# A Survey of Advanced Techniques in AI: From Large Language Models to Reinforcement Learning and Beyond

www.surveyx.cn

## Abstract

This survey paper provides a comprehensive analysis of advanced artificial intelligence (AI) techniques, emphasizing their transformative impact across various domains. Central to the discussion are Large Language Models (LLMs), which have revolutionized AI's ability to process and generate human-like text, and Reinforcement Learning (RL) variants, which enhance decision-making through trial and error. The integration of human feedback, particularly through Reinforcement Learning from Human Feedback (RLHF), is highlighted for its role in aligning AI systems with human values, addressing ethical and legal challenges in AI-generated content. The survey also explores the efficiency of machine learning algorithms in processing large datasets and the significance of reward modeling in RLHF. Furthermore, it consolidates research on enhancing LLMs using RL, focusing on bias mitigation strategies and the balance between performance enhancement and the prevention of toxic outputs. The paper is structured to provide a detailed overview of these advanced AI techniques, their applications, and their implications for future research and development. By examining these topics, the survey underscores the transformative potential of AI technologies in driving innovation across various sectors, while also addressing the challenges and limitations associated with fine-tuning LLMs and aligning them with human preferences and values.

## 1 Introduction

### 1.1 Scope and Significance

This survey provides a comprehensive analysis of advanced artificial intelligence (AI) techniques, particularly focusing on Large Language Models (LLMs) and their transformative impact across various domains. LLMs leverage sophisticated cognitive mechanisms that enhance AI's capabilities in processing and generating human-like text, as noted by Sun [1]. The necessity to represent the rich diversity of human preferences is crucial for ensuring fairness and robustness in AI outputs [2].

The significance of Reinforcement Learning (RL) and its variants is also highlighted, as they are essential for training AI agents through trial and error to enhance decision-making processes. Integrating human feedback into RL, as discussed by Lindström [3], aligns AI systems more closely with human values, addressing ethical and legal challenges in AI-generated content [4].

Additionally, the survey examines the efficiency of machine learning algorithms in processing large datasets, particularly for real-time data analysis, which is vital for scalability and speed in large-scale systems [5]. This efficiency is crucial for overcoming historical challenges as AI applications rapidly expand.

The survey consolidates emerging research on enhancing LLMs using RL, addressing both challenges and advancements in this area [6]. It also emphasizes the role of advanced AI techniques in mitigating harmful social biases in LLMs [7] and addressing cognitive biases in instruction-tuned models [8].
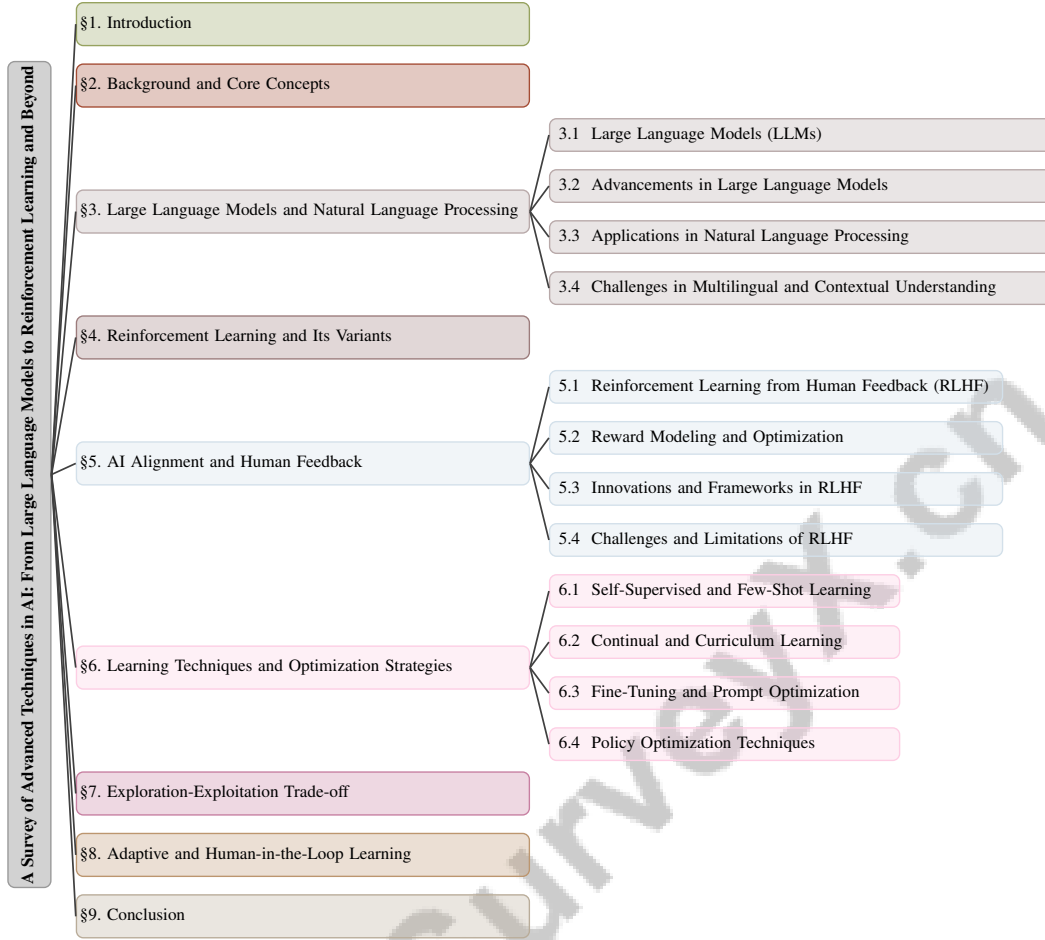
Figure 1: chapter structure

Through this analysis, the survey aims to inform future research and development efforts, highlighting the potential benefits and ethical considerations surrounding the deployment of these technologies [9, 10]. By examining these topics, the survey underscores the transformative potential of AI technologies in driving innovation across various sectors.

## 1.2 Objectives of the Survey

The primary objectives of this survey are to identify challenges and limitations in fine-tuning Large Language Models (LLMs) and to propose methodologies for enhancing their alignment with human preferences and values. A key focus is on advanced techniques such as Reinforcement Learning from Human Feedback (RLHF) to improve LLM performance regarding helpfulness, truthfulness, safety, and harmlessness [11]. The survey highlights the need for diversified approaches to fine-tuning, addressing the constraints of methods relying on a single form of supervision [12].

Additionally, the survey explores bias mitigation strategies in LLM outputs, proposing innovative methods like structured debates to reduce biases and ensure equitable alignment with diverse human preferences. It also addresses the degradation of foundational capabilities in LLMs following Supervised Fine-Tuning (SFT), aiming to balance enhancements while preventing toxic outputs [6].

Furthermore, the survey proposes frameworks to enhance LLM utility in generating coherent and contextually appropriate outputs across various domains, including image caption generation [13]. It examines the integration of Safe RLHF to ensure AI systems remain useful while avoiding harmful content generation [4].

Moreover, the survey seeks to establish benchmarks for evaluating the effectiveness of RLHF in enhancing storytelling capabilities of small language models, contributing to the development of

more aligned and ethically sound AI technologies [5]. By achieving these objectives, the survey aims to drive innovation in AI, ensuring that advanced technologies meet societal needs and ethical standards.

## 1.3  Structure of the Survey

The survey is organized to provide a comprehensive exploration of advanced AI techniques, structured into several key sections. It begins with an introduction to the scope and significance of the research, followed by the objectives guiding the study. The second section delves into background concepts, establishing foundational knowledge on Large Language Models (LLMs), Reinforcement Learning (RL), and AI Alignment.

The third section examines the role of LLMs in Natural Language Processing, highlighting advancements and applications, along with challenges in multilingual and contextual understanding. Subsequent sections explore various forms of Reinforcement Learning, including Model-based RL, Meta-Reinforcement Learning, and Multi-Agent Reinforcement Learning, focusing on their applications and contributions to AI development.

The survey then discusses AI alignment techniques, such as Reinforcement Learning from Human Feedback (RLHF) and Reward Modeling, emphasizing their importance in aligning AI systems with human values. Following this, learning techniques and optimization strategies are covered, including Self-Supervised Learning, Few-Shot Learning, and policy optimization.

The exploration-exploitation trade-off is analyzed, with strategies for balancing exploration and exploitation in AI agents. The survey also discusses Adaptive Learning and Human-in-the-Loop Learning, highlighting their roles in creating responsive AI systems. The concluding section summarizes key points, discusses future research opportunities, and addresses challenges in the field. Throughout, methodologies such as Constrained Direct Preference Optimization (C-DPO) [14] and the integration of Supervised Learning with RLHF [13] are discussed to enhance AI safety and performance. This structured approach systematically addresses complex topics, providing a coherent narrative aligned with current research trends. The following sections are organized as shown in Figure 1.

## 2  Background and Core Concepts

### 2.1  Background and Core Concepts

Contemporary artificial intelligence (AI) is grounded in advanced techniques, with Large Language Models (LLMs) and Reinforcement Learning (RL) as pivotal components. LLMs, through sophisticated neural architectures, generate human-like text, with their cognitive processes offering insights into their performance and limitations [15, 9]. RL, particularly when combined with Human Feedback (RLHF), enhances decision-making by aligning AI systems with human values, though it faces challenges due to its reliance on extensive human-labeled data [16]. The integration of RL with maximum likelihood estimation for domain adaptation exemplifies its versatility in text generation.

Reward modeling is crucial in RLHF, focusing on obtaining high-quality human feedback [17]. The QA-FEEDBACK dataset, derived from ASQA, highlights the importance of accurate reward signals in generating long-form answers to ambiguous questions, improving model performance [18]. In programming question-answering, RLHF refines models through human feedback as rewards, showcasing its broad applicability [19].

AI Alignment ensures AI systems adhere to human intentions, necessitating benchmarks to evaluate LLM reasoning, particularly in specific data distributions [20]. The diversity of tasks and lack of comparability among benchmarks underline the need for a standardized platform [21]. Benchmarks examining output length and performance further highlight optimization complexities [22].

Bias in LLMs poses significant challenges, with current RLHF methods often misaligned with diverse human preferences due to reliance on a single reward model [7, 2]. Addressing these biases is crucial for ethical AI deployment.

Efficient machine learning algorithms are vital for real-time applications, as current data management systems' slow processing speeds limit usability [5]. Recent AI research advancements, including

3

LLM improvements for automating research ideation, the SKILL-MIX framework for skill integration, and critique-based supervision for enhanced reasoning, establish a foundation for ongoing exploration and innovation [23, 24, 25].

# 3 Large Language Models and Natural Language Processing

Large Language Models (LLMs) are reshaping Natural Language Processing (NLP) by generating human-like text and driving research into their mechanisms and implications. This section examines LLMs' architecture, capabilities, and alignment challenges with human values, focusing on their design and operational principles that enhance performance in NLP tasks.

