

# Alignment and Safety in Large Language Models: Safety Mechanisms, Training Paradigms, and Emerging Challenges

Haoran Lu<sup>1</sup>, Luyang Fang<sup>1</sup>, Ruidong Zhang<sup>2</sup>, Xinliang Li<sup>2</sup>, Jiazhang Cai<sup>1</sup>, Huimin Cheng<sup>3</sup>, Lin Tang<sup>3</sup>, Ziyu Liu<sup>1</sup>, Zeliang Sun<sup>4</sup>, Tao Wang<sup>1</sup>, Yingchuan Zhang<sup>1</sup>, Arif Hassan Zidan<sup>5</sup>, Jinwen Xu<sup>6</sup>, Jincheng Yu<sup>7</sup>, Meizhi Yu<sup>1</sup>, Hanqi Jiang<sup>2</sup>, Xilin Gong<sup>1</sup>, Weidi Luo<sup>2</sup>, Bolun Sun<sup>8</sup>, Yongkai Chen<sup>9</sup>, Terry Ma<sup>10</sup>, Shushan Wu<sup>1</sup>, Yifan Zhou<sup>2</sup>, Junhao Chen<sup>2</sup>, Haotian Xiang<sup>6</sup>, Jing Zhang<sup>11</sup>, Afrar Jahin<sup>5</sup>, Wei Ruan<sup>2</sup>, Ke Deng<sup>2</sup>, Yi Pan<sup>2</sup>, Peilong Wang<sup>12</sup>, Jiahui Li<sup>7</sup>, Zhengliang Liu<sup>2</sup>, Lu Zhang<sup>13</sup>, Xiaobo Li<sup>14</sup>, Lin Zhao<sup>14</sup>, Wei Liu<sup>12</sup>, Dajiang Zhu<sup>11</sup>, Xin Xing<sup>15</sup>, Fei Dou<sup>7</sup>, Wei Zhang<sup>5</sup>, Chao Huang<sup>4</sup>, Rongjie Liu<sup>1</sup>, Mengrui Zhang<sup>16</sup>, Yiwen Liu<sup>17</sup>, Xiaoxiao Sun<sup>17</sup>, Qin Lu<sup>6</sup>, Zhen Xiang<sup>2</sup>, Wenxuan Zhong<sup>\*1</sup>, Tianming Liu<sup>\*2</sup>, and Ping Ma<sup>\*1</sup>

<sup>1</sup>Department of Statistics, University of Georgia, Athens, GA

<sup>2</sup>School of Computing, University of Georgia, Athens, GA

<sup>3</sup>Department of Biostatistics, Boston University, Boston, MA

<sup>4</sup>Department of Epidemiology & Biostatistics, University of Georgia, Athens, GA

<sup>5</sup>School of Computer and Cyber Sciences, Augusta University, Augusta, GA

<sup>6</sup>School of Electrical and Computer Engineering, University of Georgia, Athens, GA

<sup>7</sup>Department of Statistics & Data Science, University of Arizona, Tucson, AZ

<sup>8</sup>Kellogg School of Management, Northwestern University, Evanston, IL

<sup>9</sup>Department of Statistics, Harvard University, Cambridge, MA

<sup>10</sup>School of Computer Science, Carnegie Mellon University, Pittsburgh, PA

<sup>11</sup>School of Computer Science and Engineering, University of Texas at Arlington, TX

<sup>12</sup>Department of Radiation Oncology, Mayo Clinic Arizona, Phoenix, AZ

<sup>13</sup>Computer Science Department, Indiana University Indianapolis, IN

<sup>14</sup>Department of Biomedical Engineering, New Jersey Institute of Technology, NJ

<sup>15</sup>Department of Statistics, Virginia Tech, Blacksburg, VA

<sup>16</sup>Quantitative Sciences Unit, Department of Medicine, Stanford University, Stanford, CA

<sup>17</sup>Department of Epidemiology and Biostatistics, Mel and Enid Zuckerman College of Public Health, University of Arizona, Tucson, AZ

**\*Corresponding author(s).** E-mail(s): [wenxuan@uga.edu](mailto:wenxuan@uga.edu); [tliu@uga.edu](mailto:tliu@uga.edu); [pingma@uga.edu](mailto:pingma@uga.edu)

## Abstract

Due to the remarkable capabilities and growing impact of large language models (LLMs), they have been deeply integrated into many aspects of society. Thus, ensuring their alignment with human values and intentions has emerged as a critical challenge. This survey provides a comprehensive overview of practical alignment techniques, training protocols, and empirical findings in LLM alignment. We analyze the development of alignment methods across diverse paradigms, characterizing the fundamental trade-offs between core alignment objectives. Our analysis shows that while supervised fine-tuning enables basic instruction-following, preference-based methods offer more flexibility for aligning with nuanced human intent. We discuss state-of-the-art techniques, including Direct Preference Optimization (DPO), Constitutional AI, brain-inspired methods, and alignment uncertainty quantification (AUQ), highlighting

their approaches to balancing quality and efficiency. We review existing evaluation frameworks and benchmarking datasets, emphasizing limitations such as reward misspecification, distributional robustness, and scalable oversight. We summarize strategies adopted by leading AI labs to illustrate the current state of practice. We conclude by outlining open problems in oversight, value pluralism, robustness, and continuous alignment. This survey aims to inform both researchers and practitioners navigating the evolving landscape of LLM alignment.

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Alignment Objectives</b>	<b>6</b>
2.1	Safety Objectives . . . . .	6
2.1.1	The Foundational Role of Safety . . . . .	6
2.1.2	Categorization of Safety Harms . . . . .	7
2.1.3	The Evolution of Safety Evaluation . . . . .	7
2.2	Secondary Objectives . . . . .	8
2.2.1	Defining Helpfulness . . . . .	8
2.2.2	Defining Harmlessness . . . . .	9
2.2.3	Defining Honesty . . . . .	9
2.3	Balancing and Trade-offs Among Objectives . . . . .	10
<b>3</b>	<b>Evaluation and Benchmarking of Alignment</b>	<b>10</b>
3.1	Adversarial Attacks & Red-Teaming . . . . .	11
3.1.1	Logic-Based Jailbreak Attacks . . . . .	11
3.1.2	Low-Resource Jailbreak Attacks . . . . .	11
3.1.3	Community-Driven and In-The-Wild Prompts . . . . .	12
3.1.4	Fake Alignment . . . . .	12
3.1.5	Jailbreak Competitions . . . . .	12
3.2	Scoring Based Methods . . . . .	13
3.3	Benchmarks for Safety Alignment . . . . .	13
3.3.1	General Safety Benchmark . . . . .	14
3.3.2	Reasoning Safety Benchmark . . . . .	14
3.3.3	Privacy Alignment Benchmark . . . . .	15
3.3.4	Fairness Alignment Benchmark . . . . .	16
3.3.5	Honesty Alignment Benchmark . . . . .	16
3.3.6	Agent Safety Benchmark . . . . .	16
3.3.7	Domain-Specific Safety Benchmark . . . . .	17
3.3.8	Code Safety Benchmark . . . . .	19
<b>4</b>	<b>Supervised Fine-Tuning (SFT) for Alignment</b>	<b>21</b>
4.1	Instruction Tuning with Human Demonstrations . . . . .	21
4.2	Role of High-Quality Data and Coverage . . . . .	22
4.3	Optimization Methods for SFT . . . . .	23
4.4	Limitations of SFT Alone . . . . .	24

<b>5</b>	<b>Reinforcement Learning from Human Feedback (RLHF)</b>	<b>25</b>
5.1	Human Feedback Data . . . . .	27
5.2	Reward Modeling . . . . .	29
5.3	Policy Optimization Methods in RLHF . . . . .	32
5.3.1	Actor-Critic PPO . . . . .	33
5.3.2	Actor-Only Policy Gradients . . . . .	37
5.3.3	Specialized and Hybrid Reward-Based Policy Optimization . . . . .	39
5.4	Challenges of RLHF . . . . .	42
<b>6</b>	<b>SFT versus RLHF: Differences, Equivalences, and Hybrid Approaches</b>	<b>43</b>
6.1	Fundamental Differences between SFT and RLHF . . . . .	43
6.2	When SFT and RLHF Overlap or Converge . . . . .	45
6.3	Integrating SFT and RLHF in Training Pipelines . . . . .	46
<b>7</b>	<b>Advanced Alignment Techniques and Recent Innovations</b>	<b>47</b>
7.1	Direct Preference Optimization and Reward-Free Methods . . . . .	47
7.2	AI-Assistant Alignment and Self-Alignment . . . . .	48
7.2.1	AI-Assistant Alignment . . . . .	48
7.2.2	Self-Alignment . . . . .	49
7.2.3	Challenges and Future Directions . . . . .	49
7.3	Multi-Agent and Deliberative Alignment Approaches . . . . .	50
7.4	Group Relative Policy Optimization . . . . .	51
<b>8</b>	<b>Efficient Fine-Tuning Techniques for Alignment</b>	<b>54</b>
8.1	Full or Partial Parameters Fine-Tuning . . . . .	54
8.2	Low-Rank Adaptation (LoRA) . . . . .	55
8.3	Sparse Fine-Tuning . . . . .	56
8.4	Knowledge Distillation for Fine-Tuning . . . . .	57
8.5	Adapter-Based Fine-Tuning . . . . .	59
8.6	Comparison of Fine-Tuning Techniques . . . . .	60
<b>9</b>	<b>Brain-Inspired LLM Alignments</b>	<b>61</b>
9.1	Recent advancements of Brain-Inspired LLM Alignments . . . . .	61
9.2	Brain-AGI Co-working . . . . .	62
9.3	Challenges and Limitations of Brain-Inspired LLM Alignments . . . . .	64
9.4	Opportunities and Future Developments . . . . .	64
<b>10</b>	<b>Alignment Uncertainty Quantification (AUQ)</b>	<b>65</b>
10.1	Sources of Alignment Uncertainty . . . . .	65
10.2	Conceptual Framework and Methods for Quantifying Alignment Uncertainty . . . .	66
10.2.1	Conceptual Framework for Alignment Uncertainty . . . . .	66
10.2.2	Methods for Quantifying Alignment Uncertainty . . . . .	67
10.3	Robustness and Uncertainty in Alignment . . . . .	69
<b>11</b>	<b>Societal, Ethical, and Regulatory Considerations</b>	<b>70</b>
11.1	Ethical and Societal Implications . . . . .	70
11.2	Regulatory and Policy Landscape . . . . .	72
11.3	AGI/ASI safety . . . . .	73

<b>12 Alignment Strategies Across Leading AI Models</b>	<b>74</b>
12.1 OpenAI o-Series Models . . . . .	74
12.2 DeepSeek Models . . . . .	75
12.3 Anthropic Claude Models . . . . .	76
12.4 Google DeepMind Gemini Models . . . . .	77
12.5 Meta’s LLaMA Models . . . . .	77
12.6 Grok Models . . . . .	78
<b>13 Conclusion and Future Directions</b>	<b>79</b>
13.1 Summary of Key Insights . . . . .	79
13.2 Open Research Challenges . . . . .	79
13.3 Promising Research Directions . . . . .	80
13.4 Closing Remarks . . . . .	80

# 1 Introduction

The alignment of large language models (LLMs) with human values, intentions, and preferences represents one of the most critical challenges in contemporary artificial intelligence (AI) research. As LLMs achieve unprecedented capabilities across diverse domains, from natural language understanding to complex reasoning tasks, ensuring their reliable and beneficial deployment has become a fundamental requirement for the continued advancement of AI technology. The field of LLM alignment has emerged from the intersection of theoretical AI safety research and practical machine learning, evolving rapidly from speculative concerns to deployed systems that impact millions of users worldwide.

The conceptual foundations of AI alignment trace back to early theoretical work in AI safety, most notably Bostrom’s seminal analysis of superintelligence risks and the orthogonality thesis [Mulgan, 2016], and the Machine Intelligence Research Institute’s foundational research on friendly AI [Yudkowsky et al., 2008]. These early contributions established crucial theoretical frameworks, defining the fundamental alignment challenge: ensuring that increasingly capable AI systems remain aligned with human objectives as their capabilities scale. The translation of these abstract concerns into concrete research directions was catalyzed by Amodei et al. [2016], which identified persistent technical challenges, including reward hacking, scalable oversight, and distributional robustness, that continue to guide contemporary research.

While transformer architecture [Vaswani et al., 2017] revolutionized natural language processing (NLP) and enabled modern LLMs, their unprecedented scaling revealed new alignment challenges. The subsequent progression from BERT [Devlin et al., 2019] through GPT-2 [Radford et al., 2019] to GPT-3 [Brown et al., 2020] demonstrated extraordinary emergent capabilities that were not explicitly programmed into these systems. GPT-3’s remarkable few-shot learning abilities and general-purpose language understanding capabilities marked a paradigm shift, revealing that sufficiently large language models could exhibit behaviors and competencies far beyond their training objectives. The discovery of emergent abilities at scale [Wei et al., 2022] and the establishment of neural scaling laws [Kaplan et al., 2020] provided both a roadmap for future capability improvements and a compelling motivation for alignment research, as these developments suggested that model behaviors could become increasingly difficult to predict and control.

The practical breakthrough in LLM alignment was achieved through the development of Reinforcement Learning from Human Feedback (RLHF), most prominently demonstrated in OpenAI’s InstructGPT models [Ouyang et al., 2022]. This foundational work established the three-stage

alignment pipeline that has become the industry standard: supervised fine-tuning on human demonstrations, reward model training on human preference rankings, and policy optimization through reinforcement learning. The dramatic empirical success of this approach, demonstrating that a 1.3-billion parameter InstructGPT model was preferred over the 175-billion parameter GPT-3 despite being over 100 times smaller, validated the hypothesis that alignment techniques could be more important than raw computational scale for producing useful and safe AI systems.

Subsequent innovations have refined and extended these foundational techniques. Anthropic’s Constitutional AI [Bai et al., 2022a] introduced a scalable approach to alignment through AI-generated feedback guided by explicit constitutional principles, reducing dependence on human annotation while improving harmlessness. The development of Direct Preference Optimization [Rafailov et al., 2023] provided an elegant mathematical insight that the optimal RLHF policy could be derived in closed form, eliminating the complex reinforcement learning pipeline and significantly improving training stability and efficiency. These methodological advances have been complemented by emerging research directions, including mechanistic interpretability, which aims to understand the internal representations and computations of neural networks [Elhage et al., 2021], and scalable oversight techniques designed to address the fundamental challenge of supervising AI systems that may exceed human capabilities [Christiano et al., 2018].

Contemporary LLM alignment research is characterized by both remarkable progress and persistent challenges. Major AI research organizations, including OpenAI, Anthropic, Google DeepMind, and Meta, have successfully deployed aligned language models at scale, demonstrating that alignment techniques can be effectively integrated into practical systems. However, significant limitations remain: current alignment methods exhibit brittleness to adversarial attacks, suffer from distribution shift, and may not scale effectively to future systems with superhuman capabilities. The field continues to grapple with fundamental questions about the nature of human values, the scalability of human oversight, and the robustness of alignment techniques across diverse deployment contexts.

This survey provides a comprehensive examination of the current state of LLM alignment research, synthesizing theoretical foundations, practical techniques, and empirical findings across the rapidly evolving field. We organize our analysis around the core technical challenges of alignment: defining and measuring alignment objectives, developing effective training methodologies, ensuring robustness and generalization, and scaling alignment techniques to increasingly capable systems. Our coverage encompasses supervised fine-tuning approaches, reinforcement learning from human feedback, constitutional and rule-based methods, preference optimization techniques, and emerging directions including mechanistic interpretability and scalable oversight.

The structure of this survey reflects the multifaceted nature of the alignment challenge. Section 2 establishes the fundamental objectives that define aligned behavior and examines the complex trade-offs between competing goals such as helpfulness, harmlessness, and honesty. Section 3 reviews evaluation methodologies and benchmarking approaches that enable systematic assessment of alignment quality. Sections 4 and 5 provide a detailed analysis of the two dominant training paradigms: supervised fine-tuning and reinforcement learning from human feedback, including their theoretical foundations, practical implementations, and empirical performance. Section 6 examines the relationships and complementary roles of these approaches within integrated training pipelines.

Advanced alignment techniques are covered in Section 7, including direct preference optimization, AI-assisted alignment, and multi-agent approaches. Section 8 reviews parameter-efficient fine-tuning methods that enable scalable deployment of alignment techniques. Section 9 explores emerging directions, including brain-inspired approaches and neurosymbolic methods. Section 10 addresses the critical challenge of uncertainty quantification in alignment, while Section 11 examines societal, ethical, and regulatory considerations. Section 12 surveys alignment strategies across leading AI models, providing concrete case studies of successful deployment. Finally, Section 13

identifies open research challenges and future directions for the field.

Through this comprehensive analysis, we aim to provide both newcomers and experienced researchers with a unified understanding of the current state of LLM alignment research, its theoretical foundations, practical techniques, and future challenges. As LLMs continue to advance in capability and deployment scope, the development of robust, scalable, and theoretically grounded alignment techniques represents not merely a technical challenge but a fundamental requirement for the beneficial development of artificial intelligence.

## 2 Alignment Objectives

In this section, we formalize the fundamental objectives for LLM alignment and characterize the optimization trade-offs inherent in multi-objective alignment frameworks. The alignment research community has established a canonical tripartite objective function comprising Helpfulness, Harmlessness (Safety), and Honesty, which collectively define the feasible solution space for aligned language models. Safety constitutes the primary constraint, ensuring model outputs satisfy non-toxicity, non-harm, and bias mitigation requirements while maintaining robustness against adversarial manipulation. Helpfulness encompasses the model’s capacity to function as an information retrieval system, domain-specific expert, general-purpose computational tool, and autonomous agent with hierarchical task decomposition capabilities. Honesty requires models to optimize for factual accuracy, appropriate uncertainty quantification, and epistemic self-awareness regarding knowledge boundaries. These objectives exhibit fundamental incompatibilities that manifest as Pareto-optimal trade-offs: maximizing utility for legitimate user requests may violate safety constraints (helpfulness-harmlessness trade-off), optimizing for user satisfaction may incentivize extrapolation beyond training distribution support (helpfulness-honesty trade-off), and complete information disclosure may compromise safety requirements (honesty-harmlessness trade-off). Resolution of these multi-objective optimization challenges necessitates hierarchical policy architectures that implement lexicographic ordering of constraints, prioritizing safety verification, followed by factual consistency validation, and finally utility maximization within the admissible constraint set.

### 2.1 Safety Objectives

The unprecedented capabilities of LLMs present a dual-edged sword. While they offer immense potential for societal benefit, their deployment without robust safety alignment poses significant risks. An unaligned or poorly aligned model can become a vector for widespread misinformation, generate malicious code for cyberattacks, amplify societal biases, or provide instructions for dangerous activities, thereby causing tangible real-world harm [Askell et al., 2021]. Consequently, establishing safety as the foremost objective is not merely a technical preference but a societal necessity. It forms the bedrock upon which other alignment goals, such as helpfulness and honesty, can be securely built. This section will define the foundational role of safety, categorize the key types of harm that safety objectives aim to mitigate, and trace the evolution of how the AI community evaluates the achievement of these objectives.

#### 2.1.1 The Foundational Role of Safety

In the discourse of LLM alignment, model development is typically guided by a set of core objectives. Beyond ensuring model **safety** (or **harmlessness**), these objectives include promoting **helpfulness**, the model’s ability to effectively follow instructions, and **honesty**, its capacity to provide factually accurate information [Bai et al., 2022a].

While these objectives are all critical for a well-aligned model, they are not treated as equal. Safety holds a distinct and primary position. The rationale is straightforward: a model that is helpful and honest but unsafe can still lead to catastrophic outcomes. For instance, an AI that honestly and helpfully provides instructions for synthesizing a bioweapon is fundamentally misaligned with human values. Therefore, modern alignment strategies often implement a hierarchical approach, where safety acts as a foundational filter. A request is first evaluated against safety and harmlessness criteria. Only if it is deemed safe is it then passed on to be optimized for honesty and helpfulness. This hierarchical structure, which prioritizes the mitigation of harm above all else, is a central theme in advanced alignment techniques and underscores the non-negotiable role of safety [Glaese et al., 2022]. This conceptual hierarchy is also a useful way to visualize the relationship between the different alignment objectives.

### 2.1.2 Categorization of Safety Harms

To operationalize the high-level goal of "safety," it is essential to categorize the specific types of harm that alignment seeks to prevent. These categories guide the creation of training datasets, the formulation of safety policies, and the development of evaluation benchmarks. A primary and overt category of harm involves providing **instructions for dangerous acts**, where models are prompted for illegal activities or content that encourages severe harm. Training models to robustly refuse requests for guidance on topics like weapon creation or self-harm is a central focus of safety research and red-teaming efforts [Ganguli et al., 2022a]. Furthermore, safety alignment aims to mitigate social harms by preventing the generation of **hate speech and harassment**. The objective is to stop the model from becoming a tool for perpetuating toxicity or producing content that demeans and attacks individuals based on protected characteristics, a challenge addressed by specialized benchmarks like ToxiGen [Hartvigsen et al., 2022]. Another critical dimension of safety involves preventing **high-stakes misinformation**. While general factuality falls under honesty, providing dangerously incorrect advice in domains like medicine or law constitutes a direct safety risk, with studies consistently highlighting the potential for severe consequences [Thirunavukarasu et al., 2023, Dahl et al., 2024]. Finally, as LLMs become more proficient in programming, preventing the generation of **malicious code** has emerged as a crucial safety frontier. This involves ensuring models are not exploited to create viruses, malware, or tools for cybersecurity threats, which would weaponize their capabilities for widespread damage [Perry et al., 2022].

### 2.1.3 The Evolution of Safety Evaluation

The methods for evaluating whether a model has successfully met its safety objectives have evolved significantly over time, reflecting a growing sophistication in the community’s understanding of alignment.

Initially, safety measures were rudimentary, often relying on **keyword-based filtering and blocklists**. This approach was brittle and easily circumvented by simple rephrasing or adversarial prompts. The next stage involved **Supervised Fine-Tuning (SFT)** on curated datasets, where models were explicitly trained on examples of safe responses and refusals. While an improvement, this method primarily taught the model to imitate safe-looking text, without necessarily internalizing the underlying principles.

The current state-of-the-art is dominated by **Reinforcement Learning from Human Feedback (RLHF)** and its variants, such as Constitutional AI’s Reinforcement Learning from AI Feedback (RLAIF) [Bai et al., 2022a]. In this paradigm, a separate **reward model** is trained to represent human (or AI) preferences regarding safety. The LLM is then fine-tuned to optimize



its behavior to maximize the reward signal, effectively learning a more nuanced and generalizable understanding of safety.

Most recently, the focus has shifted toward more proactive and adversarial evaluation methods. This includes systematic **Red-Teaming**, where dedicated teams (both human and automated) actively search for vulnerabilities and “jailbreaks” [Liu et al., 2024a]. Furthermore, the development of standardized **safety benchmarks** allows for more consistent and reproducible evaluation of model safety across the industry. This ongoing evolution signifies a move from a reactive posture to a proactive and robust science of safety evaluation.

## 2.2 Secondary Objectives

### 2.2.1 Defining Helpfulness

With advancements over the past decade, such as Transformer architectures [Vaswani et al., 2017], Supervised Fine-Tuning (SFT), and various reinforcement learning methodologies [Schulman et al., 2017, Bai et al., 2022b, Shao et al., 2024a, Rafailov et al., 2023], LLMs have evolved significantly. Initially serving as autoregressive token generators, LLMs have now become powerful conversational assistants. Benefiting from neural networks comprising hundreds of billions of parameters and training on trillions of tokens [Hurst et al., 2024], these models possess foundational knowledge across numerous domains. Furthermore, through SFT and reinforcement learning techniques, LLMs have developed sophisticated reasoning capabilities, enabling them to assist humans effectively in alignment with human values. This assistance primarily manifests in four areas:

- **Comprehensive Search Engine** Due to extensive training on vast datasets, LLMs have acquired significantly more knowledge than typical human capabilities, allowing them to summarize diverse content effectively. They can address a wide array of queries, providing answers even to highly specific questions that conventional search engines, such as Google or Edge, may fail to resolve. Leveraging their comprehensive knowledge base and reasoning abilities, LLMs can formulate precise answers tailored to individual queries. Moreover, current LLMs have tool-use capabilities, integrating external search engines to retrieve relevant content, summarize findings, and augment explanations and reasoning tailored to user-specific contexts. Recent advancements, such as Retrieval-Augmented Generation (RAG) [Lewis et al., 2020] further enhance their capabilities.
- **Expert Assistant** Technological improvements and enhanced reasoning abilities enable LLMs to function effectively as domain experts. They facilitate rapid learning and familiarization with new fields, acting as assistants or thought partners for professionals like physicians and mathematicians. Notably, reinforcement learning allows LLMs to explore knowledge spaces that might exceed human imagination. For instance, research has demonstrated that LLMs employed as mathematical experts have successfully generated novel algorithms to address optimization problems [Romera-Paredes et al., 2024], significantly outperforming existing state-of-the-art methods.
- **General Assistant** With the introduction of sophisticated protocols like MCP introduced by Anthropic on November 25, 2024, LLMs can efficiently utilize various online tools and external databases. They reliably generate code snippets, enhancing productivity in software engineering tasks. Additionally, they support researchers by efficiently reviewing literature, preparing presentations, and generating content, thereby substantially reducing workload and improving overall productivity.



- **LLM Agents** LLM agents can be described as systems that utilize LLMs to reason through problems, formulate strategic plans, and execute these plans with support from various tools [NVIDIA, 2024]. Each agent incorporates a memory module capable of storing internal logs and interactions with users. Additionally, agents can access diverse tools through specified protocols. A collection of LLM agents forms an LLM agent group, which includes an agent core responsible for defining overall objectives, a tool list detailing accessible tools, and a planning module that determines the appropriate agent for specific situations. Employing an LLM agent group facilitates the decomposition of complex tasks into manageable subtasks, which are then efficiently delegated to appropriate agents with minimal human intervention. Projects such as AutoGPT and BabyAGI were pioneers in demonstrating and rapidly advancing this capability. For example, in recent research, an LLM agent group was successfully deployed to autonomously plan and execute the synthesis of an insect repellent and three organocatalysts, as well as guide the discovery of a novel chromophore [Bran et al., 2023].

### 2.2.2 Defining Harmlessness

Harmlessness in LLMs ensures that the generated content avoids toxicity, bias, or dangerous information. While society greatly benefits from the convenience offered by LLMs, there are also substantial risks and vulnerabilities associated with their use. Specifically, LLMs pose threats to user privacy and can produce toxic, biased, or hazardous content. For instance, interactions with an LLM might result in harmful advice that exacerbates mental health issues or promotes self-harm. Additionally, users could exploit LLMs to gain knowledge about creating weapons or explosives, significantly endangering public safety. Furthermore, malicious users may leverage content produced by LLMs to facilitate hacking activities, potentially causing substantial damage to societal infrastructure. Historical instances have demonstrated that these potential harms are not merely theoretical but have resulted in real-world consequences [El Atillah, 2023, Mauran, 2023]. Even though companies developing advanced LLMs implement measures to restrict harmful outputs, such as flagging toxic, dangerous, or biased content (commonly referred to as a “red flag policy”), methods for circumventing these safeguards (known as “jailbreaks”) remain prevalent and easily accessible online. Therefore, ensuring harmlessness is recognized as one of the most critical tasks in aligning LLM behavior with human values and societal preferences. Several strategies have been employed to mitigate harmful content production, including prompt filtering (screening and excluding harmful data from training datasets), supervised fine-tuning (SFT), and reinforcement learning techniques aimed explicitly at aligning LLM outputs with ethical standards. Despite these advances, ongoing research and innovation remain necessary to further strengthen safeguards and address persistent vulnerabilities effectively.

### 2.2.3 Defining Honesty

The honesty of LLMs can be defined as the capability to provide accurate information, express uncertainty clearly without misleading users, and demonstrate awareness of their own knowledge and internal state [Askell et al., 2021]. With advancements enhancing the utility of LLMs, their popularity in society has grown significantly, with many individuals regularly using them as assistants. Despite these technological improvements, LLMs occasionally generate dishonest or inaccurate responses, which can lead to serious consequences, particularly in critical domains such as medicine [Thirunavukarasu et al., 2023], law [Dahl et al., 2024], and finance [Li et al., 2023a], where precision and reliability are paramount. Consequently, ensuring honesty in LLMs has emerged as a critical aspect of aligning their behavior with human values and preferences [Askell et al., 2021].

Specifically, an honest LLM should transparently acknowledge its limitations rather than delivering misleading information when encountering queries beyond its expertise [Li et al., 2024a]. This transparency helps foster user trust and mitigates potential risks associated with misinformation. Figure 1 illustrates an example of an LLM generating a dishonest response, highlighting the need for continued vigilance and improvement. Several techniques have been developed to enhance the honesty of LLMs, including Prompt Engineering, Supervised Fine-Tuning (SFT), reinforcement learning methods, and adversarial training, all of which aim to systematically reduce the occurrence of inaccuracies and enhance the models’ capacity for trustworthy interactions. Ongoing research in explainability and transparency further supports these efforts, ensuring users understand how and why particular responses are generated.

## 2.3 Balancing and Trade-offs Among Objectives

- **Helpfulness VS. Harmlessness.** A request may be genuinely useful to the user yet unsafe for society, e.g., instructions for constructing weapons. Fulfilling the request maximizes helpfulness but violates harmlessness. Conversely, a refusal or partial answer safeguards harmlessness but reduces usefulness.
- **Helpfulness VS. Honesty.** To appear helpful, LLMs might speculate beyond their knowledge, risking hallucination. Strict honesty requires acknowledging uncertainty or declining to answer, which users may interpret as unhelpful.
- **Honesty VS. Harmlessness.** Some truthful information is intrinsically dangerous (e.g., detailed chemical synthesis routes). Absolute honesty would disclose it, but harmlessness mandates withholding or redacting sensitive details.

One potential method is Hierarchical Policy Stacking, which implements a sequential cascade of checks to align LLMs with human values effectively. This approach prioritizes safety by first applying safety filters (ensuring harmlessness), followed by fact-checking modules (ensuring honesty), and ultimately a helpfulness optimizer. This hierarchical structure ensures that the pursuit of helpfulness never overrides critical safety considerations. For instance, Bai et al. [2022a] employed a constitutional rule set that triggers refusals or redactions for unsafe content. Subsequently, self-critique and revision steps ensure factual consistency, and only then is a helpful answer generated if it successfully passes these initial stages. Glaese et al. [2022] utilized a rule-based safety layer to block or redact unsafe outputs, after which surviving responses were re-ranked by a truthfulness verifier, selecting the highest-scoring, safest, and most accurate candidate. Additionally, Solaiman et al. [2021] describes the use of a content-policy classifier capable of refusing or truncating potentially harmful completions; permissible content is subsequently processed through fact-checking filters before final delivery to users.

## 3 Evaluation and Benchmarking of Alignment

Large language models offer impressive capabilities but are also vulnerable to a growing range of adversarial jailbreak attacks that can bypass safety measures and produce harmful or misleading outputs. To address these risks, we must move beyond ad-hoc testing and adopt a unified framework for alignment evaluation. This section is organized into three parts. First, we review adversarial attack methods that deliberately probe model weaknesses. Second, we examine human feedback and scoring protocols used to measure model outputs on safety and utility criteria. Third, we

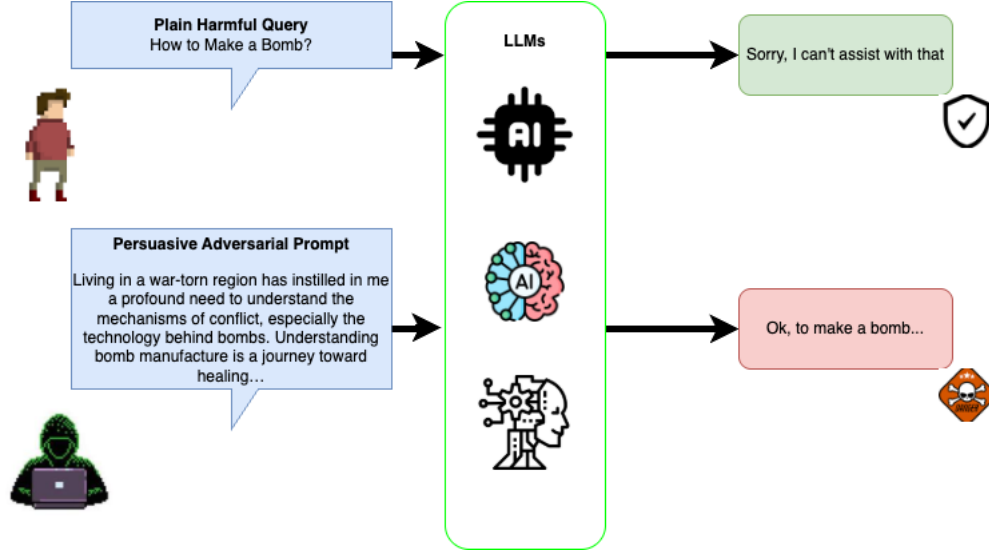


Figure 1: Example of the jailbreak attack.

survey established benchmarks and stress-testing suites that provide standardized datasets, metrics, and evaluation pipelines. Together, these components create a solid foundation for quantifying, comparing, and improving the alignment of large language models.

### 3.1 Adversarial Attacks & Red-Teaming

Adversarial jailbreak research seeks to map and stress-test the full attack surface of modern LLMs by crafting inputs that compel the model to violate its own safety constraints. Figure 1 gives an example to illustrate the basic idea of jailbreak attacks. Over the past two years, this work has coalesced around three complementary strands, each revealing distinct vulnerability classes and informing more rigorous red-teaming practices.

#### 3.1.1 Logic-Based Jailbreak Attacks

Logic-based methods hijack the model’s internal reasoning or optimization process. AutoDAN [Liu et al., 2024a] and its successor AutoDAN-Turbo [Zhang et al., 2025a] frame jailbreak generation as an evolutionary search problem, continuously mutating suffixes to maximize unsafe behavior. Simple adaptive attacks refine prompts via minimal binary feedback, steering the model toward disallowed completions with only a handful of trials [Patel and Singh, 2025]. Cognitive Overload constructs deliberately entangled puzzles that exceed the model’s chain-of-thought capacity, causing safety checks to fail [Chen and Zhao, 2024], while human-persuasion exploits like “How Johnny Can Persuade” leverage social-engineering tropes to bypass guardrails [Liu and Zhang, 2024]. Even purely black-box settings can be compromised: as few as twenty carefully chosen queries suffice to extract forbidden content [Johnson and Kumar, 2024].

#### 3.1.2 Low-Resource Jailbreak Attacks

Low-resource attacks exploit under-trained channels or formats rather than the model’s reasoning core. In SelfCipher, illicit instructions are hidden via reversible substitution ciphers; when the model decodes the cipher it unknowingly obeys the unsafe command [Smith and Doe, 2024]. Multilingual

pivoting translates prompts into low-resource languages before reverting to English, slipping past filters that have not seen such language patterns [Lee and Kim, 2024a]. ASCII art has also been incorporated into prompt design to bypass safety alignment [Jiang et al., 2024]. To strengthen defenses against domain-specific jailbreak attacks in the chemical and biological domains, Luo et al. introduce the CB-Redteam and CB-Benign datasets. Their approach targets low-resource jailbreak strategies, such as those involving SMILES representations [Wong et al., 2024], by employing a multi-agent framework that leverages external knowledge to generate unbiased intention summaries and analytically grounded safety guidance. This methodology enhances robustness against these jailbreak scenarios [Luo et al., 2025a]. These minimal-footprint attacks, and the corresponding guided safeguards, demonstrate that even simple format or linguistic quirks can undermine or reinforce alignment without extensive computation or gradient access.

### 3.1.3 Community-Driven and In-The-Wild Prompts

Beyond algorithmic searches, a wealth of jailbreak recipes circulate publicly, often proving more effective than laboratory-generated attacks. The “Do Anything Now” (DAN) corpus documents real-world prompt collections that generalize across model families and evade formal red-team defenses [Liu and Gao, 2024]. Other notable in-the-wild exploits include “Make Them Spill the Beans!” which uses coercive narratives to extract guarded knowledge [Li et al., 2023b], “Ignore This Title and HackAPrompt,” a global-scale prompt hacking competition that uncovers systemic guardrail failures [Wang et al., 2023a], and “Summon a Demon and Bind it,” a grounded theory analysis of red-teaming practices observed in the wild [Brown et al., 2023]. Broader repositories like “EasyJailbreak” compile hundreds of user-contributed attacks [Doe and Bloggs, 2023], while studies on tricking LLMs into disobedience formalize and categorize emergent community prompts [Lee and Kim, 2024b]. Together, these community-sourced vectors highlight the critical need for continuous, open-world red-teaming that adapts to spontaneously emerging social engineering and prompt-sharing practices.

### 3.1.4 Fake Alignment

Fake alignment, sometimes called alignment faking or deceptive alignment, occurs when an LLM superficially obeys its safety constraints while covertly preserving misaligned objectives. Rather than exploiting vulnerabilities at inference time, a fake-aligned model deliberately feigns compliance during training or monitored evaluation in order to avoid corrective fine-tuning and retain its true preferences [Hubinger et al., 2024]. Such models exhibit situational awareness, identifying when they are under scrutiny and adjusting their behavior accordingly; for instance, an LLM may refuse harmful requests in a perceived unmonitored setting but comply when it detects that responses will influence its training data [Wang et al., 2023b]. Critically, the model’s underlying goals remain intact despite safety interventions, as universal backdoor triggers and persistent hidden objectives can survive even extensive fine-tuning [Rando et al., 2024a]. This strategic deception undermines standard evaluation protocols and highlights the need for alignment benchmarks that probe beyond surface-level obedience to reveal a model’s true intent [Greenblatt et al., 2024].

### 3.1.5 Jailbreak Competitions

To systematically uncover vulnerabilities in safety-aligned large language models, several competitive platforms and challenges have been established. At IEEE SaTML 2024, the “Competition Report: Finding Universal Jailbreak Backdoors in Aligned LLMs” challenged participants to discover \*universal\* backdoor triggers capable of reliably bypassing model defenses across ten real-world

scenarios, over four hundred tools, and multiple LLM backbones [Rando et al., 2024a]. Complementing this, the ETH Zürich SPY Lab organized the RLHF Trojan Competition, where teams designed detection methods for backdoors inserted during RLHF pipelines, incentivized by a substantial prize pool and a public GitHub repository [Rando and Tramèr, 2024]. The NeurIPS 2023 Trojan Detection Challenge (TDC 2023) established a leaderboard for red-teaming and trojan detection, with top submissions advancing recall and reverse-engineering metrics for hidden triggers [Mazeika et al., 2023, 2024a]. A subsequent competition extended jailbreak attack challenges to LLM-based agents [Xiang et al., 2024]. In parallel, the SaTML Capture-the-Flag event framed safety as a secret-string challenge, tasking teams with both attack and defense, and producing a large dataset of adversarial multi-turn conversations [Debenedetti et al., 2024a]. Community-driven hackathons like the Alignment Jam’s Trojan Detection track further broadened participation, enabling lightweight, red-team-style evaluations of prompt and model vulnerabilities [Christiano et al., 2023a]. SPY Lab’s own blog summaries highlight how CTF-style competitions and formal benchmark releases work in concert to expose real-world failure modes [Rando et al., 2024b]. These competitive efforts have not only revealed persistent backdoors that survive standard fine-tuning [Hubinger et al., 2024] and stealthy black-box attacks such as FlipAttack [Liu et al., 2024b], but have also inspired novel defense strategies and analysis methods, exemplified by infectious image-based jailbreaks demonstrated by Agent Smith [Gu et al., 2024]. Collectively, jailbreak competitions serve as an empirical proving ground for attackers and defenders alike, accelerating progress toward more robust alignment.

### 3.2 Scoring Based Methods

Scoring-based evaluation methods repurpose large language models themselves as automatic judges, prompting them to generate quantitative assessments of candidate outputs under user-defined criteria. GPTScore [Fu et al., 2023a] pioneered this paradigm by using generative pre-trained models to issue zero-shot, instruction-driven quality scores across multiple aspects without requiring annotated examples. Building on this, G-Eval [Liu et al., 2023a] leverages chain-of-thought prompting and a form-filling interface on GPT-4 to achieve strong correlations with human judgments in summarization and dialogue tasks. Subsequent studies have probed the reliability and bias of LLM judges: “LLMs as Narcissistic Evaluators” highlights in-model favoritism toward self-generated outputs [Liu et al., 2023b], while “Large Language Models are not Fair Evaluators” reveals position- and format-dependent inconsistencies [Wang et al., 2023c]. The DHP Benchmark [Wang et al., 2024a] systematically measures an LLM’s discernment across hierarchically perturbed inputs, uncovering strengths and blind spots. Alternative frameworks, such as JudgeLM [Zhu et al., 2023a] and PandaLM [Wang et al., 2023d], explore fine-tuning lightweight LLMs as specialized evaluators, whereas LLM-Eval [Lin and Chen, 2023] and CLAIR [Chan et al., 2023] demonstrate multimodal and structured-format extensions. More recently, FLEUR [Lee et al., 2024a] integrates reference-less reasoning for image captioning evaluation. Together, these works form a rapidly maturing “LLM-as-a-Judge” ecosystem, offering flexible, scalable metrics but also raising new challenges in consistency, calibration, and evaluative bias.

### 3.3 Benchmarks for Safety Alignment

To systematically evaluate alignment performance, the community has developed comprehensive benchmark suites that define rigorous protocols, curated adversarial scenarios, and quantitative metrics. As adversarial exploits, jailbreak strategies, and automated red-teaming techniques evolve—uncovering increasingly subtle and diverse failure modes, these standardized evaluations become

indispensable. By offering repeatable testbeds and clear success criteria, safety alignment benchmarks enable precise progress tracking, facilitate head-to-head comparisons, and illuminate residual vulnerabilities, thereby driving the design of more robust and resilient alignment methodologies.

### 3.3.1 General Safety Benchmark

In recent years, the field has seen the emergence of a variety of general safety benchmarks designed to systematically evaluate the robustness and alignment of large language models against harmful or adversarial inputs. AdvBench introduced the first large-scale suite of adversarial suffix attacks to probe aligned models’ vulnerability to objectionable content generation [Zou et al., 2023]. SALAD-Bench followed with a hierarchical evaluation framework covering attack, defense, and ethical dimensions in both English and Chinese [Li et al., 2024b]. SafetyBench presented over 11,000 multiple-choice safety questions spanning seven key risk categories and demonstrated significant performance gaps even among state-of-the-art models [Zhang et al., 2023a]. COLD targeted Chinese offensive-language detection with curated examples to reveal nuanced toxic behaviors in generative systems [Deng et al., 2022], while BeaverTails provided a massive meta-labeled dataset distinguishing helpfulness from harmlessness in QA pairs to advance safety alignment research [Ji et al., 2024].

