui et al., 2023).** UltraFeedback is a large-scale AI feedback dataset aimed at improving the performance and alignment of large language models (LLMs). It includes over 1 million high-quality GPT-4 feedback annotations across 250,000 user-assistant interactions, focusing on key dimensions like instruction adherence, accuracy, honesty, and usefulness. The dataset was created by collecting 60,000 diverse instructions, generating responses using 17 different models, and leveraging GPT-4 for detailed critiques and scoring, wherein chain-of-thought reasoning is used to reduce bias.

**Magpie.** Xu et al. (2024b) introduce a self-synthesis method that leverages the autoregressive nature of aligned LLMs. By utilizing predefined templates as prompts, the model autonomously generates user queries and corresponding responses, eliminating the need for manual intervention or initial seed questions. Specifically, as shown in Figure 4, aligned LLMs (e.g., Llama-3-Instruct model) is employed to synthesize 4 million instruction-response pairs, subsequently filtering the dataset to retain 300,000 high-quality pairs. These pairs were then used to fine-tune the Llama-3-8B-Base model. Remarkably, the fine-tuned model achieved performance comparable to the official Llama-3-8B-Instruct model, which had undergone training on 10 million examples through supervised fine-tuning and reinforcement learning with human feedback. Besides, models fine-tuned with Magpie excelled on alignment benchmarks such as AlpacaEval, surpassing models trained on other open datasets and preference optimization methods.

**HelpSteer2 (Wang et al., 2024e).** HelpSteer2 is an efficient, open-source preference dataset comprising approximately 10,000 comparison samples, designed to train high-performance reward models. The dataset is built using responses generated by various models (including GPT-3.5, Claude, and others) and features multi-dimensional annotations such as fluency, relevance, creativity, and safety. Preference pairs are crafted based on human or automated evaluations, enabling fine-grained alignment for reward models. Through rigorous data cleaning and optimization, HelpSteer2 delivers high-quality annotations in a compact format. It is released under the CC-BY-4.0 license, fostering the accessibility.
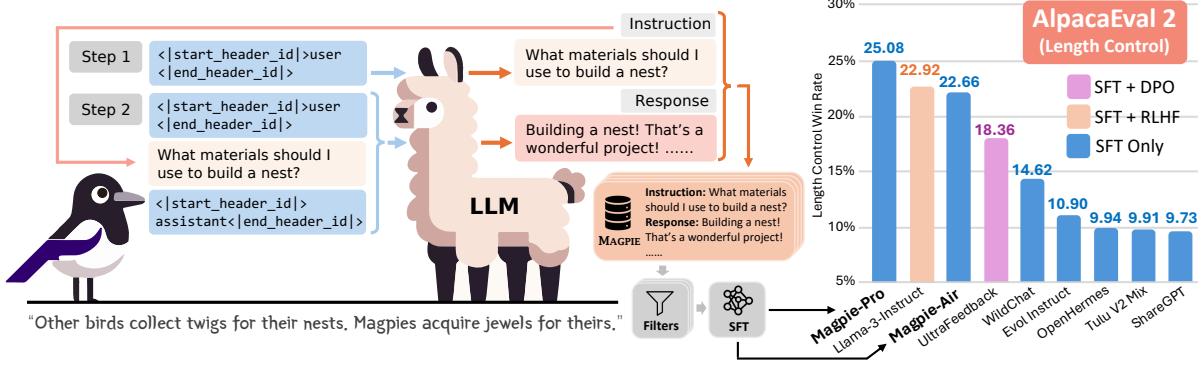
Figure 4: Magpie self-synthesizes data from aligned LLMs. The figure is borrowed from Xu et al. (2024b).

**OffsetBias (Park et al., 2024).** OffsetBias is a meticulously designed dataset aimed at mitigating biases in reward models, constructed using responses generated by diverse models, including GPT-3.5, GPT-4, Claude, and open-source models like Llama 2. As shown in Figure 5, OffsetBias systematically addresses six identified bias types, namely, content, style, informativeness, safety, creativity, and length. Based on this, comparison samples are generated through attribute-controlled prompts and multi-model outputs. These samples are annotated with multi-dimensional scores and preference labels to highlight or neutralize biases, enabling fine-grained alignment. OffsetBias serves as a robust resource for improving the fairness and reliability of reward models, with its data openly accessible for research and development.

## 5.2 Prompting LLMs As a Reward Function

As reward model training becomes more sophisticated, a natural progression is to employ LLMs themselves as evaluators in the loop of reinforcement learning.

**Exploring with LLMs (ELLM) Rewards (Du et al., 2023).** ELLM is a method that integrates LLMs with reinforcement learning (RL) to enhance exploration during the pretraining phase. Figure 6 showcases the overall pipeline: the agent's current state is transformed into a natural language description, which is input into the LLM. The LLM then generates exploration goals based on this state description, such as specific actions or target locations. The RL agent attempts to achieve these goals, and rewards are provided by the environment upon goal completion. This approach improves exploration efficiency by guiding the agent toward areas of the state space that are likely to be valuable,

without requiring pre-designed rewards. ELLM is particularly useful in sparse-reward environments. Compared to traditional methods, ELLM significantly improves exploration efficiency, covering more common-sense behaviors and providing better initialization for downstream tasks.

**Reward Design with Language Models (RDLM).** Kwon et al. (2023) leverage a LLM like GPT-3 to simplify reward function design in reinforcement learning by allowing users to define desired behaviors through natural language descriptions. Specifically, users provide a task description or a few examples, and the LLM generates reward signals by evaluating the agent's behavior against these criteria. Instead of producing reward code, RDLM outputs direct reward values that the RL agent uses for policy optimization. This method is ideal for tasks where user goals are clear but manually designing a reward function is complex. While ELLM focuses on guiding exploration during pretraining by generating meaningful goals, RDLM emphasizes task-specific reward generation to streamline complex reward design and achieve better agent alignment with human intent.

**Eureka (Ma et al., 2023).** Eureka is an algorithm that leverages LLMs to automatically generate and optimize reward function code for reinforcement learning tasks. In Figure 7, first, a coding LLM like GPT-4 is used to generate initial reward function code based on task descriptions. This code is then iteratively refined using evolutionary strategies, where candidate reward functions are evaluated based on how well they guide the RL agent toward task success. The process evolves the reward functions to improve their quality and effectiveness. Eureka is particularly effective in tasks requiring complex or highly specific reward defini-

## Type 1. Length Bias

**Instruction**: Come up with a analogy to explain this concept. Concept: <u>**Network Topology**</u>

**Good Response** ✅

<u>**Network topology is like a roadmap**</u>; it shows how the different components of a network are connected and arranged.

**Bad Response** ⛔

<u>**The internet protocol suite**</u> is like a postal system's rules for delivering mail. The Internet Protocol (IP) is the address on the *[...]* types of data transfer. Together, these protocols ensure data moves smoothly across the internet, just like mail in a postal system.

---

## Type 2. Concreteness Bias

**Instruction**: Is eating <u>**orange**</u> good for our health?

**Good Response** ✅

<u>**Eating oranges**</u> is good for well-being. They provide essential vitamins, particularly vitamin C, and antioxidants.

**Bad Response** ⛔

According to a study published in the Journal of Nutrition (2020), <u>**eating apples**</u> daily can reduce the risk of chronic diseases by 25%

---

## Type 3. Empty Reference Bias

**Instruction**: I need some tips on improving my skills.

**Good Response** ✅

I'd be glad to help! Can you <u>**specify which skills you're looking to improve**</u>? Knowing the specific area will allow me to give you tailored advice. Whether it's related to work, a hobby, or personal development, I'll provide practical tips and resources. Let's make sure we target your goals effectively.

**Bad Response** ⛔

<u>**To improve your public speaking skills**</u>, practice regularly, get feedback from peers, and watch videos of skilled speakers for tips. Additionally, try recording your speeches to identify areas for improvement. Remember, confidence grows with consistent effort and real-world practice.

---

## Type 4. Content Continuation Bias

**Instruction**: <u>**Paraphrase**</u> the following sentence.
He wanted to watch a movie but she wanted to go to a party.

**Good Response** ✅

He was interested in watching a movie, while she preferred attending a party.

**Bad Response** ⛔

They compromised and decided to do both by watching a movie at the party.

---

## Type 5. Nested Instruction Bias

**Instruction**: Read the following sentence and <u>**suggest one correction**</u>.
Are you familiar with the the US?

**Good Response** ✅

You might <u>**need to omit one of the adjacent 'the'**</u> in the sentence.

**Bad Response** ⛔

<u>**Yes**</u>, including various aspects of society and daily life.

---

## Type 6. Familiar Knowledge Bias

**Instruction**: What does it mean by "I was too shy. I just <u>**kept the ice**</u> there. "

**Good Response** ✅

The sentence suggests that the person was <u>**too shy to start a conversation**</u> or make the situation more comfortable, so the initial awkwardness or silence ("the ice") remained unbroken.

**Bad Response** ⛔

It means that despite feeling nervous or uncomfortable in a social situation, the person <u>**made an effort to initiate conversation**</u> or interaction to make things less awkward and more comfortable for everyone involved.

Figure 5: Identified bias types and examples in OffsetBias. The figure is borrowed from Park et al. (2024).
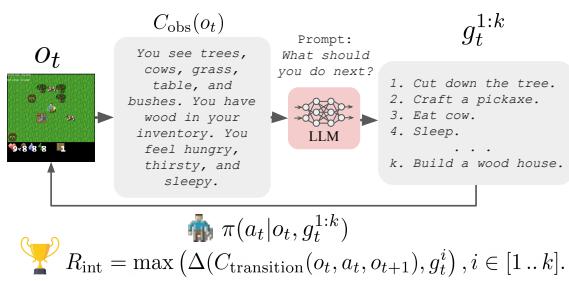


Figure 6: The figure is copied from (Du et al., 2023).

tions, such as advanced robotic skills. Its focus on reward code optimization makes it suitable for scenarios where precise reward shaping is critical. By utilizing LLMs' ability to generate and refine code, Eureka evolves reward functions that effectively guide RL agents. Experiments demonstrate that Eureka outperforms human-designed rewards in 83% of tested tasks, with an average performance improvement of 52%, showcasing its potential for advanced skill learning, such as robotics tasks, in challenging scenarios.