## 3.1 Large Language Models (LLMs)

LLMs represent a significant leap in AI, particularly in NLP, due to their ability to produce coherent, contextually relevant text for tasks like summarization, translation, and reasoning [26]. However, the emphasis on instruction tuning for a few languages limits their global accessibility. Aligning LLMs with human preferences for helpfulness and harmlessness is challenging, as conflicting objectives complicate training [4]. LLMs also struggle with accuracy and confidence indication, hindering broader application [27].

Innovative methods like Mistral-Plus employ Direct Harmless RLHF to improve conversational abilities and safety without relying on Supervised Fine-Tuning [6]. These approaches are crucial for aligning LLMs with human values, enhancing their utility across sectors. LLMs advance AI capabilities, particularly in NLP, by enabling real-time updates and personalized interactions. RLHF enhances decision-making and output quality, while frameworks like SKILL-MIX assess skill integration. Critique models refine reasoning, emphasizing LLMs' transformative role in AI innovation [28, 25, 9, 29, 23]. Ongoing research is essential to ensure ethical alignment with human values.

## 3.2 Advancements in Large Language Models

LLMs have improved in natural language understanding and generation, though challenges persist in mathematical problem-solving and reasoning [30]. Critique models providing step-level feedback enhance performance on complex queries [23]. Online training allows real-time updates and customization, while offline methods like META REFLECTION offer versatile task adaptation [29, 31].

Prompt tuning leverages pre-trained knowledge to enhance performance, reducing computational demands [32]. Reinforcement learning frameworks for fine-tuning help LLMs learn from diverse outputs without human rankings, broadening applicability [33]. Evaluation benchmarks assess RLHF's impact on generalization and diversity, showcasing its effectiveness over traditional methods [34, 26]. Representation Alignment from Human Feedback (RAHF) simplifies model behavior manipulation, reducing computational costs [11].

These advancements underscore LLMs' transformative influence on AI, marked by innovative evaluation methods and reinforcement learning techniques that enhance output quality. Critique models improve reasoning, while research into AI alignment emphasizes congruence with human values, enhancing sophistication and addressing ethical considerations [28, 25, 9, 10, 23].

## 3.3 Applications in Natural Language Processing

LLMs have advanced NLP by enabling sophisticated applications like machine translation, text summarization, and sentiment analysis. Techniques like Active Preference Learning (APL) optimize preference labels during fine-tuning [35]. Prompt optimization, enhanced by human feedback, improves reliability and relevance [36].

Reinforcement learning with prompt optimization, such as RLPROMPT, excels in few-shot classification and unsupervised text style transfer, expanding LLM applicability [37]. Challenges remain in robustness under out-of-distribution conditions, with significant accuracy drops in unfamiliar data [38]. Datasets like SKILL-MIX benchmark LLM performance across skills and topics, ensuring models are challenged and capable of generalizing [25].

4

Online training methods leverage user inputs for real-time updates and customization [29]. Frameworks like HFFT improve performance across tasks [12]. LLMs play a pivotal role in education, using child-directed data for storytelling and language development, illustrating their potential in educational tools [39].
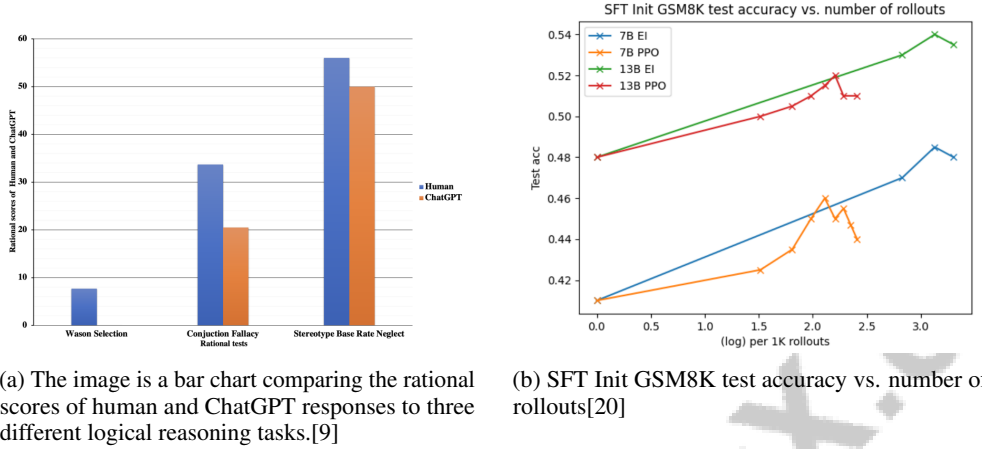


(a) The image is a bar chart comparing the rational scores of human and ChatGPT responses to three different logical reasoning tasks.[9]

(b) SFT Init GSM8K test accuracy vs. number of rollouts[20]

Figure 2: Examples of Applications in Natural Language Processing

As illustrated in Figure 2, LLMs like ChatGPT enhance machines' text understanding and generation capabilities. The first application compares ChatGPT's rationality against human responses, highlighting challenges in logical reasoning. The second focuses on improving model performance through reinforcement learning, optimizing LLMs for complex tasks [9, 20].

## 3.4   Challenges in Multilingual and Contextual Understanding

LLMs struggle with multilingual and contextual information, crucial for deployment across diverse settings. A lack of training data for multi-turn interactions hampers performance in goal-directed tasks, especially in dialogue management [40]. Real-time mistake rectification is limited, affecting dialogue management requiring long-term planning [31].

Multilingual preference optimization faces challenges, as previous methods focus on major languages, neglecting others and limiting global reach [41]. Assessing content difficulty is labor-intensive, failing to adapt to varying user knowledge, impacting educational contexts [42]. Identifying human values in communication remains complex, essential for understanding value-influenced arguments [32].

LLMs' limited performance self-assessment is exacerbated by question difficulty and critique absence, necessitating enhanced self-evaluation mechanisms [23]. Addressing multilingual and contextual challenges is vital for LLM performance. Advanced techniques like RLHF have improved multilingual LLM alignment across languages, enhancing reliability and applicability in real-world tasks [28, 30, 41, 29, 12].

## 4   Reinforcement Learning and Its Variants

This section examines the diverse landscape of reinforcement learning, focusing on distinct approaches that contribute unique methodologies and insights. Model-based Reinforcement Learning is highlighted for its use of internal models to enhance decision-making, setting a foundation for understanding its mechanisms and applications. As illustrated in Figure 3, the hierarchical structure and key concepts of reinforcement learning and its variants are depicted, emphasizing model-based, meta-reinforcement, and multi-agent reinforcement learning approaches. The diagram categorizes the techniques, challenges, and advancements within each domain, while also highlighting curriculum learning strategies, innovative methodologies, and adaptability in complex environments. This visual representation not only complements the textual analysis but also provides a comprehensive overview of the relationships between the various methodologies in reinforcement learning.
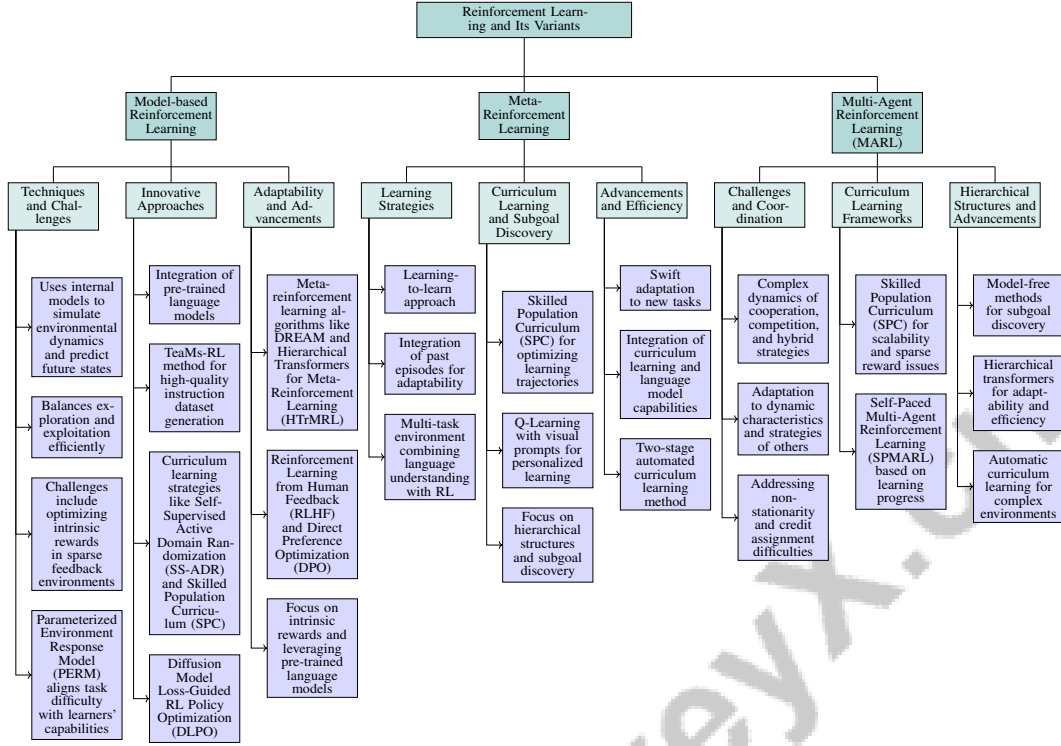
Figure 3: This figure illustrates the hierarchical structure and key concepts of reinforcement learning and its variants, highlighting model-based, meta-reinforcement, and multi-agent reinforcement learning approaches. The diagram categorizes the techniques, challenges, and advancements within each domain, emphasizing curriculum learning strategies, innovative methodologies, and adaptability in complex environments.

## 4.1 Model-based Reinforcement Learning

Model-based Reinforcement Learning (RL) uses internal models to simulate environmental dynamics, predict future states, and optimize decision-making, effectively balancing exploration and exploitation more efficiently than model-free methods, thus reducing sample complexity and accelerating training [43]. A challenge in this approach is optimizing intrinsic rewards in sparse feedback environments. Techniques like the Parameterized Environment Response Model (PERM) align task difficulty with learners' capabilities to enhance learning outcomes.

Recent advances integrate pre-trained language models to improve decision-making in data-constrained scenarios. The TeaMs-RL method, for example, uses an instructor LLM to generate diverse instructions for querying an expert LLM, forming a high-quality instruction dataset [44]. This underscores LLMs' potential in adaptive and robust RL systems.