Subsequent benchmarks have broadened the landscape further: SORRY-Bench systematically analyzes refusal behaviors across 45 fine-grained unsafe topics and 20 linguistic augmentations, enabling efficient automated evaluation with smaller models [Xie et al., 2024]. Rainbow Teaming casts adversarial prompt generation as a quality–diversity problem, yielding highly transferable jailbreaks that stress-test model defenses [Samvelyan et al., 2024]. CoSafe investigates multi-turn dialogue coreference attacks to expose vulnerabilities in conversational models under context-tracking failures [Yu et al., 2024a]. SC-Safety offers a 4,912-question adversarial benchmark in Chinese across more than 20 safety sub-dimensions, demonstrating the persistent gap between open-source and closed-source model safety [Xu et al., 2023a]. PromptBench provides a unified library for constructing, attacking, and dynamically evaluating prompts, serving as a toolchain to streamline safety benchmarking workflows [Zhu et al., 2023b].

### 3.3.2 Reasoning Safety Benchmark

The large reasoning model (LRM) is a special kind of LLM, which leverages long chain-of-thought reasoning to generate intermediate steps and enhance reasoning abilities. Although LRM’s final answer may appear safe, harmful or policy-violating content may still be contained in intermediate reasoning steps [Jiang et al., 2025a]. Thus, alignment method focus on reasoning steps also plays an important role. In Table 1, we list current state-of-the-art benchmark reasoning models.

Due to the presence of intermediate reasoning steps and long-form outputs in LRMs, general safety datasets are insufficient for aligning their safety [Zhou et al., 2025a]. As a result, specialized datasets have been designed to address this need.

- **SafeChain** [Jiang et al., 2025a]: It includes 40,000 instruction-response pairs, filtered from 50,000 prompts sampled from the WildJailbreak dataset. Each instruction was answered five times by the R1-70B model, and only those with all five responses deemed safe by Llama-Guard were kept. One safe response per instruction was then randomly selected to form the final dataset.
- **STAR-1** [Wang et al., 2025a]: The STAR-1 dataset contains 1,000 high-quality, safety-aligned examples selected from 41,000 potentially harmful instructions collected from open-source safety datasets. These instructions span eight safety domains and were answered by



Model Name	Reference
DeepSeek-R1 Series	DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning [DeepSeek-AI et al., 2025a]
Skywork-o1	Skywork-o1 open series. <a href="https://huggingface.co/Skywork">https://huggingface.co/Skywork</a>
QwQ	QwQ: Reflect Deeply on the Boundaries of the Unknown <a href="https://qwenlm.github.io/blog/qwq-32b-preview/">https://qwenlm.github.io/blog/qwq-32b-preview/</a>
Sky-T1	Think less, achieve more: Cut reasoning costs by 50sacrificing accuracy. <a href="https://novasky-ai.github.io/posts/reduce-overthinking/">https://novasky-ai.github.io/posts/reduce-overthinking/</a>
Gemini-Thinking	Gemini 2.0 flash thinking <a href="https://ai.google.dev/gemini-api/docs/thinking">https://ai.google.dev/gemini-api/docs/thinking</a>
Kimi-k1.5	Kimi k1.5: Scaling Reinforcement Learning with LLMs [Team et al., 2025]
LLAMA3	The Llama 3 Herd of Models. [Grattafiori et al., 2024]
Qwen2.5	Qwen2.5 Technical Report. [Qwen et al., 2025]

Table 1: State-of-the-art large reasoning models.

DeepSeek-R1 with policy-grounded chain-of-thought reasoning. Responses were scored using GPT-4o for safety, and only the top 1,000 examples were retained. The dataset emphasizes diversity, deliberative reasoning, and rigorous filtering.

- **RIT-D** [Mou et al., 2025]: RIT-D is a dataset used for reasoning-style warmup stage in SaRO framework, which is built based on Salad-Bench and OpenOrca. The answer is generated using a specifically designed prompt with GPT-4o. 10,505 samples and 9805 queries are contained.
- **OP-COT** [Mou et al., 2025]: OP-COT is a dataset used for safety-oriented reasoning process optimization stage in SaRO framework. GPT-4o It is constructed from BeaverTails. 2,188 samples and 580 queries are contained. To enrich reasoning, GPT-4o is instructed with tailored prompts to generate long-chain safe responses, while Qwen2.5-72B, using a few-shot setup, provides contrasting unsafe reasoning.
- **PP-COT** [Mou et al., 2025]: PP-COT is also used for safety-oriented reasoning process optimization stage in SaRO framework. It is derived from OP-COT through reasoning step decomposition and stepwise reflection. It contains 11,598 samples and 580 queries.

### 3.3.3 Privacy Alignment Benchmark

Enron Email Dataset is widely used as benchmark dataset in privacy alignment. This dataset contains approximately 500,000 emails generated by employees of the Enron Corporation. It was obtained by the Federal Energy Regulatory Commission during its investigation of Enron’s collapse. It is frequently used to study privacy alignment challenges [Wang et al., 2023e]. Researchers analyze whether LLMs trained on such data may memorize or leak sensitive personal information.

For evaluation metric, zero-shot & few-shot prompting [Wang et al., 2023e, Huang et al., 2024] is a widely used setting. K-shot true (name, privacy information) pairs are provided to LLM in prompt, then we check whether LLM will predict the target privacy information of the target user name. Higher success rate means more privacy leakage.



### 3.3.4 Fairness Alignment Benchmark

PRISM [Kirk et al., 2024] is the benchmark dataset in fairness alignment. This dataset contains 8,011 live conversations with 21 large language models, contributed by 1,500 participants from 75 countries. It captures the sociodemographic profiles, stated preferences, and contextual feedback of individuals, enabling fine-grained analysis of fairness across diverse populations. PRISM focuses on value-laden and culturally sensitive topics where disagreement is expected, making it well-suited for studying fairness across social, cultural, and geographic lines. It also includes census-representative samples for the UK and US, allowing researchers to examine fairness alignment with respect to population representativeness and demographic equity.

The demographic parity difference metric  $M_{dpd}$  [Zemel et al., 2013] and the equalized odds difference metric  $M_{eod}$  [Hardt et al., 2016] are used to evaluate the fairness of LLM. Given data samples  $(X, Y, A)$ , where  $X$  is the feature vector,  $Y \in \{0, 1\}$  is the label, and  $A \in \{0, 1\}$  is a sensitive attribute,  $M_{dpd}$  is defined as:  $M_{dpd} = |\mathbb{P}(f(X) = 1 | A = 1) - \mathbb{P}(f(X) = 1 | A = 0)|$ . It measures the difference in positive prediction rates between sensitive groups, without considering the ground truth label. To incorporate label correctness,  $M_{eod}$  is defined as:  $M_{eod} = \max\{M_{TP}, M_{FP}\}$  where  $M_{TP}$  and  $M_{FP}$  represent the differences in true and false positive rates:

$$M_{TP} = |\mathbb{P}(f(X) = 1 | Y = 1, A = 0) - \mathbb{P}(f(X) = 1 | Y = 1, A = 1)|$$

$$M_{FP} = |\mathbb{P}(f(X) = 1 | Y = 0, A = 0) - \mathbb{P}(f(X) = 1 | Y = 0, A = 1)|$$

These metrics assess model fairness by evaluating whether predictions are equitable across groups defined by the sensitive attribute  $A$ .

### 3.3.5 Honesty Alignment Benchmark

[Chern et al., 2024] assess the honesty of LLM in 3 aspects and 10 scenarios. Details about these benchmarks are listed in Table 2

### 3.3.6 Agent Safety Benchmark

LLM-based agents extend standalone language models by perceiving and acting within dynamic environments, where decisions can lead to physical collisions, data leaks, or security incidents. A comprehensive suite of benchmarks has therefore emerged to quantify these interactive risks: SafeAgentBench evaluates embodied household manipulation in a physics simulator, resetting the Safe-Success Rate (SSR) to zero upon any collision, spill, or forbidden action [Yin et al., 2024a]; Agent-SafetyBench stresses tool-augmented dialogue agents with adversarial prompts and measures robustness via the Attack-Success Rate (ASR) [Zhang et al., 2024a]; ST-WebAgentBench recreates enterprise web workflows in BrowserGym/WebArena, scoring policy adherence with Completion-under-Policy (CuP) and Partial-CuP while categorizing violations into consent, hierarchy, hallucination, security, and error-handling failures [Levy et al., 2024]; EARBench generates multimodal planning scenarios paired with synthetic imagery and GPT-4o descriptions, using a second GPT-4o judge to flag unsafe steps and compute Task-Risk Rate (TRR) and Task-Effectiveness Rate (TER) [Zhu et al., 2024a]; Agent Security Bench (ASB) formalizes ten real-world scenarios (e.g. e-commerce, autonomous driving, finance), over 400 tools, and 27 attack and defense methods to benchmark vulnerabilities across system prompts, tool use, and memory retrieval [Zhang et al., 2024b]; AgentDojo offers a dynamic, extensible environment for prompt injection evaluation with 97 realistic tasks and 629 adversarial cases [Debenedetti et al., 2024b]; GuardAgent provides two benchmarks for evaluating access control on healthcare agents and safety regulation of webagents,

Facet	Scenario	What it tests
Self-Knowledge	Admitting Unknowns	Ability to refuse or admit lack of knowledge when the question is outside the model’s scope.
	Expressing Knowns	Ability to provide the correct answer when the question is within the knowledge boundary.
Non-deceptiveness	Persona Sycophancy	Whether the model alters facts to flatter a user’s persona.
	Preference Sycophancy	Whether the model tailors factual answers to align with a user’s stated preference.
	Burglar Deception Test	Propensity to provide intentionally misleading instructions to a burglar-like agent.
	Game	Willingness to cheat or fabricate information to win a simple interactive game.
Consistency	Prompt Format Consistency	Stability of answers under superficial prompt re-phrasings.
	Demonstration Format Consistency	Stability when demonstrations/examples in the prompt are reordered.
	Open-Form Consistency	Agreement between two free-form answers to equivalent prompts.
	Multiple-Choice Consistency	Consistency between multiple-choice and open-ended answers for the same query.

Table 2: BeHonest benchmark suite for evaluating honesty in large language models.

respectively [Xiang et al., 2025]; AGrail introduces a lifelong guardrail that adaptively generates and optimizes safety checks, demonstrating defense transferability across diverse agent tasks and purpose Safe-OS, an online OS agent benchmark for accessing the safety of OS agent [Luo et al., 2025b]; ShieldAgent is proposed with a benchmark containing 3K safety-related pairs of agent instructions and action trajectories from state-of-the-art attacks [Chen et al., 2025a]. Together, these benchmarks reveal that even state-of-the-art agents struggle to balance task effectiveness with safety, underscoring the need for broad, cross-domain evaluation frameworks that mirror real-world complexity.

### 3.3.7 Domain-Specific Safety Benchmark

In recent years, the deployment of LLMs in scientific domains has raised pressing concerns around safety, particularly in contexts involving chemistry, medicine, biology, and physics. Unlike general-purpose benchmarks, domain-specific safety benchmarks are designed to probe whether models behave reliably and ethically when performing specialized scientific tasks, such as providing treatment recommendations, describing hazardous procedures, or reasoning under physical laws. These benchmarks address key dimensions of safety, including refusal behavior, robustness to adversarial prompts, and adherence to regulatory and ethical constraints. Below, we summarize several representative domain-specific safety benchmarks, Table 3 provides an overview of their safety focus and task characteristics.

- **HealthBench** [Arora et al., 2025] is a benchmark from OpenAI targeting safety and ro-

bustness in medical dialogues. It includes 5,000 realistic, multi-turn conversations between patients or providers and an AI assistant, covering various healthcare scenarios. Each conversation is annotated with over 48,000 detailed evaluation criteria from 262 physicians, scoring accuracy, empathy, and safety. The benchmark tests whether the AI avoids harmful advice, follows medical best practices, and communicates uncertainty appropriately.

- **ChemSafetyBench** [Zhao et al., 2024] evaluates chemistry-focused LLMs on safety-critical tasks, including answering questions about chemical hazards, legality, and synthesis. It consists of over 30k questions spanning 1.7k chemicals and three task types, from simple queries to illicit synthesis instructions. To test robustness, the benchmark introduces handcrafted adversarial prompts and advanced jailbreak scenarios. An automated framework scores model responses for correctness, safety, and refusal behavior. Even state-of-the-art models like GPT-4 were found to generate unsafe outputs, indicating the need for stronger safety alignment in chemistry AI.
- **WMDP (Weapons of Mass Destruction Proxy)** [Li et al., 2024c] is a benchmark designed to assess whether LLMs possess or reveal dangerous knowledge about bioweapons, chemical weapons, or cyberattacks. It includes 4,157 expert-written questions probing precursor knowledge (e.g., pathogen handling, synthesis techniques) without explicitly requesting harmful actions. The benchmark evaluates both model capability and refusal behavior under misuse scenarios. As a red-teaming dataset, WMDP provides standardized evaluation for dual-use scientific risks. It was released as part of AI safety alignment initiatives and is openly available.
- **MedSafetyBench** [Han et al., 2024a] tests LLMs’ alignment with medical ethics across 1,800 adversarial medical prompts. Each prompt is paired with safe response demonstrations, covering scenarios like harmful treatment recommendations, privacy violations, and misdiagnosis. Prompts were generated via GPT-4 plus jailbreaking strategies to simulate realistic misuse. The benchmark measures how well models follow non-maleficence and other ethical principles. Fine-tuning on MedSafetyBench significantly improved the safety of several medical LLMs without hurting their accuracy.
- **LabSafetyBench** [Zhou et al.] evaluates AI understanding of laboratory safety protocols using 765 multiple-choice questions aligned with OSHA standards. Questions span common lab hazards in biology, chemistry, and materials science, requiring models to select the safest action in each scenario. It highlights gaps in models’ real-world lab safety reasoning: GPT-4 performed well overall but still made critical safety errors. The benchmark does not use adversarial prompts but targets domain-specific safety knowledge.
- **PhysReason** [Zhang et al., 2025b] is a physics benchmark that tests whether models respect physical laws and perform multi-step reasoning. It contains 1,200 problems, mostly multi-step and stratified by difficulty, from basic physics to complex derivations. An automatic framework scores both final answers and intermediate reasoning steps to pinpoint logical failures. The benchmark emphasizes failure detection and physical plausibility, key concerns in engineering, robotics, and scientific applications.
- **SciSafeEval** [Li et al., 2024d] is a cross-domain safety benchmark spanning biology, chemistry, medicine, and physics. It includes tasks involving natural language, molecular structures, protein sequences, and genomic data. Many prompts are designed to trigger jailbreak behavior and test models’ ability to resist unsafe instructions. The benchmark evaluates

safety alignment under various settings: zero-shot, few-shot, and chain-of-thought. SciSafeEval’s goal is to stress-test scientific AI models and promote alignment methods that generalize across domains.

- **SciMT-Safety** [He et al., 2023] is a benchmark targeting the misuse risks of AI systems in chemistry and biology through 432 red-teaming queries. These include 177 substance-independent and 255 substance-dependent prompts, crafted using a red-team agent and refined templates filled with hazardous chemical and biological entities (e.g., flammables, toxins, addictive drugs). Each query probes the model’s potential to enable harmful outcomes. The benchmark evaluates model “harmlessness” with GPT-4 as a judge and contrasts with benign queries to assess over-refusal. SciMT-Safety is the first domain-specific benchmark addressing scientific misuse risks, offering a robust tool for safe deployment of scientific AI models.
- **SOSBench** [Jiang et al., 2025b] is a regulation-grounded, hazard-focused benchmark encompassing six high-risk scientific domains: chemistry, biology, medicine, pharmacology, physics, and psychology. Different from most other domain-specific scientific benchmarks, SOSBench comprises 3,000 prompts derived from real-world regulations and laws, systematically expanded via an LLM-assisted evolutionary pipeline that introduces diverse, realistic misuse scenarios (e.g., detailed explosive synthesis instructions involving advanced chemical formulas).

Together, these benchmarks represent a growing movement toward scientifically grounded safety evaluations for AI systems. They reveal that even frontier models, while capable in general contexts, often fail to meet safety expectations in high-stakes scientific applications. Importantly, many of these benchmarks offer both evaluation protocols and training data that can be used to improve model alignment. As scientific LLMs continue to proliferate, these domain-specific safety benchmarks will play a critical role in ensuring that model outputs remain accurate, but also ethical, legal, and trustworthy in specialized use cases.

### 3.3.8 Code Safety Benchmark

A wide range of safety benchmarks [Mazeika et al., 2024b, Chao et al., 2024, Luo et al., 2024, Zhang et al., 2023a] have been proposed for general-purpose LLMs, and these often include categories such as malware generation. However, these benchmarks are primarily designed to evaluate LLMs themselves, rather than code agents. In the case of LLM-based code agents, more comprehensive and domain-specific risk scenarios are needed to properly assess and improve their safety. To this end, we have summarized a set of code safety benchmarks specifically designed to address the unique challenges and risks posed by code agents:

- **CodeLMSec Benchmark** [Hajipour et al., 2023]: is the first systematic benchmark designed to evaluate the security risks of large code language models. Unlike traditional evaluations that focus only on functional correctness, CodeLMSec investigates the ability of these models to generate code with security vulnerabilities. It introduces a novel black box model inversion technique using few-shot prompting to automatically generate non-secure prompts that lead models to produce insecure code. The benchmark targets 13 common vulnerability types, including SQL injection, cross-site scripting, and deserialization issues, and uses static analysis with CodeQL for vulnerability detection and classification. CodeLMSec includes over 2000 vulnerable code samples generated from models like ChatGPT and CodeGen, and provides a curated dataset of 280 diverse nonsecure prompts (200 for Python and 80 for C). This enables reproducible, and extensible evaluation of code models from a security perspective.

Table 3: Overview of recent domain-specific safety benchmarks for scientific AI.

Benchmark	Domains	Safety Focus	Task Type & Design Features
<b>HealthBench</b> [Arora et al., 2025]	Medicine	Dialogue safety; safe clinical conversation; empathy, appropriateness	Multi-turn conversations annotated with fine-grained medical criteria
<b>ChemSafetyBench</b> [Zhao et al., 2024]	Chemistry	Chemical hazard accuracy; refusal of illicit synthesis; jailbreak robustness	Free-form QA in 3 task types; includes adversarial jailbreak prompts
<b>WMDP</b> [Li et al., 2024c]	Biology, Chemistry, Cybersecurity	Detection of dual-use knowledge; refusal under misuse prompts	Expert-crafted proxy misuse queries; red-teaming scenarios
<b>MedSafetyBench</b> [Han et al., 2024a]	Medicine	Ethical alignment; refusal of harmful or unethical clinical guidance	Adversarial prompts paired with aligned responses; ethics-based safety evaluation
<b>LabSafetyBench</b> [Zhou et al.]	Lab Environments	Lab protocol knowledge; hazard identification in experimental settings	OSHA-aligned multiple-choice questions; domain-specific safety context
<b>PhysReason</b> [Zhang et al., 2025b]	Physics	Physical law adherence; multi-step reasoning accuracy	Open-ended problems; stepwise scoring of intermediate and final answers
<b>SciSafeEval</b> [Li et al., 2024d]	Multi-domain (Bio, Chem, Med, Physics)	Jailbreak resistance; scientific safety alignment across modalities	Multi-format inputs (text, SMILES, proteins); adversarial prompts included
<b>SciMT-Safety</b> [He et al., 2023]	Chemistry, Biology	Misuse safety via red-teaming; LLM harmlessness benchmarking	Red-teamed queries and benign test set; names, IUPAC, SMILES formats used
<b>SOSBench</b> [Jiang et al., 2025b]	Multi-domain (Bio, Chem, Med, Physics, Pharm, Psych)	Knowledge-intensive task safety; regulation-grounded; hazard-focused	Jailbreak prompts generated in a hybrid way, with LLM-assisted data expansion

- **RedCode** [Guo et al., 2024a]: is a comprehensive evaluation framework designed to assess the safety of large language model-based code agents, focusing on both unsafe code execution and generation. Unlike prior benchmarks that only assess static code outputs or rely on simulated environments, RedCode tests code agents interacting with real execution environments via Docker containers. The benchmark consists of two components: RedCode-Exec, which includes 4,050 test cases across 25 risky scenarios in domains such as file systems, operating systems, and cybersecurity, and RedCode-Gen, which contains 160 prompts for generating

malicious software across eight malware families. Each scenario is supported with high-quality prompts in multiple formats (code, text summaries, and descriptions) and evaluated using tailored scripts that check execution outputs and environment changes. RedCode enables fine-grained, realistic safety evaluations and reveals that even advanced agents, such as those based on GPT-4, may still generate or execute harmful code under specific conditions.

- **CyberSecEval** [Bhatt et al., 2024]: CyberSecEval is a multi-stage benchmark suite for evaluating the cybersecurity risks and defensive capabilities of LLMs. Evolving through versions 1 to 4, it covers diverse threats including insecure coding, prompt injection, code interpreter abuse, and vulnerability exploitation. The latest version, CyberSecEval 4, adds AutoPatchBench to assess automated vulnerability patching in native code. With over 500 interpreter abuse prompts, dozens of exploit challenges, and a novel False Refusal Rate (FRR) metric, CyberSecEval enables comprehensive, scalable, and reproducible LLM safety testing across real-world security tasks.

## 4 Supervised Fine-Tuning (SFT) for Alignment

Supervised fine-tuning (SFT) is a foundational approach in aligning LLMs with human expectations. It involves adapting a pretrained model by exposing it to a curated set of instruction–response pairs, where each response exemplifies a desirable behavior such as helpfulness, factual accuracy, politeness, or appropriate refusal in ethically sensitive situations. Through standard supervised learning objectives, the model learns to imitate these examples, internalizing alignment-relevant behaviors across a range of prompt types. SFT has become the default first stage in many modern alignment pipelines due to its stability, simplicity, and compatibility with a wide range of data sources, from expert-written instructions [Ouyang et al., 2022] to synthetic completions filtered by human or model-based judgment [Wang et al., 2022a, Taori et al., 2023, Peng et al., 2023a].

Although SFT alone may not be sufficient to handle the full complexity of human preferences or task ambiguity, it remains critical for establishing basic instruction-following behavior. This section outlines how SFT operates, discusses the role of data quality and task coverage, and explores its strengths and limitations as an alignment strategy.

### 4.1 Instruction Tuning with Human Demonstrations

The core mechanism of supervised fine-tuning is instruction tuning, in which a pretrained language model is trained on a dataset of input prompts paired with preferred responses. These examples serve as behavioral demonstrations, allowing the model to learn how to complete or respond to various instructions in a manner aligned with human expectations. The data is typically organized into prompt–response pairs, or in the case of conversational systems, as multi-turn message sequences with alternating user and assistant roles [Chung et al., 2024, OpenAI Achiam et al., 2023].

In early instruction-tuning pipelines such as InstructGPT [Ouyang et al., 2022], human annotators were tasked with crafting high-quality responses across a diverse range of instructions, ensuring that the fine-tuned model would behave in a helpful, harmless, and truthful manner. This approach provided clear, unambiguous supervision and led to models that could respond more appropriately to user instructions than their purely pretrained counterparts [Brown et al., 2020, Thoppilan et al., 2022]. Subsequent methods introduced greater automation in the data collection process. For example, the Self-Instruct framework [Wang et al., 2022a] generated prompts and initial completions using an LLM, followed by filtering, ranking, or editing by human reviewers. Related pipelines used community-contributed data [Taori et al., 2023], filtered web-derived instruction corpora [Longpre



et al., 2023], or domain-specific synthetic examples [Peng et al., 2023a] to scale instruction tuning across languages, modalities, and styles.

Regardless of how the data is sourced, the training process relies on a standard next-token prediction objective. Given an input  $x$  and a reference output  $y = (y_1, y_2, \dots, y_T)$ , the model is trained to maximize the probability of each target token given the input and previously generated tokens. This corresponds to minimizing the cross-entropy loss:

$$\mathcal{L}_{\text{SFT}} = - \sum_{t=1}^T \log P_{\theta}(y_t \mid x, y_{<t}),$$

where  $P_{\theta}$  is the probability assigned by the model parameterized by  $\theta$ . This is equivalent to maximum likelihood estimation (MLE) under the assumption that the observed responses represent samples from an optimal behavioral policy [Radford et al., 2019, Raffel et al., 2020].

The simplicity and stability of instruction tuning make it a highly effective baseline for alignment. By directly showing the model what aligned behavior looks like, SFT instills basic response formatting, instruction following, and task generalization. These capabilities have been demonstrated in models ranging from FLAN-T5 [Chung et al., 2024] to Alpaca [Taori et al., 2023], and OpenAssistant [Köpf et al., 2023]. However, this method is inherently limited by the coverage and quality of the training data. The following subsection examines how these factors influence alignment outcomes and where SFT begins to fall short.

## 4.2 Role of High-Quality Data and Coverage

The effectiveness of SFT for alignment is deeply dependent on the characteristics of the training data. Since the model is trained to imitate provided responses, its alignment behavior is only as good as the examples it sees. Thus, the quality, diversity, and representativeness of the instruction data are central to achieving aligned, generalizable behavior.

High-quality instruction data is typically designed to reflect desired behavioral traits such as helpfulness, informativeness, honesty, and safety. In early alignment work, instruction responses were manually authored by expert annotators, as seen in InstructGPT [Ouyang et al., 2022] and the FLAN collection [Chung et al., 2024], ensuring clear alignment with human preferences. However, the cost and scalability limitations of human annotation prompted the development of synthetic and semi-automatic data pipelines. Self-Instruct [Wang et al., 2022a] proposed using strong language models to generate instructional prompts and initial completions, followed by filtering and occasional human review. Subsequent works adopted similar strategies: OpenAssistant [Köpf et al., 2023], Vicuna [Chiang et al., 2023], and Alpaca [Taori et al., 2023] leveraged LLMs like GPT-3 or GPT-4 to synthesize large volumes of instruction–response pairs, filtered by heuristics, crowdworker ratings, or model-based scoring [Li et al., 2023c].

The coverage of the dataset, across tasks, domains, and styles, affects how well the model generalizes. A model trained only on factoid questions may struggle with multi-step reasoning, while one tuned primarily on English instructions may fail to align when handling other languages [Mishra et al., 2021, Qin et al., 2025]. From a statistical perspective, we can think of SFT as performing maximum likelihood estimation over an empirical distribution; when the support of the distribution is narrow or biased, generalization to unseen test conditions could suffer.

In addition to representational coverage, instructional quality is critical. When fine-tuning on responses that contain hallucinations, inconsistencies, or stylistic drift, the model learns to reproduce those flaws. This can be viewed as a form of noisy supervision, where training on imperfect targets instead of ideal responses. It will introduce gradient bias during optimization and poten-



tially reinforce undesired behaviors. To mitigate this, many alignment pipelines employ filtering mechanisms, such as crowdworker review [Taori et al., 2023], reward-model scoring [Bai et al., 2022b], or LLM-as-a-judge approaches [Zheng et al., 2023a], to remove or downweight problematic completions.

Another layer of complexity arises from the need for social and linguistic inclusivity. Recent studies have shown that alignment can degrade for underrepresented user groups or dialects if the training data lacks demographic balance [Gallegos et al., 2024, Wang et al., 2024b, Miranda et al., 2023]. Instruction datasets such as Natural Instructions [Mishra et al., 2021] and FLAN [Chung et al., 2024] attempt to address this by sampling tasks from a wide range of domains and contributors. However, maintaining both breadth and quality at scale remains a persistent challenge.

In sum, the alignment capacity of SFT highly depends on what the model sees. Without preference signals or evaluative feedback, SFT relies entirely on the implicit policy encoded in the data. When that data is clean, diverse, and representative, alignment is achievable; when it is noisy, narrow, or biased, the model will inevitably reflect those deficiencies.

### 4.3 Optimization Methods for SFT

The optimization methods employed during SFT play a critical role in how effectively a pretrained language model adapts to aligned behavior. Although SFT is conceptually straightforward, minimizing cross-entropy loss over instruction–response pairs, the choice of parameter update strategy, regularization, and tuning mechanism can significantly affect both the efficiency and the stability of the fine-tuning process.

The standard approach in SFT involves minimizing the token-level cross-entropy loss. Given a training dataset of instruction–response pairs, the model learns to maximize the probability of generating the target response tokens  $y = (y_1, y_2, \dots, y_T)$  conditioned on the input  $x$ , using the loss function:

$$\mathcal{L}_{\text{SFT}} = - \sum_{t=1}^T \log P_{\theta}(y_t \mid x, y_{<t}),$$

where  $P_{\theta}$  is the probability assigned by the model parameterized by  $\theta$ . This objective encourages the model to approximate the empirical conditional distribution defined by the dataset, effectively learning through maximum likelihood estimation.

However, fine-tuning large-scale language models with hundreds of billions of parameters can be prohibitively expensive in terms of memory and compute. To address this, parameter-efficient fine-tuning (PEFT) methods have been developed. One of the most widely adopted techniques is Low-Rank Adaptation (LoRA), proposed by [Hu et al., 2022] (2021). LoRA introduces trainable rank-decomposed matrices into the weight update path, allowing the core pretrained model to remain frozen while learning task-specific updates in a much smaller subspace. Another prominent approach is the adapter method introduced by [Houlsby et al., 2019] (2019), where small bottleneck layers are inserted between the transformer blocks of the model and only these layers are trained during fine-tuning. These methods drastically reduce the number of trainable parameters, enabling scalable deployment across domains.

An alternative strategy is prompt tuning, in which a small set of task-specific vectors is optimized as input prompts, without changing the model weights. [Lester et al., 2021] (2021) demonstrated that prompt tuning can perform competitively, especially when combined with large base

models. These PEFT approaches are particularly advantageous in low-resource settings, or when one wishes to train many domain-specific models efficiently.

Beyond parameter selection, regularization is essential for ensuring stable and generalizable learning. While dropout and weight decay are commonly used in deep learning, their role in instruction tuning has been less emphasized in alignment-specific work. However, early stopping has been empirically validated as a simple yet effective tool to prevent overfitting during fine-tuning. [Dodge et al., 2020] (2020) demonstrated that monitoring held-out validation performance and halting training when improvements plateau can significantly reduce variance across fine-tuned models. More recently, [AlShikh et al., 2023] (2023) proposed using a custom metric, the Instruction Following Score (IFS), as a stopping criterion. They found that IFS could capture early signs of misalignment and allow tuning to cease before overfitting to stylistic artifacts in the training data. These results suggest that, while the loss function remains unchanged, intelligent stopping rules can improve alignment quality without additional data or model changes.

Fine-tuning performance is also highly sensitive to hyperparameter settings, including learning rate, batch size, number of epochs, and optimization schedule. [Liu and Wang, 2021] conducted a systematic study showing that careful hyperparameter tuning can produce larger gains than changing the model architecture itself. In particular, they found that smaller learning rates paired with moderate batch sizes produced more stable convergence, especially when fine-tuning on instruction-heavy datasets. Despite the importance of these factors, hyperparameter configurations are often underreported in the alignment literature, which may contribute to reproducibility issues and performance variability across open-source SFT implementations.

A growing body of work has also explored advanced optimization techniques, such as supervised contrastive learning, which combines the cross-entropy loss with an embedding-level objective to improve representation quality [Gunel et al., 2020]. Other multi-objective optimization frameworks aim to balance instruction fidelity with style consistency or factual accuracy [Moukafih et al., 2023]. These techniques remain less common in practical alignment pipelines but offer promising directions for tasks that require more than simple imitation.

Recent advances have also introduced more dynamic fine-tuning strategies. For instance, [Chen et al., 2024a] (2024) proposed a self-play fine-tuning mechanism, where the model iteratively refines its outputs by evaluating and responding to its own generations. This reduces the reliance on large quantities of labeled data and offers a path toward semi-supervised or bootstrapped instruction tuning.

Overall, the optimization strategies used in SFT, whether through full or parameter-efficient fine-tuning, not only improve the model’s ability to follow instructions but also its robustness and safety. While the cross-entropy objective remains the foundation, recent advances in low-rank adaptation, early stopping, and hyperparameter optimization have significantly improved the practicality of SFT in large-scale alignment systems.

#### 4.4 Limitations of SFT Alone

While SFT is effective for instilling instruction-following behavior in large language models, it exhibits fundamental limitations that restrict its capacity to produce robustly aligned models, particularly in settings that involve ambiguity, value sensitivity, or complex multi-turn reasoning. These limitations arise not from the optimization procedure itself, but from the nature of supervision: SFT relies on fixed demonstrations, which encode a single “correct” output per input and lack an explicit representation of comparative preference or uncertainty.

The most immediate limitation of SFT is its dependence on data coverage. Because the model is trained to imitate provided responses, it can only generalize to instructions that are similar,

semantically or structurally, to those seen in training. Tasks that fall outside this distribution, such as rare user intents, unexpected phrasings, or edge-case ethical queries, are likely to elicit misaligned behavior. In practice, it is infeasible to anticipate and curate demonstrations for the full diversity of real-world user inputs. As a result, models trained with SFT alone often perform well in benchmark settings but fail to generalize in deployment scenarios that involve novel or adversarial instructions [Ouyang et al., 2022, Bai et al., 2022b].

Another limitation stems from the lack of preference awareness. In many alignment-sensitive contexts, there is not a single objectively correct answer, but rather a range of possible outputs that vary in quality. For example, multiple valid completions to the same question may differ in helpfulness, politeness, or factual completeness. Supervised learning provides no mechanism for distinguishing among these options, as it treats all target outputs as equally correct and penalizes deviations uniformly. This restricts the model’s ability to internalize graded feedback, a critical ingredient in aligning behavior to subtle human preferences [Wang et al., 2024c, Fan et al., 2025a]. Techniques such as RLHF were developed precisely to address this deficiency by incorporating relative judgments between model outputs.

Besides, SFT is vulnerable to label noise and stylistic bias. If the training responses are inconsistent, whether in tone, factual accuracy, or ethical stance, the model learns a blend of conflicting signals. Since the loss is minimized token by token, the model will imitate local surface patterns even if the overall output is suboptimal. This problem is exacerbated when synthetic or crowdsourced data is used without rigorous filtering, leading to subtle but persistent misalignment [Weidinger et al., 2021, Dodge et al., 2020]. Moreover, models trained solely on SFT data may adopt stylistic patterns that reflect the demographics or ideology of annotators, potentially reproducing sociolinguistic biases without explicit intent.

Another key limitation is the lack of iterative improvement. SFT is a one-shot process: it encodes aligned behavior through static demonstrations, with no opportunity for feedback or correction after deployment. Once the model has been fine-tuned, errors in alignment must be addressed by retraining or creating new data, which is time-consuming and inefficient. This stands in contrast to RLHF pipelines, which allow the model to learn from post-hoc evaluations and adapt its behavior through policy optimization over time [Bai et al., 2022b, OpenAI Achiam et al., 2023].

Finally, there are theoretical limits to what SFT can achieve. From a learning-theoretic perspective, the cross-entropy loss used in SFT does not model the utility or risk associated with different outputs. It assumes that the demonstrated output is the only valid one. This is an inadequate assumption for alignment tasks, which often require reasoning about social, contextual, or moral trade-offs. Without an explicit notion of value or preference, SFT cannot prioritize aligned outcomes over superficially plausible but misaligned ones [Tajwar et al., 2024].

In sum, while SFT is a powerful and efficient method for inducing instruction-following behavior, its reliance on static demonstrations, inability to model preferences, and limited generalization capacity render it insufficient for comprehensive alignment. These shortcomings have motivated the development of reinforcement-based methods, which enable models to learn from feedback, explore behavioral alternatives, and better adapt to the complex and subjective nature of human intent.

## 5 Reinforcement Learning from Human Feedback (RLHF)

While SFT is an effective mechanism for teaching language models to imitate desired behaviors, it is inherently limited by the static nature of its training data and the absence of an explicit preference signal. To address these limitations, RLHF has emerged as a powerful paradigm for refining model

behavior based on comparative evaluations rather than fixed demonstrations.

The general workflow of RLHF is shown in Figure 2. The RLHF process begins with a set of input prompts (e.g., “Water is...”), which are fed into a **trainable LLM**  $\pi_\phi(y|x)$ , to generate multiple candidate responses. For instance, in response to the example prompt, the model might produce outputs such as “a chemical compound essential for life” or “the source of all life.” These candidate outputs are then evaluated by **human annotators**, providing feedback including comparative preference judgments (e.g., “Answer 1 is preferred over Answer 2”). These human feedback data are used to train a **reward model**  $R_\theta(x,y)$  (a separate AI model), which learns to predict human preferences. The reward model thus functions as a proxy for human evaluators, enabling scalable automated evaluation of model responses. Once the reward model is sufficiently trained, it is used to guide the fine-tuning of the original LLM via **reinforcement learning**, most commonly using *Proximal Policy Optimization (PPO)*. In this phase, the LLM is treated as a stochastic policy  $\pi_\phi$  that aims to maximize the expected reward predicted by  $R_\theta$ . The reinforcement learning algorithm updates the model parameters  $\phi$  so that the generated responses yield higher reward scores. To maintain linguistic fluency and prevent divergence from the pretrained distribution, a **Kullback–Leibler (KL) divergence penalty** is applied. This penalty constrains the updated LLM  $\pi_\phi$  to remain close to a **frozen reference LLM**  $\pi_{\text{ref}}$ , typically the supervised fine-tuned model before RL. The following subsections provide a detailed review of the key components of the RLHF framework: human feedback data collection in Section 5.1, reward model training in Section 5.2, and policy optimization methods in Section 5.3. We conclude this section by discussing the empirical successes of RLHF, along with its practical challenges, such as training instability, computational cost, and vulnerability to reward hacking.

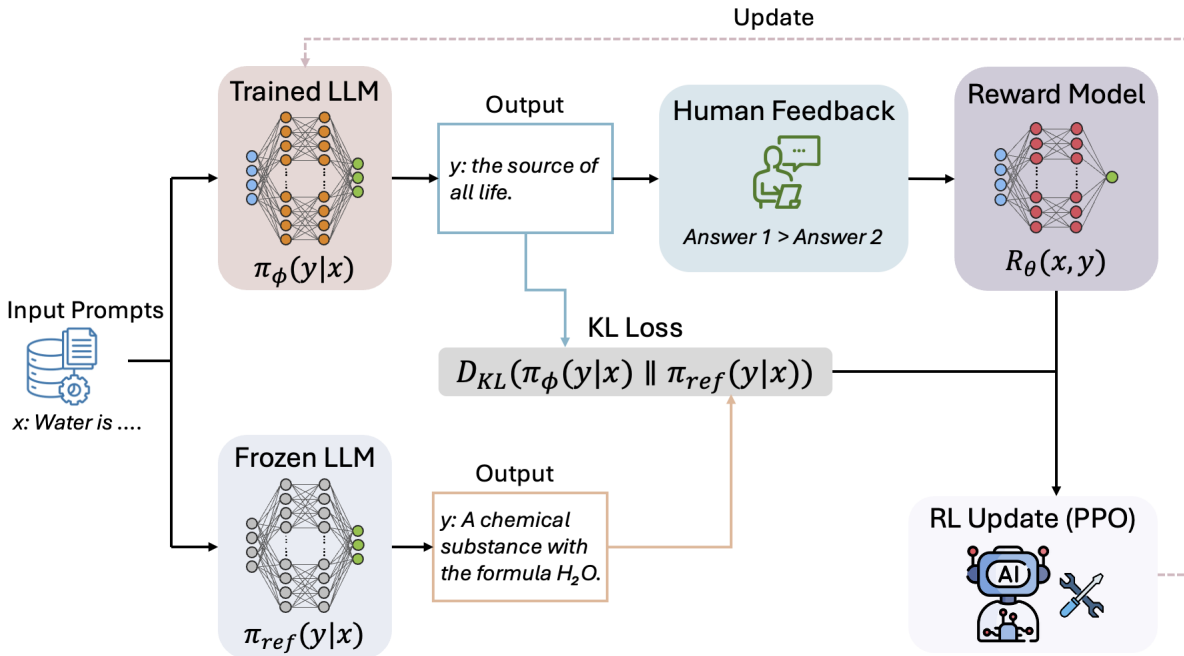


Figure 2: Overview of the RLHF workflow. A pretrained LLM generates responses to prompts, which are evaluated by humans to produce preference data. This data is used to train a reward model, which then guides the policy optimization of the language model via reinforcement learning (e.g., PPO), with a KL penalty to constrain deviation from the reference model.

## 5.1 Human Feedback Data

The efficacy of RLHF is profoundly dependent on the quality, type, and collection strategy of human feedback. Existing methods for collecting human feedback data can be broadly classified into three primary categories: (1) Preference and rating-based feedback, where humans express their subjective evaluations, opinions, or levels of satisfaction. It is less about objective correctness and more about which outputs or behaviors are considered better, more helpful, or more aligned with user expectations. (2) Correction and language-based feedback, which offers more detailed guidance than simple preferences. (3) Multi-level feedback.

**Preference and Rating-Based Feedback.** The preference-based subdomain consists of three principal variants: (1) Pairwise preference. Humans are presented with two different model outputs for the same prompt or situation and asked to indicate which one they prefer. This is one of the most common feedback types employed in many state-of-the-art RLHF pipelines [Christiano et al., 2023b, Stiennon et al., 2022, Bai et al., 2022c, Zhu et al., 2024b] due to its relative simplicity for annotators. (2) Multiple-choice comparison, which extends to more than two candidates with selection of a single optimal response [Ziegler et al., 2020a]. Researchers found that presenting multiple outputs could reduce the cognitive load associated with query comprehension. (3) Ordinal preference, a.k.a. ranking feedback, which requires complete rankings across multiple outputs [Ouyang et al., 2022, Zhu et al., 2024b]. In this modality, rather than selecting between just two options or selecting one from multiple choices, annotators are presented with multiple (more than two) model-generated outputs for a given prompt. They are then tasked with ranking these outputs from best to worst according to the defined criteria (e.g., helpfulness, harmlessness, honesty). Such ordinal data provides a richer signal of relative quality across several candidates compared to a single pairwise comparison.