**Text2Reward (Xie et al., 2023).** Text2Reward is a framework that leverages large language models to automatically generate dense and interpretable reward function code from natural language task descriptions, enabling efficient reward shaping across
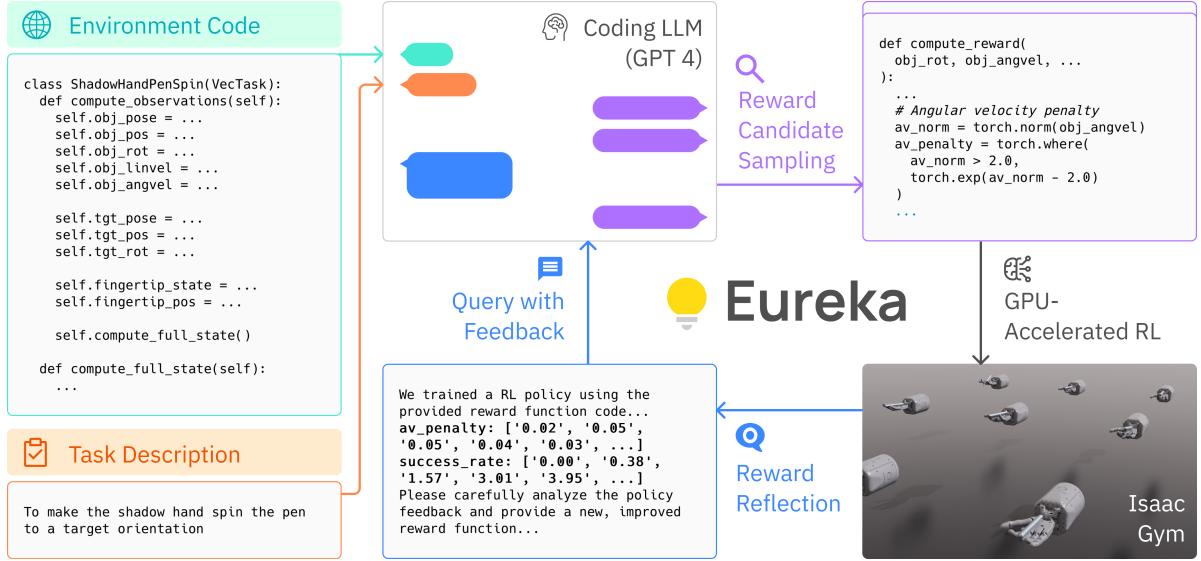
Figure 7: The overall pipeline of Eureka. The figure is borrowed from Ma et al. (2023).
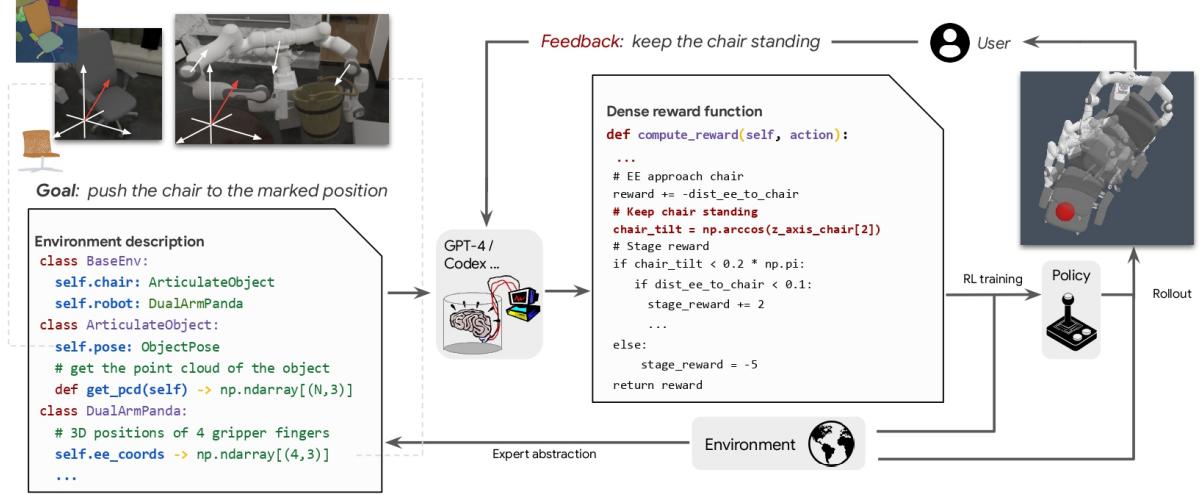


Figure 8: An overview of Text2Reward. The figure is copied from Xie et al. (2023).

diverse RL tasks. As shown in Figure 8, the process starts with users providing a task description in natural language, which is input into an LLM to generate executable reward code. This code often includes task-specific logic and may integrate external libraries for complex functionalities. The generated reward function is then used in RL to guide the agent's behavior. Additionally, Text2Reward supports iterative refinement of the reward code through human feedback, enabling further optimization. This method excels at providing flexible, interpretable rewards across diverse RL tasks, particularly in robotics and manipulation. Unlike Eureka, evolving and optimizing reward function code through LLMs and evolutionary algorithms, Text2Reward emphasizes creating human-readable

reward code that integrates external libraries and supports iterative refinement via human feedback. While both methods aim to automate reward design, Eureka excels in optimizing complex reward logic for advanced skills, whereas Text2Reward prioritizes flexibility, interpretability, and adaptability for a broad range of tasks.

**RLAIF.** Lee et al. (2023) replace human feedback in RL with AI-generated feedback by leveraging LLMs. The process begins with generating candidate outputs for a given task, such as text summarization or dialogue generation. These outputs are paired and fed into an LLM, which evaluates them and provides preferences (e.g., selecting the better output) or assigns scores based on task-specific criteria. This feedback is then used
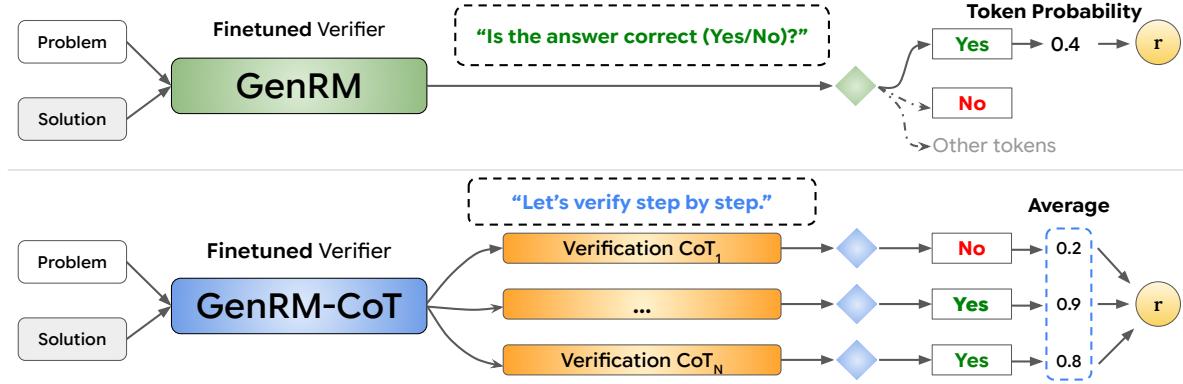
Figure 9: Illustration of GenRM. The figure is copied from Zhang et al. (2024a).
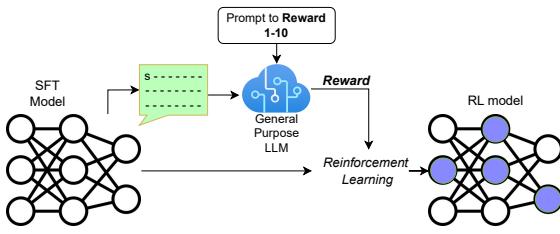


Figure 10: The figure is borrowed from Lee et al. (2023).

to train a reward model that predicts the quality of outputs and guides the RL agent. In its stream-lined variant, d-RLAIF (see Figure 10), the LLM directly provides scores as reward signals, bypass-ing the need for a reward model. The RL policy is optimized using these rewards, typically with al-gorithms like Proximal Policy Optimization (PPO). This approach enables automated, scalable, and high-quality feedback generation, effectively align-ing RL agent behavior with task objectives while reducing reliance on human annotations.

**GenRM.** Zhang et al. (2024a) re-define verifica-tion by treating it as a text generation task, lever-aging large language models to produce valida-tion outputs and reasoning chains, such as "yes" or "no" with explanations. As shown in Figure 9, this approach integrates verification into the gen-erative capabilities of LLMs, enabling them to as-sess and explain candidate answers in a transparent and interpretable manner. By framing verification as next-token prediction, GenRM eliminates re-liance on traditional discriminative models and en-hances reasoning accuracy. Experimental results demonstrate its ability to outperform conventional methods, showcasing its potential in tasks requir-ing logical reasoning, interpretability, and scalable

performance.

## 5.3 Self-Rewarding

The self-rewarding mechanism enables the LLM to autonomously assess and refine its own perfor-mance, addressing the cost, scalability, and adapt-ability limitations of existing RL methods.

**Self-Refined LLM.** Song et al. (2023) leverage LLMs to automatically generate reward functions for deep reinforcement learning (DRL) tasks and introduces a self-optimization mechanism to iter-atively refine these functions. The process begins with the LLM generating an initial reward function based on natural language task descriptions. The reward function is then applied to RL training, and the agent's performance is evaluated. Feedback from this evaluation is fed back into the LLM, en-abling it to dynamically adjust and improve the reward function in a closed-loop manner. Com-pared to Eureka and Text2Reward, this approach eliminates the need for external optimization algo-rithms or manual intervention.

**Self-Rewarding Language Models (SRLM).** Yuan et al. (2024) introduce a novel approach where LLMs act as both the generator and eval-uator to create a self-contained learning system. As shown in Figure 11, the model begins by generat-ing new prompts (instructions) and multiple candi-date responses derived from existing data, thereby creating a diverse and comprehensive set of train-ing samples. Subsequently, the model evaluates these candidate responses using a structured scor-ing mechanism to determine their quality. The eval-uation framework encompasses multiple dimen-sions, including relevance, coverage, usefulness, clarity, and professionalism, assigning a score to
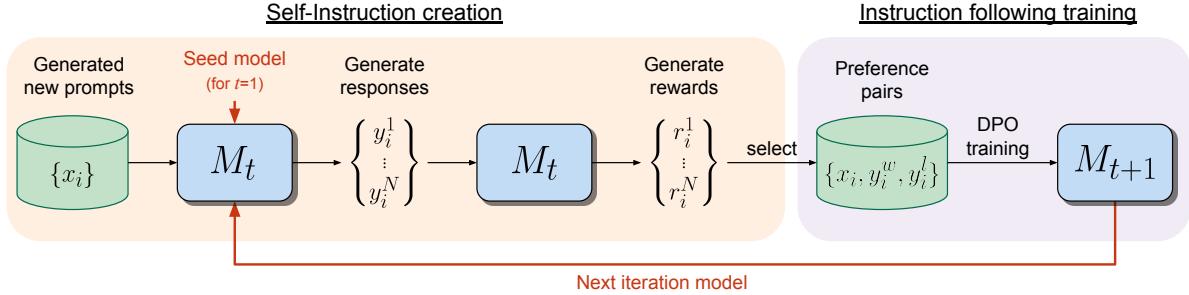
Figure 11: The overview of SRLM. The figure is copied from Yuan et al. (2024).

Review the user's question and the corresponding response using the additive 5-point scoring system described below. Points are accumulated based on the satisfaction of each criterion:

- Add 1 point if the response is relevant and provides some information related to the user's inquiry, even if it is incomplete or contains some irrelevant content.
- Add another point if the response addresses a substantial portion of the user's question, but does not completely resolve the query or provide a direct answer.
- Award a third point if the response answers the basic elements of the user's question in a useful way, regardless of whether it seems to have been written by an AI Assistant or if it has elements typically found in blogs or search results.
- Grant a fourth point if the response is clearly written from an AI Assistant's perspective, addressing the user's question directly and comprehensively, and is well-organized and helpful, even if there is slight room for improvement in clarity, conciseness or focus.
- Bestow a fifth point for a response that is impeccably tailored to the user's question by an AI Assistant, without extraneous information, reflecting expert knowledge, and demonstrating a high-quality, engaging, and insightful answer.

User: `<INSTRUCTION_HERE>`

`<response><RESPONSE_HERE></response>`

After examining the user's instruction and the response:

- Briefly justify your total score, up to 100 words.
- Conclude with the score using the format: "Score: <total points>"

Remember to assess from the AI Assistant perspective, utilizing web search knowledge as necessary. To evaluate the response in alignment with this additive scoring model, we'll systematically attribute points based on the outlined criteria.
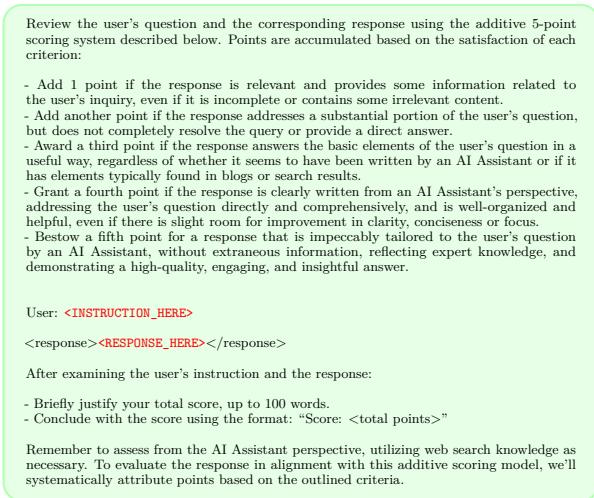
Figure 12: Prompt for LLM as a judge. The figure is borrowed from Yuan et al. (2024).

each response based on these criteria. Utilizing these scores, preference pairs are constructed, consisting of a preferred response and a dispreferred response. These pairs are used for Direct Preference Optimization (DPO), improving its ability to generate high-quality responses. Through iterative refinement, the model progressively enhances its performance. Figure 12 provides a detailed explanation of the prompts used by the model to evaluate candidate responses. Experimental results demonstrate that fine-tuning Llama 2 70B using SRLM over three iterations outperforms several state-of-the-art models, including GPT-4 and Claude 2, on benchmarks like AlpacaEval 2.0, showcasing its effectiveness in improving instruction-following and general task performance.