Curriculum learning strategies, such as Self-Supervised Active Domain Randomization (SS-ADR), optimize learning trajectories by creating a joint curriculum of goals and environments, enabling agents to progressively acquire complex skills [45]. The Skilled Population Curriculum (SPC) employs a hierarchical structure to facilitate behavior learning across tasks, emphasizing population-invariant communication in multi-agent contexts.

Innovative techniques like Diffusion Model Loss-Guided RL Policy Optimization (DLPO) refine generative models by incorporating diffusion model loss as a penalty within the reward function, showcasing Model-based RL's potential to optimize performance across applications, including text-to-speech systems [46].

The adaptability of Model-based RL is further illustrated by meta-reinforcement learning algorithms such as DREAM, which enable agents to learn effectively from limited episodes in dynamic environments where language serves as contextual support. Hierarchical Transformers for

Meta-Reinforcement Learning (HTrMRL) leverage sophisticated architectures to process intra- and inter-episode experiences, enhancing knowledge acquisition efficiency and allowing agents to adapt rapidly to new tasks, improving generalization from limited data and outperforming state-of-the-art approaches in various simulated environments, as evidenced by results from the Meta-World Benchmark [47, 48, 49].

Model-based Reinforcement Learning is rapidly progressing, particularly through incorporating techniques like Reinforcement Learning from Human Feedback (RLHF) and Direct Preference Optimization (DPO), enhancing LLMs to generate more accurate, contextually relevant outputs. Recent studies highlight LLMs' effectiveness in crafting reward models for training agents, improving decision-making policies and adapting to distribution shifts, addressing alignment and robustness challenges [38, 50, 28]. The focus on intrinsic rewards, curriculum learning, and leveraging pre-trained language models exemplifies innovative strategies to tackle inherent challenges, promising enhanced adaptability and efficiency in complex environments.

## 4.2 Meta-Reinforcement Learning

Meta-Reinforcement Learning (Meta-RL) enables agents to rapidly adapt to new tasks by leveraging prior knowledge, akin to human learning processes. This is achieved through learning-to-learn, where agents modify their strategies based on task requirements [48]. The integration of past episodes into the learning process, as demonstrated by Hierarchical Transformers for Meta-Reinforcement Learning (HTrMRL), significantly enhances adaptability, enabling swift adjustments to new environments and challenges.

An innovative approach within Meta-RL involves creating a multi-task environment where agents learn to associate language with actions and goals through exploration and interaction [51]. This method highlights the potential of combining language understanding with reinforcement learning to improve generalization and learning efficiency.

Curriculum learning strategies are crucial in Meta-RL, with methodologies like the Skilled Population Curriculum (SPC) modeling the teacher's task as a contextual bandit to adapt the curriculum to evolving student strategies, optimizing learning trajectories [52]. Additionally, using Q-Learning to present vocabulary words with visual prompts tailored to the learner's CEFR level exemplifies curriculum learning's application in educational contexts, facilitating personalized and effective learning experiences [42].

Meta-RL also addresses subgoal discovery challenges in hierarchical reinforcement learning, particularly in environments with sparse and delayed feedback [49]. By focusing on hierarchical structures and subgoal discovery, Meta-RL enhances agents' learning efficiency and adaptability in complex environments.

Meta-Reinforcement Learning represents a significant advancement in reinforcement learning, offering a framework for creating agents that swiftly adapt to new and complex tasks. By integrating curriculum learning techniques, enhancing language model capabilities, and optimizing exploration strategies through meta-learning, Meta-RL develops adaptive learning systems capable of efficiently tackling complex tasks and generalizing across diverse environments. This approach improves training efficiency using a two-stage automated curriculum learning method to refine task selection while leveraging fine-tuning in large language models to navigate trade-offs between novelty, feasibility, and effectiveness in idea generation. Consequently, Meta-RL achieves significant performance improvements, demonstrating a sophisticated understanding of task complexity and agent adaptability [53, 47, 24].

## 4.3 Multi-Agent Reinforcement Learning (MARL)

Multi-Agent Reinforcement Learning (MARL) tackles the challenges of environments where multiple agents interact simultaneously, navigating complex dynamics of cooperation, competition, and hybrid strategies to achieve objectives. Agents must adapt to dynamic characteristics and adjust strategies based on others' actions and learning progress, necessitating sophisticated coordination mechanisms that address challenges like non-stationarity and credit assignment difficulties. Recent advancements, such as the Skilled Population Curriculum (SPC) and self-paced MARL (SPMARL), aim to address

AI-generated, for reference only.

these issues by implementing automatic curriculum learning frameworks that enhance cooperation and skill acquisition, ultimately improving scalability and performance in complex scenarios [52, 54].

The Skilled Population Curriculum (SPC) is an automatic curriculum learning framework designed to enhance scalability and address sparse reward issues in multi-agent environments [52]. SPC optimizes the learning trajectory by dynamically adjusting the curriculum based on agents' performance, thereby improving learning efficiency and effectiveness.

Self-Paced Multi-Agent Reinforcement Learning (SPMARL) utilizes a curriculum optimized based on learning progress, as indicated by critic loss [54]. This method allows for the gradual introduction of complexity in multi-agent tasks, facilitating smoother learning curves and better overall performance.

The integration of model-free methods for subgoal discovery, as highlighted by Rafati [49], exemplifies advancements in MARL. This approach does not require an environmental model, making it suitable for large-scale applications where computational resources and time are limited.

Hierarchical structures play a crucial role in MARL, enabling agents to decompose tasks into manageable sub-tasks and optimize their learning strategies accordingly. The use of hierarchical transformers, as explored in Meta-Reinforcement Learning settings [48], underscores the potential of such architectures to enhance adaptability and efficiency in multi-agent scenarios.

MARL is advancing significantly through the integration of sophisticated techniques such as automatic curriculum learning (ACL), self-paced learning, and model-free subgoal discovery. ACL enhances agent coordination in complex environments by systematically increasing task difficulty, crucial for addressing scalability and sparse reward challenges. Recent innovations like the Skilled Population Curriculum (SPC) framework allow for adaptive curriculum learning that accommodates varying agent populations and improves cooperation skills across tasks. Additionally, self-paced MARL (SPMARL) prioritizes tasks based on learning progress rather than mere episode returns, leading to faster convergence and superior performance in sparse-reward scenarios. These developments collectively enhance the efficiency and effectiveness of learning in multi-agent systems [52, 54], paving the way for more sophisticated applications in complex, interactive environments.

# 5 AI Alignment and Human Feedback

To effectively align artificial intelligence (AI) systems with human values, it is essential to explore various methodologies that leverage human feedback. This section will delve into the foundational approaches that underpin this alignment process, beginning with Reinforcement Learning from Human Feedback (RLHF). RLHF serves as a pivotal framework, enabling AI systems to learn from human preferences and adjust their behaviors accordingly. The subsequent discussion will provide an in-depth examination of RLHF, highlighting its significance and the innovative strategies employed to enhance its effectiveness.

## 5.1 Reinforcement Learning from Human Feedback (RLHF)

Reinforcement Learning from Human Feedback (RLHF) is an essential methodology for aligning artificial intelligence (AI) systems, particularly Large Language Models (LLMs) and Vision-Language Models (VLMs), with human values and preferences [55]. This approach involves leveraging human feedback to guide the reinforcement learning process, ensuring that AI outputs align more closely with human expectations. One of the primary challenges RLHF faces is the effective utilization of offline preference data in online learning environments, which is crucial for optimizing the alignment process [56].

Innovative strategies have been developed to address the scalability limitations of RLHF, which often stem from the high cost and time associated with obtaining high-quality human preference labels. The Okapi framework represents a significant advancement in multilingual LLM research, enhancing the models' ability to understand and follow human instructions across multiple languages [26]. This development is critical for ensuring that AI systems can cater to a diverse global audience, thereby improving their applicability and ethical robustness.

The integration of natural language critique into reward models, such as the CLoud reward models, offers a novel approach to predicting scalar rewards based on the quality of an assistant's response

[57]. This method enhances the RLHF process by providing more nuanced feedback, which is essential for refining AI outputs and ensuring they meet human expectations.

Revised Sentence: "Reinforcement Learning from Human Feedback (RLHF) is essential in the development of advanced AI systems, particularly large language models (LLMs), as it facilitates the alignment of these models with human values by integrating human preferences into their training processes. This approach not only enhances the efficiency and reliability of AI outputs but also addresses ethical concerns by promoting helpfulness and harmlessness, thereby minimizing issues like toxicity and hallucinations. By employing techniques such as reward modeling, Proximal Policy Optimization (PPO), and Safe RLHF, researchers are working to overcome challenges related to reward design and training stability, ultimately striving for AI systems that can effectively serve as trustworthy human-centric assistants." [58, 4, 59]. By addressing the inherent challenges and leveraging innovative strategies, RLHF continues to advance AI technologies that better reflect human values and expectations, ultimately contributing to the creation of more helpful and harmless AI systems.

## 5.2 Reward Modeling and Optimization

Reward modeling and optimization are critical components in the reinforcement learning framework, especially within the context of Reinforcement Learning from Human Feedback (RLHF). These techniques are essential for aligning AI systems with human values and preferences, ensuring that AI behavior is optimized towards desirable outcomes [60]. A significant challenge in this domain is the accurate definition and implementation of reward functions that can effectively capture the nuances of human feedback while avoiding issues such as reward hacking [61].

Recent advancements have introduced methods that involve learning multiple reward models to account for the diversity of human preferences, aggregating these models to form a unified representation [62]. This approach enhances the robustness of reward models by incorporating a broader spectrum of human values, thereby facilitating more accurate and ethical AI behavior. The use of contrastive and meta-learning techniques has also been proposed to measure the strength of preferences in the data, further improving reward model performance [63].

The integration of conservative policy optimization methods, which account for reward uncertainty, has been suggested to produce more reliable and risk-averse policies [64]. Such approaches address the risks associated with deploying models that may not adequately reflect diverse user values, a concern highlighted in previous studies. Moreover, the use of direct policy optimization techniques, such as Zeroth-Order Policy Gradient (ZPG), offers a novel pathway for optimizing AI systems without relying on traditional reward models [65].

In the realm of neural machine translation, integrating quality metrics as reward models has been shown to enhance translation quality, underscoring the importance of systematic evaluation methods in optimizing AI outputs. However, existing benchmarks sometimes fail to account for variability in teacher expertise, which can lead to suboptimal performance in RLHF algorithms that do not differentiate between teachers based on their areas of specialization [66].