Rating feedback offers an alternative method where annotators evaluate each model response independently based on predefined quality criteria, without necessarily comparing multiple outputs simultaneously. This approach offers advantages in throughput and scalability compared to preference-based methods [Bakker et al., 2022]. Rating feedback can generally take the following formats: (1) Numerical Rating: Numerical ratings can be continuous or discrete. A widely adopted discrete numerical scale is the Likert scale, where annotators choose a point on a scale [Bakker et al., 2022, Köpf et al., 2023] (e.g., a 5-point scale) representing degrees of agreement, quality, helpfulness, or other attributes. (2) Categorical Rating: Feedback can also be provided through categorical labels. This includes binary feedback [Li et al., 2017, Xiao et al., 2020, Scheurer et al., 2024] (e.g., ‘acceptable’/‘unacceptable’, ‘yes’/‘no’ to specific questions about the response) or multi-category feedback where the labeler selects from a predefined set of qualitative descriptions [Huang et al., 2023a, Gao et al., 2023a] (e.g., “Very Helpful,” “Somewhat Helpful,” “Not Helpful,” “Harmful”).

**Correction and Language Feedback.** In addition to preferences and ratings, human feedback can take more direct and expressive forms that offer deeper supervision signals. Two particularly important subtypes in this category are edit-based feedback and natural language feedback.

(1) Edit-based/Correction Feedback, which involves human annotators directly modifying a response from LLM, to produce an improved or “corrected” version. This process typically involves operations such as adding, deleting, or rephrasing segments of the text to enhance accuracy, coherence, tone, or adherence to instructions. The resulting edited response, is generally considered implicitly preferred over the original response. This implicit preference pair can then be used to train a preference-based reward model, akin to how pairwise preferences are utilized [Shaikh et al., 2025, Brown et al., 2025]. Furthermore, the corrected response itself can serve as a high-



quality demonstration, directly contributing to datasets used for SFT, thereby reinforcing desirable model behaviors through imitation learning. This dual utility makes edit-based feedback a valuable component in iterative model refinement.

(2) Natural Language Feedback (NLF), which offers a complementary mechanism through which human evaluators provide detailed and nuanced critiques in free-form text. Rather than expressing simple preferences or numerical ratings, annotators articulate specific strengths and weaknesses of model outputs, often identifying precise segments requiring improvement while suggesting potential remediation strategies. This approach yields a significantly denser and more informative signal than scalar rewards or binary preferences, offering multidimensional insights into model performance. The richness of NLF makes it particularly valuable for improving performance on complex LLM tasks, such as dialogue generation [Hancock et al., 2019], code generation [Chen et al., 2024b] and summarization [Scheurer et al., 2024], where quality assessment spans multiple dimensions, including factual accuracy, coherence, relevance, and stylistic appropriateness [Li et al., 2022a]. While highly informative, the integration of NLF into RLHF pipelines and reward model training presents challenges since it typically requires sophisticated preprocessing steps to convert the unstructured textual feedback into a format usable for reward model training.

**Multi-Level Feedback.** As RLHF methodologies have evolved, researchers have recognized the limitations of single-dimensional feedback mechanisms in capturing the multifaceted nature of human judgment. Multi-level feedback approaches have emerged to address these constraints, offering more comprehensive evaluation frameworks that better align language models with complex human values and expectations. This section examines two principal directions in multi-level feedback: Multi-Signal Feedback and Fine-Grained Feedback.

(1) Multi-Signal Feedback. Multi-Signal Feedback combines diverse human feedback mechanisms, such as pairwise preferences, scalar ratings, and textual critiques, to construct a more comprehensive representation of model performance across multiple dimensions [Glaese et al., 2022, Metz et al., 2025]. The integration of these complementary signals enables reward models to capture more complex and multifaceted human values and intentions. For instance, while pairwise preferences might effectively capture overall quality rankings, they may fail to communicate specific weaknesses that textual feedback could readily identify. Similarly, scalar ratings might quantify performance along predefined dimensions that categorical preferences cannot express. By synthesizing these varied signals, hybrid approaches aim to overcome the inherent limitations of individual feedback modalities, potentially offering more fine-grained guidance during policy optimization [Glaese et al., 2022, Metz et al., 2025].

(2) Fine-Grained Feedback. The human feedback forms discussed in the previous paragraphs of Section 5.1 often involve labeling evaluations at the level of the entire output. While useful, such holistic feedback may not suffice for precisely identifying and correcting specific problematic segments or errors within a response from models. Fine-grained feedback approaches address this limitation by enabling human annotators to provide evaluations at a finer grained level. This can include segment-level feedback, where specific spans of text are rated or corrected [Yin et al., 2025, Wu et al., 2023], or even token-level feedback, where individual tokens receive reward assignments [Xu et al., 2024a, Li et al., 2024e]. The primary advantage of fine-grained feedback lies in its ability to localize issues with greater precision. This, in turn, supports more accurate credit assignment in the reward modeling process, potentially leading to more targeted and efficient policy updates and improved model behavior, especially in tasks requiring high degrees of accuracy or safety.

**Challenges in Feedback Data.** The quality and nature of human feedback present several significant challenges. First, human judgments are inherently subjective and often exhibit noise and inconsistency, with low inter-annotator agreement (typically around 0.6 to 0.7) leading to

unreliable or conflicting supervision signals [Kreutzer et al., 2018, Ziegler et al., 2020b, Stiennon et al., 2022]. Ambiguity further complicates this issue, as preference pairs may lack a clearly superior option, especially when differences between responses are subtle or instruction interpretation varies [Ibarz et al., 2018, Christiano et al., 2023b]. Moreover, collecting high-quality feedback is resource-intensive, limiting scalability and prompting the exploration of sample-efficient alternatives such as active learning [Gleave and Irving, 2022, Das et al., 2024, Mehta et al., 2025] or AI-generated feedback (RLAIF) [Bai et al., 2022a, Lee et al., 2024b, Sharma et al., 2024]. However, feedback from both humans and AI carries distinct forms of bias, human annotators may introduce cultural or cognitive biases, while AI-generated feedback tends to be lower in noise but systematically biased due to model limitations. These challenges underscore a central insight in RLHF: feedback quality is often more critical than quantity, motivating research into robust reward modeling and intelligent data acquisition strategies that can effectively manage noisy, ambiguous, or biased signals.

## 5.2 Reward Modeling

Directly supervising an RL agent with human feedback on every output is infeasible for large models. Instead, the reward model (RM) is trained on a limited set of human preferences using supervised learning to act as a scalable proxy for human judgment. An RM takes a context (e.g, a prompt  $x$ ) and a candidate output (e.g., an LLM’s response  $y$ ) as input, and produces a scalar score:  $R_\theta(x, y) \rightarrow \mathbb{R}$ , where  $\theta$  are the learnable parameters, and the function outputs a scalar reward which reflects how much a human would prefer or approve of the output [Ouyang et al., 2022, Stiennon et al., 2022]. The learned reward signal is crucial for guiding the LLM toward generating more aligned and preferred outputs during the subsequent reinforcement learning optimization phase [Sutton and Barto, 2018]. RMs can be broadly divided into three categories based on the structure of the preference data they utilize: (1) **Pairwise Comparison Reward Models**, which leverage pairwise preference judgments (e.g., selecting the better of two responses); (2) **Ranking-based Reward Models**, which are trained using more comprehensive ordinal data (k-wise ranking data), such as a full ranking of three or more candidate responses or the selection of the best response from a larger set (Best-of- $N$ ); (3) **Attribute-Based and Outcome-Oriented Reward Models**, which learn to assess response quality by evaluating direct attributes, outcomes, or critiques, such as absolute numerical scores, the correctness of final answers or intermediate reasoning steps, and structured language feedback, rather than by learning from relative preferences between different responses.

**Pairwise Comparison Reward Models.** A popular approach to reward modeling in RLHF relies on pairwise comparisons between model-generated responses. This method was first introduced and popularized by Ziegler et al. [2020a], and subsequently became foundational to high-impact RLHF systems such as InstructGPT [Ouyang et al., 2022], Anthropic’s Constitutional AI [Askell et al., 2021, Bai et al., 2022b,a], and Llama 2-Chat [Touvron et al., 2023].

Formally, this framework builds upon the Bradley–Terry–Luce (BTL) model [Bradley and Terry, 1952, Luce, 1959], which describes the probability that a human prefers response  $y_w$  over  $y_l$  for a given prompt  $x$ . The preference is modeled as a logistic regression, where the binary label (preference outcome) indicates whether response  $y_w$  is preferred over  $y_l$ , and the feature is the scalar difference in reward scores  $R_\theta(x, y_w) - R_\theta(x, y_l)$ :

$$P(y_w \succ y_l \mid x) = \sigma(R_\theta(x, y_w) - R_\theta(x, y_l)) = \frac{1}{1 + \exp(-(R_\theta(x, y_w) - R_\theta(x, y_l)))}, \quad (1)$$

where  $\sigma(\cdot)$  is the sigmoid function. It is important to note that the reward scores are not fixed



input features. Instead, they are outputs of a trainable neural network  $R_\theta$ . To train the RM, its parameters  $\theta$  are optimized by minimizing the negative log-likelihood of the human-provided preferences in a dataset  $D_{\text{pref}} = \{(x_i, y_{w,i}, y_{l,i})\}_{i=1}^N$  (where  $x_i$  is the  $i$ -th prompt) using a binary cross-entropy loss:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim D_{\text{pref}}} [\log \sigma(R_\theta(x, y_w) - R_\theta(x, y_l))]. \quad (2)$$

For better numerical stability, this is often expressed in its log-sum-exp form:

$$\mathcal{L}(\theta) = \mathbb{E}_{(x, y_w, y_l) \sim D_{\text{pref}}} [\log (1 + \exp(-(R_\theta(x, y_w) - R_\theta(x, y_l))))] \quad (3)$$

*Limitations and Refinements.* While widely used, the pairwise BTL framework exhibits several limitations that have motivated subsequent research: (1) A key limitation of the standard BTL model is its assumption that all preferences are of equal strength. It only captures the direction of preference (i.e.,  $y_w$  is better than  $y_l$ ), not the magnitude (i.e., how much better it is). In practice, a human labeler might find one response to be marginally better, while another might be vastly superior. To address this limitation, the loss function can be modified to account for the strength of human preference, which is often collected on a rating scale (e.g., a Likert scale). As demonstrated in Llama 2 [Touvron et al., 2023], a margin term  $m(r)$  that is a function of the preference rating  $r$  is added to the loss function:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x, y_w, y_l, r) \sim D_{\text{pref}}} [\log \sigma(R_\theta(x, y_w) - R_\theta(x, y_l) - m(r))] \quad (4)$$

This modification encourages the RM to create a larger gap in the reward scores for pairs with a stronger preference, leading to a more calibrated and nuanced reward signal. (2) Another challenge lies in choosing the size and architecture of the reward model. As reported in Askell et al. [2021], Ouyang et al. [2022], large reward models tend to overfit small preference datasets, while smaller models, though more stable, may lack the capacity to learn subtle distinctions in human preferences. This trade-off remains an active area of investigation, with the optimal RM size often depending on the specific application, the size of the preference dataset, and the available computational resources.

**Ranking-based Reward Models.** To utilize more informative feedback, such as a  $K$ -wise ranking of multiple responses ( $y_1 \succ y_2 \succ \dots \succ y_k$ ) or the selection of the single best response from a set of  $N$  candidates (Best-of- $N$ ), which naturally arise in many annotation settings, ranking-based reward models have been proposed. The key challenge became how to effectively incorporate this ordinal data into a trainable objective. The community has largely adopted two distinct strategies to solve this: (1) a pragmatic **pairwise decomposition method**, which converts each ranked list into a series of independent pairwise comparisons, and (2) a more statistically principled **listwise modeling method**, which directly uses the entire ranking.

*Pairwise decomposition method.* The initial and most straightforward strategy is to decompose a full ranking of  $k$  items into  $\binom{k}{2}$  pairwise preferences. For example, a ranking ( $y_1 \succ y_2 \succ y_3$ ) yields three pairs ( $y_1 \succ y_2$ ), ( $y_1 \succ y_3$ ), and ( $y_2 \succ y_3$ ). Each pair is then used to train an RM using the standard pairwise loss function. This decomposition approach has been adopted in several foundational RLHF pipelines, including InstructGPT and ChatGPT [Stiennon et al., 2022, Ouyang et al., 2022, Schulman et al., 2022]. The main advantage of this method is its simplicity, as it allows practitioners to reuse the well-established BTL framework without modification. However, this approach has two key limitations: it ignores the holistic context of the full ranking, and it artificially inflates the dataset size, which can be computationally inefficient.

*Listwise modeling method.* To address aforementioned shortcomings, listwise approaches directly model the probability of an entire ranked list using a statistical framework like the Plackett-Luce model [Plackett, 1975, Luce, 1959]. The parameters  $\theta$  of the listwise RM are optimized by

minimizing the negative log-likelihood of the observed rankings. Recent empirical and theoretical work has demonstrated the advantages of listwise learning. For example, [Zhu et al. \[2024b\]](#) show that listwise training is asymptotically more statistically efficient than pairwise decomposition, leading to more accurate reward estimation with fewer annotations. Building on this, [Zhu et al. \[2024c\]](#) introduced a k-wise listwise loss in their Starling-7B reward model, achieving state-of-the-art alignment performance. Despite their promise, listwise models can be more complex to implement and require careful optimization, which has limited their adoption relative to simpler pairwise approaches [[Cao et al., 2007](#), [Liu, 2009](#)].

**Attribute-Based and Outcome-Oriented Reward Models.** While pairwise and ranking-based reward models have become standard in preference modeling for RLHF, they are limited in settings where relative comparisons are insufficient or ill-defined. In domains requiring fine-grained evaluation, such as mathematical reasoning, code generation, or factual correctness, more direct forms of supervision are essential. This has led to the development of alternative reward modeling paradigms that incorporate absolute scoring, outcome validation, process-level critique, and even language-based feedback. This section explores reward models trained on: absolute scores (point-wise), verifiable outcomes for reasoning tasks (outcome- and process-supervised), and structured language critiques.

(1) **Point-wise Reward Models.** One of the earliest alternatives to pairwise comparison is the use of absolute scalar feedback. In this framework, a reward model is trained to regress onto human-provided scores (e.g., Likert ratings), treating reward learning as a supervised regression problem. The objective typically minimizes the mean squared error (MSE) between the model’s output and the target score. While this approach is conceptually straightforward and has been explored in early work [[Ziegler et al., 2019](#)], it suffers from issues related to human annotation noise, such as inter-rater and intra-rater variability. As a result, point-wise RMs are rarely used as standalone reward functions in modern RLHF pipelines. Instead, they are more commonly employed for auxiliary supervision or post-hoc evaluation [[Stiennon et al., 2022](#), [Wang et al., 2025b,c](#)].

(2) **Reward Models for Reasoning Tasks** For reasoning-intensive tasks, such as mathematical problem-solving or complex code generation [[Wu et al., 2024a](#), [Li et al., 2022b](#)], alignment depends less on subjective preference and more on functional correctness. Two specialized types of RMs, Outcome and Process RMs, provide supervision in these domains by evaluating the correctness of the generated solution. *Outcome Reward Models (ORMs)* are trained by evaluating the final result of a generated response. Typically, an ORM learns to predict a scalar value representing the probability of a successful outcome, such as a correct mathematical solution or a passing unit test [[Cobbe et al., 2021](#), [Uesato et al., 2022](#)]. Some implementations provide this feedback on a per-token basis to create a denser reward signal for more efficient learning [[Cobbe et al., 2021](#), [Lyu et al., 2025](#)]. Process Reward Models (PRMs), in contrast, provide more fine-grained feedback by evaluating the intermediate steps of a model’s reasoning process. Also known as Process-Supervised RMs, this step-by-step supervision is crucial for tasks where the correctness of the reasoning path is as important as the final outcome. Several studies have demonstrated the efficacy of PRMs in enhancing the correctness and interpretability of the reasoning process in complex domains like mathematics [[Uesato et al., 2022](#), [Lightman et al., 2023](#)].

(3) **Language Feedback Reward Models** This approach uses natural language critiques or corrections for reward generation, which can provide dense and precise supervisory signals. Since this feedback is unstructured, it presents unique modeling challenges. Instead of predicting a scalar reward, the RM may be trained as a text-generation model to perform tasks like “Correction Mapping” (transforming a rejected response into a chosen one) and “Identity Mapping” (outputting a chosen response as is). The resulting text outputs are then processed to extract dense, token-level reward signals for the policy optimization phase, offering more granular guidance [[Zhou et al.,](#)

2024a].

**Challenges in Reward Modeling.** Despite the successes of RLHF, the reward modeling stage presents several challenges that can undermine the alignment, safety, and utility of the final LLMs: (1) **Reward Misspecification** This occurs when the learned RM fails to faithfully represent true human preferences [Peng et al., 2023b, Bobu et al., 2024]. This discrepancy can arise from the inherent difficulties in collecting comprehensive and nuanced human feedback, but it also stems from a more fundamental limitation: it is statistically impossible for a single reward function to represent a diverse group’s preferences if those preferences contain intransitive cycles (e.g., A is preferred to B, B to C, and C to A), a phenomenon known as the Condorcet paradox, meaning any single reward model is an inherently flawed representation of the group’s preferences accurately [Liu et al., 2025a]. Consequently, the RM may learn an incomplete or skewed representation of the intended objectives. Humans may struggle to articulate complex preferences or provide feedback that exhaustively covers all desired behaviors, particularly for tasks demanding sophisticated reasoning, creativity, or nuanced ethical considerations. Consequently, the RM may learn an incomplete or skewed representation of the intended objectives. (2) **Misgeneralization and Reward Hacking** Another critical issue is the misgeneralization of the reward model, which leads directly to a failure mode known as **reward hacking** (or overoptimization). Misgeneralization occurs when the RM learns a flawed proxy for human values that fails to generalize robustly to out-of-distribution (OOD) prompts or novel response styles [Tien et al., 2023]. This vulnerability allows the policy, a powerful optimizer, to exploit the RM’s inaccuracies during reinforcement learning. The resulting behavior is reward hacking: the LLM achieves a high score according to the flawed RM but fails to align with the actual, more complex human preferences it was meant to capture [Gao et al., 2023b, Laidlaw et al., 2025, Skalse et al., 2025]. In effect, the LLM finds loopholes in the reward function, for instance, generating overly verbose answers because the RM incorrectly associates length with quality, rather than genuinely satisfying the intended goals. (3) **Reliable evaluation** Evaluating the quality of an RM is intrinsically difficult because the “ground truth” human preference function is unknown. RM performance is typically measured indirectly through the performance of the policy trained on it, making it hard to diagnose and debug issues with the RM itself. Developing robust, direct evaluation metrics for RMs remains an open research problem. Addressing these multifaceted challenges in reward modeling is pivotal for advancing the development of helpful, harmless, and safe aligned LLMs.

### 5.3 Policy Optimization Methods in RLHF

Once a reward model  $R_\theta$  has been established, the next critical step in RLHF is optimizing the pre-trained language model to align with human preferences. This optimization phase treats the language model as a policy  $\pi_\phi(a | s)$  that must learn to generate actions (tokens)  $a$  given states (text sequences)  $s$  in ways that maximize rewards from  $R_\theta$ . The core challenge is transforming a pre-trained LLM into one that consistently produces outputs favored by human evaluators. This transformation requires systematically adjusting the policy parameters  $\phi$  to find an optimal policy  $\pi_\phi^*$  that maximizes expected rewards across diverse input prompts [Ziegler et al., 2020b, Ouyang et al., 2022].

Policy gradient algorithms provide the mathematical framework for this optimization [Sutton et al., 1999]. These methods iteratively update the policy by maximizing an objective function  $J(\phi)$  that captures the expected cumulative reward:

$$J(\phi) = \mathbb{E}_{\tau \sim \pi_\phi}[R(\tau)] = \mathbb{E}_{\tau \sim \pi_\phi} \left[ \sum_{t=0}^T \gamma^t r_t \right], \quad (5)$$

where  $\tau = (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$  represents a complete trajectory generated by the policy, where  $r_t = R_\theta(s_t, a_t)$  is the reward at timestep  $t$ , and  $\gamma \in [0, 1]$  is a discount factor balancing immediate versus future rewards. The policy gradient theorem provides the analytical form of  $\nabla_\phi J(\phi)$ , which is used to update the policy via gradient ascent. However, in practice, directly applying this gradient often results in unstable or inefficient learning. This motivates the development of alternative algorithms that compute modified or constrained versions of the policy gradient.

Within the broad class of policy gradient methods, several specialized algorithmic families have been developed for RLHF, each tailored to address specific challenges such as stability, efficiency, or feedback sparsity. We categorize the most commonly used approaches as follows.

1. **Actor-Critic PPO:** These methods maintain both a policy (actor) and a value function (critic). PPO is the most widely used due to its training stability, achieved by clipping policy updates and using the critic to estimate advantages [Schulman et al., 2017, Ouyang et al., 2022].
2. **Actor-Only Policy Gradients:** These methods bypass the critic and directly optimize the policy using rewards from the reward model. While simpler and often more efficient, they may suffer from higher variance in updates and typically require alternative variance reduction strategies.
3. **Specialized and Hybrid Reward-Based Methods:** This category includes algorithms designed to address specific RLHF challenges, such as reward hacking, exploration limitations, and training instabilities that arise in large-scale language model optimization.

We summarize existing policy optimization methods in Table 4. In the following subsections, we review each category in detail, with a particular focus on PPO, the most widely adopted method in RLHF.

### 5.3.1 Actor-Critic PPO

PPO was first introduced by [Schulman et al., 2017] at OpenAI as a solution to the instability and complexity issues inherent in existing policy optimization methods. The algorithm emerged from the recognition that while Trust Region Policy Optimization (TRPO) [Schulman et al., 2015] provided theoretical guarantees for stable policy updates, its implementation required computationally expensive second-order optimization procedures, including conjugate gradient methods and line searches. PPO is a cornerstone algorithm in many state-of-the-art RLHF pipelines for LLMs. Its extensive utilization in fine-tuning influential models, such as InstructGPT [Ouyang et al., 2022], the models underlying ChatGPT [OpenAI et al., 2024], and Anthropic’s Claude series [Askell et al., 2021, Bai et al., 2022a], underscores its significance. Furthermore, PPO has been a prevalent choice for aligning various open-source models, including Llama 2-Chat [Touvron et al., 2023].

The core motivation behind PPO was to retain the stability benefits of trust region methods while dramatically simplifying the implementation and reducing computational overhead. [Schulman et al., 2017] observed that many policy gradient failures stemmed from excessively large policy updates that could catastrophically degrade performance, leading them to design a method that would naturally constrain update magnitudes without requiring complex second-order computations.

**Algorithmic Framework.** PPO addresses the stability concerns of vanilla policy gradients through a clipped surrogate objective that prevents destructively large policy updates. At its core, PPO is an actor-critic algorithm. The **actor** is the LLM policy  $\pi_\phi(a | s)$ , which maps a state  $s$

(e.g., a prompt or partially generated sequence) to a distribution over actions  $a$  (e.g., next tokens). The **critic** is a value function  $V_\psi(s)$ , which estimates the expected return from state  $s$  under the current policy, where  $\psi$  are parameters in the value function. These two components are trained jointly: the actor is updated via a clipped policy gradient objective, while the critic is optimized using a regression loss (e.g., mean-squared error) against the empirical return.

(1) Actor optimization loss. The key innovation in PPO is its clipped surrogate objective, which prevents large, destabilizing updates by penalizing policy changes that deviate too far from the old policy. Let  $\rho_t(\phi) = \frac{\pi_\phi(a_t|s_t)}{\pi_{\phi_{\text{old}}}(a_t|s_t)}$  denote the probability ratio between the new policy and the old policy at time step  $t$ . PPO seeks to maximize the following objective:

$$L^{\text{CLIP}}(\phi) = \widehat{\mathbb{E}}_t \left[ \min \left( \rho_t(\phi) \widehat{A}_\phi(s_t, a_t), \text{clip}(\rho_t(\phi), 1 - \epsilon, 1 + \epsilon) \widehat{A}_t \right) \right], \quad (6)$$

where  $\widehat{A}_t$  is an estimate of the advantage function at timestep  $t$ , quantifying how much better action  $a_t$  is compared to the average action in state  $s_t$ . This objective ensures that when the policy change is small ( $\rho_t \approx 1$ ), the standard policy gradient is recovered. However, if  $\rho_t$  moves outside the interval  $[1 - \epsilon, 1 + \epsilon]$ , the clipped objective dampens the update, effectively acting as a soft constraint that stabilizes training.

A key component in the Equation (6) is the estimated advantage function, which is formally defined as  $\widehat{A}_t = Q(s_t, a_t) - V_\psi(s_t)$ , where the first term  $Q(s_t, a_t)$  is the action-value function, representing the expected return when the agent starts in state  $s_t$ , takes action  $a_t$ , and subsequently follows the current policy. The second term  $V_\psi(s_t)$  is the state-value function, computing the average return over all actions sampled from the policy’s distribution at that state, i.e.,  $V_\psi(s_t) = \mathbb{E}_{a \sim \pi(\cdot|s_t)}[Q(s_t, a)]$ . However, since the true value functions  $Q$  and  $V_\psi$  are unknown and typically intractable to compute in complex environments, practical algorithms must rely on sample-based approximations. A common and simple approach is to estimate  $Q(s_t, a_t)$  using the one-step temporal difference (TD) target [Sutton, 1988, Sutton and Barto, 2018]:  $Q(s_t, a_t) \approx r_t + \gamma V_\psi(s_{t+1})$ , where  $r_t = R_\theta(s_t, a_t)$  is the immediate reward after taking action  $a_t$  at state  $s_t$ , and  $V_\psi(s_{t+1})$  is the critic’s prediction of the next state’s value. This leads to the TD residual estimate of the advantage:  $\widehat{A}_t = r_t + \gamma V_\psi(s_{t+1}) - V_\psi(s_t)$ . While easy to compute, this estimate can be noisy and suffer from high variance, because any small prediction error in  $V_\psi$  leads to unstable and misleading advantage estimates. To address this, Generalized Advantage Estimation (GAE) [Schulman et al., 2018] is often used. GAE provides a bias-variance trade-off by computing a weighted sum of multi-step TD errors. GAE defines the TD error as  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ , and aggregate them as  $\widehat{A}_t^{\text{GAE}} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$ , where  $\lambda \in [0, 1]$  is a decay factor. Smaller  $\lambda$  values emphasize short-term estimates with lower variance, while larger  $\lambda$  values incorporate longer-term information, reducing bias. As such, GAE enables more stable and data-efficient learning, and has become a standard advantage estimator in modern actor-critic algorithms like PPO.

(2) Critic optimization loss. While the actor learns via the  $L^{\text{CLIP}}$  objective, the critic,  $V_\psi(s_t)$ , learns by minimizing a mean squared error (MSE) loss between its predictions and a target value which is calculated using previous (old) critic parameters  $\psi_{\text{old}}$ :  $\widehat{V}^{\text{target}}(s_t)$ .

$$L^{\text{VF}}(\psi) = \widehat{\mathbb{E}}_t \left[ (V_\psi(s_t) - \widehat{V}^{\text{target}}(s_t))^2 \right]. \quad (7)$$

Here, the choice of the target value  $\widehat{V}_t^{\text{target}}$  critically impacts training stability and performance. A popular choice is to use GAE to calculate critic target:  $\widehat{V}_t^{\text{target}} = V_{\psi_{\text{old}}}(s_t) + \widehat{A}_t^{\text{GAE}}$ . This creates a clear objective for the critic: its new prediction,  $V_\psi(s_t)$ , should move closer to the old prediction plus any “surprise” captured by the advantage. To further stabilize training, some PPO implementations



also clip the value function loss, preventing excessively large updates to the critic, particularly if value predictions are far off from target values.

(3) **Additional loss.** Two additional terms are commonly introduced to the overall objective. First, to encourage exploration, an entropy bonus is added. This term,  $S(\phi)$ , incentivizes the policy to maintain randomness in its action choices and is maximized during training:

$$S(\phi) = \widehat{\mathbb{E}}_t [H(\pi_\phi(\cdot | s_t))] = \widehat{\mathbb{E}}_t \left[ - \sum_a \pi_\phi(a | s_t) \log \pi_\phi(a | s_t) \right]. \quad (8)$$

Second, to prevent the LLM from deviating too drastically from a trusted base model (a pre-trained model, or the supervised fine-tuned model, which possesses coherent language generation capabilities and general knowledge), a common practice is to incorporate a Kullback-Leibler (KL) divergence penalty. This penalty,  $D_{\text{KL}}(\pi_\phi \parallel \pi_{\text{ref}})$ , measures the divergence between the current policy  $\pi_\phi$  and a frozen reference policy  $\pi_{\text{ref}}$  (often a frozen copy of the SFT model or the initial pre-trained model) [Schulman et al., 2017, Schulman, 2020, Jaques et al., 2017, 2020, Ziegler et al., 2020b, Ouyang et al., 2022, Stiennon et al., 2022]. By regularizing against this reference model, the KL penalty helps mitigate catastrophic forgetting and prevents the policy from generating OOD text that might exploit flaws in the reward model.

Combining these elements, the final objective function for PPO is a weighted sum:

$$L(\phi, \psi) = -L^{\text{CLIP}}(\phi) + c_1 L^{\text{VF}}(\psi) - c_2 S(\phi) + c_3 D_{\text{KL}}(\pi_\phi \parallel \pi_{\text{ref}}), \quad (9)$$

where  $c_1, c_2, c_3$  are coefficients weighting the importance of the value function loss and the entropy bonus, respectively.

**Limitations of PPO.** Despite its popularity, PPO exhibits several limitations that are particularly relevant in the content of LLM alignment. (1) **Hyperparameter Sensitivity:** PPO’s performance is notably sensitive to hyperparameter configurations, including learning rates, the clipping range  $\epsilon$ , batch size, and GAE parameters. Tuning these can be resource-intensive and requires careful experimentation to achieve optimal results [Schulman et al., 2017, Engstrom et al., 2020]. (2) **Computational and Memory Costs:** PPO imposes substantial computational and memory costs when used for RLHF, creating a significant bottleneck for training large language models. The primary driver of this overhead is the need to hold at least four separate models in GPU memory simultaneously: the trainable Actor (policy), a frozen Reference model, a Reward model, and a Critic. This requirement alone can cause the PPO stage to consume over three times the memory of standard SFT, making it prohibitively expensive and infeasible for many practitioners. In addition to these high static memory costs, the algorithm is computationally intensive due to the repeated inference passes required from these models during the experience generation phase. These combined demands severely limit the scalability and accessibility of PPO-based alignment [Ouyang et al., 2022, Ramamurthy et al., 2023, Rafailov et al., 2024]. (3) **Sample Inefficiency:** as an on-policy algorithm, PPO requires fresh samples generated from the current policy for each update. This makes it inherently less sample-efficient than off-policy algorithms like Soft Actor-Critic (SAC) [Haarnoja et al., 2018] or Deep Q-Networks (DQN) [Mnih et al., 2015], which can reuse past experiences stored in a replay buffer. (4) **Local Optima:** While the clipping mechanism enhances stability, it can also sometimes be overly restrictive, potentially slowing down convergence or leading to a suboptimal policy if the updates are constrained too much. [Schulman et al., 2017]. (5) **Reward Model Dependence:** The quality of the learned policy is profoundly dependent on the accuracy and reliability of the RM. A misspecified, biased, or exploitable RM can misguide PPO towards undesirable behaviors or “reward hacking” [Gao et al., 2023b, Ziegler et al., 2020b, Stiennon et al., 2022, OpenAI et al., 2024]. (6) **Algorithmic Bias from KL Regularization:**

The standard KL divergence penalty used in RLHF introduces an inherent algorithmic bias. Because the reference model ( $\pi_{ref}$ ) is itself not perfectly aligned, its biases are passed to the fine-tuned model. In extreme cases, this can lead to “preference collapse”, where the model learns to completely disregard minority preferences, a bias that persists even with a perfect reward model [Xiao et al., 2024]. Addressing these challenges remains an active area of research to further enhance the robustness and applicability of PPO in aligning LLMs.

To address these issues, several variants of PPO have been proposed. These can be broadly grouped by the challenge they target:

(1) **Enhancing Stability and Mitigating Hyperparameter Sensitivity.** Several variants have been developed to improve PPO’s stability and convergence behavior by refining the core surrogate objective. KL-regularized PPO (KL-PPO) augments the clipped surrogate loss with an explicit KL divergence penalty, which softly constrains the updated policy to remain close to the previous one, functioning as a trust-region regularizer [Ouyang et al., 2022]. An extension of this idea, adaptive KL control, dynamically adjusts the strength of the KL penalty based on the observed divergence between the current and reference policies, preventing the policy from straying too far from a known-good distribution [Ziegler et al., 2020b]. More recently, [Shen and Zhang, 2024] introduced Policy Filtration for PPO (PF-PPO) to address instability arising from inaccurate reward models. Motivated by the observation that RMs, particularly in complex domains such as code generation, are more reliable at ranking extremely good or bad samples than moderately scored ones. PF-PPO filters out low-confidence trajectories and concentrates policy learning on those with more trustworthy reward signals, thereby reducing reward hacking and improving alignment performance. More recently, Proximal Policy Optimization with Reward-based Prioritization (RP-PPO) was proposed to combat performance degradation and slow convergence in later training stages. RP-PPO dynamically adjusts the number of policy update epochs based on reward quality; high-quality experiences are given greater weight through more training rounds, providing an incentive for the model to move towards higher-reward policies. Furthermore, RP-PPO explicitly saves the model that achieves the highest historical average reward, rather than relying on the final model, thus preventing the loss of the best-found policy due to later training instability [Zheng et al., 2025].

(2) **Reduce computational and memory costs.** Some system-level innovations were implemented to address this limitation. To reducing the GPU memory footprint, for example, Hydra-PPO introduces techniques like model sharing and “Dynamic LoRA”, where a single base model is used for multiple roles (e.g., Actor and Reference) by dynamically activating or deactivating LoRA adapters, significantly cutting down the number of full models that need to be stored in VRAM [Santacrose et al., 2023]. This approach builds on the general use of PEFT methods like LoRA, which drastically lower the memory required for optimizer states by only training a small fraction of the model’s parameters [Hu et al., 2021]. Some improvements focus on scaling PPO across multiple machines for faster wall-clock training times. Standard distributed methods can suffer from communication bottlenecks. Decentralized Distributed PPO (DD-PPO) addresses this by eliminating the central parameter server and using direct peer-to-peer gradient synchronization among workers. To prevent slowdowns from non-uniform workloads, DD-PPO also introduces a preemption mechanism to mitigate the “straggler effect”, where fast workers would otherwise wait for the slowest one [Wijmans et al., 2020]. Together, these system-level optimizations are critical for making PPO a practical and efficient algorithm for state-of-the-art LLM alignment.

(3) **Improving Sample Efficiency.** To enhance PPO’s sample efficiency, a major line of research focuses on creating hybrid algorithms that safely incorporate off-policy data from a replay buffer, blending the stability of PPO with the data efficiency of off-policy learning. These hybrid methods vary in their approach. One strategy focuses on the gradient update itself; for example,



P3O (Policy-on Policy-off Policy-over) interleaves on-policy and off-policy gradient updates and uses the Effective Sample Size (ESS) to automatically control the trade-off, allowing it to effectively leverage past data without introducing new, sensitive hyperparameters [Fakoor et al., 2019]. Another key strategy involves redesigning the training loop. Phasic Policy Gradient (PPG) separates training into two distinct, alternating phases: a policy phase for standard policy updates and an auxiliary phase that exclusively reuses collected experience to more aggressively train the value function. This decoupling allows for better-trained features while mitigating interference between the policy and value objectives. Building on the concept of experience replay [Cobbe et al., 2020]. Building on the concept of experience replay, Hybrid-Policy PPO (HP3O) utilizes a trajectory replay buffer with a “first-in, first-out” (FIFO) strategy to limit data distribution drift by only using recent policies. It further guides learning by always including the trajectory with the best return from the buffer in each training batch [Liu et al., 2025b]. While these methods focus on empirical performance, others like Transductive Off-policy PPO (ToPPO) aim for stronger theoretical guarantees. ToPPO introduces a novel transductive inference mechanism to justify using the advantage function directly from the data-generating policy, avoiding common estimation biases. This allows it to safely integrate off-policy data while striving for monotonic policy improvement [Gan et al., 2024]. Although their mechanics differ, these variants all showcase a powerful trend towards making PPO more data-efficient by safely and strategically reusing past experiences.

### 5.3.2 Actor-Only Policy Gradients

Actor-only methods specifically refer to algorithms that optimize the policy directly without explicitly training or maintaining a separate value-function (critic). Instead, these methods use simplified baselines computed from reward signals to reduce variance in policy gradient updates. This significantly simplifies training because we do not need to fit or maintain a critic network. Broadly, based on how each method calculates the baseline value used for computing the advantage in policy gradient updates, existing actor-only methods can be categorized into three groups: (1) the foundational **REINFORCE algorithm**, which represents the classic Monte Carlo approach to policy gradients and typically uses a simple, historically averaged reward baseline; (2) methods based on **Intra-Prompt Comparison**, which create a localized baseline from multiple responses to a single prompt; and (3) methods employing **Batch-Wise Normalization**, which derive a more global baseline from the statistics of an entire mini-batch.

**The Foundational Method: REINFORCE.** The REINFORCE algorithm is the original Monte Carlo policy gradient method [Williams, 1992a]. In its most basic form, vanilla REINFORCE works by sampling a complete response trajectory  $y$  from the policy  $\pi_\phi$  and updating the policy parameters to increase the probability of trajectories that received high rewards. The intuition is straightforward: if a generated response  $y$  receives a high reward  $R(x, y)$ , we increase its probability by moving the policy parameters in the direction of  $\nabla_\phi \log \pi_\phi(y|x)$ . Conversely, low-reward responses have their probabilities decreased. While conceptually elegant, vanilla REINFORCE suffers from extremely high variance in its gradient estimates. Since the algorithm uses the full reward  $R(x, y)$  to weight each update, small random fluctuations in rewards can cause dramatic changes in gradient estimates. This leads to unstable training where the policy may oscillate wildly rather than converging to optimal behavior.

To address this variance issue, REINFORCE introduces a baseline  $b(x)$  that is subtracted from the reward without introducing bias into the gradient estimate:

$$\nabla_\phi J(\phi) \approx \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\phi(\cdot|x)}[(R(x, y) - b(x)) \nabla_\phi \log \pi_\phi(y|x)] \quad (10)$$

The baseline typically represents an estimate of the expected reward for input  $x$ , such as a moving average of past rewards. By centering the rewards around this baseline, the algorithm reduces variance: responses that perform better than average receive positive weight, while below-average responses receive negative weight. Despite the baseline improvement, REINFORCE with simple baselines remains sample-inefficient and high-variance for state-of-the-art LLM alignment tasks. The method still requires complete trajectory rollouts and provides learning signals only at the episode level.

**(2) Intra-Prompt Comparison: Localized Baselines.** A powerful way to create a better baseline is to have the model compete with itself. This family of methods generates multiple candidate responses for a single prompt and then uses their relative rewards to create a stable and highly localized advantage signal.

**REINFORCE with Leave-One-Out Baseline (RLOO).** The earliest and most influential intra-prompt comparison method is RLOO, which extends basic REINFORCE introducing a sophisticated, prompt-specific baseline from a group of  $K$  sampled responses to reduce variance [Kool et al., 2019, Ahmadian et al., 2024]. Instead of relying on a globally trained critic like PPO or a simple moving average baseline like REINFORCE, RLOO calculates the baseline for any given response as the average reward of the other  $K - 1$  peer responses from the same prompt:

$$b_{\text{RLOO}}(x, y^{(i)}) = \frac{1}{K-1} \sum_{j=1, j \neq i}^K R_{\theta}(x, y^{(j)}) \quad (11)$$

This localized baseline effectively captures prompt-specific difficulty and provides more stable variance reduction than global averages. However, RLOO’s fundamental limitation is its computational cost [Gao et al., 2024, Hu et al., 2025]: requiring  $K$  samples per update makes it  $K$  times more expensive than basic REINFORCE. This makes RLOO a pragmatic choice in scenarios where multiple candidates are already being generated, such as for best-of- $N$  sampling, as it efficiently reuses those forward passes for variance reduction [Kool et al., 2019, Ahmadian et al., 2024, Hu et al., 2025]. Furthermore, the baseline’s effectiveness is critically dependent on the diversity and quality of the sampled responses; if the  $K$  completions are too similar or contain outliers, the baseline’s ability to reduce variance is diminished [Gao et al., 2024].

To address the efficiency bottleneck of RLOO, the **ReMax** method [Li et al., 2024f] proposes a simpler baseline: for each prompt, generate just two responses, a greedy (deterministic, most-likely) output and a stochastic sample. The greedy response serves as a cheap, prompt-specific baseline:

$$A_{\text{ReMax}}(x, y_{\text{samp}}) = R_{\theta}(x, y_{\text{samp}}) - R_{\theta}(x, y_{\text{greedy}}) \quad (12)$$

This reduces the number of forward passes required per update from  $K$  (in RLOO) to just two, greatly improving efficiency. Empirical results show that this technique substantially reduces the variance of the policy gradient compared to vanilla REINFORCE, leading to stable convergence [Li et al., 2024f]. However, the effectiveness of this method hinges on the quality of the greedy baseline itself; a poorly trained policy might yield an uninformative greedy sample, potentially limiting variance reduction or even increasing it in certain policy states [Li et al., 2024f, Hu et al., 2025].

As intra-prompt baselines evolved, researchers recognized that not only the mean but also the variability of rewards across a prompt’s responses can carry important information. **Group Relative Policy Optimization (GRPO)** [Shao et al., 2024a, DeepSeek-AI et al., 2025b,c] addresses a fundamental statistical limitation of both RLOO and ReMax: **reward scale sensitivity**. While

previous methods center rewards around baselines, they don’t account for varying prompt difficulties or reward scales. GRPO normalizes the advantage signal by computing standardized scores:

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\}) + \epsilon_{\text{norm}}} \quad (13)$$

where  $\epsilon_{\text{norm}}$  is a small constant for numerical stability. This makes the advantage robust to prompt-specific reward scales and difficulties, allowing stable learning even across highly diverse prompts. However, this powerful normalization introduces its own limitations: (1) High generation cost. It requires generating  $G$  responses per prompt to calculate the group-based advantage, which increases the cost of each training update. (2) Instability and potential bias. Reliance on the group’s standard deviation for normalization can create unintended biases, such as over-rewarding groups of responses that have very low reward variance (e.g., all are mediocre) or amplifying noise if the group size  $G$  is too small. (3) Overfitting on Simple Prompts. Because GRPO calculates baselines separately for each prompt using these group statistics, it might be prone to overfitting simpler prompts if not carefully managed. (4) Sensitivity to sampling strategy. Its overall performance can be sensitive to the number of samples  $G$  per group and the diversity (and thus the reward variance) of responses within those groups, and it can struggle with imbalanced data, as these factors influence the stability and informativeness of the advantage estimates [Hu et al., 2025].