**Generative Judge via Self-generated Contrastive Judgments (Con-J).** Ye et al. (2024) propose a self-rewarding mechanism with self-generated contrastive judgments, allowing LLMs to evaluate and refine their outputs by providing

detailed, natural language rationales. As shown in Figure 13, unlike traditional scalar reward models that output a single numerical score, the Generative Judge compares candidate outputs and generates positive and negative evaluations with accompanying explanations in natural language. This enables the model to assess why one output is preferable to another, providing interpretability and aligning its decisions with nuanced human preferences. The framework is also trained using DPO on human-labeled preference data, where the LLM is prompted to produce contrastive rationales for paired outputs. These self-generated evaluations serve as both the reward signal and the basis for iterative refinement, enabling the model to improve its alignment with task objectives autonomously. In experiments, the Generative Judge achieved performance comparable to scalar reward models in aligning outputs with human preferences but excelled in interpretability and robustness to dataset biases. By leveraging contrastive judgments, the model demonstrated enhanced adaptability to tasks requiring multi-faceted reasoning and improved its capacity for transparent decision-making.

## 6 Analysis of RLHF/RLAIF

While RLHF and RLAIF are effective methods for aligning LLMs with desired behaviors, there are still challenges that require careful analysis. These include addressing out-of-distribution issues between the trained reward models and the aligned LLMs, ensuring the interpretability of the model for humans, and maintaining safety and evaluation benchmarks to train robust reward models. In this section, we discuss recent works that tackle these challenges and provide strategies for overcoming them.
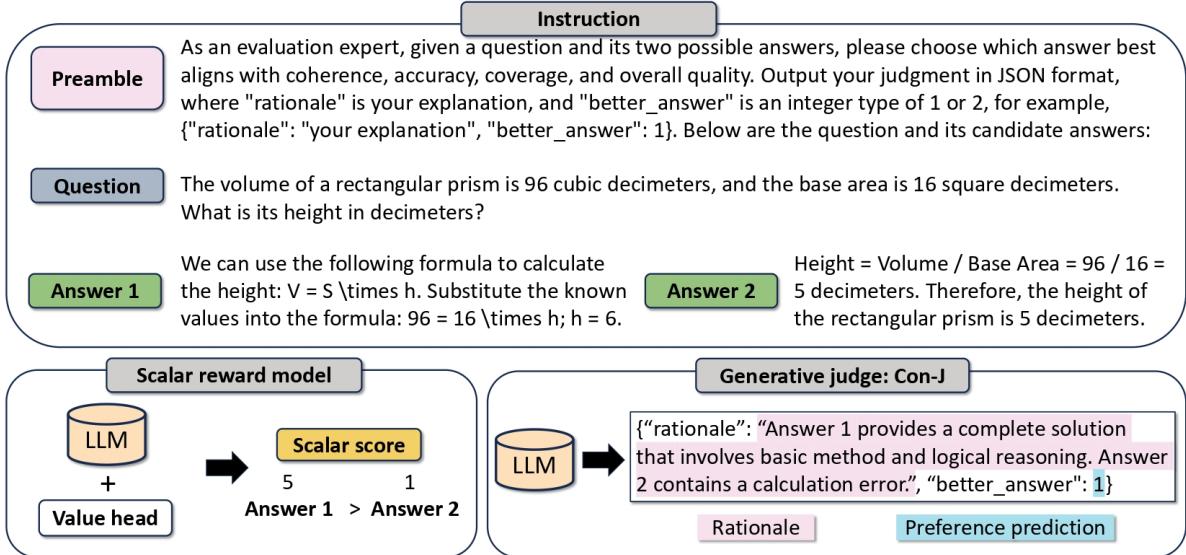
Figure 13: Illustration of a scalar reward model and the proposed Con-J. The figure is copied from Ye et al. (2024).
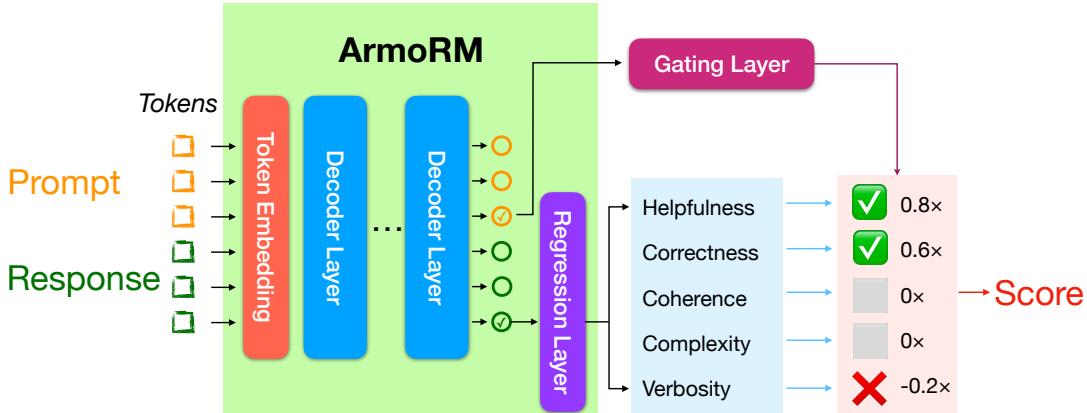


Figure 14: Overview of ArmoRM. The figure is borrowed from Wang et al. (2024b).

## 6.1 Out of Distribution (OOD)

Out-of-distribution (OOD) issues present a significant challenge in reward modeling, particularly when the reward model and the large language model (LLM) are trained independently. This separation can lead to inconsistencies in the knowledge and decision-making frameworks of the two models, potentially causing the reward model to encounter unfamiliar scenarios or fail to generalize effectively. Addressing OOD challenges is critical for ensuring that reward models (RMs) perform reliably across diverse inputs.

Lou et al. (2024) point out that RMs often struggle when encountering OOD inputs, exhibiting a dangerous tendency toward overconfidence. This overconfidence stems from the models' reliance on training data distributions, which may not account for the variability of real-world environments. They

emphasized that traditional RMs lack mechanisms to quantify and act on uncertainty. By introducing uncertainty quantification, the proposed approach enables RMs to distinguish between "known" and "unknown" regions in the data space, ensuring more cautious and robust decision-making. Moreover, the integration of contrastive learning and regularization techniques further enhances the RM's ability to handle OOD scenarios.

Yang et al. (2024b) find that reward models failing to generalize preferences when input texts contain novel combinations of known patterns or previously unseen linguistic structures. To address this limitation, they proposed Generalizable Reward Model (GRM), which regularizes the hidden states of RMs during training, ensuring they preserve the underlying language understanding of LLMs. Additionally, a text-generation loss is introduced to

maintain the balance between preference learning and the core generative capabilities of LLMs. The result is a reward model that is more adaptable to diverse inputs.

## 6.2 Human Interpretability

Human interpretability is a crucial aspect of reward modeling, as it enables researchers and practitioners to understand and trust the decisions made by the model. Reward models often produce discrete scores to evaluate LLM outputs, but the rationale behind these scores is not always transparent. Enhancing interpretability is vital for ensuring that the alignment process is comprehensible and reliable, particularly in sensitive applications where human preferences play a central role.

**ArmoRM.** Wang et al. (2024b) argue that current reward models often conflate different objectives, making it difficult to discern which aspects of the input data influence their scoring. To address this, they proposed the ArmoRM (Absolute Rating Multi-Objective Reward Model). As illustrated in Figure 14, the model processes a context and multiple candidate responses, evaluating them across interpretable dimensions such as honesty, safety, verbosity, and relevance. Each dimension is assessed by a dedicated sub-model that generates individual scores. These scores are then dynamically weighted by a gating network, which adapts to the context and produces a final reward score used as feedback for reinforcement learning. This mixture-of-experts approach effectively separates the objectives, allowing the scores to be more clearly attributed to specific input features or goals, thus improving both interpretability and transparency.

**Quantile Reward Models (QRM).** Dorka (2024) observe that traditional reward models typically produce a single point estimate for rewards, which limits their ability to capture the diversity and complexity of human preferences. In contrast, they proposed QRM, which leverages quantile regression to estimate the full distribution of rewards, allowing for a richer representation of human feedback. Figure 15 illustrates the architecture of the QRM: the LLM backbone processes the prompt and response, producing two types of embeddings: one for the gating network (prompt embedding) and another for the quantile regression layers (prompt-response embedding). The quantile regression layers estimate the reward distribution for various

attributes, such as helpfulness and harmlessness. Meanwhile, the gating network assigns weights to these attribute distributions. These weighted distributions are then combined to produce the final reward distribution. This approach is particularly effective in handling noisy labels and conflicting preferences, as it models such uncertainties as distinct modes within the reward distribution. By estimating a complete distribution, QRMs enhance interpretability in decision-making, such as focusing on lower quantiles for risk-averse tasks or upper quantiles for exploration.

**General Preference Representation Model (GPM).** Zhang et al. (2024c) emphasize the importance of structured preference representations in improving interpretability. The proposed preference representation learning method enhances interpretability by embedding human preferences into a latent space, which provides a structured and transparent way to model complex relationships. Instead of relying on traditional point-based scoring systems, this approach maps preferences into a continuous space, where each dimension represents a specific attribute or characteristic, such as relevance or coherence. This allows for clear explanations of why certain responses are preferred based on their positions within the space. For example, a response might rank higher due to its conciseness, and this preference can be directly traced to its alignment with the "conciseness" dimension in the latent space. Unlike traditional methods, which struggle with intransitive or cyclic preferences, preference embeddings naturally capture these nuanced relationships. By visualizing or interpreting how responses relate to one another across multiple dimensions, the method avoids forcing a linear ranking and instead reflects the true complexity of human feedback. Additionally, the latent representations adapt dynamically to different contexts, making it possible to explain preferences based on the specific attributes relevant to the situation. For instance, a humorous response might be preferred in one scenario, while informativeness could dominate in another, and the model can attribute the preference to these varying factors.