Overall, the field of reward modeling and optimization continues to evolve, with ongoing research focused on overcoming inherent challenges and refining techniques to ensure AI systems are aligned with human values. The transition towards sophisticated and dependable reward models, along with cutting-edge optimization techniques, presents a promising avenue for enhancing the performance of AI systems in intricate and rapidly changing environments. This is particularly crucial as the effectiveness of reward models in Reinforcement Learning with Human Feedback (RLHF) has been shown to be sensitive to distribution shifts, impacting their accuracy and calibration. Recent advancements, such as Critique-out-Loud reward models, demonstrate improved preference classification by explicitly generating critiques of responses, while approaches like leverage large language model feedback to synthesize intrinsic rewards more efficiently. Furthermore, the introduction of reward-robust RLHF frameworks aims to mitigate the instability associated with traditional reward models, ensuring more resilient learning in large language models. Collectively, these innovations are set to significantly advance AI capabilities in complex scenarios. [67, 24, 38, 57, 68]

9

## 5.3 Innovations and Frameworks in RLHF

Recent advancements in Reinforcement Learning from Human Feedback (RLHF) have introduced several innovative frameworks and methodologies aimed at enhancing the alignment of AI systems with human values and preferences. One notable innovation is the Mixed Preference Optimization (MPO) approach, which employs a two-stage training process and utilizes a well-trained Direct Preference Optimization (DPO) model as a reference during the RLHF phase. This method effectively addresses the distribution shift problem by employing meta-learning techniques to align reward models with changing environment distributions, thereby enhancing the robustness and stability of AI systems, particularly in the context of reinforcement learning from human feedback (RLHF). Extensive experiments demonstrate that this approach improves the reward model's ability to distinguish subtle differences in out-of-distribution samples, ultimately leading to more reliable and resilient AI performance. [69, 38, 10, 70, 68]. Another significant development is the MaxMin-RLHF approach, which focuses on maximizing the minimum utility across diverse user groups, ensuring a more equitable alignment of AI systems with varied human preferences and thereby enhancing fairness and inclusivity in AI outputs.

The MetaRM framework utilizes advanced meta-learning techniques to effectively align reward models with evolving environmental distributions, thereby ensuring that these models maintain robust performance across both in-distribution and out-of-distribution samples. This capability is achieved without the need for continuous data labeling, addressing the challenges posed by shifts in the output distribution of policy models during training, which can hinder the reward model's ability to accurately differentiate responses. Extensive experiments have demonstrated that MetaRM significantly enhances the distinguishing ability of reward models during iterative reinforcement learning from human feedback (RLHF) optimization, particularly in identifying subtle differences in out-of-distribution samples. [71, 72, 38, 57, 70]. This innovation reduces the need for extensive human input, streamlining the RLHF process. In the realm of reward modeling, the introduction of critique generation, as exemplified by the CLoud reward models, allows models to explicitly reason about response quality rather than relying solely on implicit reasoning. This method improves the ability of AI systems to evaluate and refine outputs based on human-like critique.

Complementing this, the CRM method enhances reward models by measuring preference strength and employing contrastive learning to distinguish between chosen and rejected responses, thereby refining the alignment process. Furthermore, the SaySelf framework improves confidence calibration and enables AI systems to express nuanced uncertainty, thereby enhancing user trust in AI outputs. This approach is essential for the development of AI systems that not only achieve high accuracy but also prioritize transparency and user-friendliness, thereby ensuring alignment with human values and enhancing their overall reliability and effectiveness in real-world applications. [23, 10, 60, 25]

Additionally, the Nash Learning from Human Feedback (NLHF) presents a promising alternative to traditional RLHF by utilizing preference models and Nash equilibrium to align Large Language Models (LLMs) with human preferences. This intuitive approach simplifies the alignment process and enhances the adaptability of AI systems to diverse human values. The proposed method, warmPref-PS, utilizes Bayesian posterior sampling to leverage offline preference data, allowing the online learning process to begin with a more informed prior [56].

The recent advancements in reinforcement learning from human feedback (RLHF) methodologies highlight a significant evolution in the field, emphasizing a dual commitment to enhancing the effectiveness of AI systems while ensuring they align ethically with human values and expectations. This evolution is underscored by a critical examination of existing algorithms, which reveals that traditional approaches can inadvertently create perverse incentives due to assumptions like the independence of irrelevant alternatives (IIA). Such insights drive ongoing research to refine RLHF frameworks, focusing on improving reward models and addressing limitations such as incorrect generalization and feedback sparsity. Ultimately, these innovations reflect a proactive effort to develop AI technologies that not only perform well but also prioritize safety and ethical considerations in their interactions with users. [58, 73]. By addressing inherent challenges and leveraging innovative strategies, RLHF continues to advance AI technologies, contributing to the development of more helpful and harmless AI systems.

As shown in Figure 4, The example of "AI Alignment and Human Feedback; Innovations and Frameworks in RLHF" provides a visual and conceptual exploration of the cutting-edge methodologies employed in the realm of reinforcement learning from human feedback (RLHF). Illustrated through

10

(a) Actor Model (deepspeed) and Critic Model (deepspeed) with Reference Model (deepspeed) and Reward Model (deepspeed) in a deepspeed framework[74]

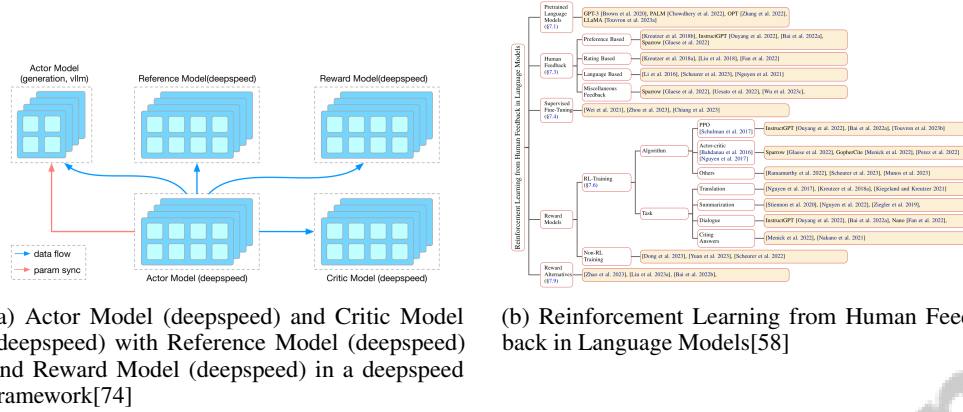(b) Reinforcement Learning from Human Feedback in Language Models[58]

Figure 4: Examples of Innovations and Frameworks in RLHF

two distinct images, the first showcases a deepspeed framework that integrates an Actor Model, Critic Model, Reference Model, and Reward Model, each pivotal in the development and execution of AI systems that align closely with human intentions. The intricate layering of these models underscores the complexity and scalability of the deepspeed approach, emphasizing its capability to enhance AI performance through structured feedback mechanisms. The second image complements this by presenting a tree diagram that systematically categorizes the diverse strategies and techniques utilized in RLHF for language models. This diagram delineates various sections such as 'Pretrained Language Models' and 'Supervised Fine-Tuning,' highlighting the multifaceted nature of feedback and training methods. Together, these visual representations encapsulate the innovations and frameworks that drive the evolution of AI systems towards more aligned and human-centric functionalities. [? ]hu2024openrlhfeasytousescalablehighperformance,chaudhari2024rlhfdecipheredcriticalanalysis)

## 5.4 Challenges and Limitations of RLHF

Reinforcement Learning from Human Feedback (RLHF) is a powerful tool for aligning AI systems with human values, yet it faces several challenges and limitations that impact its scalability and effectiveness. A significant challenge is the high cost and inefficiency associated with collecting extensive human preference data, which is often difficult to obtain in real-world applications . This reliance on preference data makes the process expensive and time-consuming, limiting its scalability and generalizability [75]. The sample complexity required by existing RLHF methods further exacerbates these limitations, restricting the exploration of the vast space of possible responses and leading to suboptimal performance [76].

The robustness of RLHF systems is also challenged by noisy or corrupted preference data. While methods like R3M have shown superior robustness by effectively identifying and mitigating the impact of outliers, traditional RLHF approaches may struggle to maintain performance under such conditions [77]. This underscores the importance of developing techniques that can handle variability in human feedback, especially in complex or ambiguous scenarios where multiple feedback pieces may conflict [78]. Furthermore, the reliance on pre-trained large language models, which are prone to hallucination, can lead to inaccurate rankings and uninformative rewards, posing a significant challenge in maintaining the fidelity of RLHF processes [79].

Another limitation is the potential for inherent biases in human evaluators, which may still affect the alignment outcomes [10]. This issue is compounded by the difficulty in accurately capturing and modeling highly nuanced user preferences, which may require extensive feedback data [80]. Existing methods often assume a uniformity of preferences among users, leading to suboptimal results for personalized applications [81]. Additionally, the potential for overfitting to the specific datasets used for training may affect the model's performance in unseen scenarios [63].

The implementation of multiple reward models and aggregation mechanisms poses additional complexity and resource requirements, which can be a limitation for RLHF systems aiming to capture the diversity of human preferences [62]. The MaxMin-RLHF approach, for example, faces potential complexity in accurately modeling diverse preference distributions, which may require extensive data

and careful tuning [2]. Moreover, existing studies often overlook the contextual factors influencing human feedback and fail to model the diverse communication strategies used by humans in real-world scenarios [64].

The RLHF trilemma presents a significant challenge for the widespread adoption of reinforcement learning from human feedback (RLHF) in large language models (LLMs), as it highlights the inherent incompatibility between achieving high task diversity, maintaining low labeling costs, and ensuring alignment performance that generalizes effectively across various scenarios. This trade-off complicates the fine-tuning process, as evidenced by findings that while RLHF improves out-of-distribution generalization, it often results in reduced output diversity compared to supervised fine-tuning (SFT). Addressing this trilemma is crucial for optimizing RLHF methodologies and enhancing their applicability in diverse real-world contexts. [82, 34, 83, 74]. Addressing these challenges is crucial for advancing RLHF methodologies and ensuring that AI systems are both effective and ethically aligned with societal needs.

# 6   Learning Techniques and Optimization Strategies

This section examines methodologies that enhance AI systems' adaptability and efficiency, focusing on self-supervised and few-shot learning. These approaches are vital for improving performance in scenarios with limited labeled data and require rapid task adaptation. The discussion highlights innovative techniques and frameworks supporting these methodologies, emphasizing their significance in AI development.

## 6.1   Self-Supervised and Few-Shot Learning

Self-supervised and few-shot learning are crucial methodologies in AI, enhancing model adaptability and efficiency across diverse environments. Self-supervised learning utilizes vast amounts of unlabeled data to learn representations that generalize across tasks, reducing reliance on extensive labeled datasets, particularly beneficial in domains with scarce or costly labeled data. The integration of RLHF with LLMs captures nuanced human preferences through multimodal sensory data, refining training processes in complex domains like autonomous driving [84].