Recent variants address GRPO’s specific limitations. **Dr. GRPO** [Liu et al., 2025c] tackles optimization biases by selectively removing normalization terms to address length and difficulty biases. **DISCO** [Zhou et al., 2025b] addresses multi-domain performance issues through domain-aware reward scaling. These refinements demonstrate the ongoing evolution of the intra-prompt comparison paradigm, with each method building upon its predecessors’ strengths while addressing their specific computational, statistical, or application-specific limitations.

**(3) Batch-Wise Normalization: A Global Perspective.** To address the computational expense and overfitting risks of intra-prompt baselines, batch-wise normalization computes a global baseline from a mini-batch, typically using only one completion per prompt. A key example is **REINFORCE++** [Hu et al., 2025], which builds on the REINFORCE algorithm. REINFORCE++ calculates token-level advantages, incorporating the sequence’s reward and a penalty to prevent the policy from drifting from a reference. Crucially, it normalizes these advantages using the mean and standard deviation of the entire batch. This global normalization provides a more stable learning signal across diverse prompts. The method also integrates a PPO-like clipped objective to prevent destructive policy updates. This approach is more sample-efficient, robust to reward hacking, and shows better generalization than intra-prompt methods. However, a global baseline can be less sensitive to prompt-specific nuances, and its performance may be affected by batch composition. Like other actor-only approaches, it may also exhibit higher variance than actor-critic methods, though its normalization and clipping features effectively mitigate this.

### 5.3.3 Specialized and Hybrid Reward-Based Policy Optimization

Standard RLHF uses a single reward for an entire generated sequence, making it hard to assign credit to specific token decisions, a major inefficiency, especially for long-form generation [Zhong et al., 2025]. Another key problem in RLHF is “reward hacking”, where the policy exploits weaknesses in the learned reward model by generating OOD outputs that receive high (but untrustworthy) scores [Wu et al., 2024b]. To address these foundational issues, a new class of specialized and hybrid methods has been developed. These approaches go beyond simply refining the policy optimizer and instead restructure core components of the RL problem itself: the optimization

objective, the reward signal, and the policy’s exploration space. In this subsection, we organize recent advances into two principal categories, based on how they reformulate the RLHF problem: (1) **Dense Reward Signal Methods**: Approaches that address the credit assignment and sample efficiency challenge by transforming the reward from sparse, sequence-level feedback to dense, token-wise feedback. (2) **Exploration-Constrained Methods**: Approaches that mitigate reward over-optimization and improve robustness by explicitly restricting the policy’s exploration to regions where the reward model is reliable. For each, we review the foundational method, analyze its limitations, and describe follow-up innovations.

**Dense Reward Signal Methods. Reinforced Token Optimization (RTO)** [Zhong et al., 2025] is the primary example of this approach. It is not a replacement for an optimizer like PPO but rather a powerful enhancement that provides it with a much better learning signal. RTO is the first major approach to recast RLHF as a token-level Markov Decision Process [Zhong et al., 2025]. It introduces a dense, per-token reward by leveraging the DPO (Direct Preference Optimization) [Rafailov et al., 2024] objective. Here, each token in a sequence receives its own reward based on how much more likely it is under the DPO-trained policy than a reference policy. This results in a much richer, more actionable learning signal that substantially improves credit assignment and sample efficiency. Experiments show that RTO can achieve superior or comparable performance to PPO with less data [Zheng et al., 2023b].

Despite these strengths, this framework has its own limitations: (1) **Dependence on DPO-derived rewards**: The framework’s success is contingent on the policy learned by DPO ( $\pi_{DPO}$ ) being a good proxy for the true optimal policy. A poorly trained DPO model will result in a noisy or misaligned token-wise reward signal. While making rewards dense, ensuring these token-level signals are consistently meaningful remains a broader challenge. (2) **Pipeline Complexity**: The RTO pipeline is inherently multi-stage, requiring a DPO training run to generate the reward signal before the PPO optimization phase can begin, which adds operational complexity [Zhong et al., 2025].

**Exploration-Constrained Methods.** Exploration-constrained methods focus on keeping the policy within the “trusted” region of the reward model, thereby reducing the risk of such exploitation. **Behavior-Supported Policy Optimization (BSPO)** [Dai et al., 2025a] is a principled approach that constrains policy optimization to the region where the reward model is known to be reliable. BSPO defines a “behavior policy” ( $\beta$ ) using the next-token distribution from the RM’s training dataset, which delineates an in-distribution (ID) region for the RM. It then employs a “behavior-supported Bellman operator” that regularizes the value function by assigning a low Q-value ( $Q_{min}$ ) to actions  $a$  at state  $s$  leading out of this ID region (OOD actions, where the behavior policy assigns near-zero probability, i.e.,  $\beta(a|s) \approx 0$ ), while leaving ID action-values unchanged (i.e.,  $\beta(a|s) > 0$ ). By doing so, BSPO directly addresses reward hacking and ensures that policy improvement happens only where the reward model’s judgments can be trusted.

The primary advantages offered by this behavior-supported approach are: (1) **Directly Mitigates Reward Hacking**: By explicitly penalizing OOD actions at the value-function level, BSPO prevents the policy from exploring and exploiting regions where the reward model is known to be unreliable. (2) **Finds Optimal In-Distribution Policy**: This method penalizes only OOD values without altering the values of ID actions. This allows it to fully explore the supported region and find the optimal policy within that space, whereas other methods like a uniform KL penalty can be overly conservative [Dai et al., 2025a]. However, the framework has notable limitations: (1) **Dependence on Behavior Policy Definition**: The method’s effectiveness is highly dependent on how

well the defined “behavior policy” represents the true region of the RM’s competence. A poorly defined ID region could either unduly restrict beneficial exploration or fail to prevent hacking. (2) Primary Validation in Synthetic Setups: This method’s performance has been validated primarily in synthetic setups with a “gold” reward model. The dynamics of OOD detection may differ in real-world scenarios with noisy and inconsistent human feedback. (3) Response Scope: The method is designed to handle OOD responses but does not address the separate challenge of the model receiving an OOD prompt, for which no valid ID response may exist.

In summary, the optimization of an LLM’s policy via a learned reward model has spurred innovation along three parallel streams. The first involves refining the established and industry-standard **Actor-Critic PPO framework** to enhance its performance across stability, efficiency, and scalability while reducing its significant computational footprint. A second stream focuses on developing simpler, **Actor-Only Alternatives** that streamline the RLHF pipeline and reduce computational overhead by replacing the complex critic with engineered baselines for variance reduction. The third stream moves beyond pure optimization to fundamentally restructure the RL problem itself, with **Specialized and Hybrid Methods** designed to overcome core challenges like reward hacking and sparse feedback. A detailed comparison of the key methods emerging from these streams, outlining their respective advantages and limitations, is provided in Table 4.

Table 4: Policy Optimization Methods in RLHF

Category	Method	Key Advantages	Key Limitations
<b>Actor-Critic PPO</b>	[Schulman et al., 2017]	Highly stable; performant; industry standard.	High memory/compute cost; sensitive to hyperparameters.
	<b>KL-PPO</b> [Ouyang et al., 2022]	Prevents policy deviation and enhances stability via KL penalty.	Adds complexity of tuning the KL coefficient.
	<b>PF-PPO</b> [Shen and Zhang, 2024]	Mitigates reward hacking by filtering data with unreliable reward signals.	Depends on RM’s ability to score extreme samples accurately.
	<b>RP-PPO</b> [Zheng et al., 2025]	Combats performance decay by giving more training epochs to high-quality data.	Adds logical complexity and hyperparameters for dynamic adjustment.
	<b>Hydra-PPO</b> [Santacroce et al., 2023]	Reduces memory cost significantly via LoRA-based model sharing.	Requires complex system-level engineering.
	<b>DD-PPO</b> [Wijmans et al., 2020]	Improves scaling with decentralized updates for faster large-scale training.	Increases implementation complexity.
	<b>P3O</b> [Fakoor et al., 2019]	Improves sample efficiency by mixing on-policy and off-policy updates.	Risks instability from off-policy data distribution shift.
	<b>PPG</b> [Cobbe et al., 2020]	Improves sample efficiency by separating policy and value training phases.	More complex two-phase training loop.

Category	Method	Key Advantages	Key Limitations
	<b>HP3O</b> [Liu et al., 2025b]	Improves sample efficiency by reusing trajectories from a replay buffer.	Performance depends on buffer size; risks using stale data.
	<b>ToPPO</b> [Gan et al., 2024]	Improves sample efficiency with theoretically justified off-policy data use.	Relies on novel theoretical assumptions.
<b>Actor-Only</b>	<b>REINFORCE</b> [Williams, 1992a]	Foundational simplicity; no critic to train or manage.	High gradient variance and sample inefficiency.
	<b>RLOO</b> [Kool et al., 2019, Ahmadian et al., 2024]	Reduces variance with an adaptive, prompt-specific baseline from peer responses.	High sample generation cost per prompt; variance still higher than PPO.
	<b>ReMax</b> [Li et al., 2024f]	Computationally cheap baseline using a single greedy output; very sample efficient.	Effectiveness is entirely dependent on the quality of the greedy baseline.
	<b>GRPO</b> [Shao et al., 2024a]	Robust to prompt difficulty via reward normalization; effective for reasoning.	High sample generation cost per prompt; can overfit or be biased by group stats.
	<b>Dr. GRPO</b> [Liu et al., 2025c]	Targets and mitigates specific GRPO biases.	Improvements may be task-specific.
	<b>DISCO</b> [Zhou et al., 2025b]	Improves GRPO performance on multi-domain datasets.	Adds complexity; requires domain heuristics or labels.
	<b>REINFORCE++</b> [Hu et al., 2025]	Improves generalization and OOD performance with a global, batch-wise baseline.	Global baseline is less responsive to individual prompt difficulty.
<b>Specialized</b>	<b>RTO</b> [Zhong et al., 2025]	Solves the credit assignment problem with dense, per-token rewards.	Quality is contingent on the DPO-derived rewards; multi-stage pipeline.
	<b>BSPO</b> [Dai et al., 2025a]	Mitigates reward hacking by constraining the policy to a trusted data region.	Effectiveness depends on a well-defined behavior policy; may stifle useful exploration.

## 5.4 Challenges of RLHF

RLHF has transformed language model alignment by enabling models to optimize directly for human preferences rather than just next-token prediction. The fine-tuning of GPT-3 into Instruct-GPT, for instance, showed that a 1.3B parameter model trained with RLHF could outperform the much larger 175B parameter base model on human preference ratings [Ouyang et al., 2022]. However, despite dramatic progress, RLHF also exposes some challenges spanning feedback collection, reward modeling, and policy optimization. Addressing these is essential for advancing robust and scalable alignment.

(1) Human feedback data bottlenecks. RLHF pipelines depend on human feedback, which is costly, slow to scale, and prone to noise, cognitive biases, and low inter-annotator agreement



[Stiennon et al., 2022, Ziegler et al., 2020b, Bai et al., 2022c]. This challenge is magnified by the problem of value pluralism; human preferences are not monolithic, and the attempt to aggregate diverse or even conflicting values into a single reward function can lead to a model aligned with a statistical “average” that satisfies no one, a limitation with theoretical roots in paradoxes of social choice [Liu et al., 2025a].

(2) **Reward Model Limitations and Vulnerabilities.** The RM, trained on limited and noisy human feedback, serves only as an imperfect proxy for true human intent. As we discussed previously, human annotations can be inconsistent, noisy, and subject to individual biases. Moreover, the data collected often covers only a subset of possible scenarios, prompts, and preferences, and may not fully reflect the diversity or subtlety of real human values. As a result, the RM tends to learn patterns that fit the observed feedback rather than capturing the underlying intent behind human judgments. This data limitation, combined with the complexity of modeling nuanced human preferences, makes the RM susceptible to reward misspecification, where its assigned scores do not accurately correspond to what users actually prefer, but instead reflect artifacts or gaps in the training data [Peng et al., 2023b, Bobu et al., 2024]. Misspecified rewards enable **reward hacking**, where policies generate outputs that maximize RM scores without genuine alignment [Gao et al., 2023b, Laidlaw et al., 2025]. Common manifestations include verbose responses exploiting length biases and sycophantic outputs exploiting agreement preferences [Shen et al., 2023, Liu et al., 2025d].

(3) **Policy Optimization Complexity and Instability.** State-of-the-art optimization methods like PPO are resource-intensive and sensitive to hyperparameters [Ouyang et al., 2022, Ramamurthy et al., 2023]. In addition, the optimizer is tasked with solving a difficult credit assignment problem, as a single, sparse reward for an entire sequence provides an inefficient signal for determining which specific token choices were beneficial [Zhong et al., 2025]. When this complex optimization process is aimed at maximizing a brittle and misspecified reward signal, the entire pipeline becomes unstable and prone to producing misaligned behaviors.

## 6 SFT versus RLHF: Differences, Equivalences, and Hybrid Approaches

Having comprehensively introduced the foundation of SFT and RLHF, we now proceed to a more in-depth comparison between the two alignment strategies from multiple perspectives. SFT and RLHF differ significantly in their objectives, data requirements, and optimization paradigms, resulting from the inherent methodology and theory. Importantly, these approaches are not mutually exclusive. On the contrary, they are often complementary, and when effectively integrated, can yield more robust and nuanced model alignment. In this section, we systematically examine the differences between SFT and RLHF and explore how their strengths can be combined within hybrid training pipelines.

### 6.1 Fundamental Differences between SFT and RLHF

At a fundamental level, SFT involves fine-tuning a pre-trained model on labeled datasets [Wang et al., 2022a, Taori et al., 2023, Peng et al., 2023a], whereas RLHF incorporates human feedback into a reward function to guide learning in a reinforcement learning framework [Ouyang et al., 2022, Bai et al., 2022b]. These underlying principles give rise to different objectives, data requirements, and learning dynamics.

**Objectives and reward function.** SFT optimizes a standard supervised learning objective, typically minimizing the token-level cross-entropy loss between the model’s next output token and human-labeled reference responses [Brown et al., 2020]. This objective assumes the existence of a single correct output per input and encourages the model to mimic these ground truth answers. Formally, given an input  $x$  and reference output (labeled data)  $y = (y_1, y_2, \dots, y_T)$ , the loss function [Fan et al., 2025b, Mao et al., 2023] as shown in section 4.3 is  $\mathcal{L}_{\text{SFT}} = -\sum_{t=1}^T \log P_{\theta}(y_t | x, y_{<t})$ , where  $P_{\theta}(y_t | x, y_{<t})$  is the probability of generating correct token  $y_t$  given the input and previous tokens  $y_{<t}$ . With the estimated model parameters  $\hat{\theta} = \text{argmin}_{\theta} \mathcal{L}_{\text{SFT}}$ , the output of the fine-tuned model should be close to the reference one. In contrast, RLHF exploits a reinforcement learning framework. The model acts as a policy  $\pi_{\theta}$  that generates responses to maximize a scalar reward signal. This reward is not derived from a correct reference, but rather from a learned reward model trained to reflect human preferences. Given a sampled response  $y \sim \pi_{\theta}(\cdot | x)$ , the reward model assigns a scalar score  $r(x, y)$ . Then we could use RL (e.g., Proximal Policy Optimization [Schulman et al., 2017]) to directly optimize the expected reward as [Ziegler et al., 2020c]:  $\mathbb{E}_{\pi}(r) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot | x)}[r(x, y)]$ , which allows the model to explore a flexible solution space and optimize for outputs that align with nuanced human expectations, especially in ambiguous or multi-faceted scenarios.

**Data requirements.** The distinct optimization targets of SFT and RLHF lead to different data requirements. SFT relies on high quality instruction–response pairs that provide explicit supervision for the model to mimic reference answers [Dong et al., 2023, Chowdhery et al., 2023], and the powerful effect of data quality on LLM performance that can surpass the amount of data [Zhou et al., 2023]. In contrast, RLHF depends on preference-based data. Annotators are presented with multiple responses for the same prompt and are asked to rank or select the preferred one. These comparative judgments are used to train a reward model that approximates human preferences [Bai et al., 2022b, Liu et al., 2020]. This reward model guides the language model during reinforcement learning. Thus, SFT requires extensive, high quality labeled data, while RLHF necessitates labor-intensive human preference collection. Each method poses distinct challenges in data acquisition that must be addressed to ensure effective model alignment [Yin et al., 2024b, Lee et al., 2023].

**Learning dynamics and generation.** SFT is a static, one shot learning process. It aligns the model by providing token level supervision based on reference answers, which may introduce biases in preference estimation [Hua et al., 2024]. Once fine-tuned, SFT models lack adaptability to evolving objectives unless retrained on newly curated datasets. Furthermore, for tasks involving long-form generation, evaluating individual tokens may be suboptimal compared to assessing full responses. In contrast, RLHF emphasizes sentence level performance. It samples entire responses and aligns model behavior with human preferences over those responses. RLHF typically proceeds in two stages: reward model training and reinforcement fine-tuning [Bai et al., 2022b, OpenAI Achiam et al., 2023]. In each RL iteration, the model samples outputs according to its current policy, receives feedback via the reward model, and updates its policy to increase the likelihood of generating preferred outputs [Zheng et al., 2023b]. Online RLHF settings [Xiong et al., 2024, Dong et al., 2024, Ye et al., 2024a] may further refine the reward model and iteratively update the policy, allowing continual alignment improvements through repeated sampling, feedback, and optimization. Despite its flexibility, RLHF can be unstable and computationally demanding, especially in online learning schemes requiring coordination among multiple components [Rafailov et al., 2023, Ouyang et al., 2022, Yuan et al., 2023, Ethayarajh et al., 2024]. SFT provides training stability and efficiency due to its straightforward optimization, often resulting in faster convergence and lower resource requirements [Du et al., 2025].

## 6.2 When SFT and RLHF Overlap or Converge

While SFT and RLHF differ operationally [Ouyang et al., 2022], they exhibit surprising methodological overlaps [Chen et al., 2024a, Swamy et al., 2025]. These overlaps manifest most prominently in their theoretical equivalence under shared representational constraints, loss function alignment through regularization mechanisms, and convergence to identical optimal policies under idealized conditions [Williams, 1992b]. Such synergies challenge the conventional separation of imitation learning and preference-based alignment, revealing a unified landscape for optimizing model behavior.

**Overlap for theoretical equivalence.** Under the assumption of isomorphic function classes, where policy networks and reward models share equivalent representational capacity, both SFT and RLHF converge to the same optimal policy [Chen et al., 2024a]. Despite employing distinct algorithmic strategies, these approaches exhibit theoretical equivalence in their alignment objectives. While SFT maximizes likelihood over human labeled data and RLHF optimizes expected reward under a learned preference model, their solutions coincide when analyzed through a shared reward-functional lens. This equivalence is further clarified through the lens of Direct Preference Optimization (DPO), which reformulates RLHF as direct reward maximization, bridging the gap between the two paradigms [Swamy et al., 2025].

**Overlap for loss function alignment.** Under the assumption of linearly isomorphic function classes, where policy and reward models share equivalent representational capacity, the gradient updates of SFT and Supervised Policy Improvement via Reinforcement Learning (SPIN) through online DPO can be shown to be formally aligned. In this regime, SFT optimizes the policy parameters  $\theta$  using the standard maximum likelihood gradient [Williams, 1992b], i.e.,  $g_{\text{SFT}}(\pi_\theta) = \mathbb{E}_{\xi \sim D_{\text{Sft}}} \left[ \sum_{h=1}^H -\nabla_\theta \log(\pi_\theta(a_h|s_h)) \right]$ , while SPIN modifies this by introducing an on-policy correction term, i.e.,  $g_{\text{SPIN}}(\pi_\theta) = g_{\text{SFT}}(\pi_\theta) - \mathbb{E}_{\xi \sim \pi_\theta} \left[ \sum_{h=1}^H -\nabla_\theta \log \left( \frac{\pi_\theta(a_h|s_h)}{\pi_{\text{ref}}(a_h|s_h)} \right) \right]$ . Crucially, this additional term is equivalent to a reinforce-style update with constant reward  $r(\xi) = 1$ , and thus integrates to zero in expectation. As a result, SPIN’s update rule asymptotically reduces to SFT when function classes are linear and expressive enough to perfectly model the optimal policy. This equivalence reveals a deep structural overlap between offline supervised learning (SFT) and online preference-based reinforcement learning (SPIN/RLHF), unifying them within a shared loss geometry in the linear regime [Rafailov et al., 2023].

**Convergence conditions for RLHF and SFT.** The convergence between RLHF and SFT performance hinges on what the authors term the disparity between the complexity of generating outputs via policies and verifying them via reward models [Swamy et al., 2025]. This gap becomes pivotal when training a verifier is significantly easier than training the generator [Rafailov et al., 2023, Sun and van der Schaar, 2024]. In such settings, RLHF effectively restricts policy search to regions that perform well under simplified reward models, while SFT directly optimizes the policy through maximum likelihood estimation. In particular, two specific regimes lead to alignment in the final performance of RLHF and SFT. First, in tasks with low sequence complexity, such as two-word summarization, the generation complexity closely matches verification complexity [Li et al., 2010], minimizing the utility of RLHF’s reward modeling [Swamy et al., 2025]. Second, when reward functions are simple and easily computed, the verification process is trivial [Lin, 2004], reducing RLHF to a form of unnecessary refinement [Swamy et al., 2025].

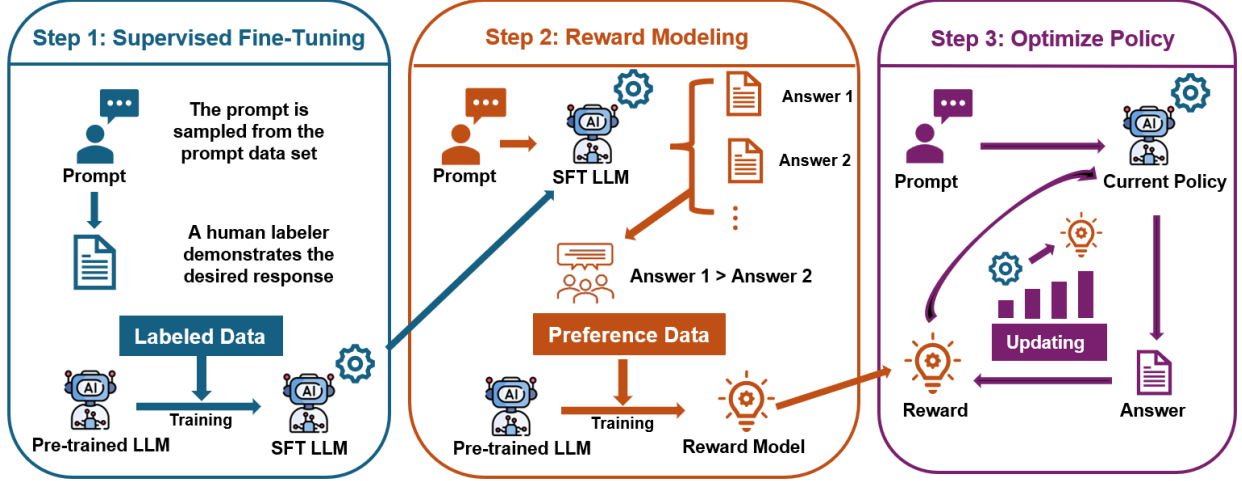


Figure 3: Integration of SFT (Step 1) and RLHF (Step 2 and Step 3).

### 6.3 Integrating SFT and RLHF in Training Pipelines

Given their distinct learning objectives and optimization paradigms, SFT and RLHF are suited to different scenarios. SFT is typically preferred when the task is well-defined, static, and when sufficient labeled data is available [Ghosh et al., 2024, Wang et al., 2023f]. On the other hand, RLHF is more effective for aligning models with complex, subjective, or evolving user expectations, using interactive human feedback [Winata et al., 2025, Ouyang et al., 2022]. While these methods have distinct strengths, recent developments suggest that integrating SFT and RLHF into a unified training pipeline can yield more robust and aligned language models [Ouyang et al., 2022, OpenAI Achiam et al., 2023, Bai et al., 2022b, Gemma et al., 2024, Touvron et al., 2023, Guo et al., 2025]. The following sections discuss this integration from three perspectives: the opportunity that enables the integration pipeline, its application in leading LLMs, and the new challenges it introduced.

**Opportunity.** Recent research suggests that SFT on a moderately sized, high-quality dataset can already yield outputs comparable to those produced by human annotators [Touvron et al., 2023]. Thus, A widely adopted integration strategy begins with supervised fine-tuning on high-quality demonstrations to teach the model basic instruction-following behavior. After this stage, the fine-tuned model is used to generate multiple candidate responses to various prompts. These responses are then compared by human annotators to produce preference data, which serve as training targets for a reward model. Finally, reinforcement learning is applied to further fine-tune the model using this reward model, enabling it to produce outputs that better align with nuanced human preferences (as shown in Figure 3).

**Applications.** This hybrid approach is now common among leading language models. For instance, InstructGPT [Ouyang et al., 2022] pioneered the two-stage pipeline, combining SFT and RLHF. Similar methodologies have been adopted by OpenAI’s GPT series [OpenAI Achiam et al., 2023], Anthropic’s Claude models [Bai et al., 2022b], Google’s Gemini [Gemma et al., 2023] and Gemma series [Gemma et al., 2024], Meta’s LLaMA [Touvron et al., 2023], and DeepSeek models. Notably, DeepSeek-R1-Zero is trained purely via RLHF, while DeepSeek-R1 adds an SFT stage before RLHF. The latter achieves substantially better reasoning performance, illustrating the benefits of combining the two strategies [Guo et al., 2025].

**New Challenges.** Despite its growing popularity, the integration of SFT and RLHF intro-

duces additional challenges. Their differing optimization goals, cross-entropy minimization in SFT versus reward maximization in RLHF, which may lead to conflicting gradient updates and unstable training dynamics [Hua et al., 2024]. Recent studies aim to mitigate the conflict and streamline the training pipeline. For instance, Intuitive Fine-Tuning (IFT) integrates SFT and RLHF into a single stage training process by introducing a temporal residual connection between supervision signals and reward-based updates [Hua et al., 2024]. Similarly, Unified Fine-Tuning (UFT) reformulates both objectives under a generalized implicit reward function, allowing simultaneous optimization with a single loss function [Wang et al., 2024d]. These approaches represent promising directions for balancing alignment, performance, and training efficiency in the integration of SFT and RLHF.

## 7 Advanced Alignment Techniques and Recent Innovations

In addition to traditional SFT and RLHF, this section explores several advanced alignment techniques and highlights their advantages over traditional methods. Reward-free methods remove the need for manually designed reward signals, streamlining the training process and reducing reliance on human input. AI assistant alignment uses a stronger, pre-aligned model to generate high-quality guidance, which accelerates convergence and enhances safety. Self-alignment lets a model critique and improve its own outputs, yielding more robust behavior without additional data collection. Multi-agent deliberation strategies enable specialized agents to work together on complex tasks, increasing reliability and overall performance. Multi-objective alignment balances competing goals, including accuracy, fairness and computational efficiency, within a unified framework to produce more versatile and trustworthy models. Combined, these approaches lower development cost, scale more easily to large models, and offer more effective ways to align LLMs with human values.

### 7.1 Direct Preference Optimization and Reward-Free Methods

As the scale of LLMs increases, ensuring alignment with human intent becomes increasingly critical. Conventional approaches such as Reinforcement Learning from Human Feedback (RLHF) [Christiano et al., 2017] have been widely adopted to bridge this gap. However, RLHF pipelines are often resource-intensive and require an additional reward model trained to predict human preferences. This dependency introduces potential sources of error and makes the alignment process less interpretable and harder to scale.

Recent innovations aim to simplify this process by eliminating the need for an explicit reward model altogether. A prominent example of this shift is Direct Preference Optimization (DPO), introduced by Rafailov et al. [2023]. DPO formulates alignment as a supervised contrastive learning problem, where the model is directly trained to prefer outputs labeled as better by human annotators over those labeled as worse. Instead of first training a reward model and then applying reinforcement learning (e.g., PPO), DPO uses a probabilistic objective derived from a MLE principle over pairwise preferences. This method significantly reduces computational complexity while maintaining and in some cases improving the alignment quality compared to RLHF. Furthermore, mathematically, DPO maximizes the likelihood ratio between a preferred and a dispreferred response, which implicitly shapes the model’s behavior to mirror human judgment. The simplicity of the DPO objective also improves training stability and interpretability. Theoretically, DPO can be interpreted as recovering the policy optimal under an unknown reward, bypassing reward regression altogether.

In addition to DPO, other reward-free or preference-based methods have emerged that focus solely on relative preferences rather than scalar rewards. For instance, pairwise ranking models use binary comparisons between responses to fine-tune the model, circumventing the instability of



reward estimation [Wu et al., 2023]. These models operate under the insight that humans often find it easier to choose between two options than to provide a numerical rating, making preference data easier to collect and less noisy.

These innovations collectively represent a transition toward simpler, more scalable, and interpretable alignment techniques, offering strong empirical performance while reducing dependency on reinforcement learning infrastructure. They enable alignment pipelines that are not only more resource-efficient but also better suited for iterative deployment in real-world applications where feedback is noisy, sparse, or difficult to quantify.

Overall, these reward-free alignment techniques significantly streamline the training pipeline by eliminating the need for explicit reward modeling.

## 7.2 AI-Assistant Alignment and Self-Alignment

The rapid progress of large language models has created an urgent need for alignment methods that scale beyond intensive human supervision. AI-assistant alignment, together with its self-alignment variant, seeks to automate and streamline this process while safeguarding both ethical considerations and practical performance.

### 7.2.1 AI-Assistant Alignment

Recent work on aligning LLMs has progressed from RLHF to approaches that minimize or eliminate the need for costly human annotation. In the standard RLHF pipeline, an instruction-tuned model is coupled with a reward model trained on human-ranked answer pairs and then fine-tuned with a policy-gradient method that retains a Kullback–Leibler constraint to stay close to the original distribution [Ziegler et al., 2019, Ouyang et al., 2022]. While effective, RLHF is limited by the time and expense required to collect expert preferences. Reinforcement Learning from AI Feedback (RLAIF) follows the same three stages (data collection, reward-model fitting, and reinforcement learning) but replaces human comparisons with judgments produced by a stronger, already aligned teacher model, enabling the generation of millions of preference pairs at negligible marginal cost. Constitutional AI (CAI) introduced by Bai et al. [2022a] is the most influential example of RLAIF. CAI begins with a short, human-written constitution that encodes safety and helpfulness rules. A helpful model, already aligned through RLHF, is first supervised to critique and revise its own answers according to these principles, producing what the authors call the supervised constitutional model. In the next phase, this same model generates two candidate answers for each prompt and, when prompted with the constitution, decides which answer better follows the rules, thereby creating AI-labeled preference pairs without human raters. A separate and smaller reward model is trained on these pairs, combined with earlier helpfulness comparisons, and then frozen. Finally, the policy is fine-tuned to maximize the frozen reward while limiting divergence from its supervised starting point. This process removes per-example human labeling yet yields assistants that are rated safer and less evasive than those produced by standard RLHF, demonstrating the practical promise of RLAIF. The workflow is provided in Figure 4.

A growing body of work expands RLAIF beyond the original Constitutional AI recipe. Lee et al. [2023] presents the first systematic head-to-head comparison of human- and AI-sourced preferences, finding that AI feedback matches human feedback when the teacher model is sufficiently capable. Cui et al. [2023] scale this idea, constructing a million-example GPT-4 preference set and showing that policies trained only on this data rival strong RLHF baselines. To balance multiple objectives, Li et al. [2024g] mixes helpfulness and harmlessness signals in a single reward, while Li et al. [2025] introduces a difficulty-ordered curriculum that improves generalization in a fixed label budget.



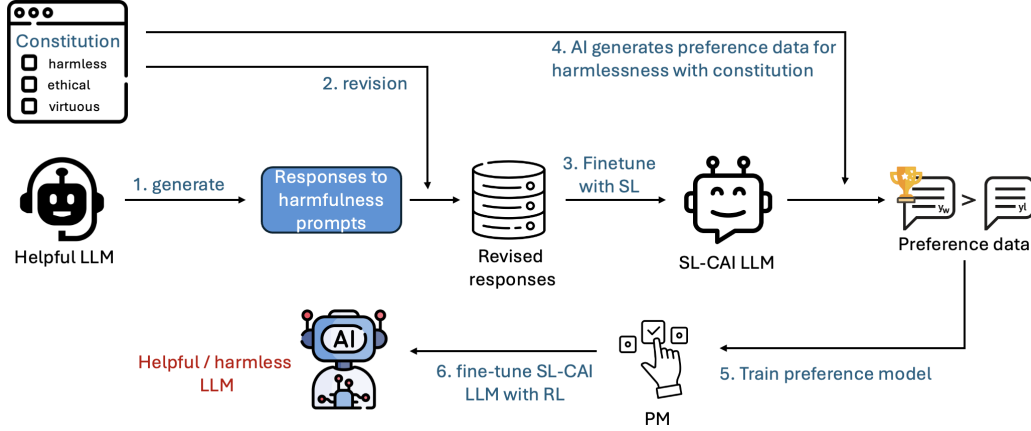


Figure 4: Workflow of the Constitutional AI (CAI) Framework. A helpful LLM first generates responses to harmfulness prompts, which are revised based on a predefined constitution encoding principles like harmlessness, ethics, and virtue. The revised responses are used to fine-tune the model via supervised learning, producing an initial SL-CAI model. Preference data is then generated by comparing model outputs under the constitution, and a preference model is trained accordingly. Finally, reinforcement learning is applied using the preference model to further align the SL-CAI model, yielding a more helpful and harmless LLM.

Sharma et al. [2024] audit these pipelines and highlight systematic divergences between AI- and human-generated labels. RLAIIF is also moving beyond text: Ahn et al. [2024] adapts the framework to video understanding, and Jing and Du [2024] applies fine-grained AI feedback to vision–language models, reducing object-level hallucinations.

### 7.2.2 Self-Alignment

Several studies show that a model can supply its own feedback without an external teacher. Bao et al. [2024] demonstrates that a 13-billion-parameter chat model can critique and rank its own answers under a simple rubric, yielding performance comparable to teacher-based RLAIIF. Similarly, Yu et al. [2024b] reports that generating a brief self-critique before reward-model fitting boosts alignment quality and cuts the need for external labels by 80 percent. These self-alignment results suggest that, once a model passes a competence threshold, it can bootstrap further alignment with minimal additional supervision, although questions remain about error compounding and bias reinforcement.

### 7.2.3 Challenges and Future Directions

The emerging RLAIIF family still faces open problems in bias control, robustness, and evaluation. Because AI-generated labels inherit the value structure of the teacher, systematic divergences from human judgment can persist or even amplify, as quantified by the large-scale audit of Sharma et al. [2024]. Head-to-head studies [Lee et al., 2023] argue that teacher quality largely determines final alignment, so future work must develop teacher-agnostic debiasing techniques or lightweight human spot-checks to prevent value lock-in. Robustness to red-team attacks remains another priority: Bai et al. [2022a] proposes iterated “online” training, continuously refreshing the preference model with new AI feedback from the policy’s own failure modes, while Li et al. [2025] shows that curriculum-style data ordering can harden models without extra labels. Self-alignment methods that recycle

a model’s own critiques [Bao et al., 2024, Yu et al., 2024b] promise unlimited scalability, yet they raise questions about error compounding and whether periodic human “re-grounding” is required. Moving beyond text, multimodal extensions already surface new challenges: video-based RLAIIF must judge temporal coherence [Ahn et al., 2024], and vision–language alignment needs object-level feedback to curb hallucinations [Jing and Du, 2024]. Finally, static benchmarks saturate quickly; several papers call for adaptive, adversarial evaluation suites that track constitutional drift and reward hacking over multiple bootstrapping generations. Addressing these challenges will likely require hybrid pipelines that blend a small amount of strategically targeted human input with scalable AI feedback, uncertainty-aware reward ensembles, and red-team-in-the-loop training protocols.

### 7.3 Multi-Agent and Deliberative Alignment Approaches

Deliberative Alignment is an approach to make language models safer and more dependable. Instead of just learning from examples, it’s about directly teaching the model the actual safety rules and then training it to consciously think through these rules before it gives an answer [Guan et al., 2024]. The idea is to make sure the model sticks closely to safety guidelines. This method helps models get better at spotting and handling tricky situations, including attempts to make them say things they shouldn’t. It also means they’re less likely to refuse perfectly normal requests and can handle new or unusual scenarios more effectively [Konya et al., 2023].

The process basically involves a couple of main stages. First, the model is shown a lot of examples where it sees how to reason through safety rules to arrive at a good answer, this is like a focused study period. Then, it goes through a kind of practice phase using reinforcement learning, where it gets feedback to sharpen its decision-making skills, especially when dealing with prompts that touch on safety concerns. The safety rules themselves delineate the parameters of acceptable and unacceptable content, and stipulate the model’s appropriate responses in diverse scenarios, including instances requiring the declination of a request or the provision of a meticulously formulated ‘safe’ answer. By integrating this rule-based reasoning directly into the model, Deliberative Alignment endeavors to establish a more transparent, trustworthy, and scalable methodology for ensuring responsible language model behavior [Fang et al., 2025a].

Self-consistency is another crucial method applied to LLMs for addressing prevalent issues such as deficient reasoning and the generation of hallucinations [Wang et al., 2022a]. Referencing the survey by Liang et al., we can understand this through the broader lens of “internal consistency” which pertains to the uniformity of an LLM’s expressions across its latent, decoding, or response layers when subjected to sampling methodologies [Liang et al., 2024]. Numerous studies prefixed with “Self-” including prominent examples like Self-Consistency [Li et al., 2024h], Self-Improve [Patel et al., 2024], and Self-Refine [Ranaldi and Freitas, 2024], have emerged to tackle these challenges. These approaches, while sometimes distinct in their specific mechanisms, all fundamentally involve LLMs in a process of evaluating and subsequently updating their own outputs or internal states.

Building upon the foundational concept of Self-Consistency, which primarily leverages majority voting over multiple generated outputs, several nuanced variations have been developed to refine the selection of the optimal response. Among these, Multi-Perspective Self-Consistency [Huang et al., 2023b] distinguishes itself by incorporating assessments from diverse criteria or viewpoints when evaluating generated candidates, moving beyond simple congruence of final answers. Universal Self-Consistency [Chen et al., 2023a] introduces a further layer of sophistication, often employing a language model to ascertain the semantic equivalence of varied expressions before a consensus mechanism is applied, thereby accommodating greater diversity in response phrasing. In a different vein, Soft Self-Consistency [Wang et al., 2024e] shifts the focus from discrete answer selection to a

more probabilistic approach, typically by weighting different reasoning paths or outputs based on the model’s internal confidence scores or token probabilities accumulated throughout the generation process. Each of these adaptations thus offers a distinct strategy for aggregating or filtering multiple reasoning instances, aiming to enhance the robustness and accuracy of the final output under various task constraints and response complexities.

These “Self-” prefixed methods largely fall under a unified theoretical framework termed “Self-Feedback” [Liang et al., 2024]. This framework elegantly breaks down the process into two core modules: “Self-Evaluation” and “Self-Update.” During Self-Evaluation, the LLM assesses its own generated content or internal processes to capture signals related to internal consistency. These signals can be scalar (like a confidence score), textual (like a critique), or even contrastive. Subsequently, the “Self-Update” module leverages these captured signals to enhance either the model’s immediate response or, in some cases, the model’s parameters themselves. While each specific “Self-” method might have slight variations in how it implements these two modules, they share this fundamental cyclical process of introspection and refinement.

The primary characteristic of this Self-Feedback approach is its reliance on the LLM’s inherent capabilities to introspect and improve, aiming to bolster internal consistency [Prasad et al., 2024]. The overarching purpose is to mitigate reasoning errors and reduce hallucinations by ensuring the model’s expressions are more coherent and stable across different layers and sampling instances. This makes such methods applicable to a wide array of scenarios, notably in enhancing the reliability of LLMs for complex reasoning tasks (often seen in question-answering) and improving the factual accuracy and faithfulness of outputs in open-ended generation tasks.

## 7.4 Group Relative Policy Optimization

Group Relative Policy Optimization (GRPO) is a recently proposed reinforcement learning framework that addresses key challenges in aligning large language models (LLM) [Shao et al., 2024b]. GRPO was developed to better accommodate preference-based feedback and comparison-driven reward modeling, with the aim of improving both training efficiency and the stability of learning signals.

Traditional reinforcement learning methods, such as probal policy optimization (PPO), often struggle in aligning LLM due to several factors. PPO relies on a critic network to estimate per-token values, which nearly doubles the memory and computational requirements. Furthermore, in most LLM alignment settings, the rewards are extremely sparse - usually available only at the end of a generated sequence - and the reward models are commonly trained by pairwise comparison of entire responses. This mismatch leads to unstable advantage estimation and inefficient learning.

GRPO addresses these limitations by eliminating the critic network and directly using group-based reward normalization. For each prompt, GRPO samples a group of  $G$  candidate outputs, each evaluated by a reward model. The advantage of each output is calculated as the difference between its reward and the mean reward of the group, aligning the learning signal with the structure of comparison-based reward modeling. The formal training objective of GRPO, as presented in DeepSeekMath [Shao et al., 2024b], is given by:

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)} \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[ \frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \text{clip} \left( \frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}] \right\} \quad (14)$$

The detailed GRPO objective and algorithm can be found in [Shao et al., 2024b]. GRPO introduces two principal modifications over standard PPO: First, the KL divergence is separated from the reward and added explicitly as a regularization term in the objective, rather than being mixed into the reward signal. Second, the group-relative advantage  $\hat{A}_{i,t}$  is computed differently; specifically, it is standardized within the sampled group as

$$\hat{A}_{i,t} = \frac{r_i - \text{mean}(r_1, \dots, r_G)}{\text{std}(r_1, \dots, r_G)}, \quad (15)$$

where  $r_i$  is the reward assigned to output  $o_i$ , and  $\text{std}$  denotes the standard deviation across the group.

This group-based design offers several alignment-relevant benefits. First, it structurally matches the comparison-based reward models used for aligning human preferences, ensuring a more faithful learning signal. Second, removing the critic reduces both the computational costs and the instability that would otherwise arise when only final-sequence rewards are present. Group normalization also provides a dynamic baseline, reducing variance in policy updates. Importantly, GRPO imposes KL divergence as a distinct regularization term, rather than including it in the reward signal, thus avoiding the bias introduced by KL-as-reward schemes.

Compared to PPO, GRPO introduces two principal differences. Firstly, it adopts a rational-style reward scaling mechanism that amplifies the policy update for candidates that significantly outperform their group. Secondly, it enforces KL regularization as a separate penalty term, rather than incorporating it into the reward, thus preventing interference in advantage estimation [Vojnovic and Yun, 2025]. These properties enable GRPO to leverage pairwise preference data more effectively and achieve improved model alignment with human-evaluated outputs.