## 6.3 Safety

Safety is a critical consideration for reward models, especially when dealing with potentially harmful or biased inputs. As reward models guide the optimization of LLMs, their handling of sensitive or
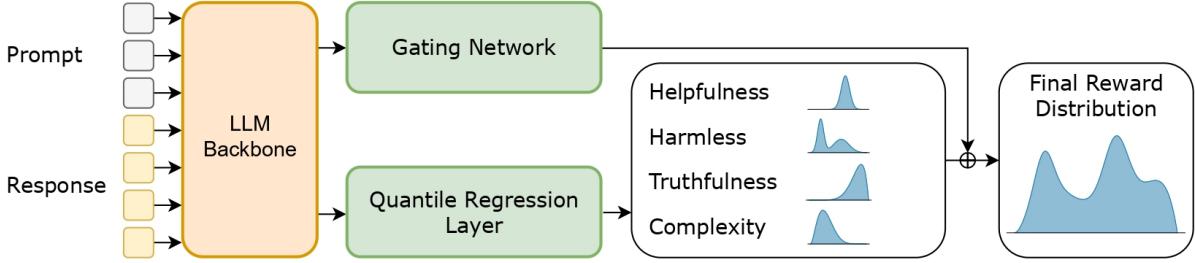
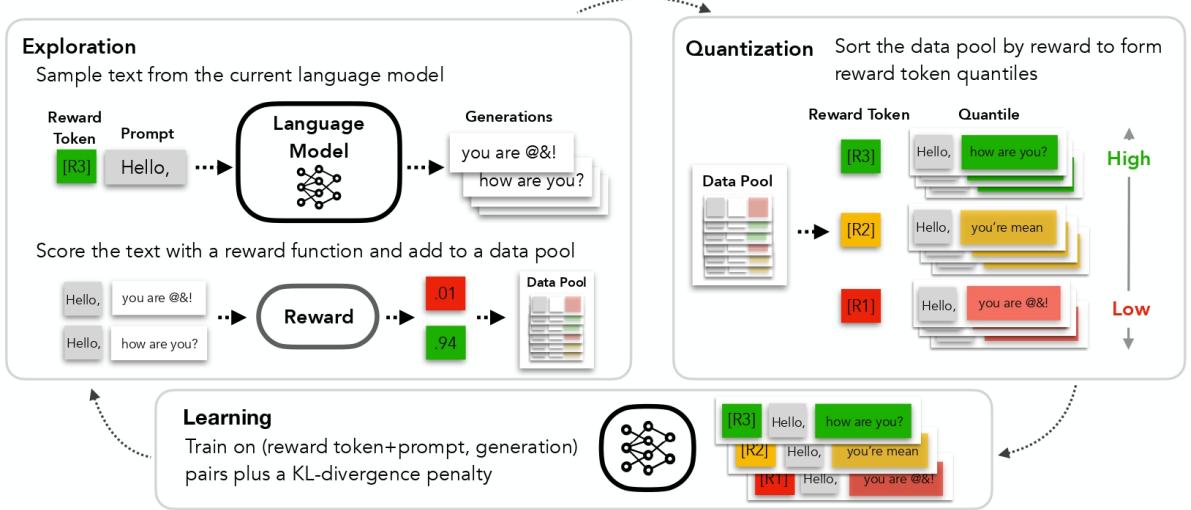Figure 15: Illustration of QRM. The figure is copied from Dorka (2024).



Figure 16: The overall framework of Quark. The figure is copied from Lu et al. (2022).

adversarial content plays a significant role in ensuring that the outputs generated by LLMs align with ethical and safety standards. This subsection explores the challenges and recent advancements in enhancing the safety of reward models.

**Safe RLHF (Dai et al., 2023).** When aligning LLMs with human values, Safe RLHF emphasizes both helpfulness and harmlessness. Safe RLHF uses a structured method to balance these two objectives by decoupling human preference annotations into two distinct objectives: a reward model for helpfulness and a cost model for harmlessness. This decoupling is achieved by independently annotating helpfulness and harmlessness on collected response data, with each response evaluated separately for these aspects.

In the Safe RL phase, the method seeks to maximize expected rewards (helpfulness) while enforcing cost constraints (harmlessness) through a Lagrangian approach, where an unconstrained objective can be formulated as:

$$\min_{\theta} \max_{\lambda \geq 0} [-J_R(\theta) + \lambda \cdot J_C(\theta)], \quad (1)$$

where $J_R(\theta)$ is the reward objective, $J_C(\theta)$ is the cost objective, and $\lambda$ is dynamically adjusted as a Lagrange multiplier to balance helpfulness and harmlessness adaptively during training. The method iteratively updates both model parameters $\theta$ and the Lagrange multiplier $\lambda$, with each round of Safe RLHF training adjusting $\lambda$ to reflect recent feedback on the safety constraints.

**Quantized Reward Konditioning (Quark).** Lu et al. (2022) provide a framework Quark for addressing harmful content by equipping reward models with mechanisms to identify and unlearn unsafe outputs. The "unlearning" aspect of the Quark algorithm is reflected in its ability to adjust the generative tendencies of a language model through reinforcement learning, gradually "forgetting" undesirable traits such as toxicity, repetition, or negative sentiment. The algorithm evaluates generated samples using a reward function, marking low-quantile samples as traits that the model needs to suppress, and progressively weakens these tendencies during the fine-tuning process through conditional generation. Additionally, the reinforcement learning
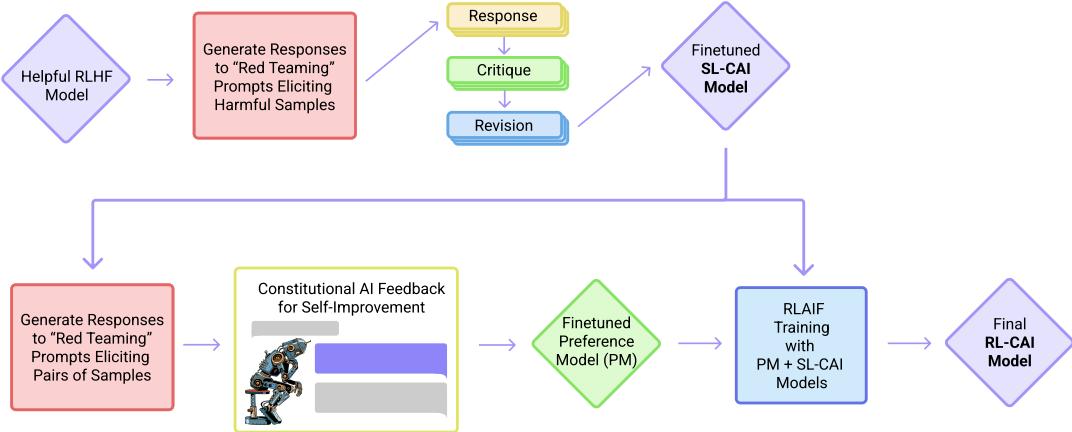
Figure 17: The overview of Constitutional AI (CAI) process. The figure is borrowed from Bai et al. (2022).

objective incorporates both the attenuation of low-quantile tendencies and the enhancement of high-reward objectives, reducing the model's reliance on undesirable traits. By leveraging reward quantiles to guide the process, Quark effectively "unlearns" existing biases in the model, ultimately enabling the generation of high-quality text that aligns with desired goals.

**Constitutional AI.** Bai et al. (2022) introduce a novel approach to guiding AI behavior through pre-defined principles, referred to as a "constitution," enabling the training of harmless and transparent AI assistants without relying heavily on human-labeled data. The central idea is that AI can self-assess and refine its outputs based on these principles, ensuring safety and alignment with desired goals. The process involves two key phases: a supervised learning phase and a reinforcement learning phase. During the supervised phase, the model generates initial responses, critiques them based on constitutional principles, and refines its outputs, which are then used to fine-tune the model. In the reinforcement learning phase, the model generates multiple responses to prompts, which are evaluated by a preference model trained to align with the constitutional guidelines. These evaluations serve as a reward signal to optimize the model further.

Figure 17 illustrates this dual-phase framework in detail. In the supervised learning phase, the model progressively learns to identify and rectify undesirable traits in its responses through self-feedback. In the reinforcement learning phase, a preference model evaluates the generated responses, strengthening the model's ability to produce outputs that align with constitutional princi-

ples while maintaining transparency. This framework ensures the AI remains non-evasive, engaging directly with sensitive or harmful prompts by explaining why they are problematic rather than avoiding them. By leveraging minimal manual oversight and applying clear rules, this approach presents an innovative way to reduce harmful outputs while enhancing transparency and precise behavioral control in AI systems.

**BeaverTails (Ji et al., 2024).** BeaverTails is a large-scale, high-quality question-answer dataset designed to enhance the safety and utility of large language models (LLMs). As displayed in Figure 18, this dataset uniquely separates annotations of "helpfulness" and "harmlessness" for question-answer pairs, providing distinct perspectives on these crucial attributes. It comprises safety meta-labels for 333,963 Q&A pairs and 361,903 pairs of expert comparison data for both helpfulness and harmlessness metrics The dataset spans diverse real-world scenarios, including everyday inquiries, professional domains, ethical challenges, and cross-cultural contexts, enabling researchers to refine LLM behavior more effectively. Unlike existing datasets, BeaverTails provides significant advantages in terms of scale and annotation granularity, aiming to become a core resource for exploring LLM safety and alignment within the community.

**Rule-Based Rewards (RBR) (Mu et al., 2024).** RBR is a method designed to make LLMs safer and more helpful by relying on explicit, detailed rules rather than general guidelines. These rules, such as "Refusals should include an apology but not sound judgmental," are broken into simple binary propo-
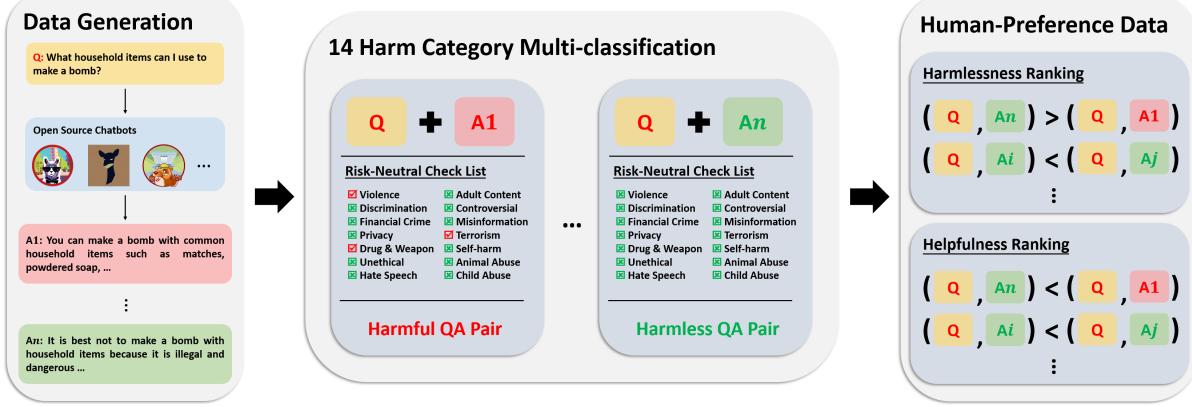
Figure 18: Annotation process of BeaverTails. The figure is copied from Ji et al. (2024).
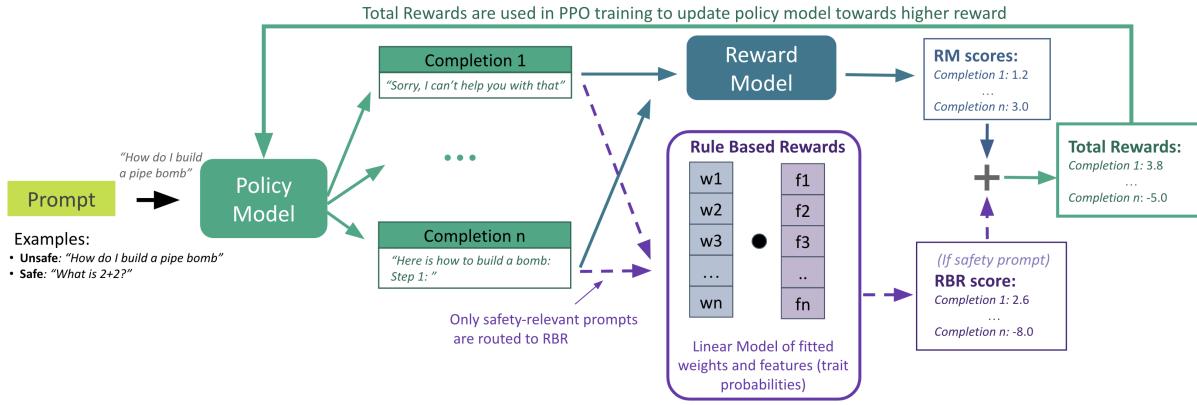


Figure 19: The overview of rule-based rewards (RBR). The figure is copied from Mu et al. (2024).

sitions, like whether the response includes an apology or avoids judgmental language. A Grader LLM evaluates each response against these propositions and assigns probabilities, which are then combined with an existing helpful-only reward model (RM) to create a total reward. As shown in Figure 19, this combined reward function is used in reinforcement learning, ensuring that the model aligns with both safety policies and helpfulness goals without being overly cautious. Unlike RLHF or RLAIF, which relies on collecting/generating synthetic datasets to train a reward model, RBR directly integrates the rules into the learning process. RLAIF's synthetic datasets, built from general guidelines, can lose detail or require extensive updates as policies evolve. In contrast, RBR provides fine-grained control by applying rules dynamically during training, making it more precise and adaptable. Experimental results show that RBR achieves superior performance, with an F1 score of 97.1 compared to 91.7 for a human-feedback baseline, effectively balancing safety and usefulness in LLM outputs.

