



BFCAl
Faculty of Computers and Artificial Intelligence
Computer Science Department

Reinforcement Learning in Generative AI: A Comprehensive Study of RLHF, Reward Modeling, and Applications in Game Design

R. Ali, R.Sharif, and O. Abdelhameed, " Reinforcement Learning in Generative AI: A Comprehensive Study of RLHF, Reward Modeling, and Applications in Game Design" supervised by Dr.W.Shalash, Faculty of Computers & Artificial Intelligence (BFCAl), Egypt.

Abstract

Generative Artificial Intelligence (AI) and Reinforcement Learning (RL) are among the most transformative advancements in modern Computer Science. RL has demonstrated remarkable success across various domains, including generative AI, games, and natural language processing (NLP). RL has proven to be a powerful paradigm for optimizing decision-making and has been integrated into generative AI for content creation, game design, and large language models (LLMs).This survey

explores the intersection of RL and generative AI, discussing its applications in content generation, optimization-driven outputs, and embedding complex characteristics into models. Additionally, we analyze the role of RL in game-based AI research, highlighting its dominance in deep learning literature and its impact on training sophisticated models surpassing human performance. Furthermore, we examine RL's growing influence in real-world applications such as healthcare, finance, robotics, and NLP, emphasizing its adaptability, state representation capabilities, and potential for improving large language models. By synthesizing recent advancements, we identify key trends, challenges, and future directions in leveraging RL for generative AI and beyond.

Keywords: reinforcement learning; deep learning; artificial intelligence; review; game industry

1.Introduction

Generative AI, including large language models (LLMs) like ChatGPT, is rapidly advancing. While these models have shown success in many domains, they still struggle with tasks requiring logical reasoning or specialized knowledge. Reinforcement learning from human feedback (RLHF) has proven effective in improving generative models, aligning their outputs with human values. However, enhancing these models while ensuring correct interaction with the environment remains a challenge due to sparse and potentially inaccurate reward signals. This survey aims to review the applications and challenges of applying reinforcement learning (RL) to generative AI, with a focus on large models. It discusses applications in fields like natural language processing, computer vision, and code generation. It also addresses future directions and open research questions to enhance the effectiveness of these models.