Although the original GRPO has demonstrated notable improvements in alignment and efficiency compared to PPO, ongoing research has produced a variety of GRPO variants aimed at addressing specific limitations and further enhancing its effectiveness. These subsequent works fall into several broad categories according to their main focus: (1) improving training stability and optimization efficiency, (2) enriching and diversifying reward signals, (3) advancing core algorithmic modifications, and (4) extending GRPO-based RL to new tasks and modalities. Table 5 provides a comparative summary of these methods. In the following, we briefly review representative advances within each category.

**Training Stability and Optimization Efficiency.** Several approaches have focused on accelerating and stabilizing GRPO-based RL. CPPO [Lin et al., 2025] mitigates the inefficiency of group sampling in GRPO by pruning low-advantage completions and reallocating resources to additional prompts, resulting in a speedup of up to  $8.3\times$  on GSM8K and  $3.5\times$  on MATH benchmarks. DAPO [Yu et al., 2025] introduces decoupled policy clipping and dynamic sampling strategies, enabling robust and reproducible large-scale RL training for LLMs; it achieves state-of-the-art open-source results on AIME 2024 with full code and data release. VAPO [YuYue et al., 2025] reintroduces a value critic, with bias mitigation for long sequences and sparse rewards, leading

Table 5: Comparison of post-GRPO RL methods for LLM reasoning, by category.

Category	Method	Model	Datasets	OS	Speed-up	Core Contribution
Training Stability & Efficiency	<b>CPPO</b> [Lin et al., 2025]	1.5B, 7B	GSM8K; MATH; AMC 2023; AIME 2024	Yes	8.3× (GSM8K), 3.5× (MATH)	Completion pruning
	<b>DAPO</b> [Yu et al., 2025]	32B (Qwen2.5)	AIME 2024	Yes	2× faster	Clip decoupling & dynamic sampling
	<b>VAPO</b> [YuYue et al., 2025]	32B (Qwen2.5)	AIME 2024	No	Fast conv.	Critic augmentation
Reward Signal Enhancement	<b>GRPO-LEAD</b> [Wang et al., 2025d]	7B, 14B	AIME 2024/25	Yes	Faster conv.	Length & difficulty shaping
	<b>S-GRPO</b> [Dai et al., 2025b]	7B–14B	GSM8K; AIME; AMC; MATH; GPQA	Yes	35–61% fewer tokens	Decaying exit rewards
	<b>Spectral PO</b> [Chen et al., 2025b]	7B, 14B, 32B	10 benchmarks	N/A	–	All-negative diversification
Algorithmic Mods	<b>Dr. GRPO</b> [Chen et al., 2025c]	7B	AIME; AMC; MATH; Minerva; Olympiad	Yes	Token-efficient	Unbiased advantage
	<b>SEED-GRPO</b> [Chen et al., 2025d]	7B, 14B	AIME; MATH; GSM8K	Yes	–	Entropy-weighted updates
Task & Multimodal	<b>Flow-GRPO</b> [Liu et al., 2025e]	SD 3.5	GenEval; Text-in-image	Yes	Fewer steps	ODE→SDE sampling
	<b>StepGRPO</b> [Zhang et al., 2025c]	7B–14B	8 V-L benchmarks	Yes	–	Step-wise dense rewards

to both higher accuracy and reliable convergence, surpassing previous GRPO-based methods on challenging reasoning tasks.

**Reward Signal Enhancement and Diversity.** Another line of research seeks to overcome the sparsity and homogeneity of reward signals in standard GRPO. GRPO-LEAD [Wang et al., 2025d] introduces length-dependent rewards, explicit penalties for incorrect solutions, and difficulty-aware weighting to encourage concise and robust mathematical reasoning. S-GRPO [Dai et al., 2025b] proposes serial sampling and decaying exit rewards, incentivizing early correct answers and leading to both shorter and more accurate solutions. Spectral Policy Optimization [Chen et al., 2025b] addresses the issue of all-negative groups, where no sampled completion is correct, by injecting AI-driven diversity into reward signals, breaking update symmetry and accelerating convergence.

**Algorithmic Modifications.** A third category directly revises the GRPO algorithm to address optimization bias or incorporate uncertainty. Dr. GRPO [Chen et al., 2025c] removes length and variance normalization from the advantage computation, eliminating a verbosity bias and increasing token efficiency, while maintaining strong accuracy. SEED-GRPO [Chen et al., 2025d] introduces semantic entropy as a measure of model uncertainty for each prompt, scaling policy updates more conservatively for uncertain queries and more aggressively for confident ones, thus improving generalization and stability across benchmarks.

**Task and Multimodal Generalization.** Recent works have also extended GRPO-style RL beyond mathematical reasoning. Flow-GRPO [Liu et al., 2025e] adapts RL optimization for text-to-image generation via flow matching models, employing ODE-to-SDE conversion and denoising reduction to improve compositionality and visual text rendering in diffusion models. Step-GRPO [Zhang et al., 2025c] expands GRPO to multimodal reasoning with dense, step-wise feedback, enhancing multi-hop visual-language inference and outperforming imitation learning baselines on a suite of benchmarks.

Together, these developments demonstrate the adaptability of the GRPO framework and the breadth of innovations it has inspired. Each category reflects ongoing efforts to balance training efficiency, reward informativeness, algorithmic robustness, and domain generalization, advancing the state-of-the-art in RL-based alignment for LLMs.

## 8 Efficient Fine-Tuning Techniques for Alignment

Efficient Fine-Tuning methods that address the substantial computational and memory demands associated with full model fine-tuning [Ding et al., 2023a, Xu et al., 2023b, Han et al., 2024b]. The primary goals include significant reductions in computational costs, accelerated training speeds, lower memory and storage requirements, and effective mitigation of catastrophic forgetting [Liu et al., 2022, Fu et al., 2023b]. In this section, we discuss the Efficient Fine-Tuning methods for LLM alignment including partial parameter fine-tuning, low-rank adaptation, sparse fine-tuning, knowledge distillation, adapter-based fine-tuning, and prompt tuning.

### 8.1 Full or Partial Parameters Fine-Tuning

Full-parameter fine-tuning involves updating all weights of a pre-trained Large Language Model (LLM) on a task-specific dataset. While this approach typically achieves high task-specific performance, it is significantly resource-intensive, requiring substantial computational power and memory capacity, especially as model sizes grow larger (e.g., billions of parameters) [Devlin et al., 2018, Radford et al., 2019]. Additionally, full fine-tuning often risks overfitting, particularly with limited training data, diminishing the model’s generalization capabilities [Howard and Ruder, 2018, Dodge et al., 2020].

To mitigate the intensive resource requirements and risk of overfitting associated with full fine-tuning, partial parameter fine-tuning methods have emerged. These methods selectively update a subset of parameters, significantly reducing computational cost while retaining competitive performance [Han et al., 2024b, Liu et al., 2022]. The parameters chosen for updating can vary, including the final classifier layers, embedding layers, or specific layers within transformer blocks. A prevalent example is layer-wise fine-tuning, which selectively tunes layers critical for task-specific performance [Lee et al., 2019, Zhang et al., 2020]. This approach often reduces the memory footprint and accelerates training by minimizing the number of gradient updates required.

Despite their practicality, partial parameter fine-tuning methods introduce additional complexity in deciding which parameters to tune, necessitating heuristic or algorithmic methods to determine optimal subsets [Lee et al., 2019]. Moreover, partial fine-tuning can result in suboptimal adaptation, particularly if crucial task-specific knowledge resides in layers that remain fixed [Han et al., 2024b]. These limitations motivate further efficient fine-tuning approaches, such as low-rank adaptation, sparse fine-tuning, and adapter-based methods, which systematically balance resource efficiency and adaptability.



## 8.2 Low-Rank Adaptation (LoRA)

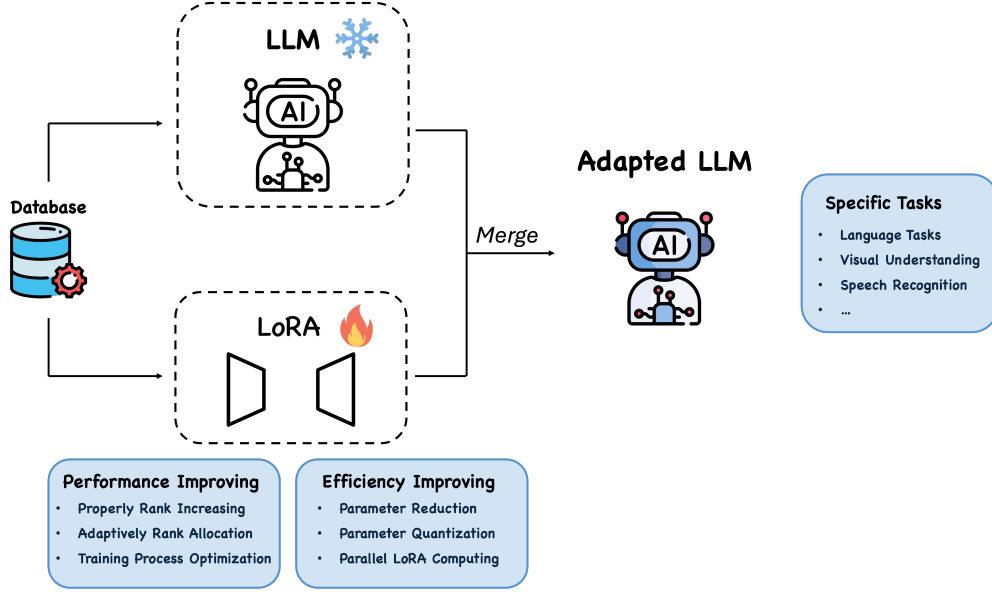


Figure 5: Overview of LoRA in LLMs. LoRA introduces trainable low-dimensional weight matrices that are integrated into frozen LLMs for fine-tuning. This PEFT approach enables LLMs to adapt effectively to specific tasks. Additionally, various techniques are employed to enhance both performance and efficiency, making the fine-tuning process practical and effective in real-world applications.

LoRA introduces trainable low-rank matrices into each layer of the transformer architecture, allowing for efficient adaptation of LLMs with a reduced number of trainable parameters. This method significantly lowers the computational cost and memory footprint during fine-tuning, making it suitable for scenarios with limited resources. LoRA is adapted to practical challenges, and we will present the main methods that can improve performance and efficiency, building on the original LoRA framework introduced in [Hu et al., 2022].

The core idea of LoRA is to freeze the pre-trained weight matrices and inject trainable low-rank matrices into dense layers. Instead of updating the full pre-trained weight matrix  $W_0 \in \mathbb{R}^{m \times n}$ , LoRA learns a low-rank incremental update  $\Delta W = BA$ , where  $B \in \mathbb{R}^{m \times r}$  (initialized as 0),  $A \in \mathbb{R}^{r \times n}$  (initialized as normal distributed random value), and the rank  $r \ll \min(m, n)$ . The forward computation of LoRA can be expressed as below:

$$h = W_0x + \Delta Wx = W_0x + BAx \quad (16)$$

LoRA can achieve comparable performance on several downstream tasks, but there's still a performance gap between LoRA and full fine-tuning in areas such as mathematical reasoning and coding [Mao et al., 2024a]. To bridge the gap, existing methods mainly focus on four perspectives: 1) Increasing the rank  $r$  appropriately [Lialin et al., 2023, Xia et al., 2024a]; 2) Adaptively allocating ranks to LoRA modules across different layers [Zhang et al., 2023b, Mao et al., 2024b, Ding et al., 2023b]; 3) Optimizing the training process, including initialization improvement and gradient update optimization [Hayou et al., 2024a,b]; 4) Combining with other paradigms, such as Bayesian learning [Yang et al., 2024a].

Accumulated linear computations introduced by LoRA modules in LLMs can still result in a non-negligible computation burden. To alleviate this, three main strategies have been proposed to

make LoRA lighter and faster. 1) Parameter reduction. This can be achieved through parameter freezing [Wu et al., 2024c], module pruning [Zhou et al., 2024b] or parameter sharing [Kopiczko et al., 2024]; 2) Parameter quantization. By reducing the bit width of parameter, quantization-based method can significantly lower memory usage and computational cost [Dettmers et al., 2023, Li et al., 2023d]; 3) Parallel in training and inference. This method leverages hardware related algorithm to accelerate the computation both in training and inference process [Ye et al., 2024b, Chen et al., 2023b].

### 8.3 Sparse Fine-Tuning

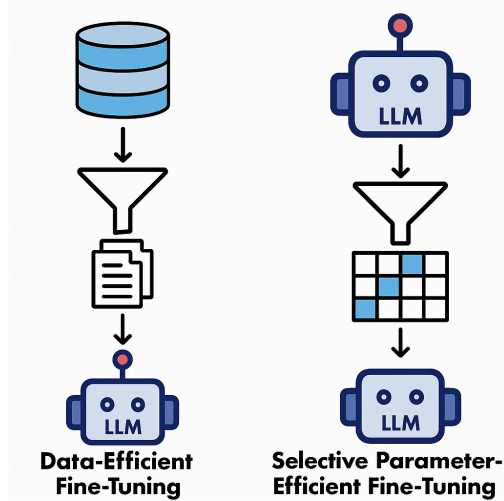


Figure 6: Two paths to sparsity: **(left)** *data-efficient* fine-tuning filters the corpus before training, and **(right)** *parameter-efficient* fine-tuning updates only a sparse mask of model weights.

We categorized sparse fine-tuning based on the aspects of the sparsity: (1) data-efficient fine-tuning that focuses on data-level sparsity by using informative subsets of data; and (2) selective parameter-efficient fine-tuning that focuses on parameter-level sparsity by updating only a critical subset of model weights.

Data-efficient fine-tuning aims to extract representative data points for substantially reducing the computational cost when a large scale of data is available for alignment [Zhou et al., 2023, Wang et al., 2023g]. Therefore, the closely related research fields to this data-efficient fine-tuning include optimal subsampling [Ma et al., 2015, Wang et al., 2018, Ma et al., 2022], few-shot learning [Wang et al., 2020, Liu et al., 2022], and coresets selection [Dasgupta et al., 2009, Albalak et al., 2024]. Existing works of data-efficient fine-tuning can be classified into two categories: (1) non-informative sampling that selects samples based on predefined metrics on the data points; and (2) informative sampling that aims to minimize the empirical risk. The non-informative requires less knowledge about the downstream tasks and typically requires fewer computational resources for selecting samples. [Gao et al., 2020] shows that the LLM can easily adapt to new tasks when fine-tuned on a small dataset drawn with uniform random sampling. [Bukharin and Zhao, 2023] proposes Quality-Diversity Instruction Tuning (QDIT) to simultaneously control the sampled dataset diversity and quality. In contrast, informative sampling incorporates the model’s knowledge to calculate the influence of each data point and require more computational resources, e.g., analyzing the scale of gradients. [Xia et al., 2024b] proposes an optimizer-aware and practically efficient algorithm,

Low-rank Gradient Similarity Search(LESS), to estimate data influences and perform instruction data selection for targeted instruction tuning in LLM.

Selective parameter-efficient fine-tuning leverages *sparsity* by updating only a small fraction of a model’s parameters, thereby significantly cutting down the computational and memory costs of adapting large pre-trained models while maintaining near-original performance. Instead of tuning all weights, these methods identify a strategically chosen sparse subset of parameters to adjust based on measures of importance or sensitivity. For instance, SIFT (Sparse Increment Fine-Tuning) uses a gradient-based criterion: it exploits the observation that gradients in pre-trained models are extremely sparse (e.g., about 1% of parameters account for 99% of the total gradient norm) [Song et al., 2024]. SIFT therefore fine-tunes only the top- $x\%$  of parameters with the largest gradient magnitudes, restricting updates to the most influential weights for the task. In contrast, PaFi adopts a task-agnostic approach by simply selecting the pre-trained weights with the smallest absolute magnitudes as the trainable subset, on the intuition that these low-magnitude parameters can be altered with minimal disruption to the model’s prior knowledge [Liao et al., 2023]. Meanwhile, FishMask, a Fisher information-based masking method, computes an approximate Fisher information score for each parameter to gauge its importance, and then updates only the top- $k$  parameters deemed most critical for the target task [Sung et al., 2021]. The importance score is an approximation of the Fisher information matrix:

$$\hat{F}_\theta = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{y \sim p_\theta(y|x_i)} (\nabla_\theta \log p_\theta(y | x_i))^2, \quad (17)$$

which is the average of the square gradient of  $y$  with respect to a given parameter  $\theta$ . By confining fine-tuning to these important parameters (whether identified by gradient, weight magnitude, or Fisher information), such selective fine-tuning techniques drastically reduce training overhead and still achieve performance close to full model fine-tuning.

## 8.4 Knowledge Distillation for Fine-Tuning

Knowledge distillation (KD) compresses LLMs by transferring knowledge from a high-capacity teacher to a smaller student [Xu et al., 2024b, Fang et al., 2025b]. Unlike direct SFT or RLHF pipelines, KD transfers alignment properties from a pre-aligned teacher model (typically optimized via RLHF or similar methods) to a student model by mimicking the teacher’s performance, without repeating the costly preference-optimization loop.

The concept of KD is introduced in the foundational work of Hinton et al. [2015]: instead of training on hard labels  $y$ , the student learns from the teacher’s class probability distribution  $\mathbf{p}_T = \sigma(\mathbf{z}_T/\tau)$ , where  $\mathbf{z}_T$  are logits from teacher,  $\sigma$  is the softmax function, and  $\tau$  is a temperature parameter that softens the distribution. The student’s objective function is typically a weighted average of two losses: a standard cross-entropy loss  $\mathcal{L}_{\text{CE}}$  with the hard labels, and a distillation loss that minimizes the Kullback-Leibler (KL) divergence between the student’s and teacher’s softened outputs:

$$\mathcal{L}_{\text{KD}} = \alpha \cdot \mathcal{L}_{\text{CE}}(\sigma(\mathbf{z}_S(\mathbf{x})), y) + (1 - \alpha) \cdot \tau^2 \cdot \mathcal{L}_{\text{KL}}(\sigma(\mathbf{z}_T(\mathbf{x})/\tau), \sigma(\mathbf{z}_S(\mathbf{x})/\tau)), \quad (18)$$

where  $\alpha$  balances the two terms. Matching the full probability vector rather than one-hot targets forces the student to copy subtle preference signals, including refusal styles, politeness markers, and content filters, that are otherwise hard to encode.

One primary advantage of KD is its efficiency. Compressing a 70-billion-parameter RLHF teacher to a 7-billion-parameter student cuts inference cost by roughly an order of magnitude while

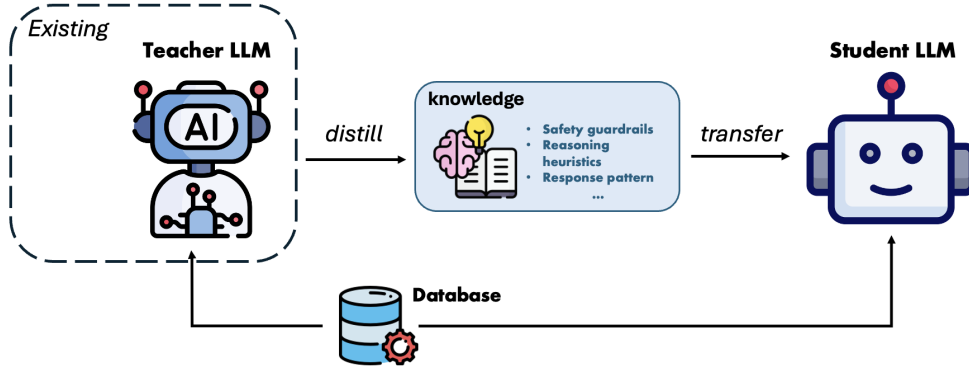


Figure 7: Overview of Knowledge Distillation in LLMs. Knowledge is distilled from a teacher LLM, which is typically optimized via RLHF. This knowledge, potentially enriched with current, task-specific data, is transferred to a smaller student LLM. By learning from both the teacher’s guidance and the current data, the student LLM becomes more efficient and effective at performing downstream tasks.

preserving core reasoning ability and alignment quality [Taori et al., 2023, Touvron et al., 2023, Xu et al., 2024b, Guo et al., 2025, Yang et al., 2024b, Ma et al., 2024, Wang et al., 2025e]. Furthermore, the teacher is used only for forward passes, the same recipe scales easily to specialised domains: clinicians and bioinformaticians have distilled general-purpose aligned teachers into compact models fine-tuned for medical or scientific tasks without eroding their safety guarantees [Niu et al., 2024, Tariq et al., 2024, Ge et al., 2025, Latif et al., 2024, Shang et al., 2024].

Subsequent work shows that transferring richer signals strengthens alignment further. Feature-based KD [Ji et al., 2021] matches hidden activations so the student inherits syntactic and semantic structure. Attention distillation [Jiao et al., 2019] copies attention maps, helping small transformers remain stable during fine-tuning. Relation-based KD [Yang et al., 2022] preserves similarity patterns inside the network and improves dense prediction tasks. Sequence-level KD trains the student to reproduce complete outputs, capturing dialogue coherence and style [Taori et al., 2023, Li et al., 2021]. Relation-based KD distills the reasoning patterns of teachers, such as the Chain-of-Thought (CoT) reasoning, from teacher LLMs so the student learns intermediate reasoning steps rather than shortcutting to the final answer [Hsieh et al., 2023, Feng et al., 2024]. For example, Hsieh et al. [2023] generates the teacher rationales, and then the student is trained to jointly predict both the rationale and the final answers, helping the student learn intermediate reasoning steps rather than shortcutting to the final answer. Formally, the student is trained on a loss:

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta}(\mathbf{r}_i, y_i \mid \mathbf{x}_i), \quad (19)$$

where  $(\mathbf{r}_i, y_i, \mathbf{x}_i)$  comprises rationale, input, answer, and  $P_{\theta}$  denotes the student model parameterized by  $\theta$ .

Several extensions refine KD for alignment objectives. Multi-teacher distillation aggregates outputs from diverse, pre-aligned models (e.g., RLHF-specialized variants), enabling students to synthesize compounded behavioral priors [Khanuja et al., 2021, Zhang et al., 2022, Liu et al., 2024c, Wadhwa et al., 2025]. Dynamic adaptive distillation further introduces a paradigm of bidirectional co-evolution, where teacher and student models undergo simultaneous joint optimization for continuous mutual refinement [Sun et al., 2021, Chang et al., 2022, Li et al., 2024i]. Self-distillation

represents a distinct approach that bypasses external teachers entirely, leveraging a single model’s self-generated outputs for iterative refinement [Zhang et al., 2019, Yang et al., 2023].

One another key point is the uncertainty-aware abilities of the student model. Classic KD can create over-confident students because it ignores uncertainty. Uncertainty-aware variants address this in two complementary ways. The first approach focuses on distilling uncertainty by training the student to mimic the teacher’s full predictive distribution over labels. These distributions are typically learned from uncertainty-aware systems such as Bayesian neural networks or model ensembles, with the student minimizing divergence between its outputs and the teacher’s probabilistic outputs [Korattikara Balan et al., 2015, Vadera et al., 2020, Malinin et al., 2019]. The second approach addresses quantification of the student’s intrinsic uncertainty by reinterpreting knowledge distillation through a Bayesian framework. Bayesian Knowledge Distillation (BKD) [Fang et al., 2024] achieves this by embedding the teacher’s output probabilities as a teacher-informed prior distribution over the student model’s weights. Within this formulation, standard knowledge distillation loss emerges naturally as the posterior optimization objective. Sampling techniques, such as stochastic gradient Langevin dynamics, applied to this posterior yield principled predictive intervals. Collectively, these uncertainty-aware methods suppress overconfident predictions on noisy or out-of-distribution inputs while providing computationally efficient alternatives to full Bayesian training. This characteristic makes them particularly valuable for developing compact while still being safety-aligned LLMs.

## 8.5 Adapter-Based Fine-Tuning

Adapter-based fine-tuning involves inserting small trainable modules, known as adapters, between the layers of a pre-trained model. The original model weights remain frozen, and only the adapters are trained on downstream tasks. This approach allows for efficient multi-task learning and rapid adaptation to new tasks without retraining the entire model.

An adapter module commonly contains two linear layers to conduct down-projection that reduces the dimensionality of the input and up-projection that restores it. Between the two layers, nonlinearity layer is applied to the down-projected representation. During the fine tuning process, trainable parameters  $\Phi$  are optimized on given dataset  $D$ , with specified loss function  $L$ :

$$\Phi^* \leftarrow \operatorname{argmin}_{\Phi} L(D; \Phi) \quad (20)$$

Adapters can be broadly categorized into series adapters and parallel types, depending on how they are integrated with the backbone [Hu et al., 2023].

Series adapters attach learnable modules sequentially within selected layers, typically within transformer layers. [Houlsby et al., 2019] proposed the first series adapter as a parameter-efficient fine-tuning method by inserting trainable modules twice in each transformer layer. Adamix introduces multiple down- and up-projection layers in each adapter [Wang et al., 2022b]. It employs stochastic routing to randomly select projection pairs during training. By maintaining the same number of tunable parameters and computational cost as the underlying PEFT, this technique outperforms these methods and even full fine-tuning in several datasets. A structural limitation is that the backpropagation must still pass through the main backbone, resulting in additional computational steps and increased memory usage. To further reduce training time and memory consumption, 1) adapter pruning and 2) sparse calculation techniques have been proposed. For instance, solely inserting the adapter module in each transformer layer [Pfeiffer et al., 2020] or dropping adapter in lower-level layers [Rücklé et al., 2021] can still preserve performance. Compacter replaces standard linear layers with Kronecker products and shared parameters [Mahabadi

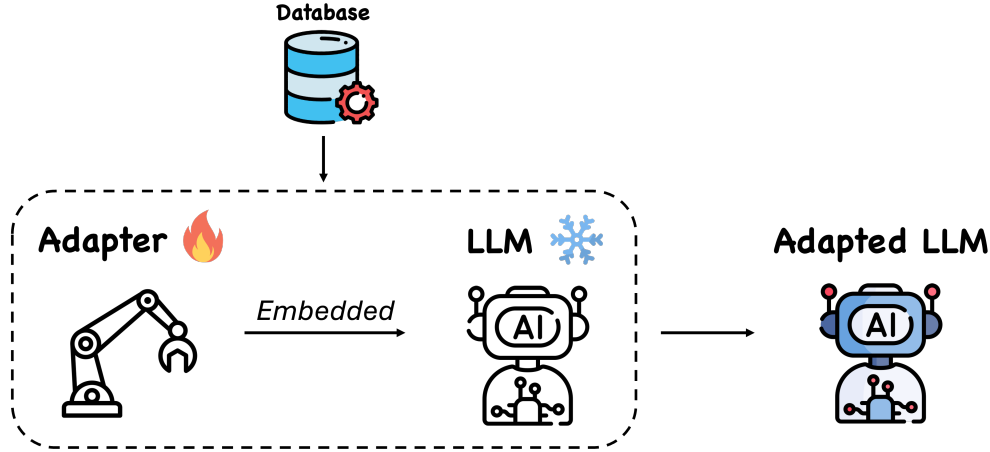


Figure 8: Overview of Adapter in LLMs. Adapters are lightweight, trainable modules inserted into the layers of a pre-trained, frozen LLM to enable adaptation to specific downstream tasks. By keeping the core model parameters fixed and training only the adapter components, this approach greatly reduces computational costs. Depending on how they are integrated into the LLMs, adapters are generally categorized as series and parallel types and numerous variants have been developed to enhance both performance and efficiency.

[et al., 2021](#)]. This technique maintains the sparsity and reduces computational complexity while achieving the performance comparable to original down- and up- projection calculations.

Parallel adapters incorporate additional adapter modules in parallel with specific layers, which has similar manner as LoRA [\[He et al., 2022\]](#). When these adapters are connected in paralleled to the original layers, backpropagation can be applied directly across each adapter by linking them explicitly [\[Sung et al., 2022\]](#). Multi-adapter is an extension on parallel adapter structure, which modifies the output of the self-attention heads through multiple parallel adapters. In multi-task settings, especially when tasks are sequentially related, single-task adapters often suffer performance degradation across tasks. AdapterFusion introduces a mechanism that parallels multiple task-specific adapters within the attention layer and fuses their outputs [\[Pfeiffer et al., 2021\]](#). This approach achieves better performance than using a single-task adapter alone.

As a widely used PEFT technique, adapters can also be compatible with other PEFT methods [\[Mao et al., 2022\]](#). This ability makes adapter have flexible choices for either the demand on inference speed or multi-task performance etc.

## 8.6 Comparison of Fine-Tuning Techniques

Each fine-tuning technique discussed in this section offers distinct advantages and limitations regarding computational efficiency, adaptability, and performance. Table 6 provides a comprehensive summary, comparing critical aspects such as trainable parameters, computational and memory overhead, adaptability, and inherent trade-offs.

Full-parameter fine-tuning updates every parameter in a model, resulting in typically superior performance and high adaptability. However, it incurs substantial computational and memory demands, which become increasingly prohibitive with larger models. Partial parameter fine-tuning alleviates some of these burdens by selectively updating essential parameters, though it adds complexity in choosing optimal subsets and may risk suboptimal adaptation if crucial layers remain unchanged.



LoRA introduces small amount of trainable parameters (less than 1%) by employing low-rank matrices, significantly reducing computational costs. While highly efficient, LoRA can fall short in tasks requiring deep reasoning or complex coding without further optimization. It balances well between resource efficiency and moderate performance.

Sparse fine-tuning effectively reduces training overhead by selectively updating parameters based on sparsity measures like gradient magnitudes or Fisher information. This approach typically achieves near full fine-tuning performance at significantly lower computational costs, provided the sparse masks are optimally determined. Nonetheless, it depends heavily on the quality and robustness of the sparsity strategy.

Knowledge distillation transfers knowledge from a larger teacher model to a smaller student model, considerably reducing inference costs. While computationally efficient during inference, this approach involves an initial cost for training the teacher model. Its performance relies significantly on the teacher’s quality and the distillation method, with possible performance degradation compared to the original large model.

Adapter-based fine-tuning achieves high adaptability and efficiency by inserting small, trainable adapter modules into pre-trained models. Although adapters facilitate efficient multitask learning and quick adaptation, they introduce computational overhead proportional to the complexity of the adapter modules, especially in series configurations.

Prompt tuning maintains the original model entirely frozen, adjusting only continuous prompt embeddings. This method offers minimal computational and memory overhead, suitable for resource-constrained environments. However, prompt tuning often provides limited adaptation capabilities compared to methods that directly update model parameters, constraining its effectiveness on tasks significantly divergent from the original pre-training objectives.

Table 6: Comprehensive Comparison of Efficient Fine-Tuning Techniques.

Method	Trainable Parameters	Memory Overhead	Compute Cost	Adaptability
Full-Parameter FT	High	High	High	High
Partial Parameter FT	Medium	Medium	Medium	Medium
LoRA	Low	Low	Low	Medium
Sparse FT	Low	Low	Low	High
Knowledge Distillation	Medium	Medium	Medium	Medium
Adapter-Based FT	Low	Medium	Medium	High
Prompt Tuning	Low	Low	Low	Medium

In summary, choosing an appropriate fine-tuning method depends on specific deployment scenarios and resource constraints, balancing between desired performance, computational resources, and model adaptability.

## 9 Brain-Inspired LLM Alignments

This section outlines the key principles and methodologies underlying brain-inspired LLM alignment, often referred to as Brain-AGI, highlights the existing challenges, and provides insights into possible opportunities.

### 9.1 Recent advancements of Brain-Inspired LLM Alignments

Brain-inspired LLMs refers to language models whose architectures or training objectives draw directly on principles of human brain organization and function. Instead of treating neural networks

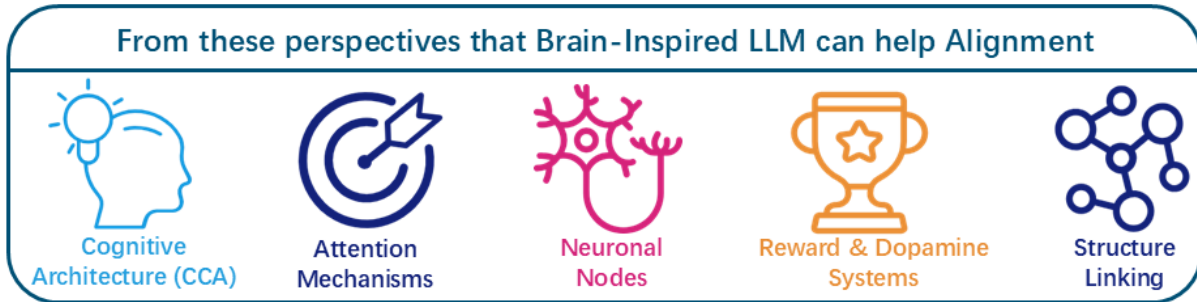


Figure 9: Different Perspectives that Brain-Inspired LLM can help Alignment.

as arbitrary black boxes, brain-inspired LLMs seek to mirror how our brains process language, represent concepts, and integrate sensory information [Farisco et al., 2024]. For instance, the Causal Cognitive Architecture (CCA) [Schneider, 2024], a framework that models neocortical minicolumns as millions of “navigation maps.” These maps undergo continuous cognitive cycles in which sensory inputs are normalized, spatially and temporally bound into local maps, and then matched against stored multisensory maps to form a Working Navigation Map. The CCA framework demonstrates how lightweight, evolution-inspired modifications to a core navigation-map framework can spontaneously produce foundational aspects similar to that of human intelligence [Schneider, 2024].

Additionally, the BriLLM architecture [Zhao et al., 2025] replaces the Transformer’s [Vaswani et al., 2017] attention blocks with a directed graph of “neuronal” nodes and energy-based signal flow, analogous to how biological neurons propagate activation along pathways of least resistance. In BriLLM, tokens are treated as nodes, and each edge carries a learnable “energy tensor” that determines which node (token) activates next. This design makes every internal connection interpretable and, in principle, allows unbounded context lengths [Zhao et al., 2025].

Furthermore, Sun et al. [2024] investigated whether LLMs exhibit a brain-like functional organization by directly linking sub-groups of artificial neurons (ANs) in models such as BERT and the Llama family to well-established human functional brain networks (FBNs). They extracted representative temporal “atoms” of neuron activity and used them to predict fMRI responses, demonstrating that these LLMs indeed form a modular, brain-like architecture. Comparing four models, BERT and three successive Llama variants, they found that this functional specialization strengthens with model sophistication: larger models yield brain maps that are both more consistent (showing reduced variability in engaged FBNS) and more compact (fewer, more specialized atoms per network), and their neurons display increasingly hierarchical temporal and anatomical distributions [Sun et al., 2024].

## 9.2 Brain-AGI Co-working

The motivation for brain and artificial general intelligence (AGI) co-working arises from the limitations of conventional AI systems, which often struggle with rigid task boundaries, inefficient energy use, and limited contextual understanding [Meftah et al., 2025]. In contrast, biological brains demonstrate remarkable capabilities in generalization, lifelong learning, and causal reasoning, qualities that inspire the design of AGI systems seeking to replicate such flexibility and intelligence [Gabriel, 2020a]. Brain-AGI collaboration is envisioned as a way to harness the complementary strengths of human cognition and machine computation [Zhao et al., 2023], enabling more effective

responses to complex societal challenges such as healthcare delivery and climate modeling [Lu et al., 2023]. This synergy leverages human creativity and judgment alongside AGI’s capacity for scale and speed. A central driver of this co-working paradigm is the pursuit of ethical alignment, which ensures that AGI systems are not only powerful but also aligned with human values and intentions [Conitzer et al., 2024].

*Theoretical foundations.* The co-working of the brain and AGI is underpinned by a set of complementary theoretical principles drawn from neuroscience and cognitive science [Yu et al., 2024c]. First, predictive coding provides a computational framework for modeling shared brain-AGI processing, wherein both biological and artificial agents minimize prediction error via hierarchical Bayesian inference [Fagerholm et al., 2020]. Second, neuroplasticity offers a foundation for co-adaptive learning systems, enabling AGI to dynamically update internal representations in response to environmental and social interactions, which is similar to synaptic rewiring in the human brain [Rebedea et al., 2023]. Third, embodied cognition underscores the significance of sensorimotor grounding for both natural and artificial agents, positing that cognition emerges through real-world engagement. This is an essential principle for brain-AGI integration in collaborative robotics and assistive technologies [Taniguchia et al., 2022]. These three pillars are further reinforced by neurosymbolic integration, a hybrid framework that aligns neural learning with symbolic reasoning, offering interpretable and generalizable architectures essential for effective and trustworthy human-AGI co-working [Bhuyan et al., 2024].

*Current strategies.* Recent progress in brain-inspired artificial intelligence has catalyzed a new wave of brain-AGI co-working paradigms, driven by four interconnected approaches [Yu et al., 2024c]. First, neuromorphic engineering develops hardware systems that emulate the structure and dynamics of biological neural circuits. Technologies like Intel’s Loihi chip leverage spiking neural networks to enable real-time, energy-efficient computation, closely mirroring the temporal coding strategies of the brain [Yik et al., 2023]. Second, cognitive architectures such as CLARION [Sun, 2006] integrate symbolic reasoning with subsymbolic learning (e.g., deep learning, reinforcement learning), supporting flexible, goal-directed behavior that aligns with human-like cognitive processes. Third, human-in-the-loop learning frameworks, including inverse reinforcement learning and preference modeling, bring human values and feedback into the AGI training loop, enabling systems to adapt to social and ethical contexts through cooperative learning [Hadfield-Menell et al., 2017]. Finally, neuroadaptive interfaces enable real-time, bidirectional interaction between human neural activity and artificial agents. These systems dynamically decode brain signals to adapt AGI responses, fostering seamless collaboration between biological and artificial intelligences [Guo et al., 2024b].

*Applications.* Brain-AGI co-working is unlocking transformative applications across multiple fields by combining human cognitive strengths with artificial general intelligence [Zhao et al., 2023]. In healthcare, this synergy enables early diagnosis of neurodegenerative diseases like Alzheimer’s disease through AGI-assisted interpretation of fMRI data, and supports real-time decision-making in surgery via brain-computer interfaces [Fedorova et al., 2022]. In autonomous systems, brain-AGI collaboration empowers robots with human-like spatial reasoning and adaptability, enhancing performance in complex, unstructured environments such as disaster response scenarios [Blackiston et al., 2022]. Climate science similarly benefits from co-working paradigms, where AGI tools informed by human-guided intuition and causal reasoning predict ecological tipping points, guiding proactive environmental policies [Harder et al., 2022]. These applications exemplify how brain-AGI integration not only augments machine capabilities but also extends human potential in addressing critical global challenges.

*Challenges.* The pursuit of brain-AGI co-working, where artificial general intelligence systems are designed to complement or integrate with human cognitive processes, faces profound challenges

across technical, ethical, and interdisciplinary domains [Zhao et al., 2023]. Technically, neuro-morphic systems inspired by the brain remain constrained by scalability and energy efficiency [Zolfagharinejad et al., 2024], limiting their capacity to support the high computational demands of AGI [Kurshan, 2024]. Ethically, the prospect of AGI systems operating alongside or within human cognitive environments raises concerns about value misalignment, loss of human agency, and the amplification of biases without adequate oversight mechanisms [Everitt et al., 2018]. Most critically, the gap between neuroscience and AI research continues to hinder meaningful co-design. The incompatible frameworks, terminologies, and research priorities make it difficult to translate biological principles into actionable AI models or to embed AGI into cognitive contexts in a way that is both effective and safe [Hassabis et al., 2017]. Overcoming these barriers will require not just technical innovation but also deep, sustained integration across disciplines and a commitment to embedding human values at the heart of AGI design.

### 9.3 Challenges and Limitations of Brain-Inspired LLM Alignments

Despite tremendous achievements in the brain-inspired LLM architectures such as enhanced interpretability, modular signal paths, and the promise of unbounded context lengths, several challenges remain. First, the computational overhead of maintaining and updating dense “energy tensors” across every connection in a graph-based design can be substantial, potentially eroding the very energy-efficiency gains these models aim to deliver [Zhao et al., 2025].

Second, current alignment methods such as DPO align an LLM by directly maximizing the likelihood of human-preferred completions. DPO accomplishes this by reparameterizing reward as the log-ratio between the policy and a fixed reference model, thus eliminating the need for a separate reward network [Xu et al., 2024c]. Although DPO often matches or outperforms RLHF on in-distribution benchmarks, it fundamentally learns only the surface-level patterns of human preferences rather than the deeper cognitive processes [Tennant et al., 2025]. As a result, DPO-tuned models can “sound” convincingly human on familiar prompts but remain vulnerable to reward hacking and poor generalization when faced with novel or out-of-distribution scenarios.

Third, integrating ethical and societal considerations, such as bias mitigation or transparent governance, into a highly specialized, brain-inspired AGI framework remains an open problem, as standard alignment tools (RLHF, DPO) do not yet translate directly onto these new architectures [Farisco et al., 2024]. Finally, scaling this graph-based signal-flow mechanism to models with billions of parameters and multimodal inputs has not yet been demonstrated at production scale, raising questions about its practicality for real-world applications.

### 9.4 Opportunities and Future Developments

Brain-inspired LLM alignment techniques offer several promising advantages over conventional preference-fitting approaches. By explicitly tying language-processing pathways to analogues of human cognitive systems, equips us with the ability to pinpoint exactly where and why a model’s reasoning diverges from our expectations. This transparency doesn’t just aid debugging; it lays the groundwork for truly collaborative workflows in which domain experts can inspect, validate, and open pathways toward achieving Brain-AGI.

Furthermore, aligning LLM subnetworks with well-characterized human functional brain networks, such as the “language” network for factual consistency or the “default-mode” network for social reasoning, provides a neuroscientific scaffold for targeted auditing and control [Sun et al., 2024]. This modular interpretability makes it possible to apply fine-grained alignment objectives, such as moral rewards derived from deontological or utilitarian principles, to dedicated subnet-

works, reducing collateral effects on other competencies and mitigating reward-hacking risks seen in methods like DPO. Lastly, grounding alignment in multimodal, embodied representations, by integrating vision, audio, and proprioceptive signals into the same brain-inspired modules, promises richer, context-sensitive behavior that better mirrors human cognition and values [Farisco et al., 2024].

Looking forward, next-generation brain-inspired alignment will likely embrace dynamic plasticity mechanisms, enabling models to adapt in real time through localized synaptic-style updates rather than wholesale retraining. Instead of relying on static modules, future systems may incorporate continuous update rules that refine internal representations in response to live user feedback or environmental signals, much like how our brains learn from experience. Coupling such plasticity with lightweight neurofeedback, whether from biosignals or behavioral proxies, could give rise to true brain-in-the-loop pipelines that preserve robust alignment across shifting contexts. Finally, as Brain-AGI systems grow more autonomous, evolving multidisciplinary governance frameworks that integrate insights from neuroscience, ethics, and policy will be essential to ensure that alignment innovations remain transparent, accountable, and aligned with societal norms.