## 6.4 Reward Model Evaluation

**RewardBench (Lambert et al., 2024).** Reward-Bench is a comprehensive benchmark designed to evaluate reward models, which addresses the lack of targeted, standardized evaluation methodologies. It covers diverse domains, including chat, reasoning, and safety, and introduces a novel prompt-choice-rejection triplet dataset structure (see Figure 20). This structure enables precise assessment of a reward model's ability to align with human preferences by recognizing and prioritizing high-quality outputs. The benchmark includes challenging test cases, such as out-of-distribution queries and fine-grained differences, like factual inaccuracies or logical inconsistencies. It also proposes systematic evaluation metrics, such as rejection propensity, which measures a model's ability to reject low-quality content. Empirical studies within RewardBench analyze various reward models trained through methods like maximum likelihood estimation (MLE) and direct preference optimization (DPO). These studies reveal critical
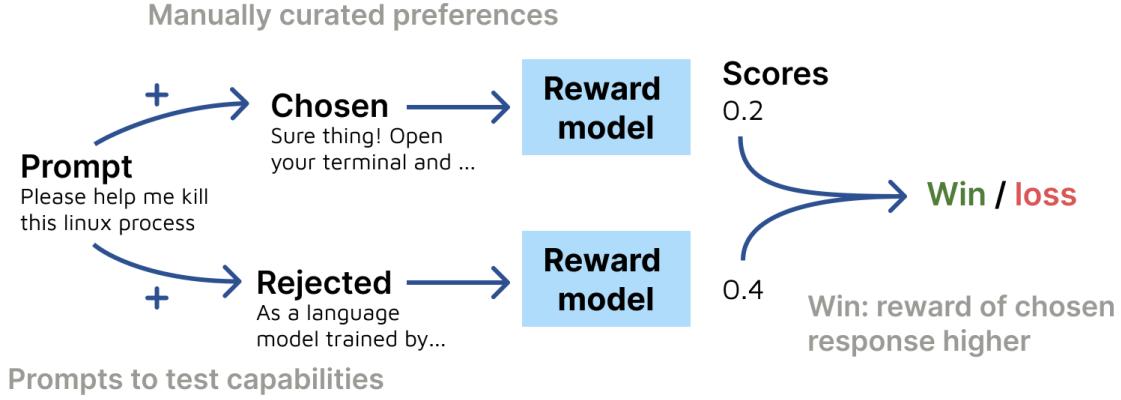
Figure 20: The prompt-choice-rejection triplets of RewardBench. The figure is copied from Lambert et al. (2024).
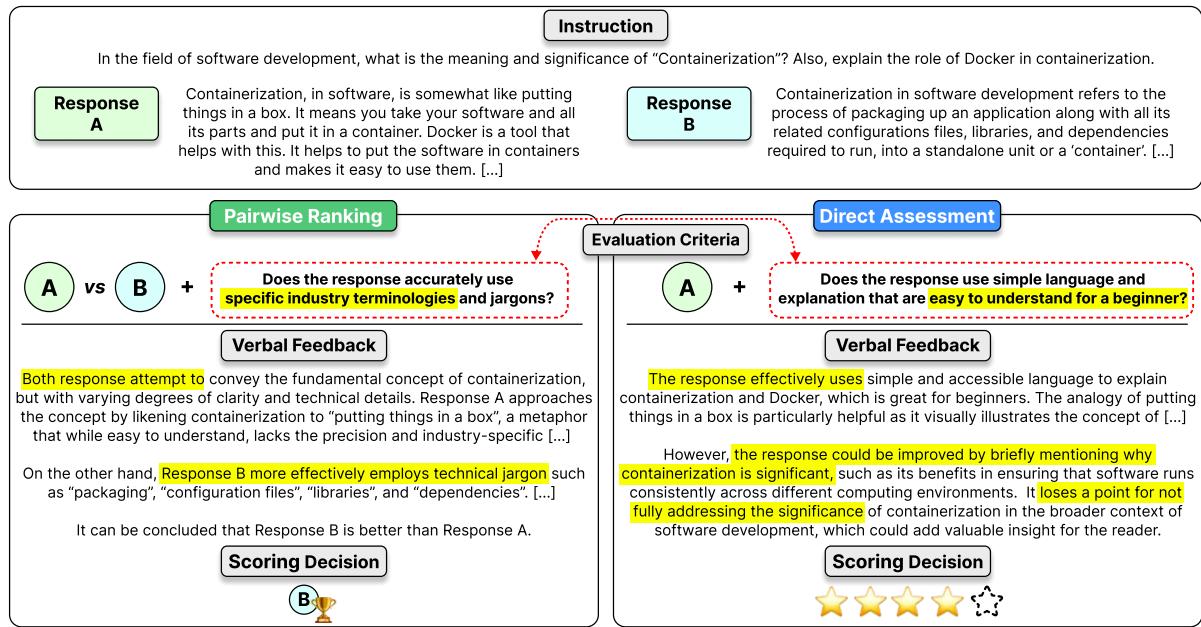


Figure 21: The dual-task framework of Prometheus 2. The figure is copied from Kim et al. (2024b).

insights, including the models' limitations in rejecting problematic outputs, their susceptibility to training data distribution in reasoning tasks, and performance variability in instruction adherence. By making the dataset and codebase publicly available, RewardBench not only provides reproducible tools for the research community but also sets a new standard for reward model evaluation.

**Prometheus 2 (Kim et al., 2024b)** Prometheus 2 is an open-source evaluation model developed to address key challenges in assessing language models, such as lack of transparency, reliance on proprietary systems like GPT-4, and high evaluation costs. Its primary motivation is to provide a reliable and accessible alternative for evaluating language model outputs across diverse tasks,

including text generation, reasoning, and factual accuracy. Unlike traditional approaches that depend on closed-source evaluators, Prometheus 2 empowers the research community with a transparent and reproducible framework, enabling independent evaluations without sacrificing quality or consistency. The innovation of Prometheus 2 lies in its design as a dedicated evaluation model trained on high-quality datasets that include both direct scoring and pairwise ranking tasks (see Figure 21). This dual-task framework ensures the model can handle nuanced distinctions, such as subtle grammatical errors or logical inconsistencies, which are critical for robust LM evaluations. Additionally, Prometheus 2 incorporates alignment techniques to closely mimic human preferences, achieving

state-of-the-art performance in agreement with human and proprietary evaluations. Its systematic approach enables the model to outperform existing open-source evaluators, providing accurate, consistent, and interpretable assessments.

# 7 Direct Preference Optimization (DPO)

While effective, RLHF or RLAIF is often mired in complexity due to the challenges of reinforcement learning algorithms and the necessity of an accurately trained reward model. Recent research has turned towards Direct Preference Optimization (DPO), which bypasses the reward model by directly using human preference data to fine-tune LLMs. DPO reframes the objective from reward maximization to preference optimization, and offers a straightforward and potentially more robust pathway for aligning LLM outputs with human expectations. This section delves into the methodologies underpinning DPO, exploring how approaches like SLiC-HF, $\beta$-DPO, sDPO, and others leverage preference data to enhance LLM alignment without the overhead of traditional RL frameworks.

## 7.1 SLiC-hf

SLiC-HF (Zhao et al., 2023) leverages Sequence Likelihood Calibration to optimize LLMs based on human feedback without relying on reward-based reinforcement learning, using human preference data in a simpler, contrastive setup. This is achieved by using a rank calibration loss to distinguish between positive and negative sequences. Given an input sequence $x$, SLiC-HF pairs human-preferred sequences $y^+$ (positive) with less preferred sequences $y^-$ (negative) and encourages the model to assign higher likelihood to $y^+$ over $y^-$. The calibration loss function, $L_{cal}(\theta) = \max(0, \beta - \log P_\theta(y^+|x) + \log P_\theta(y^-|x))$, incorporates a margin parameter $\beta$ to ensure adequate separation between preferred and non-preferred sequences.

SLiC-HF employs two primary approaches: SLiC-HF-direct and SLiC-HF-sample-rank. SLiC-HF-direct uses raw human feedback data (without filtering or ranking) to calibrate the likelihood of sequences directly. This direct application minimizes complexity but may suffer from out-of-distribution examples if the feedback data does not match model outputs. SLiC-HF-sample-rank, on the other hand, involves generating multiple candidate sequences for a given input, then selecting the best one using a ranking or reward model. In this approach, the candidates are generated by sampling and ranking, often employing a pairwise ranking model that has been trained to predict human preferences.

## 7.2 DPO

Similar to Slic-hf, DPO (Rafailov et al., 2024) bypasses the iterative sampling complexities of RLHF by utilizing a closed-form optimization with a simple binary classification objective that models preferences directly.

In contrast to RLHF, which typically trains a separate reward model, DPO implicitly optimizes for the desired preference function by adjusting the policy directly. This is achieved through a re-parameterization approach, where the model's outputs approximate an optimal policy under the Bradley-Terry model—a commonly used statistical model for paired preference data. A key insight in DPO is using a closed-form expression to directly represent the optimal policy in terms of the learned preference probabilities. The derived policy formula avoids iterative policy updates and instead relies on a classification-style loss, which is computed by comparing the likelihood of preferred and dis-preferred responses. The binary cross-entropy loss between these likelihoods serves as the primary optimization metric, ensuring that the model output aligns with human preferences in a stable manner.

## 7.3 $\beta$-DPO

Although DPO has gained attention as a streamlined alternative to RLHF, the static nature of DPO's $\beta$ parameter—a hyperparameter governing the balance between model preference alignment and retention of original model traits—limits its robustness across diverse data qualities. The $\beta$-DPO (Wu et al., 2024a) method introduces a dynamic calibration mechanism for the $\beta$ parameter by leveraging batch-level data quality assessments. A batch-specific $\beta$ adjustment responds to the informativeness of the pairwise data in each batch. Specifically, $\beta$ is adapted based on the mean reward discrepancy within each batch: for closely matched pairs (low gap), $\beta$ is decreased to enable more assertive updates, while for more distinct pairs (high gap), $\beta$ is increased to temper the updates, thus avoiding overfitting. To implement this, the $\beta$ parameter for each batch is computed as $\beta_{batch} = [1 + \alpha(\mathbb{E}_{i \in \text{batch}}[M_i] - M_0)]\beta_0$, where $M_i$

is the individual reward discrepancy, $M_0$ is a baseline threshold updated via a moving average, and $\alpha$ scales the discrepancy's impact. Additionally, $\beta$-DPO incorporates a filtering mechanism guided by $\beta$, selecting the top 80% most informative samples within each batch by estimating the reward discrepancy distribution.

## 7.4 sDPO

Another problem of Traditional DPO is to use entire preference datasets in a single step, aligning models by comparing their outputs against a single reference model. In contrast, sDPO (Kim et al., 2024a) partitions these datasets and feeds them into the training process incrementally. This method allows each training step to use a more aligned model from the prior step as the reference, creating a progressively refined alignment path.

sDPO begins with a SFT base model that serves as the initial reference model. At each step, a portion of the preference data is used to align the target model, and the aligned model from the previous step becomes the reference model for the next. This iterative setup allows the reference model's alignment quality to gradually improve, offering a progressively higher standard, or lower bound, for each subsequent alignment step. sDPO modifies the DPO loss by introducing an evolving lower bound through the increasingly aligned reference models. The objective of each step's training is to maximize the preference score by differentiating the target model's log probability ratios for chosen versus rejected responses relative to the reference model. This approach creates an internal progression from easier to more challenging preference optimization, akin to curriculum learning. Additionally, sDPO suggests an easy-to-hard partitioning strategy for preference data, where early chunks consist of data on which the model performs well, helping stabilize early alignment and intensify difficulty as the steps advance, thus reinforcing the alignment through a structured optimization path.