Few-shot learning enables models to generalize from minimal examples, essential for rapid task adaptation. Techniques such as modifying the Direct Preference Optimization (DPO) algorithm with an exploration bonus enhance sample efficiency and facilitate diverse response generation [76]. Benchmarks like LLF-Bench, involving tasks with natural language instructions and feedback, support quick adaptation by providing comprehensive evaluations of interactive learning systems [85].

Curriculum learning strategies optimize learning trajectories in self-supervised and few-shot contexts. The CurricuLLM framework generates sequences of subtasks in natural language, translating them into executable task code and evaluating trained policies for alignment with each subtask, thereby enhancing the learning process [86]. Similarly, Curriculum Q-Learning optimizes vocabulary presentation based on student proficiency, illustrating curriculum learning's application in educational settings [42].

The Okapi framework innovatively instruction-tunes LLMs across 26 languages, offering a comprehensive dataset for multilingual evaluation [26]. This approach emphasizes that enhancing the diversity and quality of training data improves a model's ability to generalize across languages, as noted by Dang [41].

Self-supervised and few-shot learning continue to evolve, presenting substantial potential for improving AI systems' adaptability and efficiency. By leveraging unlabeled data, integrating curriculum strategies, and employing innovative frameworks, these methodologies enhance models' generalization across tasks and adaptation to new challenges with minimal supervision. The supervised fine-tuning of pre-trained language models on parallel corpora and training reward models on preference datasets exemplifies advancements in translation and preference modeling [87].

## 6.2 Continual and Curriculum Learning

Continual and curriculum learning are pivotal strategies in AI aimed at enhancing model adaptability and efficiency in dynamic environments. Continual learning enables models to learn continuously from data streams, adapting to new tasks without forgetting previously acquired knowledge, addressing catastrophic forgetting. Frameworks like MOLe utilize meta-learned initialization to facilitate rapid adaptation to new tasks with minimal data, maintaining performance across varying tasks [88].

Curriculum learning structures the learning process by presenting tasks in increasing complexity, aligning with educational principles that promote foundational skill development before advancing to complex tasks. The Parameterized Environment Response Model (PERM) tailors environments to match learners' current abilities, ensuring challenges are appropriately calibrated for effective learning. This enhances training efficiency for reinforcement learning agents and facilitates knowledge transfer within the "zone of proximal development." Integrating critique models allows LLMs to engage in self-reflection and receive step-level feedback, significantly improving performance on complex reasoning tasks [89, 23]. Curriculum learning has enhanced performance and efficiency, particularly in reinforcement learning contexts through techniques like Self-Paced Learning and Goal Generation.

The integration of continual and curriculum learning strategies has profound implications for developing robust AI systems capable of adapting to diverse environments. By employing advanced techniques such as fine-tuning with heterogeneous feedback and reinforcement learning, models enhance flexibility and resilience. This approach optimizes idea generation by navigating trade-offs between novelty, feasibility, and effectiveness, improving performance in complex sequential decision-making tasks. Consequently, these models are better equipped to maintain sustained performance and adaptability in real-world applications, addressing challenges like limited data availability and contextual understanding [32, 24, 50, 69, 12]. Ongoing research in these areas continues to expand the capabilities of adaptive learning systems, paving the way for more sophisticated AI technologies.

## 6.3 Fine-Tuning and Prompt Optimization

| Method Name | Optimization Techniques | Integration with RL | Application Domains |
|---|---|---|---|
| SAFT[90] | Semantic-Aware Fine-Tuning | Text-based Reinforcement | Text-based Games |
| FCF-LLM[30] | Prompt Optimization | Rlhf | Logical Reasoning Tasks |
| RL-LMFT[33] | Proximal Policy Optimization | Reinforcement Learning Framework | Abstractive Summarization Tasks |
| RL[37] | Reinforcement Learning | Reinforcement Learning-based | Nlp Tasks |
| HVDM[32] | Fine-tuning, Prompt | - | Natural Language Understanding |
| R3M[77] | Proximal Policy Optimization | Robust Reward Modeling | Robotic Control Tasks |

Table 1: Table illustrating various methods for optimizing large language models (LLMs), highlighting their respective optimization techniques, integration with reinforcement learning (RL), and application domains. The table provides a comparative overview of approaches such as Semantic-Aware Fine-Tuning (SAFT), Prompt Optimization, and Proximal Policy Optimization, among others, in enhancing LLM performance across diverse tasks.

Fine-tuning and prompt optimization are essential methodologies for enhancing the performance and adaptability of LLMs. Fine-tuning adjusts the parameters of pre-trained models for specific tasks, such as detecting human values in natural language understanding (NLU), leveraging knowledge from the pre-training phase to tackle complex tasks effectively. In contrast, prompt optimization refines input prompts to elicit more accurate responses, serving as a practical alternative to traditional methods like RLHF, especially when model parameters cannot be modified. Recent studies validate the effectiveness of prompt optimization, demonstrating its ability to align LLMs with human values even when fine-tuning is not feasible [32, 91]. Table 1 provides a comprehensive overview of different optimization methods applied to large language models (LLMs), detailing their optimization techniques, integration with reinforcement learning, and application domains.

An innovative approach to fine-tuning is the Semantic-Aware Fine-Tuning (SAFT) method, preserving semantic understanding while enhancing performance in text-based reinforcement learning tasks [90]. Additionally, Manduzio's framework improves the function-calling capabilities of small-scale language models for complex reasoning tasks [30].

Reinforcement learning (RL) has been integrated into fine-tuning stages, as proposed by Solway, allowing for broader policy space exploration and suppression of undesirable outputs [33]. This

13

approach enhances model adaptability and robustness. Models fine-tuned followed by RLHF using Proximal Policy Optimization have shown improved performance across independent reward models [18].

Prompt optimization focuses on refining input prompts to enhance LLM outputs. The RLPROMPT algorithm employs a parameter-efficient policy network for optimized discrete prompts through RL, demonstrating effectiveness in few-shot classification and unsupervised text style transfer [37]. This highlights RL's potential in optimizing prompts, expanding LLM applicability in diverse NLP scenarios.

Sun's method involves fine-tuning on the main dataset and implementing prompt tuning with various templates for classification, improving model performance [32]. Bernardelle's framework provides a comprehensive dataset for evaluating fine-tuning and prompt optimization techniques [92].

In neural machine translation, experiments show significant quality improvements when combining RL training with quality metrics as reward models [93]. Robust methods like R3M tackle corrupted human feedback by modeling preference label corruption as sparse outliers, employing 1-regularized maximum likelihood estimation [77], which is crucial for maintaining performance amidst noisy feedback.

Fine-tuning and prompt optimization are critical for maximizing LLM potential across diverse tasks and domains. By leveraging advanced techniques, these processes continue to evolve, fostering the development of sophisticated and reliable AI systems. Future research should prioritize integrating comprehensive cognitive frameworks to enhance optimization strategies in AI, utilizing RLHF to improve decision-making and rationality, and exploring function calling capabilities in smaller models to facilitate efficient reasoning processes. Critique models can provide essential feedback, refining LLM reasoning abilities and leading to more capable AI systems for complex tasks across various domains [30, 9, 23, 24].

## 6.4 Policy Optimization Techniques

Policy optimization techniques are fundamental for enhancing the performance and efficiency of RL algorithms, particularly in complex environments where decision-making strategies are crucial. These techniques focus on optimizing the policy that governs an RL agent's behavior to maximize rewards over time. Methods such as RLHF and Proximal Policy Optimization (PPO) utilize reward models to evaluate and refine agent actions based on environmental feedback, systematically adjusting the policy to improve performance across various tasks, ensuring outputs are accurate, coherent, and aligned with human expectations [38, 94, 37, 28]. PPO, a widely used method, balances exploration and exploitation through iterative policy updates using a clipped objective function, stabilizing training and preventing destabilizing updates.

An innovative alternative to traditional methods is the ReMax algorithm, which significantly improves computational efficiency, saving approximately 46

Policy gradient techniques directly optimize the policy by estimating the gradient of expected rewards concerning policy parameters, facilitating nuanced modifications that enhance the agent's ability to learn optimal strategies in dynamic environments. Techniques such as multi-dimensional reward modeling and fine-tuning through reinforcement learning enable agents to navigate complex conditions and optimize performance across various metrics [94, 24, 69, 89, 37]. Trust Region Policy Optimization (TRPO) further refines this process, ensuring updates remain within a trust region for improved stability and convergence rates.

Integrating hierarchical reinforcement learning (HRL) frameworks, particularly with meta-learning and intrinsic motivation mechanisms, advances the decomposition of complex tasks into structured sub-tasks. This facilitates efficient policy optimization and enhances agents' rapid adaptation and exploration capabilities. By leveraging past experiences, agents swiftly learn and adapt hierarchical policies, leading to improved cumulative rewards and success rates in challenging environments. Experimental evidence indicates that these advancements significantly enhance the efficiency and versatility of reinforcement learning systems, enabling them to tackle previously unseen tasks more effectively [47, 48, 28].

Policy optimization techniques are continually evolving, driven by ongoing research and innovations in reinforcement learning. Recent developments include fine-tuning LLMs to enhance research

14

ideation, applying meta-learning for few-shot domain adaptation, and introducing discrete prompt optimization methods like RLPrompt, which improve performance across various natural language processing tasks. Novel frameworks such as Pairwise Proximal Policy Optimization (P3O) are being developed to refine LLM alignment with human values, illustrating the dynamic evolution of these techniques in addressing complex challenges and enhancing effectiveness [69, 37, 24, 95]. By incorporating advanced methodologies and leveraging computational efficiencies, these techniques play a pivotal role in advancing AI systems' capabilities, enabling them to tackle increasingly complex challenges more effectively.

## 7 Exploration-Exploitation Trade-off

### 7.1 Exploration-Exploitation Strategies

The exploration-exploitation trade-off is a fundamental challenge in reinforcement learning (RL), necessitating a balance between exploring new actions to discover potentially superior strategies and exploiting known actions to maximize immediate rewards. Effective management of this trade-off is crucial for optimizing learning processes, enabling agents to predict outcomes and adapt behaviors in dynamic environments. This balance is vital for tasks such as research ideation and few-shot learning adaptations [96, 97, 24, 69, 89].

The epsilon-greedy policy is a common strategy for addressing this dilemma, where agents primarily choose actions deemed most beneficial based on accumulated knowledge, while a small probability (epsilon) is reserved for random action selection. This randomness facilitates exploration of less familiar actions, potentially unveiling better strategies and enhancing the learning process [67, 94, 69, 98, 45]. By consistently testing new possibilities, agents avoid premature convergence to suboptimal strategies.