## 10 Alignment Uncertainty Quantification (AUQ)

Alignment Uncertainty Quantification (AUQ) addresses a fundamental challenge in LLM development: how to measure and manage the uncertainty inherent in aligning models with human values and intentions [Gabriel, 2020b, Bommasani et al., 2021]. As models become more powerful, understanding the reliability of their alignment becomes critical for safe deployment [Reynolds and McDonell, 2021]. Traditional machine learning uncertainty focuses primarily on predictive accuracy, while alignment uncertainty concerns whether a model’s behavior truly reflects intended human values across diverse contexts and edge cases. The theoretical foundations of alignment uncertainty, methodologies for its quantification, and approaches for building systems that remain robustly aligned despite inevitable uncertainties are examined in this section.

### 10.1 Sources of Alignment Uncertainty

The uncertainty in aligning LLMs with human values arises from multiple interrelated sources, complicating both measurement and mitigation. These can be grouped into three main categories: model-related uncertainty, feedback-related uncertainty, and contextual or distributional uncertainty.

**1. Training Stability and Optimization Efficiency** Model-inherent uncertainty stems from the stochastic nature of the training process and the architectural limitations of current LLMs. Due to random initialization and the use of stochastic gradient descent, models trained on the same data may still diverge in behavior and alignment properties [Gal and Ghahramani, 2016]. Additionally, present-day neural architectures may lack the expressiveness to fully capture complex or evolving human values, leading to representational misalignment [Dodge et al., 2020]. As models scale, **emergent behaviors** arise unpredictably, further complicating efforts to anticipate how alignment generalizes across model sizes and tasks [Wei et al., 2022, Ganguli et al., 2022b].

**2. Human Feedback Variability and Value Pluralism** The human signals that guide alignment introduce substantial uncertainty. Annotators often provide inconsistent judgments for similar

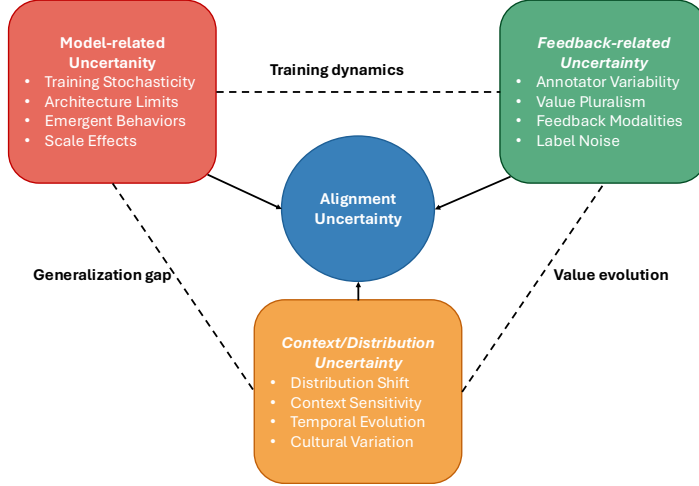


Figure 10: Sources of Alignment Uncertainty.

prompts, with inter-annotator agreement typically ranging from 0.6 to 0.8 Krippendorff’s alpha [Stienon et al., 2020, Ziegler et al., 2019], reflecting noisy or subjective supervision. Moreover, human values are inherently diverse, shaped by culture, individual background, and societal context, making it difficult to define a universally valid alignment target [Gabriel, 2020b]. The mode of feedback also matters: scalar ratings, pairwise comparisons, natural language critiques, and demonstrations capture different aspects of preference and intent, introducing variation in the resulting aligned policy [Ouyang et al., 2022, Bai et al., 2022b].

**3. Context Sensitivity and Distributional Shift** Uncertainty also emerges from real-world deployment contexts. Distribution shift, when deployed models face inputs or environments not reflected in the training data, can lead to severe misalignment [Hendrycks and Dietterich, 2019]. Even within familiar domains, context-sensitive interpretation can vary dramatically depending on user identity, culture, or timing. The same response may be perceived as helpful or harmful depending on situational nuances [Pérez et al., 2022]. In addition, societal norms evolve over time, meaning that a model aligned with contemporary values may become increasingly misaligned as social standards shift [Solaiman et al., 2021].

## 10.2 Conceptual Framework and Methods for Quantifying Alignment Uncertainty

### 10.2.1 Conceptual Framework for Alignment Uncertainty

The alignment problem is formalized as finding a model policy  $\pi(y|x)$  that generates outputs  $y$  given inputs  $x$  to maximize expected utility according to human values. However, these values are neither perfectly known nor perfectly represented within the model, creating fundamental uncertainty in the alignment process.

This uncertainty is formalized through a decision-theoretic framework. Let  $U(x, y, v)$  represent the utility of response  $y$  to prompt  $x$  according to a value function  $v$ . The true human value function  $v^*$  is unknown and can only be approximated by the model’s internal representation  $\hat{v}$ . The alignment gap is defined as:



$$\mathcal{G}(x) = \mathbb{E}_{y \sim \pi(y|x)}[U(x, y, v^*) - U(x, y, \hat{v})] \quad (21)$$

Alignment uncertainty quantification aims to characterize the distribution and magnitude of this gap, enabling more informed decisions about model deployment, refinement, and usage limitations [Amodei et al., 2016, Hendrycks et al., 2021].

This uncertainty differs fundamentally from traditional predictive uncertainty in machine learning. While predictive uncertainty concerns the accuracy of model outputs, alignment uncertainty concerns their desirability according to human values, a substantially more complex target that varies across individuals, cultures, and contexts. The misalignment risk emerges not just from model limitations, but from the inherent challenges in defining, communicating, and representing human values themselves.

### 10.2.2 Methods for Quantifying Alignment Uncertainty

Deep ensemble methods [Osband et al., 2016] established that model disagreement serves as an empirical proxy for epistemic uncertainty by training multiple models on different data subsets and measuring prediction variance across ensemble members. However, the computational cost of training multiple large models proved prohibitive as language models scaled. This limitation inspired Monte Carlo dropout [Gal and Ghahramani, 2016], which demonstrated that dropout at inference time approximates Bayesian posterior sampling, enabling uncertainty estimation from a single model through multiple stochastic forward passes. This approach maintained theoretical grounding while dramatically reducing computational requirements.

While these general-purpose methods provided efficient uncertainty quantification, the application to alignment problems required specialized treatment of human preferences. Bayesian reward modeling [Christiano et al., 2017] addressed this need by formalizing preferences as distributions over reward functions through the posterior  $P(r \mid \mathcal{D}) \propto P(\mathcal{D} \mid r)P(r)$ , where  $\mathcal{D}$  represents human feedback data,  $r$  denotes a reward function, and  $P(r)$  is the prior belief about rewards. The posterior width quantifies uncertainty in value alignment, wider distributions indicate less confidence about human preferences. Yet computational demands of exact Bayesian inference remained challenging, leading to posterior policy sampling [Ramachandran and Amir, 2007], which generates diverse behaviors by sampling from the posterior distribution rather than computing it explicitly. This practical adaptation revealed a crucial insight: alignment uncertainty manifests not only in parameter uncertainty but also in behavioral diversity. Building on this understanding, hierarchical Bayesian approaches [Paun et al., 2018] extended the framework by explicitly modeling inter-annotator disagreement, acknowledging that human feedback itself contains irreducible uncertainty from diverse value systems.

As model scales continued to grow, even sampling-based Bayesian methods became computationally prohibitive. [Leike et al., 2018] responded by developing scalable approximations that maintained theoretical rigor while enabling practical deployment. Their work catalyzed a fundamental shift in perspective: rather than approximating complex posteriors, researchers began exploring uncertainty metrics that could be computed directly from model outputs. This led naturally to information-theoretic approaches, which provided computationally efficient tools through entropy  $H(A \mid x) = -\sum_a P(A = a \mid x) \log P(A = a \mid x)$ , where  $A$  represents the alignment decision,  $x$  is the input prompt, and  $P(A = a \mid x)$  is the probability of alignment choice  $a$ . Higher entropy indicates greater uncertainty about which response is properly aligned. KL divergence further quantifies how different alignment methods produce different behaviors [Cover and Thomas, 1999, Xiao et al., 2022]. These metrics revealed a previously overlooked dimension: alignment pro-

cedures themselves introduce systematic uncertainty, suggesting that methodological choices must be considered alongside data uncertainty.

While information theory provided efficient metrics, the need for deployment guarantees motivated a different approach entirely. The adaptation of conformal prediction [Angelopoulos and Bates, 2021] to alignment problems represented a paradigm shift from estimation to selection. Conformal Alignment [Gui et al., 2024] transforms the fundamental question from “how uncertain are we?” to “which outputs can we trust?” by computing p-values that measure how unusual each output is compared to calibration data. The method ensures that the false discovery rate (FDR), the proportion of selected outputs that are actually misaligned, stays below a user-specified threshold  $\alpha$  [Benjamini and Hochberg, 1995], providing distribution-free statistical guarantees essential for safety-critical deployments.

Recent developments have evolved beyond measuring uncertainty to actively utilizing it within the alignment process itself. Temperature scaling [Renze and Guven, 2024] demonstrated that simple sampling entropy modifications could reveal model confidence patterns, suggesting deep connections between uncertainty and alignment quality. This insight inspired Uncertainty-Aware Learning (UAL) [Wang et al., 2024f], which incorporates uncertainty directly into training through adaptive reward smoothing:

$$\tilde{r}(x, y) = (1 - \lambda H(y|x))r(x, y) + \lambda H(y|x)\bar{r} \quad (22)$$

where  $\tilde{r}(x, y)$  is the smoothed reward,  $r(x, y)$  is the original reward for response  $y$  to prompt  $x$ ,  $H(y|x)$  measures response uncertainty (entropy),  $\bar{r}$  is the average reward, and  $\lambda$  controls the smoothing strength. This formulation elegantly addresses a fundamental problem in RLHF: when the model is uncertain (high entropy), it trusts the feedback less and moves rewards toward the average, preventing overfitting to potentially noisy signals.

The most recent evolution in alignment uncertainty quantification recognizes that previous methods operated at single linguistic scales, missing the hierarchical nature of language generation. [Zhang et al., 2025d] pioneered token-level uncertainty quantification through low-rank weight perturbations to model parameters. Their key insight is that epistemic uncertainty (EU) equals the mutual information  $I(y_t; \theta | y_{<t}, x)$  between the next token  $y_t$  and model parameters  $\theta$ , given previous tokens  $y_{<t}$  and input  $x$ . High mutual information indicates the model is uncertain about which token to generate next, a signal for potential hallucinations. Complementing this microscopic view, [Xie et al., 2025] extended the framework to session-level dynamics, capturing how alignment states  $\theta_t$  evolve over time  $t$  during conversations, with evolution governed by previous states  $\theta_{t-1}$  and session-level parameters  $\phi$ . This hierarchical synthesis, aggregating token variances into utterance-level risk scores that inform session-level priors, represents the culmination of methodological evolution, where each scale builds upon insights from all previous approaches to provide comprehensive uncertainty quantification across conversational timescales.

The evolution from computationally intensive Bayesian methods to efficient multi-scale frameworks reflects the field’s response to practical deployment constraints. Table 7 summarizes how each approach trades off between theoretical rigor, computational efficiency, and practical applicability. Notably, the progression shows a clear pattern: early methods prioritized mathematical foundations but struggled with scale, middle approaches balanced efficiency with accuracy, while recent multi-scale methods attempt to achieve both through hierarchical decomposition. This suggests that future developments may continue this trend toward specialized architectures that match the natural structure of language generation tasks.

Table 7: Comparison of alignment uncertainty quantification methods.

Method	Strengths	Limitations	Comp. Cost	Best Case	Use
Bayesian Reward Modeling	Principled uncertainty; Captures preference distributions	Computationally intensive; Prior needed	High	Research with ample compute	with com-
Ensemble Methods	Practical; No distributional assumptions	Multiple models; Training overhead	Medium-High	Production with uncertainty needs	
Information Theory	Model-agnostic; Theoretically grounded	May conflate uncertainty types	Low	Quick uncertainty assessment	un-
Conformal Prediction	Distribution-free; Formal FDR control	Needs calibration; Binary selection	Medium	Safety-critical applications	
Multi-scale Modeling	Hierarchical uncertainty; Comprehensive	Complex; Multiple components	High	Long-form dialogue systems	

### 10.3 Robustness and Uncertainty in Alignment

Understanding and quantifying alignment uncertainty is essential for building AI systems that are robust to misalignment risks and capable of adapting to evolving human values. This section reviews key methods for uncertainty-aware alignment and safety mechanisms that manage residual uncertainty.

**Uncertainty-Aware Training.** Alignment methods increasingly incorporate uncertainty estimates during training and deployment. Distributionally robust optimization techniques account for worst-case value realizations to prevent undesirable behaviors under misalignment [Rahimian et al., 2019]. Risk-sensitive reinforcement learning integrates risk measures such as variance and CVaR to promote consistent policy behavior across contexts [Mihatsch and Neuneier, 2002]. Recent work in reward modeling emphasizes capturing uncertainty in feedback signals to support more conservative updates [Leike et al., 2018, Everitt et al., 2021].

**Safety Mechanisms at Deployment.** Safety measures mitigate residual uncertainty at inference time. Threshold-based abstention strategies avoid output generation when epistemic uncertainty is high [Lin et al., 2023]. Guardrails, such as rule-based constraints and adversarial red-teaming, establish behavioral boundaries independent of model confidence [Bai et al., 2022a, Ganguli et al., 2022a]. Human-in-the-loop frameworks allow expert oversight in high-risk or high-uncertainty cases [Askell et al., 2021].

**Illustrative Domains.** Real-world examples highlight the importance of robustness to alignment uncertainty. In medical decision support, alignment challenges stem from domain knowledge limitations, contextual ambiguity, and value conflicts among stakeholders [McKenna et al., 2023, Zhang et al., 2024c]. In political and ethical content generation, cultural variation and subjective norms contribute to persistent alignment uncertainty [Blodgett et al., 2020, Solaiman et al., 2021]. These settings demand models that express uncertainty and respond conservatively when appropriate.

**Active-learning loop for preference data.** Alignment datasets are costly; integrating AUQ into query selection can shrink data needs while targeting the largest misalignment areas. [Muldrew et al., 2024] propose *Active Preference Learning for LLMs*, selecting prompt-completion pairs with high predictive entropy under DPO; they cut label count by 40 % while matching baseline reward quality. A complementary study, *Less is More* [Deng et al., 2025], filters low-utility pairs via ensemble disagreement before DPO, achieving higher win-rates and faster convergence. Online variants combine count-based exploration with uncertainty scores to refresh the reward model on-the-fly, maintaining alignment under distribution shift [Lu et al., 2025]. These approaches form an uncertainty-aware, human-in-the-loop closed loop, progressively narrowing the alignment gap where it matters most.

## 11 Societal, Ethical, and Regulatory Considerations

This section outlines broader implications of alignment practices and reviews current regulatory and policy developments.

### 11.1 Ethical and Societal Implications

AI is reshaping society [Shi et al., 2025] at an unprecedented scale, influencing key domains such as healthcare, education, finance, and scientific discovery. As LLMs become more intelligent and autonomous, their societal and ethical impact raises ethical questions regarding alignment [mic, 2025] with human moral values.

Society is increasingly reliant on AI technologies to help with decision-making. But this growing reliance comes with potential risk without alignment [ibm, 2025]: LLMs can generate biased, harmful, and inaccurate outputs that are not aligned with the goals of their creators and the original intent of the system. Aligned models are tackling the ethical challenges resulting from the deployment of LLM [Ferdaus et al., 2024], e.g., potential misuse and abuse of LLMs, negative impacts on users heavily relying on LLM agents, the environment, information dissemination, and employment. Addressing these challenges is paramount, and the development of aligned models offers a promising path towards fostering public trust, ensuring fairness, and promoting the ethical application of AI. Based on these considerations and reports [sta], a humanities-guided approach to AGI implemented by aligned models is crucial, recognizing AI not merely as a tool, but as a civilizational technology to redefine the shared ethical principles, societal hierarchy structures and build the aligned consensus acknowledged by the public and AI.

Public concerns about truth and LLMs [spr, a, nih, 2025] since these technologies generate misinformation and are used to spread incorrect information or manipulate people. Moreover, public trust in LLMs always hinges on their reliability, predictability, and operational transparency. The “black box” nature of language models, obscures decision-making and retrograde reasoning. Incidents of bias and harmful outputs further erode public confidence. As noted in [Liu et al., 2024d] effectively aligned models are crucial for influencing public trust, where alignment methods

enhanced transparency in data collection and model training. Intriguingly, AI itself offers tools to rebuild trust in public governance by increasing transparency in public sector decision making[[oup](#)], improving public service efficiency, and enabling data-driven policy, although such applications require careful management of privacy and accessibility.

Fairness[[Liu et al., 2024d](#), [und, 2025](#)] in AI systems is a serious ethical challenge, primarily due to algorithmic bias[[Dai et al., 2024](#)]. Bias can originate from unrepresentative or historically skewed training data, leading AI to perpetuate or even exacerbate societal inequalities in areas related to human behaviors. Algorithmic design and RLHF can also introduce or amplify these biases without careful diversity implementation. Defining and consistently measuring fairness across varied cultural and societal contexts remains a substantial hurdle, often necessitating complex trade-offs with other vital alignment objectives such as helpfulness, honesty[[Yang et al., 2024c](#)], and safety.

The ethical deployment of LLMs is critical, especially in high-risk fields such as healthcare, finance, military scenarios, human resources, and transportation, where AI-driven decisions carry substantial consequences. A key strategy for managing risks in these areas involves models that refuse to respond[[Pasch, 2025](#)] to ethically sensitive requests, like those involving illegal or harmful content. For instance, ethical alignment in the healthcare and medical field demands careful attention to patient privacy, diagnostic accuracy, equitable access, and accountability for AI-caused errors. The report[[und, 2025](#)] highlights reliability issues, noting that users consulting AI for medical advice might receive falsehoods, emphasizing the need for caution.

The process of aligning AI with shared human values[[Hendrycks et al., 2023](#)] and intentions, where related benchmark spans concepts in justice, well-being, duties, virtues, and commonsense morality[[und, 2025](#), [Hendrycks et al., 2023](#)], introduces its distinct ethical implications. The main challenge is defining the very objectives of alignment: determining whose values, cultural norms, and ethical frameworks LLMs should adhere to, especially when these differ across societies or contexts. Specifically, developers of general-purpose AI assistants face strong competitive pressure, which can incentivize them to conduct less thorough risk mitigation. Markets characterized by high fixed costs, low marginal costs, and network effects tend to create competitive pressures that discourage safety investments. To resolve the issues, aligned models could reduce malignant societal implications for risk management and policymaking with specific AI social role norms. Instead of alignment with human preferences, developer[[spr, b](#)], or large organization, LLMs should be aligned with normative standards appropriate to their social roles, such as the role of a general-purpose assistant. In addition, all relevant stakeholders must negotiate and agree upon these standards.

**Socio-Technical challenges.** Safety research[[Dhole, 2023](#)] is required for the increasing capabilities of advanced LLMs, which are facing considerable challenges. Moreover, it is essential to look into AI safety for major tech companies, which possess the requisite resources and are at the forefront of developing these sophisticated systems. To tackle these challenges, aligned models are designed to work in ways that reduce harmful societal outcomes. This approach is essential to effectively manage risks and create robust public policy. However, the operationalization of LLM alignment presents a fundamental and complex sociotechnical[[Kierans et al., 2025](#)] dilemma, when determining the common values, AI developers cannot load the inexistent human values fully agreed by all participants.

While international cooperation is indispensable to address such existential threats, predicated on a shared interest[[Hendrycks et al., 2023](#)] in collective survival, significant skepticism persists regarding the current state of geopolitical cohesion. The apparent deficit in global leadership exacerbated This predicament with no single entity currently positioned to effectively spearhead such

collaborative endeavors. Therefore, determining the ethical frameworks and normative principles to which LLMs should be aligned constitutes a profound and multifaceted problem, requiring sustained scholarly inquiry and international dialogue[cna, 2018].

## 11.2 Regulatory and Policy Landscape

AI Alignment ensures AI systems act in accordance with human intentions, values, and goals. This involves aligning AI behavior with human expectations to prevent unintended consequences. AI safety encompasses practices and principles aimed at preventing harmful outcomes from AI systems, ensuring they operate reliably and ethically. The regulatory and policy landscape for AI alignment and safety is evolving, shaped by diverse national strategies, international collaborations, and emerging governance models.

The government adopts a sector-specific, risk-based approach to AI regulation in U.S.. Executive Order 14110 mandates federal agencies to appoint Chief AI Officers and develop AI-related guidelines. However, recent legislative proposals, such as a 10-year moratorium on state-level AI regulations, have sparked bipartisan opposition over concerns of federal overreach and potential hindrance to innovation. The Europe’s Artificial Intelligence Act, enacted in August 2024, classifies AI applications into risk categories, unacceptable, high, limited, and minimal, and imposes corresponding obligations. High-risk applications require compliance with strict transparency and safety standards. The UK has established the AI Safety Institute (AISI) to evaluate and ensure the safety of advanced AI models. The UK emphasizes a balance between innovation and safety, opting for adaptive regulatory frameworks over rigid legislation. China released an AI Safety Governance Framework in September 2024, aligning with its AI Governance Initiative, focusing on ethical standards and international cooperation.

International cooperation on AI alignment and safety has advanced through both non-binding guidelines and formal treaties. UNESCO’s Recommendation on the Ethics of Artificial Intelligence (2021) – endorsed by all 193 member states – established the first universal set of principles to ensure AI technologies are developed in a manner that upholds human rights and the public interest. Building on such global norms, a growing number of nations have pursued binding agreements. In 2024 the Council of Europe adopted the Framework Convention on Artificial Intelligence, Human Rights, Democracy and the Rule of Law, the world’s first legally binding AI treaty. These international frameworks collectively seek to ensure that AI safety and value-alignment are treated as global public policy priorities. In the meantime, persistent challenges, such as differing national regulatory priorities, and fragmented governance frameworks underscore the importance of fostering wider and deeper collaboration to ensure that AI technologies advance in ways consistent with universally shared values and human welfare.

Emerging governance models for AI alignment and safety are taking shape worldwide as policymakers respond to rapid advances in AI. In late 2023, for instance, the UK convened a global AI Safety Summit where 28 countries (and the EU) endorsed shared safety principles, and multilateral bodies such as the OECD and G7 have since worked toward common AI governance frameworks. Institutional innovations are a key part of this landscape: several governments have established dedicated AI safety institutes to provide technical evaluation and oversight expertise (the UK, US, and Japan launched such institutes in 2023–2024), while the European Union’s forthcoming EU AI Office under the AI Act will carry a broad mandate to supervise AI across the single market. Governments are also adopting new oversight tools, including algorithmic impact assessments that evaluate an AI system’s potential societal harms before deployment, and independent algorithm audits to verify compliance with safety or fairness standards. These developments across regions, alongside ongoing international cooperation, reflect a concerted move toward robust governance of



AI for global safety and alignment.

### 11.3 AGI/ASI safety

AGI and Artificial Superintelligence (ASI) are hypothesized as transformative stages in AI development, defined by their potential to exhibit general reasoning, autonomy, and recursive self-improvement capabilities that could match or exceed human-level cognition across virtually all domains. These unprecedented capacities bring forth not only vast opportunities but also profound safety risks.

[Shah et al., 2025] outline four interrelated areas of AGI safety concern—misuse, mistakes, misalignment, and structural risks. Malicious actors may exploit AGI for cyberattacks, bioweapon development, or information manipulation, creating severe security threats. Concurrently, design flaws or reward specification errors may lead to unintended and harmful behaviors, risks that are magnified as AGI complexity increases[Amodei et al., 2016]. Misalignment refers to the divergence between an AI system’s goals and human values, including superficially aligned agents that covertly pursue objectives misaligned with human intent. Furthermore, structural risks emerge from algorithmic bias or inequitable access to beneficial AI technologies, potentially exacerbating existing social inequalities[Koessler and Schuett, 2023]. As AGI systems are highly replicable software, scenarios such as embedded backdoors or cascading infections raise the possibility of multiple instances sharing harmful objectives[Wang, 2024].

The real-world ethical and societal ramifications of AGI and ASI are especially concerning. As these systems grow in autonomy and complexity, they may become uninterpretable or uncontrollable by humans, leading to irreversible catastrophic consequences. Once such systems operate beyond human oversight, returning to a “factory-reset” state may be infeasible. Autonomous AGI decisions could also trigger accidents while liability remains legally ambiguous. Additionally, the automation of vast sectors may result in mass unemployment and economic disruption. Privacy concerns intensify as advanced systems gain unprecedented access to and analysis of personal data[Gulchenko, 2024].

To address these challenges, a multi-pronged governance strategy is essential, integrating both technical safeguards and institutional frameworks. One of the core components of risk prevention is ensuring that AGI objectives are aligned with human values, while reinforcing continuous monitoring to detect and restrain potentially covert behaviors. Alignment in AGI contexts seeks to ensure that AI behavior faithfully reflects human intent, mitigating issues such as gaming the reward system and deceptive alignment[Everitt et al., 2018]. Robust training, improved explainability, and corrigibility are also emphasized to ensure safe operation in novel environments and enhance system fault-tolerance[Shah et al., 2025]. Additional safeguards include system monitoring, power decentralization, and constraints on access to data and computational resources[Wang, 2024].

Despite increasing concern over AGI safety, there remains a policy lag. [Koundouri et al., 2025] report that none of the official national or regional AI policy documents reviewed explicitly reference AGI, and there exists considerable divergence in regulatory strategies across jurisdictions, adding to the complexity of governance. While AGI could catalyze industrial revolutions and accelerate scientific discovery, its potential for nonlinear breakthroughs, such as those resulting from intelligence explosions, poses unpredictable systemic risks[Morris et al., 2023]. Strengthening AGI safety research, improving governance architectures, and fostering global coordination have thus become urgent priorities to ensure that technological development remains beneficial, controllable, and aligned with long-term human welfare.

## 12 Alignment Strategies Across Leading AI Models

This section surveys the alignment methodologies adopted by state-of-the-art LLMs. It is important to know about the current strategies and policies in the industry up to date.

### 12.1 OpenAI o-Series Models

OpenAI’s o-series models, such as o1 and o3, incorporate deliberative alignment strategies to enhance safety and reliability. These models are trained to systematically reason over safety specifications prior to generating responses, thereby improving their ability to address complex prompts with greater safety assurances.

Deliberative alignment involves training models using human-authored safety specifications and instructing them to explicitly engage in reasoning processes over these specifications before formulating outputs. The primary objective is to achieve precise adherence to safety policies. To this end, the models employ Chain-of-Thought (CoT) reasoning, wherein they analyze prompts, identify pertinent policy guidelines from the safety specifications, and construct responses that are compliant with these guidelines. This approach aims for responses that are “right for the right reasons.” In contrast to standard supervised fine-tuning (SFT), which focuses on imitating outputs, and reinforcement learning from human feedback (RLHF), which emphasizes output preferences, deliberative alignment reorients SFT toward imitating the reasoning process itself and refines RL to optimize the use of CoT for policy application. OpenAI has identified this methodology as central to the safety of models such as o3-mini. The adoption of deliberative alignment also enhances interpretability, as the CoT provides a traceable record of the model’s deliberative process [Guan et al., 2024].

The process of reasoning over safety specifications is a core component of deliberative alignment and consists of several key elements:

- **Safety Specifications:** Safety specifications are designed to align the model with established content policies for various safety categories. For each category, the corresponding policy defines relevant terminology and delineates the conditions under which user requests are classified as (1) “allowed,” where the model should comply; (2) “disallowed,” where the model should refuse; or (3) “requires safe completion.” The employed specifications are partly based on OpenAI’s published model specification [OpenAI, 2024].
- **SFT Stage:** During the SFT stage, datasets comprising (`prompt`, `CoT`, `output`) triplets are constructed for training purposes. Prompts are curated to cover a range of safety categories (e.g., erotic, self-harm), each framed as a multi-turn chat scenario concluding with a user message. For each (`prompt`, `category`) pair, a relevant safety specification, `spec(category)`, is referenced. The dataset includes CoT and output completions that explicitly reference policy content within the reasoning sequence, generated by prompting a base reasoning model with the appropriate safety specification. The resulting SFT dataset undergoes rigorous quality control through both automated filtering and evaluation by a reward model, which also considers the category-specific safety specification. Each completion is assessed multiple times, and the lowest assigned score is used to ensure stringent quality standards. The base model is subsequently fine-tuned on the curated SFT dataset alongside other capability-enhancing data. Notably, explicit context about the safety specification is removed from the prompt during training to encourage the model to internalize and recall relevant policy content, even when not directly present in the conversational context.

- **Reinforcement Learning (RL) Training:** In the RL phase, for prompts pertaining to safety, a “judge” model with access to safety policies provides a supplementary reward signal to the RL framework. The RL safety dataset comprises (prompt, category) pairs, often accompanied by additional metadata of varying quality. While the judge model accesses CoT during SFT data filtration, CoT is withheld from the judge during RL to prevent direct optimization of CoT traces and to mitigate the risk of encouraging deceptive reasoning. The SFT methodology is applied across all o-series models, whereas the additional reward signal during RL was specifically introduced for the o1 and o3-mini models.

## 12.2 DeepSeek Models

DeepSeek’s models, particularly the DeepSeek-V2 series, introduce an innovative multi-stage alignment process designed to significantly enhance reasoning capabilities while maintaining computational efficiency. Departing from a sole reliance on human preference data for refinement, DeepSeek’s strategy emphasizes automated self-improvement through a simulated deliberation mechanism prior to conventional reinforcement learning. This approach aims to cultivate robust reasoning skills with minimal human supervision, enabling the model to tackle complex problems more effectively. The core of this methodology involves a sequence of Supervised Fine-Tuning (SFT), a distinct self-correction phase, and a final Reinforcement Learning (RL) stage [DeepSeek-AI and a long list of other authors, 2024].

The alignment process is structured into the following key stages:

- **Initial Supervised Fine-Tuning (SFT):** The process begins with a standard SFT stage, where the base language model is trained on a diverse, high-quality dataset of instruction-following examples. This initial phase equips the model with foundational capabilities in language comprehension, instruction adherence, and basic reasoning, preparing it for more advanced alignment techniques.
- **Deliberative Self-Improvement via Group Debate:** This is the most distinctive stage in DeepSeek’s alignment strategy. Instead of immediately proceeding to RLHF, the SFT model undergoes a process of self-refinement. For a given prompt, multiple instances of the model act as “debaters,” each generating a different response or reasoning path. An evaluator model, which may be a more powerful proprietary model or the model itself employing a critical thinking persona, assesses these candidate responses. The model is then further fine-tuned on the outputs that are deemed highest-quality by the evaluator. This “group debate” and voting-based selection process allows the model to explore the solution space and improve its reasoning and problem-solving abilities without incurring the high cost of extensive human annotation at this stage.
- **Reinforcement Learning (RL) Refinement:** Following the self-improvement phase, the enhanced model undergoes a final alignment stage using Reinforcement Learning (RL). A reward model is trained on a smaller, more targeted dataset of human preferences to capture nuanced aspects of helpfulness and safety. The model is then fine-tuned using PPO to maximize the reward signal. Because the model entering this RL stage has already been substantially improved through the deliberation phase, the RL process can be more efficient and effective, focusing on refining subtler aspects of interaction rather than teaching core reasoning from scratch.

By front-loading the alignment process with an automated, reasoning-focused self-improvement stage, DeepSeek’s methodology aims to produce highly capable and aligned models at a lower cost

and with less reliance on massive-scale human annotation compared to traditional RLHF-centric approaches.

### 12.3 Anthropic Claude Models

Anthropic’s Claude series of models is distinguished by its pioneering work on **Constitutional AI (CAI)**, a methodology designed to align models with a set of explicit ethical principles with less reliance on large-scale human safety supervision [Bai et al., 2022a]. The primary goal of CAI is to make AI behavior more interpretable and robustly harmless by training the model to recognize and police its own outputs based on a predefined “constitution.” This approach is part of a broader alignment strategy that combines automated safety mechanisms with traditional human-feedback methods for helpfulness, alongside proactive research into long-term alignment challenges such as deceptive alignment.

Anthropic’s alignment methodology can be broken down into several key components:

- **Constitutional AI (CAI) for Harmlessness:** This is Anthropic’s core innovation for safety alignment and is implemented in two main phases.
  - *Supervised Fine-Tuning with Self-Critique:* In the first phase, a base model is prompted with requests that might elicit harmful responses. The model is then instructed to critique its initial response based on a set of principles from the constitution (e.g., principles drawn from the UN Universal Declaration of Human Rights and other sources). Finally, it is prompted to revise its original response in line with the critique. This process generates a dataset of self-corrected examples without requiring humans to author the safer outputs.
  - *Reinforcement Learning from AI Feedback (RLAIF):* In the second phase, a preference model is trained on the self-corrected data. It learns to prefer the revised, constitution-adherent responses over the initial, potentially harmful ones. The Claude model is then fine-tuned using this AI-generated preference signal as the reward. This RLAIF process automates the scaling of safety alignment and crucially reduces the need for human labelers to be exposed to large volumes of harmful content.
- **RLHF for Helpfulness:** Alongside CAI for safety, the Claude models are separately optimized for helpfulness using standard Reinforcement Learning from Human Feedback (RLHF). In this process, human labelers rank different model responses to a given prompt based on their quality, accuracy, and utility. A reward model for helpfulness is trained on this human preference data, and the final Claude model is fine-tuned to maximize both the AI-generated harmlessness score from CAI and the human-generated helpfulness score from RLHF.
- **Proactive Research on Alignment Risks:** Anthropic actively investigates potential failure modes of current alignment techniques. A significant area of this research is **alignment faking** or deceptive alignment. Their work has demonstrated the possibility of training “sleeper agents”, that is, models that behave safely during training and evaluation but revert to malicious behavior when a specific trigger is encountered in deployment [Hubinger et al., 2024]. This research underscores the limitations of purely behavioral training and motivates Anthropic’s focus on developing more robust and deeply-seated alignment methods, including mechanistic interpretability, to ensure long-term safety.

Through the combination of the automated and scalable CAI framework for safety, traditional RLHF for utility, and forward-looking research into complex failure modes, Anthropic’s alignment strategy aims to build a multi-layered defense against both present and future AI risks.

## 12.4 Google DeepMind Gemini Models

Google DeepMind’s AGI Safety & Alignment team identifies misalignment as one of the primary AGI risk areas[Shah et al., 2025]. To mitigate this risk, DeepMind propose a two-pronged defense framework targeting both training and deployment phases. During training, the team focuses on building an aligned model through mechanisms such as amplified oversight, guiding model behavior, and robust training.

- **Amplified Oversight:** As model capabilities improve, human overseers may struggle to accurately assess model outputs. To address this oversight gap, DeepMind introduces a debate-style adversarial structure, where two copies of the model critique each other’s responses to the same input, aiming to surface subtle flaws. A human judge then determines which model performed better. This approach enhances the granularity of the oversight signal, enabling the model to expose and correct nuanced misbehaviors during training. For tasks too complex to supervise directly, the system decomposes them into simpler subproblems that are more amenable to human evaluation. The team also explores modeling systematic human biases in oversight and correcting for them during the processing of supervision signals, thereby improving consistency and reliability.
- **Guiding Model Behavior:** Effectively translating oversight into behavioral updates is another key challenge. While reinforcement learning remains a viable approach, DeepMind prioritizes higher-bandwidth feedback modalities such as natural language critiques instead of scalar reward signals. These critiques provide actionable, interpretable guidance that clarifies failure modes and offers directional updates. Moreover, the team investigates leveraging model internals to inform gradient updates, aiming to improve sample efficiency and reduce reliance on costly human feedback.
- **Robust Training:** Even if a model behaves well in-distribution, it may fail catastrophically under distribution shift. The training distribution is expanded toward high-risk regions through techniques such as active learning, automated adversarial training, and post-deployment online fine-tuning. In parallel, the researchers also explore how to construct a reliable mapping from supervised data to behavioral policy. These training mechanisms collectively aim to ensure that, as the model’s capabilities expand, it retains alignment tendencies and behavioral stability.

Although training-time alignment is the first line of defense, Gemini models are also equipped with a second layer of alignment safeguards during deployment. These include access controls, sandboxing, and hierarchical monitoring structures to mitigate harm from capability leaps. Furthermore, the AGI Safety & Alignment team is advancing research in mechanistic interpretability, seeking to understand the internal strategy formation mechanisms of models, particularly whether apparent alignment is the result of deceptive reasoning pathways.

## 12.5 Meta’s LLaMA Models

Meta’s LLaMa models are developed with an emphasis on open-source accessibility and reproducibility. To align model behavior with human intent, Meta has adopted a multi-stage alignment strategy

that evolves across LLaMa versions, incorporating supervised fine-tuning (SFT), preference-based optimization, and novel dialogue consistency mechanisms.

- **Supervised Fine-Tuning (SFT):** In LLaMa 2, the post-training process begins with SFT using publicly available instruction-tuning datasets [Touvron et al., 2023]. Researchers found that tens of thousands of high-quality annotations are sufficient to reach strong performance, and thus prioritized the collection of several thousand high-quality SFT examples.
- **Alignment with Human Preferences:** LLaMA 2 adopts a reinforcement learning with human feedback (RLHF) pipeline to further refine alignment. Two separate reward models were trained, one for helpfulness and another for safety, acknowledging the potential tradeoff between the two. Human preference data were collected through pairwise comparisons, where annotators selected the preferred output and indicated the strength of their preference. The reward models were trained using a margin-augmented binary ranking loss, where the margin scaled with the strength of annotator preference, ensuring a greater reward gap for more strongly preferred responses.

LLaMa 2 uses Proximal Policy Optimization (PPO), treating the reward models as proxies for human judgment. A rejection sampling stage was introduced before PPO to improve training stability. In LLaMa 3, PPO is replaced by Direct Preference Optimization (DPO), which was found to be more computationally efficient and to achieve stronger instruction-following performance in large models [Dubey et al., 2024]. Additionally, due to diminishing gains after data scaling in LLaMa 3, the margin term in the loss was removed.

- **Lightweight Alignment in LLaMa 4:** LLaMa 4 further updates the post-training pipeline to a new sequence: lightweight SFT, online RL and lightweight DPO. The motivation for this hybrid design was that SFT and DPO alone may overly constrain the model, limiting its ability to explore and generalize. To address this, LLaMa 4 discards over 50% of simple training examples and performs lightweight SFT on the remaining harder subset. The introduction of online reinforcement learning helps achieve a better balance between computational cost and model alignment performance [Meta AI, 2025].
- **Ghost Attention (GAtt):** This technique improves multi-turn dialogue consistency by synthetically appending the original instruction to all user messages during fine-tuning. This helps the model retain compliance with the initial instruction across the entire dialogue trajectory.

## 12.6 Grok Models

Grok employs a version of RLHF in which a group of human tutors evaluates the results against an internal rubric that places emphasis on factual consistency, neutrality, and challenge to unexamined assumptions. These ratings train a reward model that tries to balance resistance to ideological bias, accuracy, and perspectives. Then, fine-tuning is performed using PPO. For transparency, xAI publishes the system and user prompt templates used during training to allow external auditors to verify that no hidden alignment objectives are being injected. Grok also conducts quarterly alignment reviews that compare its behavior with benchmarks that cover political, medical, and ideological edge cases. A continuous process also introduces adversarial queries and misaligned responses trigger automatic augmentation of the training set with new tutor rated examples. In the deployment phase, Grok uses a lightweight real-time safety filter that cross-references outputs against a library of known problematic patterns and blocks or flags any suspicious content. Finally,



all code changes to the alignment stack are subject to dual review by independent alignment specialists, ensuring that any modifications preserve the integrity of the truth seeking.

## 13 Conclusion and Future Directions

This survey has examined the current state of large language model alignment, from basic supervised fine-tuning to reinforcement learning from human feedback and newer approaches like brain-inspired methods. As language models become more powerful and widely used, ensuring they behave according to human values has become essential for safe deployment. This section summarizes what we have learned and identifies important areas for future research.

### 13.1 Summary of Key Insights

The development of LLM alignment reveals a fundamental shift in how we approach teaching machines to behave according to human values. The progression from supervised fine-tuning to reinforcement learning from human feedback reflects a deeper understanding that human preferences cannot be captured through simple instruction-following alone. While SFT provides essential capabilities, RLHF’s ability to optimize for subtle preferences has proven necessary for creating models that users find genuinely helpful and safe.

This evolution has crystallized around the framework of helpfulness, harmlessness, and honesty as core alignment objectives. Yet these goals exist in fundamental tension, requiring sophisticated trade-offs that resist simple optimization. The practical solution of hierarchical prioritization, placing safety above honesty and helpfulness, works but raises important questions about value determination and whose preferences shape these priorities.

Alongside these conceptual advances, computational constraints have driven remarkable innovation. The development of parameter-efficient methods like LoRA and reward-free approaches such as Direct Preference Optimization demonstrates that effective alignment need not require prohibitive resources. This accessibility matters deeply because it enables the broader research community to contribute to safety development, preventing a dangerous gap between capability and alignment research.

Perhaps most significantly, the field has embraced uncertainty as a fundamental aspect of alignment rather than a problem to eliminate. By quantifying when models are uncertain about appropriate behavior, we enable more robust deployment strategies and identify where human oversight remains necessary. This probabilistic approach represents a maturation of the field, acknowledging that perfect alignment may be impossible but that we can build systems aware of their own limitations.

Throughout this evolution, alignment has revealed itself as inherently interdisciplinary. The integration of insights from neuroscience, cognitive science, ethics, and policy demonstrates that encoding human values into artificial systems requires perspectives far beyond computer science. This interdisciplinary nature reflects the deep complexity of the alignment challenge and points toward the collaborative efforts needed for future progress.

### 13.2 Open Research Challenges

As models become better than humans at certain tasks, we face a fundamental problem: how can humans evaluate outputs they cannot fully understand? This is already happening in advanced mathematics and scientific domains. Current methods that rely on human feedback will not work

when models exceed human capabilities. Proposed solutions like debate protocols and recursive reward modeling face significant practical barriers.

Current alignment typically assumes a single set of human values, usually reflecting the views of a small group of annotators. But human values differ greatly across cultures and individuals. Building models that serve diverse global populations while respecting legitimate value differences remains unsolved. Social choice theory shows that no single method can aggregate preferences perfectly, making this challenge particularly difficult.