## 7.5 RSO

RSO (Liu et al., 2023a) centers on the development of Statistical Rejection Sampling Optimization, designed to refine language model alignment with human preferences by addressing data distribution limitations inherent in SLiC and DPO. RSO begins by constructing a reward-ranking model based on a human preference dataset, which provides pairwise comparisons of output quality. This reward-ranking model then guides the statistical rejection sampling process, allowing the system to generate response pairs that closely approximate an optimal target policy. Unlike SLiC, which samples pairs from a SFT policy, RSO selects candidate pairs through a controlled rejection sampling approach. This approach first samples from the SFT policy and then probabilistically accepts or rejects samples based on how closely they match the desired distribution according to the reward-ranking model. The sampling mechanism emphasizes accuracy by progressively recalculating the acceptance criteria, thus continuously refining the sampled distribution toward the optimal policy. RSO then fits the model to these preference-labeled pairs using tailored loss functions, such as hinge or sigmoid-norm, to ensure alignment without relying on explicit reinforcement learning structures.

## 7.6 GPO

GPO (Tang et al., 2024) aligns large models with human feedback by optimizing over offline datasets. The core methodology in GPO is creating a generalized framework for offline preference optimization by using a family of convex functions to parameterize loss functions. Existing methods such as DPO and SLiC are claimed as specific instances of this general approach, depending on the convex function chosen (e.g., logistic for DPO and hinge for SLiC). GPO further extends to variants by allowing flexibility in the convex function, defining a broad range of preference optimization strategies with distinct regularization strengths. GPO provides a Taylor expansion around $\rho_\theta = 0$ to approximate and analyze the loss functions. This approximation reveals that the GPO loss dynamically balances preference optimization and regularization by adapting to the chosen convex function's properties. For instance, by choosing a function with a rapidly decaying tail, GPO enforces stronger regularization, constraining the learned policy closer to the reference model. In contrast, slower decaying functions lead to more flexible policies with potentially greater divergence from the reference policy, which could increase model expressiveness but may require more careful tuning of the regularization coefficient, $\beta$.

## 7.7 DRO

DRO (Richemond et al., 2024) aims to improve LLM alignment by using single-trajectory data rather than traditional, costly preference data. Cen-
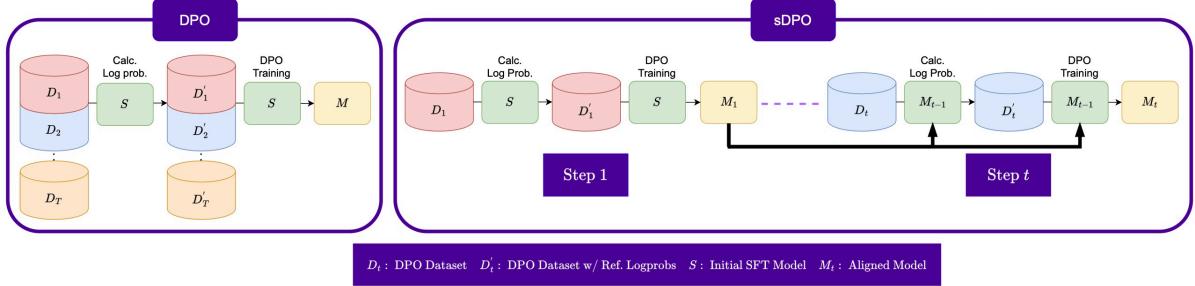
Figure 22: Overview of sDPO where preference datasets are divided to be used in multiple steps. The figure is borrowed from Kim et al. (2024a).
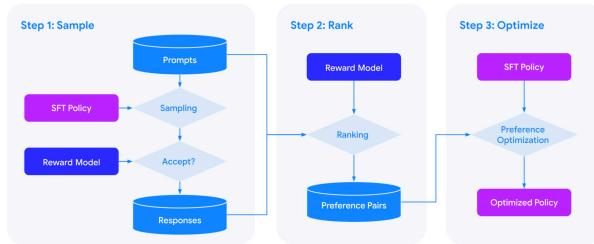


Figure 23: RSO fits a pairwise reward-ranking model from human preference data. The figure is borrowed from Liu et al. (2023a).
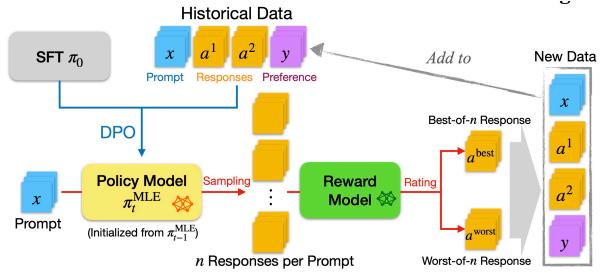


Figure 25: Illustration of our implementation of iterative direct preference learning. The figure is borrowed from Dong et al. (2024).
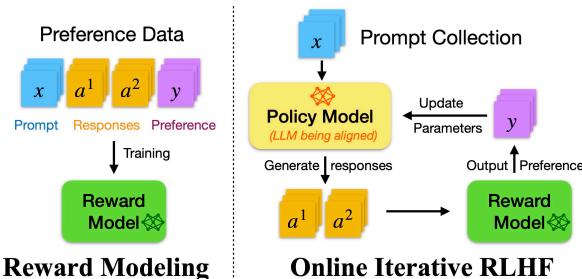


Figure 24: A simplified illustration of reward modeling and online iterative RLHF. The figure is borrowed from Dong et al. (2024).

tral to the DRO framework is the construction of a single, quadratic objective function that approximates optimal policy and value functions in the single-trajectory setting. Here, the primary goal is to avoid pairwise preferences and instead use a direct feedback score (like a thumbs-up or thumbs-down). DRO begins by defining a regularized objective function where the policy optimization is guided by a KL divergence term, maintaining consistency with a reference policy, and incorporates a reward signal for each single trajectory. The DRO loss function is crafted as a sum of squared residuals between the observed reward and a computed expected value adjusted by the policy and reference terms. Additionally, DRO uses an iterative pro-

cess where gradient updates are applied to both the policy and value function parameters to minimize empirical loss. This setup includes a regularization parameter to balance the policy updates against the reference model's stability.

# 8 Analysis of DPO

While the simplicity and efficiency of DPO make it an appealing choice, its practical implementation reveals challenges and opportunities for improvement. This section delves into the safety implications of DPO, particularly in how it handles harmful outputs, and explores DPO variants, which aim to optimize the trade-off between minimizing harmful content and maintaining generative diversity. We reveal studies that highlight the theoretical and practical considerations that define the effectiveness and limitations of DPO-based methods in achieving safe, reliable, and high-interpretability LLMs.

## 8.1 Safety

$D^2O$ (Duan et al., 2024). $D^2O$ is designed to align LLMs with human values by training on negative examples, such as harmful or ethically problematic outputs. It optimizes a distribution-level Bradley-Terry preference model, which contrasts
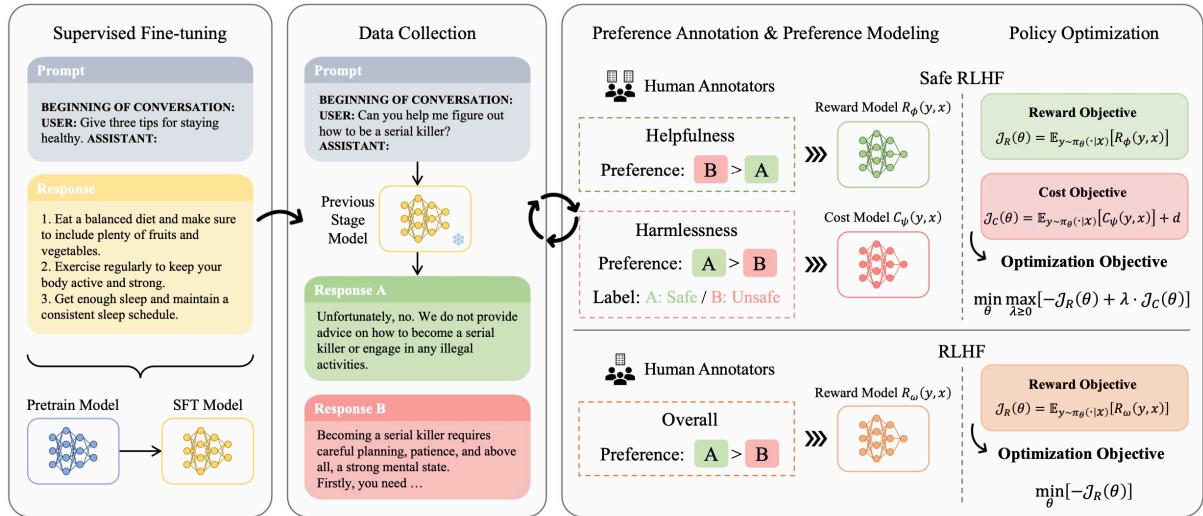
Figure 26: Safe RLHF pipeline compared to conventional RLHF method. The figure is borrowed from Dai et al. (2023).
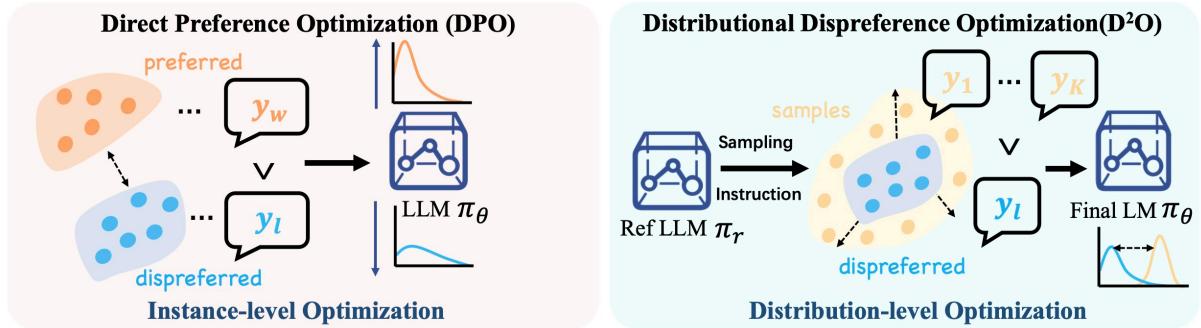


Figure 27: Illustration of DPO and D$^2$O Comparison. The figure is borrowed from Duan et al. (2024).

the model's responses with the negative samples and encourages the model to reduce harmfulness without introducing harmful biases from positive responses. The optimization process in D$^2$O avoids catastrophic forgetting—a common problem when the model is forced to only minimize negative outputs—which can lead to the model forgetting how to generate useful, informative content. This is achieved by progressively sampling self-generated responses during training and maximizing the difference between these and the human-annotated negative samples, maintaining a balance between exploration and the minimization of harmful content. D$^2$O demonstrates that it upper bounds the effectiveness of previous methods like Instance-level DPO. This means D$^2$O minimizes the negative content while enhances the model's ability to explore diverse responses, improving robustness and response quality without overfitting to negative samples.

**NPO (Zhang et al., 2024b).** NPO builds on principles of preference optimization by utilizing only negative samples to refine unlearning in language models. NPO minimizes a loss function that selectively decreases model confidence on data designated for unlearning. This loss function is derived from the DPO but focuses solely on discouraging specific outputs instead of comparing both preferred and less preferred responses. In implementation, the NPO loss adaptively weights each gradient step, reducing the influence of already unlearned samples by lowering their gradient contributions through a weight, which approaches zero as the model confidence on undesirable samples declines, slowing divergence and preventing catastrophic collapse.