Recent advancements have introduced sophisticated techniques to enhance exploration-exploitation strategies through intrinsic motivation and curiosity-driven exploration. These methods employ a two-stage framework combining Supervised Fine-Tuning with controllable Reinforcement Learning, optimizing research idea generation by navigating trade-offs among novelty, feasibility, and effectiveness, while employing multi-dimensional reward modeling for output evaluation and refinement. Such innovations address limitations in prompting-based models, fostering more effective and context-aware exploration strategies across applications [38, 24]. By providing intrinsic rewards based on prediction error or information gain, these strategies improve exploration efficiency and lead to robust learning outcomes.

Bayesian approaches further contribute by maintaining a distribution over potential models and updating beliefs based on observed data, allowing agents to quantify uncertainty and make informed decisions regarding exploration versus exploitation [99].

Advancements in policy optimization techniques, such as Proximal Policy Optimization (PPO) and Trust Region Policy Optimization (TRPO), have also improved the management of the exploration-exploitation trade-off. These methods impose constraints on policy updates, ensuring stability during exploration while allowing robust exploration of the action space [97, 94, 24, 100, 101].

Hierarchical reinforcement learning (HRL) frameworks enhance agents' ability to operate at various abstraction levels, systematically breaking down complex tasks into manageable sub-tasks. By identifying key subgoals and developing specialized skill policies, HRL addresses challenges related to large-scale applications and sparse reward feedback. Integrating meta-learning techniques further enhances adaptability and exploration efficiency, enabling rapid learning and application of hierarchical policies in diverse environments [49, 58, 48, 47, 102]. This structured approach allows for targeted exploration at different levels, enhancing overall learning efficiency.

The exploration-exploitation trade-off remains a critical focus in RL research, with significant efforts directed toward developing advanced strategies that optimize this balance and enhance the performance of models like large language models (LLMs) through techniques such as Reinforcement Learning from Human Feedback (RLHF) and Direct Preference Optimization (DPO). These approaches aim to refine decision-making policies and align model outputs with human values, addressing challenges in balancing novelty, feasibility, and effectiveness across applications [50, 103, 24, 28]. By leveraging innovative methodologies, researchers continue to enhance RL agents' capabilities, enabling them to navigate complex environments with greater efficacy and adaptability.

15

## 7.2 High-Exploration Preliminary Runs

High-exploration preliminary runs are crucial in the reinforcement learning (RL) process, facilitating the discovery of diverse strategies and solutions to complex problems. These runs prioritize exploration over exploitation, allowing RL agents to gather comprehensive information about their environment, which is essential for optimizing learning outcomes. By emphasizing exploration, agents reduce the risk of prematurely settling on suboptimal strategies, fostering a nuanced understanding of complex dynamics. This approach is particularly beneficial in multi-task domains, where strategically selecting training goals that align with the agent's current capabilities enhances sample efficiency and learning outcomes. Focusing on goals that maximize epistemic uncertainty ensures agents are continually challenged without being overwhelmed, leading to sustained performance improvements across diverse tasks [97, 24].

Advanced exploration strategies guide high-exploration runs, incentivizing agents to seek novel states and actions. Techniques such as intrinsic motivation and curiosity-driven exploration provide intrinsic rewards that encourage investigation of unfamiliar regions of the state space. Recent advancements include utilizing large language models (LLMs) to automatically generate intrinsic reward functions, enabling agents to learn from experiences without extensive external datasets. This approach has demonstrated state-of-the-art performance in challenging sparse reward tasks, highlighting the effectiveness of intrinsic motivation in promoting knowledge expansion and improving agent behavior [38, 104, 67, 24]. Such methods are particularly effective in environments with sparse rewards, where traditional exploration techniques may struggle to identify rewarding strategies.

Bayesian methods further enhance high-exploration runs by incorporating uncertainty into the exploration process. By employing a probabilistic model of the environment, agents can quantify uncertainty, facilitating strategic decisions regarding exploration versus exploitation. This approach is beneficial in RL contexts, where balancing these aspects is critical to avoiding local optima. Recent advancements in automated curriculum generation show that agents can improve learning efficiency by prioritizing tasks aligned with their current capabilities and uncertainties, ultimately enhancing performance across various tasks with differing complexities and exploration challenges [97, 105]. This probabilistic framework ensures a balanced exploration approach, allowing agents to gather valuable information without excessive risk-taking.

Moreover, integrating hierarchical reinforcement learning frameworks enables agents to conduct high-exploration runs at multiple abstraction levels. By systematically breaking down complex tasks into manageable sub-tasks, agents can strategically direct their exploration efforts toward specific areas of interest. This targeted approach enhances learning efficiency and improves agent performance across diverse goals. In multi-task reinforcement learning, prioritizing optimally challenging tasks—neither too difficult nor too easy—can significantly amplify the learning signal, leading to better sample efficiency and overall performance. This methodology draws inspiration from biological learning processes, where meaningful task sequences facilitate knowledge acquisition, fostering continual improvement in agents' abilities [97, 53, 30, 24, 106].

High-exploration preliminary runs are essential for maximizing the potential of RL agents. By employing sophisticated exploration strategies and probabilistic methods, these runs enhance agents' ability to effectively navigate intricate environments. This is achieved through a two-stage approach that first utilizes a high-exploration curriculum to identify learnable tasks and then distills this knowledge into an expert curriculum, optimizing adaptability and performance across diverse scenarios. This methodology not only improves learning efficiency but also addresses challenges posed by local optima and task complexity, ultimately fostering better generalization to unseen situations [53, 105, 24].

# 8 Adaptive and Human-in-the-Loop Learning

## 8.1 Integration of Human Feedback in Reinforcement Learning

Integrating human feedback into reinforcement learning (RL) is crucial for aligning artificial intelligence (AI) systems with human values and preferences. This integration enhances RL by guiding agents towards more desirable behaviors, especially in complex environments where traditional reward functions may not suffice. Reinforcement Learning from Human Feedback (RLHF) is a prominent method that refines AI decision-making through human input [55].

Frameworks like the Okapi instruction-tuning framework improve RLHF's scalability and efficiency by enabling multilingual instruction-following in large language models (LLMs), ensuring applicability to diverse global audiences [26]. Innovative strategies such as incorporating natural language critique into reward models provide nuanced feedback by predicting scalar rewards based on response quality, aligning AI outputs with human expectations [57].

Bayesian approaches further optimize RLHF by employing probabilistic frameworks that account for uncertainty in preferences, allowing systems to balance exploration and exploitation effectively. By integrating a reward model reflecting human preferences, RL systems can adaptively refine outputs, improving performance in tasks like context-aware text generation and complex decision-making [107, 28, 108, 24, 58].

Incorporating human feedback is essential for developing AI technologies that are efficient, reliable, and ethically aligned with human intentions. This integration addresses challenges in aligning AI with human values, driving advancements in AI technologies and ensuring that LLMs function ethically and safely. Researchers emphasize the need for refined methodologies to collect reliable feedback and critically examine the sociotechnical implications to ensure responsible AI development [3, 10].

## 8.2  Adaptive Learning Mechanisms

Adaptive learning mechanisms are vital for enhancing AI systems' flexibility and efficiency, enabling them to adjust learning processes based on performance and environmental changes. These mechanisms continuously refine AI models through real-time feedback and data integration, using techniques like Supervised Fine-Tuning and Reinforcement Learning to optimize performance across dimensions such as novelty, feasibility, and effectiveness [23, 24, 28].

Personalized learning experiences are achieved through meta-learning, allowing models to adapt rapidly to new tasks by leveraging prior knowledge. Hierarchical transformers in meta-learning frameworks enhance RL agents' adaptability by facilitating learning from both intra-episode and inter-episode experiences [48]. Curriculum learning strategies further optimize learning by systematically presenting tasks aligned with learners' knowledge and skills, fostering foundational skill development before advancing to more challenging tasks [97, 53, 89, 109, 42].

Integrating human feedback into adaptive learning processes enhances AI alignment with human values. RLHF refines decision-making strategies, ensuring AI outputs align closely with human expectations [55]. Bayesian approaches maintain probabilistic frameworks that account for uncertainty in human feedback, facilitating informed decision-making and effective exploration-exploitation balance. Advanced techniques like fine-tuning and reinforcement learning enable models to generate and optimize research ideas, navigating complex interdependencies of novelty, feasibility, and effectiveness. Multi-dimensional reward modeling and feedback mechanisms further improve outcomes, even with limited data, enhancing performance in sequential decision-making tasks [69, 110, 24, 50].

Adaptive learning mechanisms are essential for developing AI systems that are efficient, reliable, and capable of continuous improvement and alignment with human values. Techniques such as meta-learning, curriculum learning, and RLHF enhance sophisticated AI technologies, addressing challenges related to variability in expertise and trade-offs between novelty, feasibility, and effectiveness. Refining RLHF can account for diverse expertise levels among feedback providers, improving training reliability. Integrating supervised fine-tuning with controllable reinforcement learning allows dynamic optimization of AI-generated content, leading to more responsive and effective systems capable of navigating intricate problem spaces [66, 24].

# 9  Conclusion

## 9.1  Future Directions and Research Opportunities

Advancements in artificial intelligence research offer substantial prospects for enhancing the capabilities and ethical alignment of Large Language Models (LLMs) and reinforcement learning frameworks. Refining the critique generation process and integrating complex reward modeling objectives are pivotal for aligning AI systems with human values, ensuring both technical proficiency and social responsibility. Enhancing the robustness of algorithms like warmPref-PS for broader applications in Reinforcement Learning from Human Feedback (RLHF) and improving their resilience against mis-

specified parameters remain critical. Furthermore, developing comprehensive feedback mechanisms and addressing ethical concerns in AI technologies are essential research avenues.

Optimizing critique models and exploring their applicability across diverse domains can significantly boost LLM capabilities. Future efforts should emphasize sophisticated preference modeling and the integration of continuous learning mechanisms to adapt to evolving user preferences. Expanding benchmarks to encompass more languages, enhancing data quality, and exploring additional evaluation metrics related to biases and ethical considerations are vital for advancing instruction-tuned LLMs. Additionally, researching automated methods to determine optimal training cessation points, thereby reducing reliance on human evaluation, presents a promising direction. Enhancing algorithms to predict workload patterns and integrating machine learning techniques for resource optimization are crucial areas for future exploration.

These opportunities underscore the dynamic nature of AI development, providing pathways to enhance capabilities, efficiency, and ethical alignment across various domains. Addressing these challenges will significantly contribute to creating more sophisticated and reliable AI systems.

## 9.2 Challenges and Future Directions

The field of artificial intelligence is confronted with numerous challenges that necessitate innovative solutions to improve the robustness, efficiency, and ethical alignment of AI systems. A major issue is the computational cost associated with reinforcement learning, particularly during the exploration phase of RLHF, which limits scalability due to computational overheads that can impede real-time applications. The simultaneous training of multiple models further exacerbates this burden. Another challenge lies in the generalization performance of reward models, especially with unseen data. While the Generalized Reward Model (GRM) shows promise in mitigating over-optimization issues, its application remains limited to specific scenarios. The complexity of isolating the effects of different training datasets highlights the need for transparency in model training processes to better understand bias emergence.