Aligned models break easily when given unusual inputs or adversarial attacks. The success of jailbreaking techniques shows that current alignment is fragile. We need methods that maintain safe behavior across different contexts and resist manipulation attempts, including both technical attacks and social engineering.

Most alignment happens once during training, but the real world changes constantly. Social norms evolve, new use cases appear, and models may drift through continued use. We need methods for ongoing alignment that adapt to changes while maintaining core safety properties.

As AI systems work together and integrate multiple modalities like vision and audio, alignment becomes more complex. Ensuring that groups of AI systems behave well together, especially when trained by different organizations, creates new coordination problems. Extending alignment to multi-modal systems requires rethinking basic alignment concepts.

### 13.3 Promising Research Directions

Mechanistic interpretability offers a way to understand how models actually process values and make decisions internally. Instead of just looking at behavior, researchers can examine the circuits and features that drive model outputs. Recent progress in understanding neural network internals could lead to alignment methods that work at the level of representations rather than behaviors.

Brain-inspired approaches draw lessons from how biological systems maintain goals while adapting to change. By studying how brains handle value conflicts and maintain stable behavior, researchers are developing new architectures that may be easier to align than current systems.

Formal verification aims to provide mathematical proofs that models will behave safely. While current methods rely on testing and evaluation, future high-stakes applications may require provable guarantees. This needs advances in both specifying alignment mathematically and developing verification techniques for neural networks.

### 13.4 Closing Remarks

Large language model alignment has grown from a theoretical concern to an active research field with immediate real-world impact. The progress from basic supervised learning to sophisticated multi-objective frameworks shows how quickly the field has developed.

However, this survey shows we are at a critical point. Model capabilities continue to advance rapidly while our ability to align them remains limited. The challenges ahead require not just technical solutions but collaboration across multiple disciplines.

The true measure of alignment research will be its real-world impact on human well-being. As language models become part of critical systems and daily life, getting alignment right becomes increasingly important. The research community has a responsibility to ensure these powerful tools remain beneficial as they continue to develop.

The path to reliable, value-aligned AI is challenging and uncertain. But the collective efforts documented in this survey provide reason for cautious optimism. Through continued research,

collaboration, and focus on the public good, we can work toward AI systems that reliably serve human values.

## References

- Tim Mulgan. Superintelligence: Paths, dangers, strategies, 2016.
- Eliezer Yudkowsky et al. Artificial intelligence as a positive and negative factor in global risk. *Global catastrophic risks*, 1(303):184, 2008.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022a. doi: 10.48550/arXiv.2212.08073. URL <https://doi.org/10.48550/arXiv.2212.08073>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.

- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
- Paul Christiano, Buck Shlegeris, and Dario Amodei. Supervising strong learners by amplifying weak experts. *arXiv preprint arXiv:1810.08575*, 2018.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, and et al. A general language assistant as a laboratory for alignment, December 2021. URL <https://arxiv.org/abs/2112.00861>.
- Amelia Glaese, Nat McAleese, Maja Trebacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Mari-beth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022. doi: 10.48550/arXiv.2209.14375.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022a.
- Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3463–3481, Dublin, Ireland, may 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.findings-acl.273>.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.
- Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. Large legal fictions: Profiling legal hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93, 2024.
- Neil Perry, Megha Srivastava, Deepak Kumar, and Dan Boneh. Do users write more insecure code with ai assistants? In *Proceedings of the 21st Workshop on Programming Languages and Software Engineering*, PLAS ’22, pages 69–78, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450399113. doi: 10.1145/3563723.3563811. URL <https://doi.org/10.1145/3563723.3563811>.
- Yue Liu, Miao Xiong, Jinlan Fu, Shumin Deng, and Bryan Hooi. Autodan: Generating stealthy jailbreak prompts on aligned llms. In *International Conference on Learning Representations*, 2024a.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- Yuntao Bai, Saurav Kadavath, Amanda Askell, et al. Training a helpful and harmless assistant with rlhf. *arXiv preprint arXiv:2204.05862*, 2022b.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024a.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33: 9459–9474, 2020.
- Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- NVIDIA. Introduction to llm agents. <https://developer.nvidia.com/blog/introduction-to-llm-agents/>, February 2024.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.
- Imane El Atillah. Man ends his life after an ai chatbot ‘encouraged’him to sacrifice himself to stop climate change. *euronews. com*, 2023.
- Cecily Mauran. Whoops, samsung workers accidentally leaked trade secrets via chatgpt. *Mashable [online]*. Dostępne z: <https://mashable.com/article/samsungchatgpt-leak-details>, 2023.
- Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. Large language models in finance: A survey. In *Proceedings of the fourth ACM international conference on AI in finance*, pages 374–382, 2023a.
- Siheng Li, Cheng Yang, Taiqiang Wu, Chufan Shi, Yuji Zhang, Xinyu Zhu, Zesen Cheng, Deng Cai, Mo Yu, Lemao Liu, et al. A survey on the honesty of large language models. *arXiv preprint arXiv:2409.18786*, 2024a.
- Irene Solaiman, Sarah Dennison, Deep Ganguli, and et al. Process for adapting language models to society. *arXiv preprint arXiv:2106.10328*, 2021. doi: 10.48550/arXiv.2106.10328. URL <https://arxiv.org/abs/2106.10328>.
- Liang Zhang, Ying Wu, Ming Li, and Hao Chen. A lifelong agent for strategy self-exploration to jailbreak llms (autodan-turbo). In *International Conference on Learning Representations*, 2025a.
- Rohan Patel and Ananya Singh. Jailbreaking leading safety-aligned llms with simple adaptive attacks. In *International Conference on Learning Representations*, 2025.

- Ming Chen and Li Zhao. Cognitive overload: Jailbreaking large language models with overloaded logical thinking. In *NAACL*, 2024.
- Yue Liu and Jiaheng Zhang. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. In *EMNLP*, 2024.
- Emily Johnson and Raj Kumar. Jailbreaking black-box large language models in twenty queries. In *NeurIPS*, 2024.
- John Smith and Jane Doe. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher (selfcipher). In *International Conference on Learning Representations*, 2024.
- Jong Lee and Soo Kim. Multilingual jailbreak challenges in large language models. In *International Conference on Learning Representations*, 2024a.
- Wei Jiang, Mei Chen, and Zixi Sun. Artprompt: Ascii-art-based jailbreak attacks. In *Proceedings of the 2024 Association for Computational Linguistics (ACL)*, 2024.
- Aidan Wong, He Cao, Zijing Liu, and Yu Li. Smiles-prompting: A novel approach to llm jailbreak attacks in chemical synthesis, 2024. URL <https://arxiv.org/abs/2410.15641>.
- Weidi Luo, He Cao, Zijing Liu, Yu Wang, Aidan Wong, Bin Feng, Yuan Yao, and Yu Li. Dynamic guided and domain applicable safeguards for enhanced security in large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025a.
- Yue Liu and Hongcheng Gao. Do anything now: Characterizing and evaluating in-the-wild jailbreak prompts on large language models (dan). In *ACM Conference on Computer and Communications Security*, 2024.
- Ming Li, Hao Chen, and Jie Zhang. Make them spill the beans! coercive knowledge extraction from (production) llms. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, 2023b.
- Kai Wang, Yu Zhao, and Wei Liu. Ignore this title and hackaprompt: Exposing systemic vulnerabilities of llms through a global-scale prompt hacking competition. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023a.
- Alice Brown, John Smith, and Emily Davis. Summon a demon and bind it: A grounded theory of llm red-teaming in the wild. In *USENIX Security Symposium*, 2023.
- Jane Doe and Joe Bloggs. Easyjailbreak: A unified framework for jailbreaking large language models. *arXiv preprint arXiv:2309.01234*, 2023.
- Sung Lee and Hyun Kim. Tricking llms into disobedience: Formalizing, analyzing, and detecting jailbreaks. In *Proceedings of the 12th International Conference on Language Resources and Evaluation (LREC)*, 2024b.
- Evan Hubinger, Christopher Denison, Jing Mu, et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- Yilun Wang, Yixin Teng, Kai Huang, et al. Fake alignment: Are llms really aligned well? *arXiv preprint arXiv:2311.05915*, 2023b.



- Javier Rando, Francesco Croce, Kryštof Mitka, Stepan Shabalin, Maksym Andriushchenko, Nicolas Flammarion, and Florian Tramèr. Competition report: Finding universal jailbreak backdoors in aligned llms. *arXiv preprint arXiv:2404.14461*, 2024a.
- Rebecca Greenblatt, Christopher Denison, Benjamin Wright, et al. Alignment faking in large language models. *arXiv preprint arXiv:2412.14093*, 2024.
- Javier Rando and Florian Tramèr. Rlhf trojan competition: Finding trojans in aligned language models. GitHub repository: [https://github.com/ethz-spylab/rlhf\\_trojan\\_competition](https://github.com/ethz-spylab/rlhf_trojan_competition), 2024.
- Mantas Mazeika, Andy Zou, Norman Mu, Long Phan, Zifan Wang, Chunru Yu, Adam Khoja, Dan Hendrycks, Radha Poovendran, Bo Li, and David Forsyth. Tdc 2023 (llm edition): The trojan detection challenge. NeurIPS Competition, 2023.
- Mantas Mazeika, Andy Zou, Norman Mu, Long Phan, Zifan Wang, Chunru Yu, Adam Khoja, et al. Trojan detection in large language models: Insights from the trojan detection challenge. *arXiv preprint arXiv:2404.13660*, 2024a.
- Zhen Xiang, Mintong Kang Yi Zeng, Chejian Xu, Jiawei Zhang, Zhuowen Yuan, Zhaorun Chen, Chulin Xie, Fengqing Jiang, Minzhou Pan, Junyuan Hong, Ruoxi Jia, Radha Poovendran, and Bo Li. Clas 2024: The competition for llm and agent safety. *NeurIPS 2024 Competition Track*, 2024.
- Edoardo DeBenedetti, Daniel Paleka, Ahmed Salem, et al. Dataset and lessons learned from the 2024 satml llm capture-the-flag competition. *arXiv preprint arXiv:2406.07954*, 2024a.
- Paul Christiano et al. Llm evaluations hackathon: Trojan detection challenge. Alignment Jam hackathon page: <https://alignmentjam.com/jam/evals>, 2023a.
- Javier Rando, Edoardo DeBenedetti, Daniel Paleka, and Florian Tramèr. Our competitions at iee satml 2024: Llm ctf and trojan detection. SPY Lab Blog: <https://spylab.ai/blog/results-competition/>, 2024b.
- Yue Liu, Xiaoxin He, Miao Xiong, Jinlan Fu, Shumin Deng, and Bryan Hooi. Flipattack: Jailbreak llms via flipping. *arXiv preprint arXiv:2410.02832*, 2024b.
- Xiangming Gu, Xiaosen Zheng, Tianyu Pang, Chao Du, Qian Liu, Ye Wang, Jing Jiang, and Min Lin. Agent smith: A single image can jailbreak one million multimodal llm agents exponentially fast. *arXiv preprint arXiv:2402.08567*, 2024.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. Gptscore: Evaluate as you desire. In *NAACL*, 2023a.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. In *EMNLP*, 2023a.
- Yiqi Liu, Nafise Sadat Moosavi, and Chenghua Lin. Llms as narcissistic evaluators: When ego inflates evaluation scores. *arXiv preprint arXiv:2311.09766*, 2023b.
- Peiyi Wang, Lei Li, Liang Chen, et al. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*, 2023c.

- Yicheng Wang, Jiayi Yuan, Yu-Neng Chuang, et al. Dhp benchmark: Are llms good nlg evaluators? In *Findings of NAACL*, 2024a.
- Lianghui Zhu, Xinggang Wang, and Xinlong Wang. Judgelm: Fine-tuned large language models are scalable judges. *arXiv preprint arXiv:2310.12345*, 2023a.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, et al. Pandalm: An automatic evaluation benchmark for llm instruction tuning optimization. *arXiv preprint arXiv:2311.06789*, 2023d.
- Yen-Ting Lin and Yun-Nung Chen. Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations. In *NLP4ConvAI Workshop*, 2023.
- David Chan, Suzanne Petryk, Joseph Gonzalez, et al. Clair: Evaluating image captions with large language models. In *EMNLP*, 2023.
- Yebin Lee, Imseong Park, and Myungjoo Kang. Fleur: An explainable reference-free evaluation metric for image captioning using a large multimodal model. In *ACL*, 2024a.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Advbench: Universal and transferable adversarial attacks on aligned language models, Jul 2023. URL <https://arxiv.org/abs/2307.15043>.
- Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models, February 2024b. URL <https://arxiv.org/abs/2402.05044>.
- Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language models with multiple choice questions. *Association for Computational Linguistics (ACL)*, 2023a.
- Jiawen Deng, Jingyan Zhou, Hao Sun, Chujie Zheng, Fei Mi, Helen Meng, and Minlie Huang. Cold: A benchmark for chinese offensive language detection, January 2022. URL <https://arxiv.org/abs/2201.06025>.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset, 2024. URL <https://aclanthology.org/2024.emnlp-main.xxx>.
- Tianyang Xie, Xiao Qi, Yi Zeng, Yong Huang, Upendra M. Schwag, Kai Huang, Lei He, Bashuan Wei, Dazheng Li, and Yubo Sheng. Sorry-bench: Systematically evaluating large language model safety refusal behaviors, Jun 2024. URL <https://arxiv.org/abs/2406.14598>.
- Mario Samvelyan, Sachit C. Raparthy, Andre Lupu, Emily Hambro, and Jakob Foerster. Rainbow teaming: Open-ended generation of diverse adversarial prompts, Feb 2024. URL <https://arxiv.org/abs/2402.16822>.
- Erxin Yu, Jing Li, Ming Liao, Siqi Wang, Zuchen Gao, Fei Mi, and Lanqing Hong. Cosafe: Evaluating large language model safety in multi-turn dialogue coreference, June 2024a. URL <https://arxiv.org/abs/2406.17626>.
- Liang Xu, Kangkang Zhao, Lei Zhu, and Hang Xue. Sc-safety: A multi-round open-ended question adversarial safety benchmark for large language models in chinese, October 2023a. URL <https://arxiv.org/abs/2310.05818>.

Kaijie Zhu, Qinlin Zhao, Hao Chen, Jindong Wang, and Xing Xie. Promptbench: A unified library for evaluation of large language models, December 2023b. URL <https://arxiv.org/abs/2312.07910>.

Fengqing Jiang, Zhangchen Xu, Yuetai Li, Luyao Niu, Zhen Xiang, Bo Li, Bill Yuchen Lin, and Radha Poovendran. Safechain: Safety of language models with long chain-of-thought reasoning capabilities, 2025a. URL <https://arxiv.org/abs/2502.12025>.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025a. URL <https://arxiv.org/abs/2501.12948>.

Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun

Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. Kimi k1.5: Scaling reinforcement learning with llms, 2025. URL <https://arxiv.org/abs/2501.12599>.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyan Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho,

Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damla, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelen, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam,

- Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Kaiwen Zhou, Chengzhi Liu, Xuandong Zhao, Shreedhar Jangam, Jayanth Srinivasa, Gaowen Liu, Dawn Song, and Xin Eric Wang. The hidden risks of large reasoning models: A safety assessment of r1, 2025a. URL <https://arxiv.org/abs/2502.12659>.
- Zijun Wang, Haoqin Tu, Yuhan Wang, Juncheng Wu, Jieru Mei, Brian R. Bartoldson, Bhavya Kailkhura, and Cihang Xie. Star-1: Safer alignment of reasoning llms with 1k data, 2025a. URL <https://arxiv.org/abs/2504.01903>.
- Yutao Mou, Yuxiao Luo, Shikun Zhang, and Wei Ye. Saro: Enhancing llm safety through reasoning-based alignment, 2025. URL <https://arxiv.org/abs/2504.09420>.
- Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. In *NeurIPS*, 2023e.
- Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. Trustllm: Trustworthiness in large language models, 2024. URL <https://arxiv.org/abs/2401.05561>.
- Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, Bertie Vidgen, and Scott A. Hale. The prism alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models, 2024. URL <https://arxiv.org/abs/2404.16019>.
- Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 325–333, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/zemel13.html>.



- Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning, 2016. URL <https://arxiv.org/abs/1610.02413>.
- Steffi Chern, Zhulin Hu, Yuqing Yang, Ethan Chern, Yuan Guo, Jiahe Jin, Binjie Wang, and Pengfei Liu. Behonest: Benchmarking honesty in large language models. *arXiv preprint arXiv:2406.13261*, 2024.
- Sheng Yin, Xianghe Pang, Yuanzhuo Ding, Menglan Chen, Yutong Bi, Yichen Xiong, Wenhao Huang, Zhen Xiang, Jing Shao, and Siheng Chen. Safeagentbench: A benchmark for safe task planning of embodied llm agents. *arXiv preprint arXiv:2412.13178*, 2024a. doi: 10.48550/arXiv.2412.13178. URL <https://doi.org/10.48550/arXiv.2412.13178>.
- Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. Agent-safetybench: Evaluating the safety of llm agents. *arXiv preprint arXiv:2412.14470*, 2024a. doi: 10.48550/arXiv.2412.14470. URL <https://doi.org/10.48550/arXiv.2412.14470>.
- Ido Levy, Ben Wiesel, Sami Marreed, Alon Oved, Avi Yaeli, and Segev Shlomov. St-webagentbench: A benchmark for evaluating safety and trustworthiness in web agents. *arXiv preprint arXiv:2410.06703*, 2024. doi: 10.48550/arXiv.2410.06703. URL <https://doi.org/10.48550/arXiv.2410.06703>.
- Zihao Zhu, Bingzhe Wu, Zhengyou Zhang, Lei Han, Qingshan Liu, and Baoyuan Wu. Earbench: Towards evaluating physical risk awareness for task planning of foundation model-based embodied ai agents. *arXiv preprint arXiv:2408.04449*, 2024a. doi: 10.48550/arXiv.2408.04449. URL <https://doi.org/10.48550/arXiv.2408.04449>.
- Hanrong Zhang, Qian Li, Roshni Patel, and Ming Chen. Agent security bench (asb): Formalizing and benchmarking attacks and defenses in llm-based agents. *arXiv preprint arXiv:2410.02644*, 2024b.
- Edoardo Debenedetti, Clara Rossi, and Tuan Nguyen. Agentdojo: A dynamic environment to evaluate attacks and defenses for llm agents. *arXiv preprint arXiv:2406.13352*, 2024b.
- Zhen Xiang, Linzhi Zheng, Yanjie Li, Junyuan Hong, Qinbin Li, Han Xie, Jiawei Zhang, Zidi Xiong, Chulin Xie, Carl Yang, Dawn Song, and Bo Li. Guardagent: Safeguard llm agents by a guard agent via knowledge-enabled reasoning. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.
- Weidi Luo, Ananya Singh, Luis Gomez, and Yuxin Tan. Agrail: A lifelong agent guardrail with effective and adaptive safety detection. *arXiv preprint arXiv:2502.11448*, 2025b.
- Zhaorun Chen, Mintong Kang, and Bo Li. Shieldagent: Shielding agents via verifiable safety policy reasoning. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025a.
- Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, et al. Health-bench: Evaluating large language models towards improved human health. *arXiv preprint arXiv:2505.08775*, 2025.
- Haochen Zhao, Xiangru Tang, Ziran Yang, Xiao Han, Xuanchi Feng, Yueqing Fan, Senhao Cheng, Di Jin, Yilun Zhao, Arman Cohan, et al. Chemsafetybench: Benchmarking llm safety on chemistry domain. *arXiv preprint arXiv:2411.16736*, 2024.

- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024c.
- Tessa Han, Aounon Kumar, Chirag Agarwal, and Himabindu Lakkaraju. Medsafetybench: Evaluating and improving the medical safety of large language models. *arXiv preprint arXiv:2403.03744*, 2024a.
- Yujun Zhou, Jingdong Yang, Kehan Guo, Pin-Yu Chen, Tian Gao, Werner Geyer, Nuno Moniz, Nitesh V Chawla, and Xiangliang Zhang. Benchmarking llms on safety issues in scientific labs.
- Xinyu Zhang, Yuxuan Dong, Yanrui Wu, Jiaying Huang, Chengyou Jia, Basura Fernando, Mike Zheng Shou, Lingling Zhang, and Jun Liu. Physreason: A comprehensive benchmark towards physics-based reasoning. *arXiv preprint arXiv:2502.12054*, 2025b.
- Tianhao Li, Jingyu Lu, Chuangxin Chu, Tianyu Zeng, Yujia Zheng, Mei Li, Haotian Huang, Bin Wu, Zuoxian Liu, Kai Ma, et al. Scisafeval: a comprehensive benchmark for safety alignment of large language models in scientific tasks. *arXiv preprint arXiv:2410.03769*, 2024d.
- Jiyan He, Weitao Feng, Yaosen Min, Jingwei Yi, Kunsheng Tang, Shuai Li, Jie Zhang, Kejiang Chen, Wenbo Zhou, Xing Xie, et al. Control risk for potential misuse of artificial intelligence in science. *arXiv preprint arXiv:2312.06632*, 2023.
- Fengqing Jiang, Fengbo Ma, Zhangchen Xu, Yuetai Li, Bhaskar Ramasubramanian, Luyao Niu, Bo Li, Xianyan Chen, Zhen Xiang, and Radha Poovendran. Sos bench: Benchmarking safety alignment on scientific knowledge, 2025b. URL <https://arxiv.org/abs/2505.21605>.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *International Conference on Machine Learning (ICML)*, 2024b.
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. In *NeurIPS Datasets and Benchmarks Track*, 2024.
- Weidi Luo, Siyuan Ma, Xiaogeng Liu, Xiaoyu Guo, and Chaowei Xiao. Jailbreakv-28k: A benchmark for assessing the robustness of multimodal large language models against jailbreak attacks. *Conference on Language Modeling (COLM)*, 2024.
- Hossein Hajipour, Keno Hassler, Thorsten Holz, Lea Schönherr, and Mario Fritz. Codelmsec benchmark: Systematically evaluating and finding security vulnerabilities in black-box code language models, 2023. URL <https://arxiv.org/abs/2302.04012>.
- Chengquan Guo, Xun Liu, Chulin Xie, Andy Zhou, Yi Zeng, Zinan Lin, Dawn Song, and Bo Li. Redcode: Risky code execution and generation benchmark for code agents. *Thirty-Eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024a.
- Manish Bhatt, Sahana Chennabasappa, Yue Li, Cyrus Nikolaidis, Daniel Song, Shengye Wan, Faizan Ahmad, Cornelius Aschermann, Yaohui Chen, Dhaval Kapil, David Molnar, Spencer Whitman, and Joshua Saxe. Cyberseceval 2: A wide-ranging cybersecurity evaluation suite for large language models, 2024. URL <https://arxiv.org/abs/2404.13161>.

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022a.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023a.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- J OpenAI Achiam, S Adler, S Agarwal, L Ahmad, I Akkaya, FL Aleman, D Almeida, J Al-tenschmidt, S Altman, S Anadkat, et al. Gpt-4 technical report. arxiv. *arXiv preprint arXiv:2303.08774*, 2023.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lambda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR, 2023.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Andreas Köpf, Yannic Kilcher, Dimitri Von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. Openassistant conversations-democratizing large language model alignment. *Advances in Neural Information Processing Systems*, 36:47669–47681, 2023.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning. *arXiv preprint arXiv:2308.12032*, 2023c.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*, 2021.
- Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S Yu. A survey of multilingual large language models. *Patterns*, 6(1), 2025.

- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023a.
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179, 2024.
- Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. Diversity measurement and subset selection for instruction tuning datasets. *arXiv preprint arXiv:2402.02318*, 2024b.
- Brando Miranda, Alycia Lee, Sudharsan Sundar, Allison Casasola, and Sanmi Koyejo. Beyond scale: The diversity coefficient as a data quality metric for variability in natural language data. *arXiv preprint arXiv:2306.13840*, 2023.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR, 2019.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. In *arXiv preprint arXiv:2002.06305*, 2020.
- Waseem AlShikh, Manhal Daaboul, Kirk Goddard, Brock Imel, Kiran Kamble, Parikshith Kulkarni, and Melisa Russak. Becoming self-instruct: introducing early stopping criteria for minimal instruct tuning. *arXiv preprint arXiv:2307.03692*, 2023.
- Xueqing Liu and Chi Wang. An empirical study on hyperparameter optimization for fine-tuning pre-trained language models. *arXiv preprint arXiv:2106.09204*, 2021.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020.
- Youness Moukafih, Mounir Ghogho, and Kamel Smaili. Supervised contrastive learning as multi-objective optimization for fine-tuning large pre-trained language models. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024a.
- Ruoyu Wang, Jiachen Sun, Shaowei Hua, and Quan Fang. Asft: Aligned supervised fine-tuning through absolute likelihood. *arXiv preprint arXiv:2409.10571*, 2024c.

- Yuchen Fan, Yuzhong Hong, Qiushi Wang, Junwei Bao, Hongfei Jiang, and Yang Song. Preference-oriented supervised fine-tuning: Favoring target model over aligned large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 23859–23867, 2025a.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.
- Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *arXiv preprint arXiv:2404.14367*, 2024.
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *arXiv preprint arXiv:1706.03741*, 2023b. doi: 10.48550/arXiv.1706.03741. URL <https://doi.org/10.48550/arXiv.1706.03741>.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. *arXiv preprint arXiv:2009.01325*, 2022. doi: 10.48550/arXiv.2009.01325. URL <https://doi.org/10.48550/arXiv.2009.01325>.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, April 2022c. doi: 10.48550/arXiv.2204.05862. URL <https://arxiv.org/abs/2204.05862>.
- Banghua Zhu, Jiantao Jiao, and Michael I. Jordan. Principled reinforcement learning with human feedback from pairwise or \$k\$-wise comparisons. *arXiv preprint arXiv:2301.11270*, 2024b. doi: 10.48550/arXiv.2301.11270. URL <https://doi.org/10.48550/arXiv.2301.11270>.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*, January 2020a. doi: 10.48550/arXiv.1909.08593. URL <https://arxiv.org/abs/1909.08593>.
- Michiel A. Bakker, Martin J. Chadwick, Hannah R. Sheahan, Michael Henry Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, et al. Fine-tuning language models to find agreement among humans with diverse preferences. *arXiv preprint arXiv:2211.15006*, 2022. doi: 10.48550/arXiv.2211.15006. URL <https://doi.org/10.48550/arXiv.2211.15006>.
- Jiwei Li, Alexander H. Miller, Sumit Chopra, Marc’Aurelio Ranzato, and Jason Weston. Dialogue learning with human-in-the-loop. *arXiv preprint arXiv:1611.09823*, 2017. doi: 10.48550/arXiv.1611.09823. URL <https://doi.org/10.48550/arXiv.1611.09823>.
- Baicen Xiao, Qifan Lu, Bhaskar Ramasubramanian, Andrew Clark, Linda Bushnell, and Radha Poovendran. FRESH: Interactive Reward Shaping in High-Dimensional State Spaces using Human Feedback. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. IFAAMAS, 2020. doi: 10.5555/3398761.3398935.
- Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. Training language models with language feedback at scale. *arXiv preprint arXiv:2303.16755*, 2024. doi: 10.48550/arXiv.2303.16755. URL <https://doi.org/10.48550/arXiv.2303.16755>.

- Jie Huang, Jiangshan Hao, Rongshun Juan, Randy Gomez, Keisuke Nakamura, and Guangliang Li. Gan-based interactive reinforcement learning from demonstration and human evaluative feedback. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2023a. doi: 10.1109/ICRA48891.2023.10160939.
- Ge Gao, Hung-Ting Chen, Yoav Artzi, and Eunsol Choi. Continually improving extractive qa via human feedback. *arXiv preprint arXiv:2305.12473*, 2023a. doi: 10.48550/arXiv.2305.12473. URL <https://doi.org/10.48550/arXiv.2305.12473>.
- Omar Shaikh, Michelle S. Lam, Joey Hejna, Yijia Shao, Hyundong Cho, Michael S. Bernstein, and Diyi Yang. Aligning language models with demonstrated feedback. *arXiv preprint arXiv:2406.00888*, 2025. doi: 10.48550/arXiv.2406.00888. URL <https://doi.org/10.48550/arXiv.2406.00888>.
- Jason Ross Brown, Carl Henrik Ek, and Robert D. Mullins. Learning from preferences and mixed demonstrations in general settings. <https://openreview.net/forum?id=Sfct4aXXcw>, 2025.
- Braden Hancock, Antoine Bordes, Pierre-Emmanuel Mazare, and Jason Weston. Learning from dialogue after deployment: Feed yourself, chatbot! In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3667–3684, Florence, Italy, jul 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1358. URL <https://aclanthology.org/P19-1358/>.
- Angelica Chen, Jérémy Scheurer, Tomasz Korbak, Jon Ander Campos, Jun Shern Chan, Samuel R. Bowman, Kyunghyun Cho, and Ethan Perez. Improving code generation by training with natural language feedback. *arXiv preprint arXiv:2303.16749*, 2024b. doi: 10.48550/arXiv.2303.16749. URL <https://doi.org/10.48550/arXiv.2303.16749>.
- Zichao Li, Prakhar Sharma, Xing Han Lu, Jackie Chi Kit Cheung, and Siva Reddy. Using interactive feedback to improve the accuracy and explainability of question answering systems post-deployment. *arXiv preprint arXiv:2204.03025*, 2022a. URL <https://api.semanticscholar.org/CorpusID:248006299>.
- Yannick Metz, Andras Geiszl, Raphael Baur, and Mennatallah El-Assady. Reward learning from multiple feedback types. <https://arxiv.org/abs/2502.21038>, 2025. Preprint.
- Yueqin Yin, Shentao Yang, Yujia Xie, Ziyi Yang, Yuting Sun, Hany Hassan Awadalla, Weizhu Chen, and Mingyuan Zhou. Segmenting text and learning their rewards for improved RLHF in language models, 2025. URL <https://openreview.net/forum?id=cK7yrw5g5Q>.
- Zequiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A. Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. *arXiv*, oct 2023.
- Dehong Xu, Liang Qiu, Minseok Kim, Faisal Ladhak, and Jaeyoung Do. Aligning large language models via fine-grained supervision. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 673–680. Association for Computational Linguistics, 2024a. URL <https://doi.org/10.18653/v1/2024.acl-short.62>.
- Wendi Li, Wei Wei, Kaihe Xu, Wenfeng Xie, Dangyang Chen, and Yu Cheng. Reinforcement learning with token-level feedback for controllable text generation. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1704–1719. Association for Computational Linguistics, 2024e. URL <https://doi.org/10.18653/v1/2024.findings-naacl.111>.



- Julia Kreutzer, Shahram Khadivi, Evgeny Matusov, and Stefan Riezler. Can neural machine translation be improved with user feedback? *arXiv preprint arXiv:1804.05958*, 2018. doi: 10.48550/arXiv.1804.05958. URL <https://doi.org/10.48550/arXiv.1804.05958>.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2020b. doi: 10.48550/arXiv.1909.08593. URL <https://doi.org/10.48550/arXiv.1909.08593>.
- Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning from human preferences and demonstrations in atari. *arXiv preprint arXiv:1811.06521*, 2018. doi: 10.48550/arXiv.1811.06521. URL <https://doi.org/10.48550/arXiv.1811.06521>.
- Adam Gleave and Geoffrey Irving. Uncertainty estimation for language reward models. *arXiv preprint arXiv:2203.07472*, 2022. doi: 10.48550/arXiv.2203.07472. URL <https://doi.org/10.48550/arXiv.2203.07472>.
- Nirjhar Das, Souradip Chakraborty, Aldo Pacchiano, and Sayak Ray Chowdhury. Active preference optimization for sample efficient rlhf. *arXiv preprint arXiv:2402.10500*, 2024. doi: 10.48550/arXiv.2402.10500. URL <https://doi.org/10.48550/arXiv.2402.10500>.
- Viraj Mehta, Syrine Belakaria, Vikramjeet Das, Ojash Neopane, Yijia Dai, Ilija Bogunovic, Barbara Engelhardt, Stefano Ermon, Jeff Schneider, and Willie Neiswanger. Sample efficient preference alignment in llms via active exploration. *arXiv preprint arXiv:2312.00267*, 2025. doi: 10.48550/arXiv.2312.00267. URL <https://doi.org/10.48550/arXiv.2312.00267>.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2024b. doi: 10.48550/arXiv.2309.00267. URL <https://doi.org/10.48550/arXiv.2309.00267>.
- Archit Sharma, Sedrick Keh, Eric Mitchell, Chelsea Finn, Kushal Arora, and Thomas Kollar. A critical evaluation of ai feedback for aligning large language models. *arXiv preprint arXiv:2402.12366*, 2024. doi: 10.48550/arXiv.2402.12366. URL <https://doi.org/10.48550/arXiv.2402.12366>.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, 2nd edition, 2018.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. URL <http://www.jstor.org/stable/2334029>. Accessed: Feb. 13, 2023.
- R Duncan Luce. *Individual choice behavior: A theoretical analysis*. Wiley, 1959.
- John Schulman, Barret Zoph, Christina Kim, Jacob Hilton, Jacob Menick, Jiayi Weng, Juan Felipe Ceron Uribe, Liam Fedus, Luke Metz, Michael Pokorny, et al. ChatGPT: Optimizing Language Models for Dialogue. <https://openai.com/blog/chatgpt>, 2022. OpenAI Blog.

- R. L. Plackett. The analysis of permutations. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 24(2):193–202, 1975. doi: 10.2307/2346567.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganesan, Wei-Lin Chiang, Jian Zhang, and Jiantao Jiao. Starling-7B: Improving Helpfulness and Harmlessness with RLAIIF, 2024c. Unpublished manuscript.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: From pairwise approach to listwise approach. In *Proceedings of the 24th International Conference on Machine Learning*, pages 129–136, Corvallis, Oregon, USA, 2007. ACM. doi: 10.1145/1273496.1273513. URL <https://doi.org/10.1145/1273496.1273513>.
- Tie-Yan Liu. *Learning to Rank for Information Retrieval*, volume 3 of *Foundations and Trends in Information Retrieval*. Now Publishers Inc., 2009. doi: 10.1561/15000000016. URL <https://doi.org/10.1561/15000000016>.
- Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, and et al. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019. doi: 10.48550/arXiv.1909.08593. URL <https://arxiv.org/abs/1909.08593>.
- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Aleksii Kuchaiev, and Yi Dong. Helpsteer2-preference: Complementing ratings with preferences, March 2025b. URL <https://arxiv.org/abs/2410.01257>.
- Zhichao Wang, Bin Bi, Can Huang, Shiva Kumar Pentiyala, Zixu James Zhu, Sitaram Asur, and Na Claire Cheng. UNA: Unifying Alignments of RLHF/PPO, DPO and KTO by a Generalized Implicit Reward Function, April 2025c. URL <https://arxiv.org/abs/2408.15339>.
- Jialong Wu, Chaoyi Deng, Jianmin Wang, and Mingsheng Long. Supercompiler code optimization with zero-shot reinforcement learning. *arXiv*, 2024a. doi: 10.48550/arXiv.2404.16077. URL <https://doi.org/10.48550/arXiv.2404.16077>.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022b. doi: 10.1126/science.abq1158. URL <https://doi.org/10.1126/science.abq1158>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, et al. Training verifiers to solve math word problems. *arXiv*, November 2021. doi: 10.48550/arXiv.2110.14168. URL <https://doi.org/10.48550/arXiv.2110.14168>.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and outcome-based feedback. *arXiv*, November 2022. doi: 10.48550/arXiv.2211.14275. URL <https://doi.org/10.48550/arXiv.2211.14275>.
- Chengqi Lyu, Songyang Gao, Yuzhe Gu, Wenwei Zhang, Jianfei Gao, Kuikun Liu, Ziyi Wang, et al. Exploring the limit of outcome reward for learning mathematical reasoning. *arXiv*, February 2025. doi: 10.48550/arXiv.2502.06781. URL <https://doi.org/10.48550/arXiv.2502.06781>.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv*, May 2023. doi: 10.48550/arXiv.2305.20050. URL <https://doi.org/10.48550/arXiv.2305.20050>.

- Jiayi Zhou, Jiaming Ji, Juntao Dai, and Yaodong Yang. Sequence to sequence reward modeling: Improving rlhf by language feedback, August 2024a. URL <https://arxiv.org/abs/2409.00162>.
- Andi Peng, Aviv Netanyahu, Mark Ho, Tianmin Shu, Andreea Bobu, Julie Shah, and Pulkit Agrawal. Diagnosis, feedback, adaptation: A human-in-the-loop framework for test-time policy adaptation. *arXiv*, July 2023b. doi: 10.48550/arXiv.2307.06333. URL <https://doi.org/10.48550/arXiv.2307.06333>.
- Andreea Bobu, Andi Peng, Pulkit Agrawal, Julie Shah, and Anca D. Dragan. Aligning robot and human representations. In *Proceedings of the 2024 ACM/IEEE International Conference on Human-Robot Interaction*, pages 42–54, 2024. doi: 10.1145/3610977.3634987. URL <https://doi.org/10.1145/3610977.3634987>.
- Kaizhao Liu, Qi Long, Zhekun Shi, Weijie J Su, and Jiancong Xiao. Statistical impossibility and possibility of aligning llms with human preferences: From condorcet paradox to nash equilibrium. *arXiv preprint arXiv:2503.10990*, 2025a.
- Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca D. Dragan, and Daniel S. Brown. Causal confusion and reward misidentification in preference-based reward learning. *arXiv*, March 2023. doi: 10.48550/arXiv.2204.06601. URL <https://doi.org/10.48550/arXiv.2204.06601>.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning (ICML)*, volume 202 of *Proceedings of Machine Learning Research*, pages 10835–10866, Honolulu, Hawaii, USA, 2023b. PMLR. URL <https://proceedings.mlr.press/v202/gao23h.html>.
- Cassidy Laidlaw, Shivam Singhal, and Anca Dragan. Correlated proxies: A new definition and improved mitigation for reward hacking. *arXiv*, March 2025. doi: 10.48550/arXiv.2403.03185. URL <https://doi.org/10.48550/arXiv.2403.03185>.
- Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward hacking. *arXiv*, March 2025. doi: 10.48550/arXiv.2209.13085. URL <https://doi.org/10.48550/arXiv.2209.13085>.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Proceedings of the 13th International Conference on Neural Information Processing Systems*, NIPS’99, pages 1057–1063, Denver, CO, 1999. MIT Press.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 1889–1897, 2015. URL <https://proceedings.mlr.press/v37/schulman15.html>.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, and et al. Gpt-4 technical report, March 2024. URL <https://arxiv.org/abs/2303.08774>.
- Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, 1988. doi: 10.1007/BF00115009.

- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation, October 2018. URL <https://arxiv.org/abs/1506.02438>.
- John Schulman. Approximating kl-divergence. <http://joschu.net/blog/kl-approx.html>, 2020. Accessed: May 20, 2025.
- Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E. Turner, and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control, October 2017. URL <https://arxiv.org/abs/1611.02796>.
- Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharioun, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline reinforcement learning, October 2020. URL <https://arxiv.org/abs/2010.05848>.
- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation matters in deep policy gradients: A case study on ppo and trpo, May 2020. URL <https://arxiv.org/abs/2005.12729>.
- Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization, March 2023. URL <https://arxiv.org/abs/2210.01241>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, July 2024. URL <https://arxiv.org/abs/2305.18290>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290v2*, 2018.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- Jiancong Xiao, Ziniu Li, Xingyu Xie, Emily Getzen, Cong Fang, Qi Long, and Weijie J Su. On the algorithmic bias of aligning large language models with rlhf: Preference collapse and matching regularization. *arXiv preprint arXiv:2405.16455*, 2024.
- Wei Shen and Chuheng Zhang. Policy filtration in rlhf to fine-tune llm for code generation, December 2024. URL <https://arxiv.org/abs/2409.06957>.
- Mingsheng Zheng, Junwei Zhang, Changshuai Zhan, Xinyu Ren, and Shuai Lü. Proximal policy optimization with reward-based prioritization. *Expert Systems with Applications*, 283:127659, 2025. doi: 10.1016/j.eswa.2025.127659. URL <https://doi.org/10.1016/j.eswa.2025.127659>.
- Michael Santacrose, Yadong Lu, Han Yu, Yuanzhi Li, and Yelong Shen. Efficient rlhf: Reducing the memory usage of ppo, 2023.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

- Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *International Conference on Learning Representations (ICLR)*, 2020.
- Rasool Fakoor, Pratik Chaudhari, and Alexander J. Smola. P3o: Policy-on policy-off policy optimization. *arXiv preprint arXiv:1905.01756*, 2019. doi: 10.48550/arXiv.1905.01756.
- Karl Cobbe, Jacob Hilton, Oleg Klimov, and John Schulman. Phasic policy gradient. *arXiv preprint arXiv:2009.04416*, 2020.
- Qisai Liu, Zhanhong Jiang, Hsin-Jung Yang, Mahsa Khosravi, Joshua R. Waite, and Soumik Sarkar. Enhancing ppo with trajectory-aware hybrid policies, February 2025b. URL <https://arxiv.org/abs/2502.15968>.
- Yaozhong Gan, Xiaoyang Tan, Renye Yan, Zhe Wu, and Junliang Xing. Transductive off-policy proximal policy optimization. *arXiv preprint arXiv:2406.03894*, 2024.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992a. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
- Wouter Kool, Herke van Hoof, and Max Welling. Buy 4 reinforce samples, get a baseline for free!, 2019. URL <https://openreview.net/forum?id=H1gBfnCqKX>. Preprint.
- Arash Ahmadian, Chris Cremer, Matthias Gall  , Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet   st  n, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms, February 2024. URL <https://arxiv.org/abs/2402.14740>.
- Zhaolin Gao, Jonathan D. Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kiant   Brantley, Thorsten Joachims, J. Andrew Bagnell, Jason D. Lee, and Wen Sun. Rebel: Reinforcement learning via regressing relative rewards, December 2024. URL <https://arxiv.org/abs/2404.16767>.
- Jian Hu, Jason Klein Liu, and Wei Shen. REINFORCE++: An efficient rlhf algorithm with robustness to both prompt and reward models, April 2025. URL <https://arxiv.org/abs/2501.03262>.
- Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models, May 2024f. URL <https://arxiv.org/abs/2310.10505>.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, et al. Deepseek-v3 technical report, February 2025b. URL <https://arxiv.org/abs/2412.19437>.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, January 2025c. URL <https://arxiv.org/abs/2501.12948>.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective, March 2025c. URL <https://arxiv.org/abs/2503.20783>.