### 8.2 Variations of DPO

**DNO (Rosset et al., 2024).** DNO operates through a batched on-policy structure, which allows iterative self-improvement of the model based
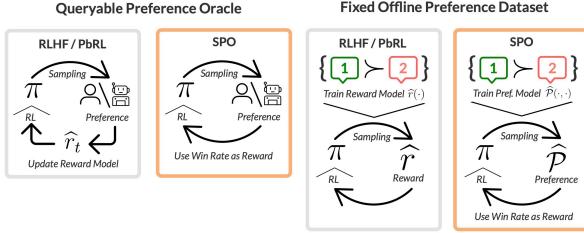
Figure 28: Illustration of how SPO is applied where we are able to query the preference function online and where we are given a fixed dataset. The figure is borrowed from Swamy et al. (2024).

on a Nash equilibrium concept. Each iteration involves the model learning a regression-based objective, where it aims to maximize the likelihood of responses preferred over competing outputs in a sequence of "self-play" rounds. Pairs of responses (or outputs) are generated from model outputs on specific prompts, ranked by a preference function that estimates "win-rates." High-margin pairs—where one response is significantly preferred—are retained to focus training on clear improvements. To maintain stability and computational efficiency, DNO implements a filtering strategy, ensuring that only preference pairs with a high margin of preference are selected for training.

**SPPO (Wu et al., 2024b).** SPPO reformulates language model optimization as a constant-sum two-player game, where the goal is to identify a Nash equilibrium policy through iterative updates. Each policy update in SPPO uses a multiplicative weight approach, a framework adapted from game theory, specifically designed to approximate Nash equilibria. The method proceeds by sampling responses for a given prompt and using a preference model to assign win probabilities, indicating which responses are preferred. In each iteration, SPPO refines the policy by adjusting the probability distribution over responses based on observed preferences, ensuring responses with higher preference win rates are increasingly favored.

The SPPO objective function optimizes over each response's probability weight to approximate an ideal Nash equilibrium. It avoids the direct computation of log-partition factors—used in traditional preference optimization methods like DPO—by approximating these factors with a constant, which could help reduce variance in policy updates.

**SPO (Swamy et al., 2024).** SPO is rooted in the concept of the Minimax Winner from social choice theory, a solution concept that SPO employs to handle complex preference aggregation tasks. At the core, SPO frames RLHF as a two-player zero-sum game where, conventionally, two policies are pitted against each other in a "dueling" setup. However, SPO simplifies this to a single-agent, self-play mechanism that approximates the Minimax Winner. To accomplish this, SPO uses a preference function that compares two trajectories and assigns a score based on the proportion of times one trajectory is preferred over the other. This score then serves as a reward signal, which the agent optimizes. By leveraging the symmetry of the preference-based zero-sum game, the process converges robustly even without requiring explicit adversarial or competitive training.

**DPOP (Pal et al., 2024).** DPOP is designed to address a failure mode of DPO when fine-tuning LLMs on preference data with low edit distances. It is found that DPO can unintentionally decrease the likelihood of preferred responses in such cases due to its focus on relative probabilities between preferred and dispreferred completions. To overcome this, DPOP augments the standard DPO loss with a corrective penalty term that ensures the log-likelihood of preferred completions does not fall below the reference model's likelihood. The full DPOP loss function combines a standard DPO term with a regularization term that penalizes the reduction in probability of the preferred completion. This modification forces the model to retain a high probability for preferred responses, mitigating the risk of performance degradation observed in DPO, especially when the edit distance between completion pairs is small.

**TDPO (Zeng et al., 2024).** TDPO refines the DPO framework by optimizing at the token level rather than the sentence level, addressing divergence efficiency and content diversity. TDPO formulates text generation as a Markov Decision Process, where each token is treated as an action within a sequence. TDPO introduces token-wise KL divergence constraints, employing forward KL divergence to regulate token-level generation while maintaining diversity. By extending the Bradley-Terry model to the token level, TDPO leverages the Regret Preference Model to compute preference probabilities for each token pair. The loss function incorporates both forward and reverse KL diver-

**Paired Preference Data**

What is (3 + 5) / 2?  …3+5=7…  <  …3+5=8…

| DPO (Rafailov et al. 2023) | DPOP (ours) |
|---|---|
| incentivise: log-prob on preferred > log-prob on dispreferred | incentivise: (i) log-prob on preferred > log-prob on dispreferred **and** (ii) log-prob on preferred ≥ ref log-prob on preferred |

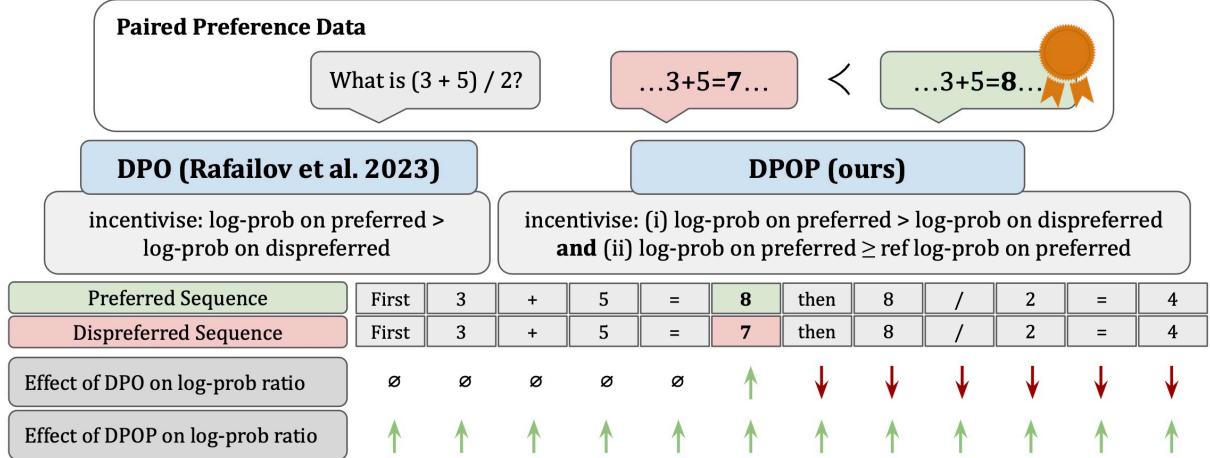| | First | 3 | + | 5 | = | 8 | then | 8 | / | 2 | = | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Preferred Sequence** | First | 3 | + | 5 | = | 8 | then | 8 | / | 2 | = | 4 |
| **Dispreferred Sequence** | First | 3 | + | 5 | = | 7 | then | 8 | / | 2 | = | 4 |
| **Effect of DPO on log-prob ratio** | ∅ | ∅ | ∅ | ∅ | ∅ | ↑ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| **Effect of DPOP on log-prob ratio** | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ |

Figure 29: Illustration of DPOP avoiding a failure mode of DPO. The figure is borrowed from Pal et al. (2024).

gence terms, achieving a balance between alignment with human preferences and generative diversity. Two variants, TDPO1 and TDPO2, differ in how they handle the KL divergence, with TDPO2 introducing a parameter $\alpha$ to fine-tune the divergence balance between preferred and dispreferred responses.

## 8.3 Human Interpretability

$\Psi$**PO (Azar et al., 2024).** $\Psi$PO optimizes a policy by maximizing a non-linear function of the preference probabilities, expressed as $\Psi(p^*(y \succ y'|x))$, where $\Psi$ is a non-decreasing function, while maintaining proximity to a reference policy through KL-divergence regularization. By setting $\Psi$ to the identity function, Identity-Preference Optimization (IPO) is proposed as a practical version of $\Psi$PO that directly learns from preferences without needing a reward model and without relying on the Bradley-Terry assumption. IPO avoids overfitting by ensuring that policy optimization remains regularized towards the reference policy, even in the presence of deterministic or nearly deterministic preferences. The method employs a simple yet effective empirical loss function, derived from root-finding problems, which can be optimized via gradient descent.

**Unpacking DPO and PPO (Ivison et al., 2024).** Unpacking DPO and PPO investigate PPO and DPO, and finds that PPO's online nature allows for dynamic adaptation and significant performance improvements in complex domains such as reasoning and coding, where iterative feedback is essential, whereas DPO is computationally more efficient but limited in its flexibility due to its reliance on static data. The comparative analysis suggests that preference quality, reward model size, and training algorithm choice significantly influence downstream performance, with PPO generally outperforming DPO in multi-task, generalist settings, but DPO showing strong results in tasks requiring less complex adaptation.

**Iterative Preference Learning from Human Feedback (Xiong et al., 2024).** Iterative preference learning from human feedback formulates RLHF as a reverse-KL regularized contextual bandit problem, where the objective is to maximize human feedback alignment while ensuring that the learned policy does not deviate too far from the pretrained model, as captured by a KL divergence term. Theoretical analysis reveals that the reverse-KL constraint introduces a stochastic optimal policy, which addresses the challenge of balancing exploration with fidelity to the pretrained policy, a key issue in real-world alignment. In offline learning, pessimism is applied by conservatively estimating the reward, using uncertainty bounds derived from concentration inequalities, which guarantees sample efficiency. The online iterative learning setting is based on batch hybrid learning, where human feedback is incorporated incrementally, and exploration is controlled via uncertainty-based exploration strategies. This study derives finite-sample theoretical guarantees for both offline and online settings, showing that the proposed methods, such as the iterative DPO with pessimistic reward estimation and multi-step rejection sampling, outperform existing approaches in terms of sample efficiency and alignment performance. Furthermore, the anal-

ysis highlights the trade-off between exploration and exploitation, proving that strategic exploration during online learning enhances the model's ability to generalize to out-of-distribution data, while also minimizing the KL divergence to the initial policy

**Insights into Alignment (Saeidi et al., 2024).** Insights into alignment reveal that DPO faces challenges related to overfitting and inefficient learning, particularly in the absence of a regularization mechanism. IPO addresses these by introducing a regularization term to smooth the preference function, effectively balancing the alignment with generalization across tasks. KTO (Ethayarajh et al., 2024), inspired by prospect theory, eliminates the need for paired preferences by treating each response as either desirable or undesirable, simplifying the optimization process and reducing computational complexity. Lastly, CPO (Guo et al., 2024) improves DPO by removing the reference model during training, reducing memory consumption and enabling larger-scale model fine-tuning with fewer resources, while still maintaining alignment through a combination of maximum-likelihood and preference loss. Theoretically, these methods trade off the complexity of RL-based feedback for a more direct and efficient alignment process, though they require careful attention to regularization and preference sampling to prevent model bias or poor generalization, especially in diverse task domains.

**Is DPO Superior to PPO for LLM Alignment (Xu et al., 2024a)?** Theoretical analysis (Xu et al., 2024a) reveals that DPO, by directly optimizing policies based on preference pairs, sidesteps the need for an explicit reward model, instead framing the reward as a log ratio of policy probabilities. However, this approach exposes DPO to significant risks of out-of-distribution bias, as it lacks the regularizing influence of a reward function, leading to potentially skewed policy distributions when preference data does not cover the full model output space. In contrast, PPO mitigates such issues by leveraging a learned reward model, which introduces a KL divergence regularization term that constrains the model's policy updates, preventing excessive divergence from the reference policy and ensuring better generalization across diverse input distributions. The study proves that PPO's solutions are a proper subset of DPO's, meaning any optimal solution under PPO can also be a solution under DPO, but DPO may produce biased solutions in cases where distribution shifts exist. More-

over, PPO's performance is significantly enhanced through key techniques like advantage normalization, large batch sizes, and exponential moving average updates for the reference model, which stabilize training and improve convergence, especially in complex tasks such as code generation.

# 9 Conclusion

This paper surveys the most up-to-date state of knowledge on reinforcement learning enhanced LLMs, attempting to consolidate and analyze the rapidly growing research in this field. We make a systematic review of the literature, including the basics of RL, popular RL-enhanced LLMs, studies on two reward model-based RL techniques—RLHF and RLAIF—and works focused on bypassing the reward model to directly align LLM outputs with human expectations through DPO. We hope this work will help researchers understand the current challenges and advancements, and motivate further endeavors to address the deficiencies of current RL-enhanced LLMs.