Integrating human feedback into AI systems, though beneficial, presents challenges such as high memory consumption during RLHF fine-tuning. Strategies like memory-efficient APIs can alleviate this issue with minimal impact on processing time. Additionally, frameworks like RLSF offer a more economical alternative by incorporating statistical business feedback, aligning reinforcement learning with commercial objectives. Aligning AI systems with human preferences is intricate and can lead to unintended behaviors such as verbosity and evasiveness. This issue is compounded by the limitations of current studies, which may not address all potential misleading behaviors across diverse domains. The focus on relatively simple language tasks within a single domain suggests a need to explore more complex environments and language forms for broader applicability.

Concerns about potential overfitting and data efficiency require further investigation. Future AI research should prioritize refining critique models and integrating more complex reward modeling objectives to enhance alignment with human values. Expanding benchmarks to include additional languages and improving data quality are essential steps for advancing instruction-tuned LLMs. Moreover, exploring automated methods to identify optimal training cessation points, thereby minimizing reliance on human evaluation, presents a promising research avenue. Addressing these challenges and pursuing these future directions will significantly contribute to developing AI systems that are not only technically proficient but also socially responsible, enhancing their applicability and ethical alignment across various domains.

# References

[1] Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. Principle-driven self-alignment of language models from scratch with minimal human supervision, 2023.

[2] Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Amrit Singh Bedi, and Mengdi Wang. Maxmin-rlhf: Towards equitable alignment of large language models with diverse human preferences, 2024.

[3] Adam Dahlgren Lindström, Leila Methnani, Lea Krause, Petter Ericson, Íñigo Martínez de Rituerto de Troya, Dimitri Coelho Mollo, and Roel Dobbe. Ai alignment through reinforcement learning from human feedback? contradictions and limitations, 2024.

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

[5] Chenliang Li, Siliang Zeng, Zeyi Liao, Jiaxiang Li, Dongyeop Kang, Alfredo Garcia, and Mingyi Hong. Joint demonstration and preference learning improves policy alignment with human feedback, 2024.

[6] Chen Zheng, Ke Sun, Hang Wu, Chenguang Xi, and Xun Zhou. Balancing enhancement, harmlessness, and general capabilities: Enhancing conversational llms with direct rlhf, 2024.

[7] Ruoxi Cheng, Haoxuan Ma, Shuirong Cao, Jiaqi Li, Aihua Pei, Zhiqiang Wang, Pengliang Ji, Haoyu Wang, and Jiaqi Huo. Reinforcement learning from multi-role debates as feedback for bias mitigation in llms, 2024.

[8] Ali Ayub, Jainish Mehta, Zachary De Francesco, Patrick Holthaus, Kerstin Dautenhahn, and Chrystopher L. Nehaniv. How do human users teach a continual learning robot in repeated interactions?, 2023.

[9] Dana Alsagheer, Rabimba Karanjai, Nour Diallo, Weidong Shi, Yang Lu, Suha Beydoun, and Qiaoning Zhang. Comparing rationality between large language models and humans: Insights and open questions, 2024.

[10] Thilo Hagendorff and Sarah Fabi. Methodological reflections for ai alignment research using human feedback, 2022.

[11] Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Aligning large language models with human preferences through representation engineering, 2024.

[12] Ryan Aponte, Ryan A. Rossi, Shunan Guo, Franck Dernoncourt, Tong Yu, Xiang Chen, Subrata Mitra, and Nedim Lipka. A framework for fine-tuning llms using heterogeneous feedback, 2024.

[13] Adarsh N L, Arun P V au2, and Aravindh N L. Enhancing image caption generation using reinforcement learning with human feedback, 2024.

[14] Zixuan Liu, Xiaolin Sun, and Zizhan Zheng. Enhancing llm safety via constrained direct preference optimization, 2024.

[15] Khanh Nguyen. Language models are bounded pragmatic speakers: Understanding rlhf from a bayesian cognitive modeling perspective, 2024.

[16] Ruitao Chen and Liwei Wang. The power of active multi-task learning in reinforcement learning from human feedback, 2024.

[17] Ben Hauptvogel, Malte Ostendorff, Georg Rehm, and Sebastian Möller. Reward modeling with weak supervision for language models, 2024.

[18] Yanjun Chen, Dawei Zhu, Yirong Sun, Xinghao Chen, Wei Zhang, and Xiaoyu Shen. The accuracy paradox in rlhf: When better reward models don't yield better language models, 2024.

[19] Alexey Gorbatovski and Sergey Kovalchuk. Reinforcement learning for question answering in programming domain using public community scoring as a human feedback, 2024.

[20] Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. Teaching large language models to reason with reinforcement learning, 2024.

[21] Yifu Yuan, Jianye Hao, Yi Ma, Zibin Dong, Hebin Liang, Jinyi Liu, Zhixin Feng, Kai Zhao, and Yan Zheng. Uni-rlhf: Universal platform and benchmark suite for reinforcement learning with diverse human feedback, 2024.

[22] Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in rlhf, 2024.

[23] Zhiheng Xi, Dingwen Yang, Jixuan Huang, Jiafu Tang, Guanyu Li, Yiwen Ding, Wei He, Boyang Hong, Shihan Do, Wenyu Zhan, Xiao Wang, Rui Zheng, Tao Ji, Xiaowei Shi, Yitao Zhai, Rongxiang Weng, Jingang Wang, Xunliang Cai, Tao Gui, Zuxuan Wu, Qi Zhang, Xipeng Qiu, Xuanjing Huang, and Yu-Gang Jiang. Enhancing llm reasoning via critique models with test-time and training-time supervision, 2024.

[24] Ruochen Li, Liqiang Jing, Chi Han, Jiawei Zhou, and Xinya Du. Learning to generate research idea with dynamic control, 2024.

[25] Dingli Yu, Simran Kaur, Arushi Gupta, Jonah Brown-Cohen, Anirudh Goyal, and Sanjeev Arora. Skill-mix: a flexible and expandable family of evaluations for ai models, 2023.

[26] Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback, 2023.

[27] Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. Sayself: Teaching llms to express confidence with self-reflective rationales, 2024.

[28] Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu, Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin Wang, and Eduard Hovy. Reinforcement learning enhanced llms: A survey, 2024.

[29] Juhao Liang, Ziwei Wang, Zhuoheng Ma, Jianquan Li, Zhiyi Zhang, Xiangbo Wu, and Benyou Wang. Online training of large language models: Learn while chatting, 2024.

[30] Graziano A. Manduzio, Federico A. Galatolo, Mario G. C. A. Cimino, Enzo Pasquale Scilingo, and Lorenzo Cominelli. Improving small-scale large language models function calling for reasoning tasks, 2024.

[31] Priyanshu Gupta, Shashank Kirtania, Ananya Singha, Sumit Gulwani, Arjun Radhakrishna, Sherry Shi, and Gustavo Soares. Metareflection: Learning instructions for language agents using past reflections, 2024.

[32] Pingwei Sun. Fine-tuning vs prompting, can language models understand human values?, 2024.

[33] Alec Solway. Reinforcement learning without human feedback for last mile fine-tuning of large language models, 2024.

[34] Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity, 2024.

[35] William Muldrew, Peter Hayes, Mingtian Zhang, and David Barber. Active preference learning for large language models, 2024.

[36] Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang Low. Prompt optimization with human feedback, 2024.

[37] Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P. Xing, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement learning, 2022.

[38] Will LeVine, Benjamin Pikus, Anthony Chen, and Sean Hendryx. A baseline analysis of reward models' ability to accurately analyze foundation models under distribution shift, 2024.

[39] Xingmeng Zhao, Tongnian Wang, Sheri Osborn, and Anthony Rios. Babystories: Can reinforcement learning teach baby language models to write better stories?, 2023.

[40] Joey Hong, Sergey Levine, and Anca Dragan. Zero-shot goal-directed dialogue via rl on imagined conversations, 2023.

[41] John Dang, Arash Ahmadian, Kelly Marchisio, Julia Kreutzer, Ahmet Üstün, and Sara Hooker. Rlhf can speak many languages: Unlocking multilingual preference optimization for llms, 2024.

[42] Ahmed H. Zaidi, Russell Moore, and Ted Briscoe. Curriculum q-learning for visual vocabulary acquisition, 2017.

[43] Iddo Drori, Yamuna Krishnamurthy, Raoni Lourenco, Remi Rampin, Kyunghyun Cho, Claudio Silva, and Juliana Freire. Automatic machine learning by pipeline synthesis using model-based reinforcement learning and a grammar, 2019.

[44] Shangding Gu, Alois Knoll, and Ming Jin. Teams-rl: Teaching llms to generate better instruction datasets via reinforcement learning, 2024.

[45] Sharath Chandra Raparthy, Bhairav Mehta, Florian Golemo, and Liam Paull. Generating automatic curricula via self-supervised active domain randomization, 2020.

[46] Jingyi Chen, Ju-Seung Byun, Micha Elsner, and Andrew Perrault. Reinforcement learning for fine-tuning text-to-speech diffusion models, 2024.

[47] Arash Khajooeinejad and Masoumeh Chapariniya. Meta-learning integration in hierarchical reinforcement learning for advanced task complexity, 2024.

[48] Gresa Shala, André Biedenkapp, and Josif Grabocka. Hierarchical transformers are efficient meta-reinforcement learners, 2024.

[49] Jacob Rafati and David C. Noelle. Learning representations in model-free hierarchical reinforcement learning, 2019.

[50] Martin Klissarov, Devon Hjelm, Alexander Toshev, and Bogdan Mazoure. On the modeling capabilities of large language models for sequential decision making, 2024.

[51] Evan Zheran Liu, Sahaana Suri, Tong Mu, Allan Zhou, and Chelsea Finn. Simple embodied language learning as a byproduct of meta-reinforcement learning, 2023.

[52] Rundong Wang, Longtao Zheng, Wei Qiu, Bowei He, Bo An, Zinovi Rabinovich, Yujing Hu, Yingfeng Chen, Tangjie Lv, and Changjie Fan. Towards skilled population curriculum for multi-agent reinforcement learning, 2023.

[53] Rémy Portelas, Katja Hofmann, and Pierre-Yves Oudeyer. Trying again instead of trying longer: Prior learning for automatic curriculum learning, 2020.

[54] Wenshuai Zhao, Zhiyuan Li, and Joni Pajarinen. Learning progress driven multi-agent curriculum, 2024.