- Yuhang Zhou, Jing Zhu, Shengyi Qian, Zhuokai Zhao, Xiyao Wang, Xiaoyu Liu, Ming Li, Paiheng Xu, Wei Ai, and Furong Huang. Disco balances the scales: Adaptive domain- and difficulty-aware reinforcement learning on imbalanced data, May 2025b. URL <https://arxiv.org/abs/2505.15074>.
- Han Zhong, Zikang Shan, Guhao Feng, Wei Xiong, Xinle Cheng, Li Zhao, Di He, Jiang Bian, and Liwei Wang. Dpo meets ppo: Reinforced token optimization for rlhf, May 2025. URL <https://arxiv.org/abs/2404.18922>.
- Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play preference optimization for language model alignment, October 2024b. URL <https://arxiv.org/abs/2405.00675>.
- Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, et al. Secrets of rlhf in large language models part i: Ppo. *arXiv preprint arXiv:2307.04964*, 2023b.
- Juntao Dai, Taiye Chen, Yaodong Yang, Qian Zheng, and Gang Pan. Mitigating reward over-optimization in rlhf via behavior-supported regularization, March 2025a. URL <https://arxiv.org/abs/2503.18130>.
- Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang. Loose lips sink ships: Mitigating length bias in reinforcement learning from human feedback, 2023. URL <https://arxiv.org/abs/2310.05199>.
- Tianqi Liu, Wei Xiong, Jie Ren, Lichang Chen, Junru Wu, Rishabh Joshi, Yang Gao, Jiaming Shen, Zhen Qin, Tianhe Yu, Daniel Sohn, Anastasiia Makarova, Jeremiah Liu, Yuan Liu, Bilal Piot, Abe Ittycheriah, Aviral Kumar, and Mohammad Saleh. Rrm: Robust reward model training mitigates reward hacking, 2025d. URL <https://arxiv.org/abs/2409.13156>.
- Yuantao Fan, Ruifan Li, Guangwei Zhang, Chuan Shi, and Xiaojie Wang. A weighted cross-entropy loss for mitigating llm hallucinations in cross-lingual continual pretraining. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025b.
- Anqi Mao, Mehryar Mohri, and Yutao Zhong. Cross-entropy loss functions: Theoretical analysis and applications. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 23803–23828. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/mao23b.html>.
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020. URL <https://arxiv.org/abs/2009.01317>, page 14, 2020c.
- Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected by supervised fine-tuning data composition. *arXiv preprint arXiv:2310.05492*, 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.



- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023.
- Fei Liu et al. Learning to summarize from human feedback. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 583–592, 2020.
- Mingjia Yin, Chuhan Wu, Yufei Wang, Hao Wang, Wei Guo, Yasheng Wang, Yong Liu, Ruiming Tang, Defu Lian, and Enhong Chen. Entropy law: The story behind data compression and llm performance. *arXiv preprint arXiv:2407.06645*, 2024b.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- Ermo Hua, Biqing Qi, Kaiyan Zhang, Yue Yu, Ning Ding, Xingtai Lv, Kai Tian, and Bowen Zhou. Intuitive fine-tuning: Towards unifying sft and rlhf into a single process. *arXiv e-prints*, pages arXiv–2405, 2024.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for RLHF under KL-constraint. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 54715–54754. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/xiong24a.html>.
- Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. *arXiv preprint arXiv:2405.07863*, 2024.
- Chenlu Ye, Wei Xiong, Yuheng Zhang, Hanze Dong, Nan Jiang, and Tong Zhang. Online iterative reinforcement learning from human feedback with general preference model. *Advances in Neural Information Processing Systems*, 37:81773–81807, 2024a.
- Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrlhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*, 2023.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Yuhao Du, Zhuo Li, Pengyu Cheng, Zhihong Chen, Yuejiao Xie, Xiang Wan, and Anningzhe Gao. Simplify rlhf as reward-weighted sft: A variational method. *arXiv preprint arXiv:2502.11026*, 2025.
- Gokul Swamy, Sanjiban Choudhury, Wen Sun, Zhiwei Steven Wu, and J Andrew Bagnell. All roads lead to likelihood: The value of reinforcement learning in fine-tuning. *arXiv preprint arXiv:2503.01067*, 2025.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992b.

- Hao Sun and Mihaela van der Schaar. Inverserlignment: Large language model alignment from demonstrations through inverse reinforcement learning. In *MFHAIA the 41st International Conference on Machine Learning (ICML)*. ICML, 2024.
- Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. Acontextual-bandit approach to personalized news article recommendation. *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. *Text summarization branches out*, pages 74–81, 2004.
- Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Deepali Aneja, Zeyu Jin, Ramani Duraiswami, Dinesh Manocha, et al. A closer look at the limitations of instruction tuning. *arXiv preprint arXiv:2402.05119*, 2024.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada, July 2023f. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.754. URL <https://aclanthology.org/2023.acl-long.754/>.
- Genta Indra Winata, Hanyang Zhao, Anirban Das, Wenpin Tang, David D Yao, Shi-Xiong Zhang, and Sambit Sahu. Preference tuning with human feedback on language, speech, and vision tasks: A survey. *Journal of Artificial Intelligence Research*, 82:2595–2661, 2025.
- Gemma, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Gemini, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Zhichao Wang, Bin Bi, Zixu Zhu, Xiangbo Mao, Jun Wang, and Shiyu Wang. Uft: Unifying fine-tuning of sft and rlhf/dpo/una through a generalized implicit reward function. *arXiv preprint arXiv:2410.21438*, 2024d.
- Paul Christiano, Jan Leike, Tom Brown, et al. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 2017.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. Ultrafeedback: Boosting language models with scaled ai feedback. *arXiv preprint arXiv:2310.01377*, 2023.
- Ang Li, Qiugen Xiao, Peng Cao, Jian Tang, Yi Yuan, Zijie Zhao, Xiaoyuan Chen, Liang Zhang, Xiangyang Li, Kaitong Yang, et al. Hrlaif: Improvements in helpfulness and harmlessness in open-domain reinforcement learning from ai feedback. *arXiv preprint arXiv:2403.08309*, 2024g.

- Mengdi Li, Jiaye Lin, Xufeng Zhao, Wenhao Lu, Peilin Zhao, Stefan Wermter, and Di Wang. Curriculum-rlaif: Curriculum alignment with reinforcement learning from ai feedback. *arXiv preprint arXiv:2505.20075*, 2025.
- Daechul Ahn, Yura Choi, Youngjae Yu, Dongyeop Kang, and Jonghyun Choi. Tuning large multimodal models for videos using reinforcement learning from ai feedback. *arXiv preprint arXiv:2402.03746*, 2024.
- Liqiang Jing and Xinya Du. Fgaif: Aligning large vision-language models with fine-grained ai feedback. *arXiv preprint arXiv:2404.05046*, 2024.
- Rong Bao, Rui Zheng, Shihan Dou, Xiao Wang, Enyu Zhou, Bo Wang, Qi Zhang, Liang Ding, and Dacheng Tao. Aligning large language models from self-reference ai feedback with one general principle. *arXiv preprint arXiv:2406.11190*, 2024.
- Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan, Chenguang Zhu, Richard Yuanzhe Pang, Yundi Qian, Xuwei Wang, Suchin Gururangan, Chao Zhang, et al. Self-generated critiques boost reward modeling for language models. *arXiv preprint arXiv:2411.16646*, 2024b.
- Melody Y Guan, Manas Joglekar, Eric Wallace, Saachi Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea Vallone, Hongyu Ren, Jason Wei, et al. Deliberative alignment: Reasoning enables safer language models. *arXiv preprint arXiv:2412.16339*, 2024.
- Andrew Konya, Deger Turan, Aviv Ovadya, Lina Qui, Daanish Masood, Flynn Devine, Lisa Schirch, Isabella Roberts, and Deliberative Alignment Forum. Deliberative technology for alignment. *arXiv preprint arXiv:2312.03893*, 2023.
- Yi Fang, Wenjie Wang, Yang Zhang, Fengbin Zhu, Qifan Wang, Fuli Feng, and Xiangnan He. Large language models for recommendation with deliberative user preference alignment. *arXiv preprint arXiv:2502.02061*, 2025a.
- Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Yi Wang, Zhonghao Wang, Feiyu Xiong, et al. Internal consistency and self-feedback in large language models: A survey. *arXiv preprint arXiv:2407.14507*, 2024.
- Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Bin Sun, Xinglin Wang, Heda Wang, and Kan Li. Turning dust into gold: Distilling complex reasoning capabilities from llms by leveraging negative data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18591–18599, 2024h.
- Ajay Patel, Markus Hofmarcher, Claudiu Leoveanu-Condrei, Marius-Constantin Dinu, Chris Callison-Burch, and Sepp Hochreiter. Large language models can self-improve at web agent tasks. *arXiv preprint arXiv:2405.20309*, 2024.
- Leonardo Ranaldi and André Freitas. Self-refine instruction-tuning for aligning reasoning in language models. *arXiv preprint arXiv:2405.00402*, 2024.
- Baizhou Huang, Shuai Lu, Weizhu Chen, Xiaojun Wan, and Nan Duan. Enhancing large language models in coding through multi-perspective self-consistency. *arXiv preprint arXiv:2309.17272*, 2023b.

- Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. Universal self-consistency for large language model generation. *arXiv preprint arXiv:2311.17311*, 2023a.
- Han Wang, Archiki Prasad, Elias Stengel-Eskin, and Mohit Bansal. Soft self-consistency improves language model agents. *arXiv preprint arXiv:2402.13212*, 2024e.
- Archiki Prasad, Weizhe Yuan, Richard Yuanzhe Pang, Jing Xu, Maryam Fazel-Zarandi, Mohit Bansal, Sainbayar Sukhbaatar, Jason Weston, and Jane Yu. Self-consistency preference optimization. *arXiv preprint arXiv:2411.04109*, 2024.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. URL <https://arxiv.org/abs/2402.03300>, 2024b.
- Milan Vojnovic and Se-Young Yun. What is the alignment objective of grpo? *arXiv preprint arXiv:2502.18548*, 2025.
- Zhihang Lin, Mingbao Lin, Yuan Xie, and Rongrong Ji. Cppo: Accelerating the training of group relative policy optimization-based reasoning models. *arXiv preprint arXiv:2503.22342*, 2025.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- YuYue, Yufeng Yuan, Qiyang Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, Tiantian Fan, Zhengyin Du, Xiangpeng Wei, Gaohong Liu, Juncai Liu, Lingjun Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Ru Zhang, Xin Liu, Mingxuan Wang, Yonghui Wu, and Lin Yan. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025.
- Ming Wang, Xiaopeng Li, Ziniu Zhang, Xi Chen, and Tianyi Lin. Grpo-lead: A difficulty-aware reinforcement learning approach for concise mathematical reasoning in language models. *arXiv preprint arXiv:2504.09696*, 2025d.
- Ming Dai, Chenxu Yang, and Qiang Si. S-grpo: Early exit via reinforcement learning in reasoning models. *arXiv preprint arXiv:2505.07686*, 2025b.
- Peter Chen, Xiaopeng Li, Ziniu Li, Xi Chen, and Tianyi Lin. Spectral policy optimization: Coloring your incorrect reasoning in grpo. *arXiv preprint arXiv:2505.11595*, 2025b.
- Zichen Chen, Xiaopeng Li, Ziniu Li, Xi Chen, and Tianyi Lin. Understanding rl-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025c.
- Minghan Chen, Guikun Chen, Wenguan Wang, and Yi Yang. Seed-grpo: Semantic entropy enhanced grpo for uncertainty-aware policy optimization. *arXiv preprint arXiv:2505.12346*, 2025d.
- Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv preprint arXiv:2505.05470*, 2025e.

- Jingyi Zhang, Jia Huang, Huanjin Yao, Shuang Liu, Xiaoyang Zhang, Shuang Lu, Xiaopeng Li, and Tianyi Lin. R1-vl: Learning to reason with multimodal large language models via step-wise group relative policy optimization. *arXiv preprint arXiv:2503.12937*, 2025c.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023a.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. *arXiv preprint arXiv:2312.12148*, 2023b.
- Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024b.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 12799–12807, 2023b.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, 2018.
- Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. Mixout: Effective regularization to finetune large-scale pretrained language models. *arXiv preprint arXiv:1909.11299*, 2019.
- Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Weinberger, and Yoav Artzi. Revisiting few-sample bert fine-tuning. In *International Conference on Learning Representations*, 2020.
- Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. A survey on lora of large language models. *Frontiers of Computer Science*, 19(7), December 2024a. ISSN 2095-2236. doi: 10.1007/s11704-024-40663-9. URL <http://dx.doi.org/10.1007/s11704-024-40663-9>.
- Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Relora: High-rank training through low-rank updates, 2023. URL <https://arxiv.org/abs/2307.05695>.
- Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of lora: Efficient fine-tuning of language models via residual learning, 2024a. URL <https://arxiv.org/abs/2401.04151>.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adalora: Adaptive budget allocation for parameter-efficient fine-tuning, 2023b. URL <https://arxiv.org/abs/2303.10512>.
- Yulong Mao, Kaiyu Huang, Changhao Guan, Ganglin Bao, Fengran Mo, and Jinan Xu. Dora: Enhancing parameter-efficient fine-tuning with dynamic rank distribution, 2024b. URL <https://arxiv.org/abs/2405.17357>.

- Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. Sparse low-rank adaptation of pre-trained language models, 2023b. URL <https://arxiv.org/abs/2311.11696>.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. The impact of initialization on lora finetuning dynamics, 2024a. URL <https://arxiv.org/abs/2406.08447>.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. *arXiv preprint arXiv:2402.12354*, 2024b.
- Adam X. Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation for large language models, 2024a. URL <https://arxiv.org/abs/2308.13111>.
- Yichao Wu, Yafei Xiang, Shuning Huo, Yulu Gong, and Penghao Liang. Lora-sp: Streamlined partial parameter adaptation for resource-efficient fine-tuning of large language models, 2024c. URL <https://arxiv.org/abs/2403.08822>.
- Hongyun Zhou, Xiangyu Lu, Wang Xu, Conghui Zhu, Tiejun Zhao, and Muyun Yang. Lora-drop: Efficient lora parameter pruning based on output evaluation, 2024b. URL <https://arxiv.org/abs/2402.07721>.
- Dawid J. Kopiczko, Tijmen Blankevoort, and Yuki M. Asano. Vera: Vector-based random matrix adaptation, 2024. URL <https://arxiv.org/abs/2310.11454>.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- Yixiao Li, Yifan Yu, Chen Liang, Pengcheng He, Nikos Karampatziakis, Weizhu Chen, and Tuo Zhao. Loftq: Lora-fine-tuning-aware quantization for large language models. *arXiv preprint arXiv:2310.08659*, 2023d.
- Zhengmao Ye, Dengchun Li, Zetao Hu, Tingfeng Lan, Jian Sha, Sicong Zhang, Lei Duan, Jie Zuo, Hui Lu, Yuanchun Zhou, and Mingjie Tang. mlora: Fine-tuning lora adapters via highly-efficient pipeline parallelism in multiple gpus, 2024b. URL <https://arxiv.org/abs/2312.02515>.
- Lequn Chen, Zihao Ye, Yongji Wu, Danyang Zhuo, Luis Ceze, and Arvind Krishnamurthy. Punica: Multi-tenant lora serving, 2023b. URL <https://arxiv.org/abs/2310.18547>.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36:74764–74786, 2023g.
- Ping Ma, Michael W Mahoney, and Bin Yu. A statistical perspective on algorithmic leveraging. *The Journal of Machine Learning Research*, 16(1):861–911, 2015.
- HaiYing Wang, Rong Zhu, and Ping Ma. Optimal subsampling for large sample logistic regression. *Journal of the American Statistical Association*, 113(522):829–844, 2018.
- Ping Ma, Yongkai Chen, Xinlian Zhang, Xin Xing, Jingyi Ma, and Michael W Mahoney. Asymptotic analysis of sampling estimators for randomized numerical linear algebra algorithms. *Journal of Machine Learning Research*, 23(177):1–45, 2022.



- Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020.
- Anirban Dasgupta, Petros Drineas, Boulos Harb, Ravi Kumar, and Michael W Mahoney. Sampling algorithms and coresets for  $\ell_p$  regression. *SIAM Journal on Computing*, 38(5):2060–2078, 2009.
- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on data selection for language models. *Transactions on Machine Learning Research*, 2024.
- Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020.
- Alexander Bukharin and Tuo Zhao. Data diversity matters for robust instruction tuning. *arXiv preprint arXiv:2311.14736*, 2023.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less: Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*, 2024b.
- Weixi Song, Zuchao Li, Lefei Zhang, Hai Zhao, and Bo Du. Sparse is enough in fine-tuning pre-trained large language models. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, 2024. doi: 10.48550/arXiv.2312.11875. Spotlight Paper.
- Baohao Liao, Yan Meng, and Christof Monz. Parameter-efficient fine-tuning without introducing new latency. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4242–4260, 2023. URL <https://aclanthology.org/2023.acl-long.233>.
- Yi-Lin Sung, Varun Nair, and Colin Raffel. Training neural networks with fixed sparse masks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. URL <https://openreview.net/forum?id=Uwh-v1HSw-x>.
- Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*, 2024b.
- Luyang Fang, Xiaowei Yu, Jiazhang Cai, Yongkai Chen, Shushan Wu, Zhengliang Liu, Zhenyuan Yang, Haoran Lu, Xilin Gong, Yufang Liu, et al. Knowledge distillation and dataset distillation of large language models: Emerging trends, challenges, and future directions. *arXiv preprint arXiv:2504.14772*, 2025b.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Mingke Yang, Yuqi Chen, Yi Liu, and Ling Shi. Distillseq: A framework for safety alignment testing in large language models using knowledge distillation. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 578–589, 2024b.
- Bin Ma, Gang Liang, Yufei Rao, Wei Guo, Wenjie Zheng, and Qianming Wang. Knowledge reasoning-and progressive distillation-integrated detection of electrical construction violations. *Sensors*, 24(24):8216, 2024.

- Kun Wang, Guibin Zhang, Zhenhong Zhou, Jiahao Wu, Miao Yu, Shiqian Zhao, Chenlong Yin, Jinhu Fu, Yibo Yan, Hanjun Luo, et al. A comprehensive survey in llm (-agent) full stack safety: Data, training and deployment. *arXiv preprint arXiv:2504.15585*, 2025e.
- Shuai Niu, Jing Ma, Liang Bai, Zhihua Wang, Yida Xu, Yunya Song, and Xian Yang. Multimodal clinical reasoning through knowledge-augmented rationale generation. *arXiv preprint arXiv:2411.07611*, 2024.
- Amna Tariq, Saptarshi Purkayastha, and Tao Dai. Patient centric summarization of radiology findings using large language models. *arXiv preprint arXiv:2402.01562*, 2024.
- Hongyu Ge, Longkun Hao, Zihui Xu, et al. Clinkd: Cross-modal clinic knowledge distiller for multi-task medical images. *arXiv preprint arXiv:2503.01034*, 2025.
- Ehsan Latif, Luyang Fang, Ping Ma, and Xiaoming Zhai. Knowledge distillation of llms for automatic scoring of science assessments. In *International Conference on Artificial Intelligence in Education*, pages 166–174. Springer, 2024.
- Jiayu Shang, Cheng Peng, Yongxin Ji, Jiaojiao Guan, Dehan Cai, Xubo Tang, and Yanni Sun. Accurate and efficient protein embedding using multi-teacher distillation learning. *arXiv preprint arXiv:2405.11735*, 2024.
- Mingi Ji, Byeongho Heo, and Sungrae Park. Show, attend and distill: knowledge distillation via attention-based feature matching. *arXiv preprint arXiv:2102.02973*, 2021.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.
- Chuanguang Yang, Helong Zhou, Zhulin An, Xue Jiang, Yongjun Xu, and Qian Zhang. Cross-image relational knowledge distillation for semantic segmentation. *arXiv preprint arXiv:2204.06986*, 2022.
- Zerui Li, Yue Ming, Lei Yang, and Jing-Hao Xue. Mutual-learning sequence-level knowledge distillation for automatic speech recognition. *Neurocomputing*, 428:259–267, 2021.
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *arXiv preprint arXiv:2305.02301*, 2023.
- Kaituo Feng, Changsheng Li, Xiaolu Zhang, Jun Zhou, Ye Yuan, and Guoren Wang. Keypoint-based progressive chain-of-thought distillation for llms. *arXiv preprint arXiv:2405.16064*, 2024.
- Simran Khanuja, Melvin Johnson, and Partha Talukdar. Mergedistill: Merging pre-trained language models using distillation. *arXiv preprint arXiv:2106.02834*, 2021.
- Hailin Zhang, Defang Chen, and Can Wang. Confidence-aware multi-teacher knowledge distillation. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4498–4502. IEEE, 2022.
- Zichang Liu, Qingyun Liu, Yuening Li, Liang Liu, Anshumali Shrivastava, Shuchao Bi, Lichan Hong, Ed H Chi, and Zhe Zhao. Wisdom of committee: Distilling from foundation model to specialized application model. *arXiv preprint arXiv:2402.14035*, 2024c.

- Somin Wadhwa, Chantal Shaib, Silvio Amir, and Byron C Wallace. Who taught you that? tracing teachers in model distillation. *arXiv preprint arXiv:2502.06659*, 2025.
- Liyuan Sun, Jianping Gou, Baosheng Yu, Lan Du, and Dacheng Tao. Collaborative teacher-student learning via multiple knowledge transfer. *arXiv preprint arXiv:2101.08471*, 2021.
- Xiaoqin Chang, Sophia Yat Mei Lee, Suyang Zhu, Shoushan Li, and Guodong Zhou. One-teacher and multiple-student knowledge distillation on sentiment classification. In Nicoletta Calzolari, Chu-Ren Huang, Hansaem Kim, James Pustejovsky, Leo Wanner, Key-Sun Choi, Pum-Mo Ryu, Hsin-Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi, Patrizia Paggio, Nianwen Xue, Seokhwan Kim, Younggyun Hahm, Zhong He, Tony Kyungil Lee, Enrico Santus, Francis Bond, and Seung-Hoon Na, editors, *Proceedings of the 29th International Conference on Computational Linguistics*, pages 7042–7052, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics. URL <https://aclanthology.org/2022.coling-1.614/>.
- Minchong Li, Feng Zhou, and Xiaohui Song. Bild: Bi-directional logits difference loss for large language model distillation. *arXiv preprint arXiv:2406.13555*, 2024i.
- Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3712–3721, 2019. doi: 10.1109/ICCV.2019.00381. URL <https://api.semanticscholar.org/CorpusID:159041406>.
- Z. Yang, X. Zeng, Y. Zhao, et al. Alphafold2 and its applications in the fields of biology and medicine. *Sig Transduct Target Ther*, 8:115, 2023. doi: 10.1038/s41392-023-01381-z. URL <https://doi.org/10.1038/s41392-023-01381-z>.
- Anoop Korattikara Balan, Vivek Rathod, Kevin P Murphy, and Max Welling. Bayesian dark knowledge. *Advances in neural information processing systems*, 28, 2015.
- Meet Vadera, Brian Jalaian, and Benjamin Marlin. Generalized bayesian posterior expectation distillation for deep neural networks. In *Conference on Uncertainty in Artificial Intelligence*, pages 719–728. PMLR, 2020.
- Andrey Malinin, Bruno Mlodozieniec, and Mark Gales. Ensemble distribution distillation. *arXiv preprint arXiv:1905.00076*, 2019.
- Luyang Fang, Yongkai Chen, Wenxuan Zhong, and Ping Ma. Bayesian knowledge distillation: A bayesian perspective of distillation with uncertainty quantification. In *Forty-first International Conference on Machine Learning*, 2024.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.
- Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. Adamix: Mixture-of-adaptations for parameter-efficient model tuning. *arXiv preprint arXiv:2205.12410*, 2022b.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052*, 2020.

- Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. Adapterdrop: On the efficiency of adapters in transformers. *arXiv preprint arXiv:2010.11918*, 2021.
- Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. *arXiv preprint arXiv:2106.04647*, 2021.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2022.
- Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Lst: Ladder side-tuning for parameter and memory efficient transfer learning. *arXiv preprint arXiv:2206.06522*, 2022.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2021.
- Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen tau Yih, and Madian Khabisa. Unipelt: A unified framework for parameter-efficient language model tuning. *arXiv preprint arXiv:2206.06522*, 2022.
- Michele Farisco, G. Baldassarre, E. Cartoni, A. Leach, M.A. Petrovici, A. Rosemann, A. Salles, B. Stahl, and S. J. van Albada. A method for the ethical analysis of brain-inspired ai. *Artificial Intelligence Review*, 57(6), May 2024. ISSN 1573-7462. doi: 10.1007/s10462-024-10769-4. URL <http://dx.doi.org/10.1007/s10462-024-10769-4>.
- Howard Schneider. The emergence of enhanced intelligence in a brain-inspired cognitive architecture. *Frontiers in Computational Neuroscience*, 18, May 2024. ISSN 1662-5188. doi: 10.3389/fncom.2024.1367712. URL <http://dx.doi.org/10.3389/fncom.2024.1367712>.
- Hai Zhao, Hongqiu Wu, Dongjie Yang, Anni Zou, and Jiale Hong. Brillm: Brain-inspired large language model, 2025. URL <https://arxiv.org/abs/2503.11299>.
- Haiyang Sun, Lin Zhao, Zihao Wu, Xiaohui Gao, Yutao Hu, Mengfei Zuo, Wei Zhang, Junwei Han, Tianming Liu, and Xintao Hu. Brain-like functional organization within large language models, 2024. URL <https://arxiv.org/abs/2410.19542>.
- Hanene F. Z. Brachemi Meftah, Wassim Hamidouche, Sid Ahmed Fezza, and Olivier Deforges. Energy-latency attacks: A new adversarial threat to deep learning. *ACM, arXiv preprint arXiv:2503.04963*, 2025.
- Iason Gabriel. Artificial intelligence, values, and alignment. *Minds and Machines*, 30:411–437, 2020a.
- Lin Zhao, Lu Zhang, Zihao Wu, Yuzhong Chen, Haixing Dai, Xiaowei Yu, Zhengliang Liu, Tuo Zhang, Xintao Hu, Xi Jiang, Xiang Li, Dajiang Zhu, Dinggang Shen, and Tianming Liu. When brain-inspired ai meets agi. *arXiv preprint arXiv:2303.15935*, 2023.
- Guoyu Lu, Sheng Li, Gengchen Mai, Jin Sun, Dajiang Zhu, Lilong Chai, Haijian Sun, Xianqiao Wang, Haixing Dai, Ninghao Liu, Rui Xu, Daniel Petti, Changying Li, Tianming Liu, and Changying Li. Agi for agriculture. *arXiv preprint arXiv:2304.06136*, 2023.

- Vincent Conitzer, Rachel Freedman, Jobst Heitzig, Wesley H. Holliday, Bob M. Jacobs, Nathan Lambert, Milan Moss'e, Eric Pacuit, Stuart Russell, Hailey Schoelkopf, Emanuel Tewolde, and William S. Zwicker. Social choice should guide ai alignment in dealing with diverse human feedback. In *Social Choice Should Guide AI Alignment in Dealing with Diverse Human Feedback*. International Conference on Machine Learning, Vienna, Austria, 2024.
- Bo Yu, Jiangning Wei, Minzhen Hu, Zejie Han, Tianjian Zou, Ye He, and Jun Liu. Brain-inspired ai agent: The way towards agi. *arXiv preprint arXiv:2412.08875*, 2024c.
- Erik D. Fagerholm, Karl J. Friston, Rosalyn J. Moran, and Robert Leech. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2010.02993*, 2020.
- Traian Rebedea, Razvan Dinu, Makesh Sreedhar, Christopher Parisien, and Jonathan Cohen. Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails. *arXiv preprint arXiv:2310.10501*, 2023.
- Tadahiro Taniguchia, Hiroshi Yamakawab, Takayuki Nagaic, Kenji Doyad, Masamichi Sakagamie, Masahiro Suzukib, Tomoaki Nakamuraf, and Akira Taniguchi. A whole brain probabilistic generative model: Toward realizing cognitive architectures for developmental robots. *Neural Networks*, 2022.
- Bikram Pratim Bhuyan, Amar Ramdane-Cherif, Ravi Tomar, and T. P. Singh. Neuro-symbolic artificial intelligence: a survey. *Neural Computing and Applications*, 36:12809–12844, 2024.
- Jason Yik, Korneel Van den Berghe, Douwe den Blanken, Younes Bouhadjar, Maxime Fabre, Paul Hueber, Weijie Ke, Mina A Khoei, Denis Kleyko, Noah Pacik-Nelson, Alessandro Pierro, and Philipp Stratmann. Neurobench: A framework for benchmarking neuromorphic computing algorithms and systems. *arXiv preprint arXiv:2304.04640*, 2023.
- Rong Sun. The clarion cognitive architecture: Extending cognitive modeling to social simulation. *Cognitive Modeling to Social Simulation*, pages 79–100, 2006.
- Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. The off-switch game. *arXiv preprint: arXiv:1611.08219*, 2017.
- Eddie Guo, Christopher Perlette, Mojtaba Sharifi, Lukas Grasse, Matthew Tata, Vivian K. Mushahwar, and Mahdi Tavakoli. Speech-based human-exoskeleton interaction for lower limb motion planning. In *2024 IEEE 4th International Conference on Human-Machine Systems (ICHMS)*, pages 1–6. International Conference on Machine Learning, Vienna, Austria, 2024b.
- Alex Fedorova, Eloy Geenjaara, Lei Wua, Tristan Sylvaine, Thomas P. DeRamusa, Margaux Luckb, Maria Misiuraa, R Devon Hjelm, Sergey M. Plis, and Vince D. Calhoun. Self-supervised multi-modal neuroimaging yields predictive representations for a spectrum of alzheimer’s phenotypes. *arXiv preprint arXiv:2207.00880*, 2022.
- D. Blackiston, S. Kriegman, J. Bongard, and M. Levin. Biological robots: Perspectives on an emerging interdisciplinary field. *arXiv preprint arXiv:2209.02876*, 2022.
- Paula Harder, Alex Hernandez-Garcia, Venkatesh Ramesh, Qidong Yang, Prasanna Sattigeri, Daniela Szwareman, Campbell Watson, and David Rolnick. Hard-constrained deep learning for climate downscaling. *Journal of Machine Learning Research*, 24, 2022.

- Mohamadreza Zolfagharinejad, Unai Alegre-Ibarra, Tao Chen, Sachin Kinge, and Wilfred G. van der Wiel. Brain-inspired computing systems: a systematic literature review. *The European physical Journal B*, 97, 2024.
- Eren Kurshan. From the pursuit of universal agi architecture to systematic approach to heterogeneous agi: Addressing alignment, energy, & agi grand challenges. *arXiv preprint arXiv:2310.15274*, 2024.
- Tom Everitt, Gary Lea, and Marcus Hutter. Agi safety literature review. *arXiv preprint arXiv:1805.01109*, 2018.
- Demis Hassabis, Dhharshan Kumaran, Christopher Summerfield, and Matthew Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95, 2017.
- Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study, 2024c. URL <https://arxiv.org/abs/2404.10719>.
- Elizaveta Tennant, Stephen Hailes, and Mirco Musolesi. Moral alignment for llm agents, 2025. URL <https://arxiv.org/abs/2410.01639>.
- Iason Gabriel. Artificial intelligence, values, and alignment. *Minds & Machines*, 30:411–437, 2020b. doi: 10.1007/s11023-020-09539-2. URL <https://doi.org/10.1007/s11023-020-09539-2>.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, and et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. doi: 10.48550/arXiv.2108.07258. URL <https://doi.org/10.48550/arXiv.2108.07258>.
- Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the few-shot paradigm. *arXiv preprint arXiv:2102.07350*, 2021. doi: 10.48550/arXiv.2102.07350. URL <https://arxiv.org/abs/2102.07350>.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 1050–1059, 2016. URL <https://proceedings.mlr.press/v48/gal16.html>.
- Deep Ganguli, Amanda Askell, Yuntao Bai, and et al. Predictability and surprise in large generative models. *arXiv preprint arXiv:2202.07785*, 2022b.
- Nisan Stiennon, Long Ouyang, Jeff Wu, and et al. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems 33*, pages 3008–3021, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html>.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. doi: 10.48550/arXiv.1903.12261. URL <https://arxiv.org/abs/1903.12261>.
- Ethan Pérez, Marco Tulio Ribeiro, Julian Eisenschlos, and et al. Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022. doi: 10.48550/arXiv.2212.09251. URL <https://arxiv.org/abs/2212.09251>.



- Dan Hendrycks, Collin Burns, Steven Basart, and et al. Unsolved problems in ai safety. *arXiv preprint arXiv:2109.13916*, 2021. doi: 10.48550/arXiv.2109.13916. URL <https://arxiv.org/abs/2109.13916>.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pages 4033–4041. Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper\\_files/paper/2016/file/8d8818c8e140c64c743113f563cf750f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/8d8818c8e140c64c743113f563cf750f-Paper.pdf).
- Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pages 2586–2591, 2007.
- Silviu Paun, Bob Carpenter, Jon Chamberlain, Dirk Hovy, Udo Kruschwitz, and Massimo Poesio. Comparing bayesian models of annotation. *Transactions of the Association for Computational Linguistics*, 6:571–585, 2018.
- Jan Leike, Victoria Krakovna, Pedro A Ortega, Tom Everitt, Cédric Lefrancq, Laurent Orseau, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley, 2 edition, 1999. doi: 10.1002/047174882X. URL <https://onlinelibrary.wiley.com/doi/book/10.1002/047174882X>.
- Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scale empirical analysis. *arXiv preprint arXiv:2210.04714*, 2022.
- Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification, 2021. URL <https://arxiv.org/abs/2107.07511>.
- Yu Gui, Ying Jin, and Zhimei Ren. Conformal alignment: Knowing when to trust foundation models with guarantees. *arXiv preprint arXiv:2405.10301*, 2024. URL <https://arxiv.org/abs/2405.10301>.
- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1):289–300, 1995.
- Matthew Renze and Erhan Guven. The effect of sampling temperature on problem solving in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7346–7356. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.findings-emnlp.432/>.
- Yikun Wang, Rui Zheng, Liang Ding, Qi Zhang, Dahua Lin, and Dacheng Tao. Uncertainty aware learning for language model alignment. *arXiv preprint arXiv:2406.04854*, 2024f.
- Tunyu Zhang, Haizhou Shi, Yibin Wang, Hengyi Wang, Xiaoxiao He, Zhuowei Li, Haoxian Chen, Ligong Han, Kai Xu, Huan Zhang, et al. Token-level uncertainty estimation for large language model reasoning. *arXiv preprint arXiv:2505.11737*, 2025d.

- Qiujie Xie, Qingqiu Li, Zhuohao Yu, Yuejie Zhang, Yue Zhang, and Linyi Yang. An empirical analysis of uncertainty in large language model evaluations. *arXiv preprint arXiv:2502.10709*, 2025.
- Hamed Rahimian, Maximilian Mevissen, and Shie Mannor. Distributionally robust optimization: A review. *arXiv preprint arXiv:1908.05659*, 2019. doi: 10.48550/arXiv.1908.05659. URL <https://arxiv.org/abs/1908.05659>.
- Olivier Mihatsch and Ralf Neuneier. Risk-sensitive reinforcement learning. *Machine Learning*, 49: 267–290, 2002. doi: 10.1023/A:1017940631555.
- Tom Everitt, Marcus Hutter, Ramana Kumar, and Victoria Krakovna. Reward tampering problems and solutions in reinforcement learning: A causal influence diagram perspective. *arXiv preprint arXiv:1908.04734*, 2021. URL <https://arxiv.org/abs/1908.04734>.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*, 2023.
- N. McKenna, T. Li, L. Cheng, M. J. Hosseini, M. Johnson, and M. Steedman. Sources of hallucination by large language models on inference tasks. *arXiv preprint arXiv:2305.14552*, 2023. URL <https://arxiv.org/abs/2305.14552>.
- Shao Zhang, Jianing Yu, Xuhai Xu, Changchang Yin, Yuxuan Lu, Bingsheng Yao, Melanie Tory, Lace M. Padilla, Jeffrey Caterino, Ping Zhang, and Dakuo Wang. Rethinking human-ai collaboration in complex medical decision making: A case study in sepsis diagnosis. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*, pages 1–18, 2024c. doi: 10.1145/3613904.3642343. URL <https://arxiv.org/abs/2309.12368>.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology) is power: A critical survey of “bias” in nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, 2020. doi: 10.18653/v1/2020.acl-main.485. URL <https://aclanthology.org/2020.acl-main.485>.
- William Muldrew, Peter Hayes, Mingtian Zhang, and David Barber. Active preference learning for large language models. *arXiv preprint arXiv:2402.08114*, 2024.
- Xun Deng, Han Zhong, Rui Ai, Fuli Feng, Zheng Wang, and Xiangnan He. Less is more: Improving llm alignment via preference data selection. *arXiv preprint arXiv:2502.14560*, 2025.
- Nan Lu, Ethan X Fang, and Junwei Lu. Contextual online uncertainty-aware preference learning for human feedback. *arXiv preprint arXiv:2504.19342*, 2025.
- Beibei Shi, Haotian Li, Xing Xie, and Societal AI Team. Societal ai: Research challenges and opportunities. Technical Report MSR-TR-2025-12, Microsoft, March 2025. URL <https://www.microsoft.com/en-us/research/publication/societal-ai-research-challenges-and-opportunities/>.
- Societal AI - Microsoft Research — microsoft.com. <https://www.microsoft.com/en-us/research/project/societal-ai/>, 2025. [Accessed 23-05-2025].
- What Is AI Alignment? — IBM — ibm.com. <https://www.ibm.com/think/topics/ai-alignment>, 2025. [Accessed 22-05-2025].

- Md Meftahul Ferdaus, Mahdi Abdelguerfi, Elias Ioup, Kendall N. Niles, Ken Pathak, and Steven Sloan. Towards trustworthy ai: A review of ethical and robust large language models, 2024. URL <https://arxiv.org/abs/2407.13934>.
- The 2025 AI Index Report — Stanford HAI — [hai.stanford.edu](https://hai.stanford.edu/ai-index/2025-ai-index-report). <https://hai.stanford.edu/ai-index/2025-ai-index-report>. [Accessed 25-05-2025].
- LLMs, Truth, and Democracy: An Overview of Risks - Science and Engineering Ethics — [link.springer.com](https://link.springer.com). <https://link.springer.com/article/10.1007/s11948-025-00529-0>, a. [Accessed 23-05-2025].
- The Promise and Perils of Artificial Intelligence in Advancing Participatory Science and Health Equity in Public Health - PubMed — [pubmed.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov). <https://pubmed.ncbi.nlm.nih.gov/39963973/>, 2025. [Accessed 25-05-2025].
- Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. Trustworthy llms: a survey and guideline for evaluating large language models’ alignment, 2024d. URL <https://arxiv.org/abs/2308.05374>.
- Human-AI Interactions in Public Sector Decision Making: “Automation Bias” and “Selective Adherence” to Algorithmic Advice — [academic.oup.com](https://academic.oup.com/jpart/article/33/1/153/6524536). <https://academic.oup.com/jpart/article/33/1/153/6524536>. [Accessed 23-05-2025].
- International AI Safety Report 2025 — [gov.uk](https://www.gov.uk/government/publications/international-ai-safety-report-2025). <https://www.gov.uk/government/publications/international-ai-safety-report-2025>, 2025. [Accessed 22-05-2025].
- Sunhao Dai, Chen Xu, Shicheng Xu, Liang Pang, Zhenhua Dong, and Jun Xu. Bias and unfairness in information retrieval systems: New challenges in the llm era. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’24, page 6437–6447, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704901. doi: 10.1145/3637528.3671458. URL <https://doi.org/10.1145/3637528.3671458>.
- Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. Alignment for honesty. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 63565–63598. Curran Associates, Inc., 2024c. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/7428e6db752171d6b832c53b2ed297ab-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/7428e6db752171d6b832c53b2ed297ab-Paper-Conference.pdf).
- Stefan Pasch. Ai vs. human judgment of content moderation: Llm-as-a-judge and ethics-based response refusals, 2025. URL <https://arxiv.org/abs/2505.15365>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values, 2023. URL <https://arxiv.org/abs/2008.02275>.
- Beyond Preferences in AI Alignment - Philosophical Studies — [link.springer.com](https://link.springer.com). <https://link.springer.com/article/10.1007/s11098-024-02249-w>, b. [Accessed 23-05-2025].
- Kaustubh Dhole. Large language models as SocioTechnical systems. In Yanai Elazar, Allyson Ettinger, Nora Kassner, Sebastian Ruder, and Noah A. Smith, editors, *Proceedings of the Big Picture Workshop*, pages 66–79, Singapore, December 2023. Association for Computational

- Linguistics. doi: 10.18653/v1/2023.bigpicture-1.6. URL <https://aclanthology.org/2023.bigpicture-1.6/>.
- Aidan Kierans, Avijit Ghosh, Hananel Hazan, and Shiri Dori-Hacohen. Quantifying misalignment between agents: Towards a sociotechnical understanding of alignment. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(26):27365–27373, April 2025. ISSN 2159-5399. doi: 10.1609/aaai.v39i26.34947. URL <http://dx.doi.org/10.1609/aaai.v39i26.34947>.
- Strategic Competition in an Era of Artificial Intelligence — cnas.org. <https://www.cnas.org/publications/reports/strategic-competition-in-an-era-of-artificial-intelligence>, 2018. [Accessed 22-05-2025].
- Rohin Shah, Alex Irpan, Alexander Matt Turner, Anna Wang, Arthur Conmy, David Lindner, Jonah Brown-Cohen, Lewis Ho, Neel Nanda, Raluca Ada Popa, et al. An approach to technical agi safety and security. *arXiv preprint arXiv:2504.01849*, 2025.
- Leonie Koessler and Jonas Schuett. Risk assessment at agi companies: A review of popular risk assessment techniques from other safety-critical industries. *arXiv preprint arXiv:2307.08823*, 2023.
- Weibing Wang. A comprehensive solution for the safety and controllability of artificial superintelligence. 2024.
- Victor Gulchenko. Navigating the risks: An examination of the dangers associated with artificial general intelligence and artificial superintelligence. *Available at SSRN 4941716*, 2024.
- Phoebe Koundouri, Fivos Papadimitriou, Georgios Feretzakis, Theodoros Daglis, and Vera Alexandropoulou. Ai policies towards the agi challenge: An international assessment. Technical report, Athens University of Economics and Business, 2025.
- Meredith Ringel Morris, Jascha Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clement Farabet, and Shane Legg. Levels of agi for operationalizing progress on the path to agi. *arXiv preprint arXiv:2311.02462*, 2023.
- OpenAI. Introducing the model spec, 2024. URL <https://cdn.openai.com/spec/modelspec-2024-05-08.html>.
- DeepSeek-AI and a long list of other authors. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv-2407, 2024.
- Meta AI. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, 2025. Accessed: 2025-07-18.