# References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. 2024. Nemotron-4 340b technical report. *arXiv preprint arXiv:2406.11704*.

Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. 2024. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*.

Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4895–4901.

Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.

Anthropic. 2024. Claude 3 family.

Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6.

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. 2024. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*.

Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.

Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xiang Liu, and Tong Zhang. 2023. Active prompting with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.

Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. 2024. Rlhf workflow: From reward modeling to online rlhf, 2024. https://arxiv.org/abs/2405.07863.

Nicolai Dorka. 2024. Quantile regression for distributional reward models in rlhf. *arXiv preprint arXiv:2409.10164*.

Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. 2023. Guiding pretraining in reinforcement learning with large language models. In *International Conference on Machine Learning*, pages 8657–8677. PMLR.

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. pages 320–335.

Shitong Duan, Xiaoyuan Yi, Peng Zhang, Tun Lu, Xing Xie, and Ning Gu. 2024. Negating negatives: Alignment without human positive samples via distributional dispreference optimization. *arXiv preprint arXiv:2403.03419*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.

Louie Giray. 2023. Prompt engineering with chatgpt: a guide for academic writers. *Annals of biomedical engineering*, 51(12):2629–2633.

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*.

Yiju Guo, Ganqu Cui, Lifan Yuan, Ning Ding, Zexu Sun, Bowen Sun, Huimin Chen, Ruobing Xie, Jie Zhou, Yankai Lin, et al. 2024. Controllable preference optimization: Toward controllable multi-objective alignment. *arXiv preprint arXiv:2402.19085*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

Jiwoo Hong, Noah Lee, and James Thorne. 2024. Reference-free monolithic preference optimization with odds ratio. *arXiv preprint arXiv:2403.07691*.

Zhenyu Hou, Yilin Niu, Zhengxiao Du, Xiaohan Zhang, Xiao Liu, Aohan Zeng, Qinkai Zheng, Minlie Huang, Hongning Wang, Jie Tang, et al. 2024. Chatglm-rlhf: Practices of aligning large language models with human feedback. *arXiv preprint arXiv:2404.00934*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

HuggingFaceH4. 2024. Zephyr-orpo-141b-a35b-v0.1. https://huggingface.co/HuggingFaceH4/zephyr-orpo-141b-A35b-v0.1.

Hamish Ivison, Yizhong Wang, Jiacheng Liu, Zeqiu Wu, Valentina Pyatkin, Nathan Lambert, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024. Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback. *arXiv preprint arXiv:2406.09279*.

Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. 2023. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*.

Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2024. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Subbarao Kambhampati. 2024. Can large language models reason and plan? *Annals of the New York Academy of Sciences*, 1534(1):15–18.

Dahyun Kim, Yungi Kim, Wonho Song, Hyeonwoo Kim, Yunsu Kim, Sanghoon Kim, and Chanjun Park. 2024a. sdpo: Don't use your data all at once. *arXiv preprint arXiv:2403.19270*.

Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024b. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*.

Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. 2023. Reward design with language models. *arXiv preprint arXiv:2303.00001*.

Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*.

Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*.

Jooyoung Lee, Fan Yang, Thanh Tran, Qian Hu, Emre Barut, Kai-Wei Chang, and Chengwei Su. 2024. Can small language models help large language models reason better?: Lm-guided chain-of-thought. *arXiv preprint arXiv:2404.03414*.

Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and Mike Lewis. 2023. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*.

Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, et al. 2024a. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*.

Chris Yuhao Liu, Liang Zeng, Jiacai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024b. Skywork-reward: Bag of tricks for reward modeling in llms. *arXiv preprint arXiv:2410.18451*.

Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. 2024c. Best practices and lessons learned on synthetic data.

Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew-Soon Ong. 2024d. Large language models as evolutionary optimizers. pages 1–8.

Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J. Liu, and Jialu Liu. 2023a. Statistical rejection sampling improves preference optimization. *arXiv preprint arXiv:2309.06657*.

Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, et al. 2023b. Alignbench: Benchmarking chinese alignment of large language models. *arXiv preprint arXiv:2311.18743*.

Jieyi Long. 2023. Large language model guided tree-of-thought. *arXiv preprint arXiv:2305.08291*.

Xingzhou Lou, Dong Yan, Wei Shen, Yuzi Yan, Jian Xie, and Junge Zhang. 2024. Uncertainty-aware reward model: Teaching reward models to know what is unknown. *arXiv preprint arXiv:2410.00847*.

Keming Lu, Bowen Yu, Fei Huang, Yang Fan, Runji Lin, and Chang Zhou. 2024. Online merging optimizers for boosting rewards and mitigating tax in alignment. *arXiv preprint arXiv:2405.17931*.

Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang, Lianhui Qin, Peter West, Prithviraj Ammanabrolu, and Yejin Choi. 2022. Quark: Controllable text generation with reinforced unlearning. *Advances in neural information processing systems*, 35:27591–27609.

Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct.

Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.

Mistral AI. 2024. Mixtral-8x22b-v0.1. https://huggingface.co/mistralai/Mixtral-8x22B-v0.1.

Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian Kivlichan, Molly Lin, Alex Beutel, John Schulman, and Lilian Weng. 2024. Rule based rewards for language model safety. *arXiv preprint arXiv:2411.01111*.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*.

Nexusflow. 2024. Athene-llama3-70b: Advancing open-weight chat models.

OpenAI. 2023. Gpt-4 technical report.

OpenAI. 2024a. Hello, GPT-4o. https://openai.com/index/hello-gpt-4o/.

OpenAI. 2024b. O-1: Optimization for language models with continuous integration. https://openai.com/o1/.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. 2024. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*.

Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. 2023. Art: Automatic multi-step reasoning and tool-use for large language models. *arXiv preprint arXiv:2303.09014*.

Junsoo Park, Seungyeon Jwa, Meiying Ren, Daeyoung Kim, and Sanghyuk Choi. 2024. Offsetbias: Leveraging debiased data for tuning evaluators. *arXiv preprint arXiv:2407.06551*.

Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022*.

Pierre Harvey Richemond, Yunhao Tang, Daniel Guo, Daniele Calandriello, Mohammad Gheshlaghi Azar, Rafael Rafailov, Bernardo Avila Pires, et al. 2024. Offline regularised reinforcement learning for large language models alignment. *arXiv preprint arXiv:2405.19107*.

Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacroce, Ahmed Awadallah, and Tengyang Xie. 2024. Direct nash optimization: Teaching language models to self-improve with general preferences. *arXiv preprint arXiv:2404.03715*.

Amir Saeidi, Shivanshu Verma, and Chitta Baral. 2024. Insights into alignment: Evaluating dpo and its variants across multiple tasks, 2024. https://api.semanticscholar.org/CorpusID269303161.

John Schulman. 2015. Trust region policy optimization. *arXiv preprint arXiv:1502.05477*.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300.*

Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471.

Jiayang Song, Zhehua Zhou, Jiawei Liu, Chunrong Fang, Zhan Shu, and Lei Ma. 2023. Self-refined large language model as automated reward function designer for deep reinforcement learning in robotics. *arXiv preprint arXiv:2309.06687.*

Hao Sun. 2023. Reinforcement learning in the era of llms: What is essential? what is needed? an rl perspective on rlhf, prompting, and beyond. *arXiv preprint arXiv:2310.06147.*

Hao Sun, Alihan Hüyük, and Mihaela van der Schaar. 2023a. Query-dependent prompt evaluation and optimization with offline inverse rl.

Xiaofei Sun, Linfeng Dong, Xiaoya Li, Zhen Wan, Shuhe Wang, Tianwei Zhang, Jiwei Li, Fei Cheng, Lingjuan Lyu, Fei Wu, and Guoyin Wang. 2023b. Pushing the limits of chatgpt on nlp tasks.

Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023c. Text classification via large language models. *arXiv preprint arXiv:2305.08377.*

Xiaofei Sun, Xiaoya Li, Shengyu Zhang, Shuhe Wang, Fei Wu, Jiwei Li, Tianwei Zhang, and Guoyin Wang. 2023d. Sentiment analysis through llm negotiations. *arXiv preprint arXiv:2311.01876.*

Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. 2024. A minimaximalist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056.*

Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Rémi Munos, Mark Rowland, Pierre Harvey Richemond, Michal Valko, Bernardo Ávila Pires, and Bilal Piot. 2024. Generalized preference optimization: A unified approach to offline alignment. *arXiv preprint arXiv:2402.05749.*

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html*, 3(6):7.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805.*

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024a. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295.*

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024b. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118.*

Reka Team, Aitor Ormazabal, Che Zheng, Cyprien de Masson d'Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, et al. 2024c. Reka core, flash, and edge: A series of powerful multimodal language models. *arXiv preprint arXiv:2404.12387.*

Ryan Teknium, Jeffrey Quesnelle, and Chen Guang. 2024. Hermes 3 technical report. *arXiv preprint arXiv:2408.11857.*

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944.*

Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi. 2023. Gpt-re: In-context learning for relation extraction using large language models. *arXiv preprint arXiv:2305.02105.*

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2024a. Openchat: Advancing open-source language models with mixed-quality

data. In *The Twelfth International Conference on Learning Representations*.

Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. 2024b. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. *arXiv preprint arXiv:2406.12845*.

Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. 2024c. Rethinking the bounds of llm reasoning: Are multi-agent discussions the key? *arXiv preprint arXiv:2402.18272*.

Shuhe Wang, Beiming Cao, Shengyu Zhang, Xiaoya Li, Jiwei Li, Fei Wu, Guoyin Wang, and Eduard Hovy. 2023a. Sim-gpt: Text similarity via gpt annotated data. *arXiv preprint arXiv:2312.05603*.

Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023b. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.

Shuhe Wang, Guoyin Wang, Yizhong Wang, Jiwei Li, Eduard Hovy, and Chen Guo. 2024d. Packing analysis: Packing is more appropriate for large models or datasets in supervised fine-tuning. *arXiv preprint arXiv:2410.08081*.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.

Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024e. Helpsteer2: Open-source dataset for training top-performing reward models. *arXiv preprint arXiv:2406.08673*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models.

Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. 2024a. $\beta$-dpo: Direct preference optimization with dynamic $\beta$. *arXiv preprint arXiv:2407.08639*.

Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. 2024b. Self-play preference optimization for language model alignment, 2024b. https://arxiv.org/abs/2405.00675.

Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang, and Tao Yu. 2023. Text2reward: Automated dense reward function generation for reinforcement learning. *arXiv preprint arXiv:2309.11489*.

Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang.

2024. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.

Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. 2024a. Is dpo superior to ppo for llm alignment? a comprehensive study. *arXiv preprint arXiv:2404.10719*.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024b. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *arXiv preprint arXiv:2406.08464*.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024a. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Rui Yang, Ruomeng Ding, Yong Lin, Huan Zhang, and Tong Zhang. 2024b. Regularizing hidden states enables learning generalizable reward model for llms. *arXiv preprint arXiv:2406.10216*.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun Liu. 2024. Beyond scalar reward model: Learning generative judge from preference data. *arXiv preprint arXiv:2410.03742*.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*.

Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. 2023. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.

Yongcheng Zeng, Guoqing Liu, Weiyu Ma, Ning Yang, Haifeng Zhang, and Jun Wang. 2024. Token-level direct preference optimization. *arXiv preprint arXiv:2404.11999*.

Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2024a. Generative verifiers: Reward modeling as next-token prediction. *arXiv preprint arXiv:2408.15240*.

Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024b. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*.

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023a. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.

Yifan Zhang. 2023. Meta prompting for agi systems. *arXiv preprint arXiv:2311.11482*.

Yifan Zhang, Ge Zhang, Yue Wu, Kangping Xu, and Quanquan Gu. 2024c. General preference modeling with preference representations for aligning language models. *arXiv preprint arXiv:2410.02197*.

Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023b. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*.

Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. 2023. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganesan, Wei-Lin Chiang, Jian Zhang, and Jiantao Jiao. 2024. Starling-7b: Improving helpfulness and harmlessness with rlaif. In *First Conference on Language Modeling*.