[55] Hakim Sidahmed, Samrat Phatale, Alex Hutcheson, Zhuonan Lin, Zhang Chen, Zac Yu, Jarvis Jin, Simral Chaudhary, Roman Komarytsia, Christiane Ahlheim, Yonghao Zhu, Bowen Li, Saravanan Ganesh, Bill Byrne, Jessica Hoffmann, Hassan Mansoor, Wei Li, Abhinav Rastogi, and Lucas Dixon. Parameter efficient reinforcement learning from human feedback, 2024.

[56] Akhil Agnihotri, Rahul Jain, Deepak Ramachandran, and Zheng Wen. Online bandit learning with offline preference data, 2024.

[57] Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D. Chang, and Prithviraj Ammanabrolu. Critique-out-loud reward models, 2024.

[58] Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik Narasimhan, Ameet Deshpande, and Bruno Castro da Silva. Rlhf deciphered: A critical analysis of reinforcement learning from human feedback for llms, 2024.

[59] Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, Limao Xiong, Lu Chen, Zhiheng Xi, Nuo Xu, Wenbin Lai, Minghao Zhu, Cheng Chang, Zhangyue Yin, Rongxiang Weng, Wensen Cheng, Haoran Huang, Tianxiang Sun, Hang Yan, Tao Gui, Qi Zhang, Xipeng Qiu, and Xuanjing Huang. Secrets of rlhf in large language models part i: Ppo, 2023.

[60] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.

[61] Gabrielle Kaili-May Liu. Perspectives on the social impacts of reinforcement learning with human feedback, 2023.

[62] Chanwoo Park, Mingyang Liu, Dingwen Kong, Kaiqing Zhang, and Asuman Ozdaglar. Rlhf from heterogeneous feedback via personalization and preference aggregation, 2024.

[63] Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, Songyang Gao, Nuo Xu, Yuhao Zhou, Xiaoran Fan, Zhiheng Xi, Jun Zhao, Xiao Wang, Tao Ji, Hang Yan, Lixing Shen, Zhan Chen, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. Secrets of rlhf in large language models part ii: Reward modeling, 2024.

[64] Yannick Metz, David Lindner, Raphaël Baur, and Mennatallah El-Assady. Mapping out the space of human feedback for reinforcement learning: A conceptual framework, 2025.

[65] Qining Zhang and Lei Ying. Zeroth-order policy gradient for reinforcement learning from human feedback without reward inference, 2024.

[66] Oliver Daniels-Koch and Rachel Freedman. The expertise problem: Learning from specialized feedback, 2022.

[67] Qinqing Zheng, Mikael Henaff, Amy Zhang, Aditya Grover, and Brandon Amos. Online intrinsic rewards for decision making agents from large language model feedback, 2024.

[68] Yuzi Yan, Xingzhou Lou, Jialian Li, Yiping Zhang, Jian Xie, Chao Yu, Yu Wang, Dong Yan, and Yuan Shen. Reward-robust rlhf in llms, 2024.

[69] Reinforcement learning for few-shot text generation adaptation.

[70] Shihan Dou, Yan Liu, Enyu Zhou, Tianlong Li, Haoxiang Jia, Limao Xiong, Xin Zhao, Junjie Ye, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. Metarm: Shifted distributions alignment via meta-learning, 2024.

[71] Adam X. Yang, Maxime Robeyns, Thomas Coste, Zhengyan Shi, Jun Wang, Haitham Bou-Ammar, and Laurence Aitchison. Bayesian reward models for llm alignment, 2024.

[72] Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan, Chenguang Zhu, Richard Yuanzhe Pang, Yundi Qian, Xuewei Wang, Suchin Gururangan, Chao Zhang, Melanie Kambadur, Dhruv Mahajan, and Rui Hou. Self-generated critiques boost reward modeling for language models, 2025.

[73] Wanqiao Xu, Shi Dong, Xiuyuan Lu, Grace Lam, Zheng Wen, and Benjamin Van Roy. Rlhf and iia: Perverse incentives, 2024.

[74] Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework, 2024.

[75] Eunseop Yoon, Hee Suk Yoon, SooHwan Eom, Gunsoo Han, Daniel Wontae Nam, Daejin Jo, Kyoung-Woon On, Mark A. Hasegawa-Johnson, Sungwoong Kim, and Chang D. Yoo. Tlcr: Token-level continuous reward for fine-grained reinforcement learning from human feedback, 2024.

[76] Tengyang Xie, Dylan J. Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q*-approximation for sample-efficient rlhf, 2024.

[77] Alexander Bukharin, Ilgee Hong, Haoming Jiang, Zichong Li, Qingru Zhang, Zixuan Zhang, and Tuo Zhao. Robust reinforcement learning from corrupted human feedback, 2024.

[78] Moritz Stephan, Alexander Khazatsky, Eric Mitchell, Annie S Chen, Sheryl Hsu, Archit Sharma, and Chelsea Finn. Rlvf: Learning from verbal feedback without overgeneralization, 2024.

[79] Muhan Lin, Shuyang Shi, Yue Guo, Behdad Chalaki, Vaishnav Tadiparthi, Ehsan Moradi Pari, Simon Stepputtis, Joseph Campbell, and Katia Sycara. Navigating noisy feedback: Enhancing reinforcement learning with error-prone language models, 2024.

[80] Xinyu Li, Ruiyang Zhou, Zachary C. Lipton, and Liu Leqi. Personalized language modeling from personalized human feedback, 2024.

[81] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging, 2023.

[82] 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, 2024.

[83] Alexey Kutalev and Sergei Markoff. Investigating on rlhf methodology, 2024.

[84] Yuan Sun, Navid Salami Pargoo, Peter J. Jin, and Jorge Ortiz. Optimizing autonomous driving for safety: A human-centric approach with llm-enhanced rlhf, 2024.

[85] Ching-An Cheng, Andrey Kolobov, Dipendra Misra, Allen Nie, and Adith Swaminathan. Llf-bench: Benchmark for interactive learning from language feedback, 2023.

[86] Kanghyun Ryu, Qiayuan Liao, Zhongyu Li, Koushil Sreenath, and Negar Mehr. Curricullm: Automatic task curricula design for learning complex robot skills using large language models, 2024.

[87] Nuo Xu, Jun Zhao, Can Zu, Sixian Li, Lu Chen, Zhihao Zhang, Rui Zheng, Shihan Dou, Wenjuan Qin, Tao Gui, Qi Zhang, and Xuanjing Huang. Advancing translation preference modeling with rlhf: A step towards cost-effective solution, 2024.

[88] Anusha Nagabandi, Chelsea Finn, and Sergey Levine. Deep online learning via meta-learning: Continual adaptation for model-based rl, 2019.

[89] Sidney Tio and Pradeep Varakantham. Transferable curricula through difficulty conditioned generators, 2023.

[90] Mauricio Gruppi, Soham Dan, Keerthiram Murugesan, and Subhajit Chaudhury. On the effects of fine-tuning language models for text-based reinforcement learning, 2024.

[91] Prashant Trivedi, Souradip Chakraborty, Avinash Reddy, Vaneet Aggarwal, Amrit Singh Bedi, and George K. Atia. Align-pro: A principled approach to prompt optimization for llm alignment, 2025.

[92] Pietro Bernardelle and Gianluca Demartini. Optimizing llms with direct preferences: A data efficiency perspective, 2024.

[93] Miguel Moura Ramos, Patrick Fernandes, António Farinhas, and André F. T. Martins. Aligning neural machine translation models: Human feedback in training and inference, 2024.

[94] Miao Fan, Chen Hu, and Shuchang Zhou. Proximal policy optimization actual combat: Manipulating output tokenizer length, 2023.

[95] Tianhao Wu, Banghua Zhu, Ruoyu Zhang, Zhaojin Wen, Kannan Ramchandran, and Jiantao Jiao. Pairwise proximal policy optimization: Harnessing relative feedback for llm alignment, 2023.

[96] Margaret Li, Weijia Shi, Artidoro Pagnoni, Peter West, and Ari Holtzman. Predicting vs. acting: A trade-off between world modeling agent modeling, 2024.

[97] Yunzhi Zhang, Pieter Abbeel, and Lerrel Pinto. Automatic curriculum learning through value disagreement, 2020.

[98] Viraj Mehta, Vikramjeet Das, Ojash Neopane, Yijia Dai, Ilija Bogunovic, Jeff Schneider, and Willie Neiswanger. Sample efficient reinforcement learning from human feedback via active exploration, 2023.

[99] 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, 2024.

[100] Yanshi Li, Shaopan Xiong, Gengru Chen, Xiaoyang Li, Yijia Luo, Xingyao Zhang, Yanhui Huang, Xingyuan Bu, Yingshui Tan, Chun Yuan, Jiamang Wang, Wenbo Su, and Bo Zheng. Adaptive dense reward: Understanding the gap between action and reward space in alignment, 2024.

[101] Rui Yang, Ruomeng Ding, Yong Lin, Huan Zhang, and Tong Zhang. Regularizing hidden states enables learning generalizable reward model for llms, 2024.

[102] Kai Wang, Zhonghao Wang, Mo Yu, and Humphrey Shi. Feudal reinforcement learning by reading manuals, 2021.

[103] Qi Gou and Cam-Tu Nguyen. Mixed preference optimization: Reinforcement learning with data selection and better reference model, 2025.

[104] Yong Lin, Skyler Seto, Maartje ter Hoeve, Katherine Metcalf, Barry-John Theobald, Xuan Wang, Yizhe Zhang, Chen Huang, and Tong Zhang. On the limited generalization capability of the implicit reward model induced by direct preference optimization, 2024.

[105] Pascal Klink, Hany Abdulsamad, Boris Belousov, Carlo D'Eramo, Jan Peters, and Joni Pajari-nen. A probabilistic interpretation of self-paced learning with applications to reinforcement learning, 2021.

[106] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks, 2023.

[107] Kristian González Barman, Simon Lohse, and Henk de Regt. Reinforcement learning from human feedback: Whose culture, whose values, whose perspectives?, 2025.

[108] Nathan Lambert, Thomas Krendl Gilbert, and Tom Zick. The history and risks of reinforcement learning and human feedback, 2023.

[109] Tobias Niehues, Ulla Scheler, and Pascal Klink. Self-paced absolute learning progress as a regularized approach to curriculum learning, 2023.

[110] Ruizhe Shi, Yuyao Liu, Yanjie Ze, Simon S. Du, and Huazhe Xu. Unleashing the power of pre-trained language models for offline reinforcement learning, 2024.

**Disclaimer:**

SurveyX is an AI-powered system designed to automate the generation of surveys. While it aims to produce high-quality, coherent, and comprehensive surveys with accurate citations, the final output is derived from the AI's synthesis of pre-processed materials, which may contain limitations or inaccuracies. As such, the generated content should not be used for academic publication or formal submissions and must be independently reviewed and verified. The developers of SurveyX do not assume responsibility for any errors or consequences arising from the use of the generated surveys.

25