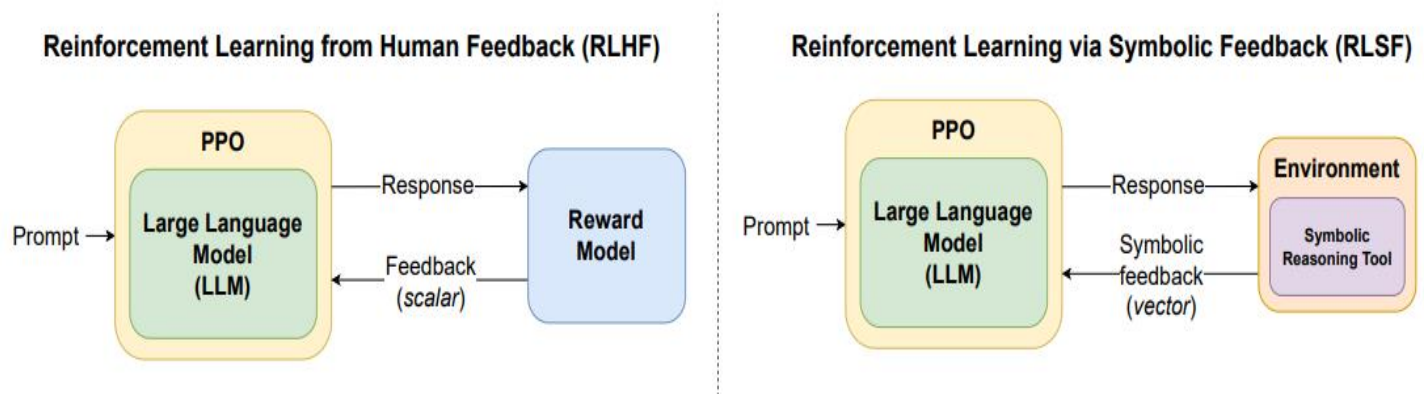


Figure 1: Contrasting RLHF with RLSF: The image depicts two distinct fine-tuning paradigms. (Left) RLHF operates within an environment governed by a black-box reward model, typically offering scalar feedback. (Right) By contrast, the environment in RLSF leverages sound symbolic reasoning tools and also provides fine-grained token-level vector (dense) feedback that is, in turn, based on poly-sized certificates produced by these symbolic tools.

2. Background & Related Work

2.1 Overview of RL and LLMs

Reinforcement Learning (RL) is an area of machine learning designed to solve decision-making problems, where an agent learns by interacting with an environment. Unlike supervised or unsupervised learning, RL requires the agent to autonomously determine the best sequence of actions to maximize a goal, known as a cumulative reward. The RL framework is based on a **Markov Decision Process (MDP)**, which includes:

- **States (S)**: Describes the current situation of the agent.
- **Actions (A)**: The choices the agent can make.
- **Transition function (T)**: The probability of transitioning from one state to another after taking an action.
- **Reward function (R)**: The feedback the agent receives after taking an action.

The goal of the agent is to find an **optimal policy (π)** that maps states to actions, maximizing long-term rewards. RL methods can be divided into three main categories: **dynamic programming**, **Monte Carlo methods**, and **temporal**

difference methods. In simpler environments, the agent can store policies in a lookup table, but more complex tasks use function approximators like **neural networks** to model the policy.

In **Natural Language Processing (NLP)**, **Large Language Models (LLMs)**, such as **BERT**, **GPT**, **PaLM**, and **LaMDA**, have become crucial tools. These models aim to predict word sequences, leveraging the chain rule of probability to break down the probability of word sequences into conditional probabilities. LLMs use **transformer-based architectures**, where the mechanism of **attention** weighs the importance of each word in a sentence. The term "large" refers to the immense number of parameters these models contain, enabling them to perform complex tasks.

LLMs can generate text, answer queries, translate languages, summarize content, and even generate creative writing, such as poetry. Their deep learning architecture, specifically the **transformer network**, has revolutionized NLP, leading to the widespread adoption of LLMs in various applications.

Connection between RL and LLMs Incorporating **Reinforcement Learning** with LLMs enhances the adaptability and efficiency of NLP models. RL techniques, such as **Reinforcement Learning from Human Feedback (RLHF)**, can refine LLMs by providing feedback loops that optimize model behavior in dynamic, real-world tasks, like conversational AI and content generation. These methods enable LLMs to learn from both labeled and unstructured data, creating more effective and human-like interactions.

2.2 State-of-Art Review Studies

Reinforcement Learning (RL) and Natural Language Processing (NLP) have rapidly evolved, attracting considerable research attention due to their broad range of applications. Both fields have produced a plethora of survey studies aimed at synthesizing and evaluating state-of-the-art research.

RL, since its inception, has captured the interest of researchers in computer science, robotics, and control. As a result, numerous surveys have been published covering a wide spectrum of RL topics. These range from general overviews of RL [63, 5] to more focused reviews on specific techniques such as Offline RL [98], Meta-Reinforcement Learning [11], and RL on graphs [84]. Other surveys address RL applications in fields like healthcare [140], robotics [48], and generative AI [20].

Additionally, studies have explored RL in dynamically changing environments [91] and complex systems like Multi-Agent Deep RL [51, 40]. With the rise of Large Language Models (LLMs), researchers have also started to publish surveys dedicated to integrating RL with LLMs, such as those focusing on Reinforcement Learning from Human Feedback (RLHF) [116].

A similar trend is seen in NLP, particularly since the introduction of deep learning techniques. Several surveys have analyzed the progression of NLP concepts and methods, including works on pretrained models [100], graphs in NLP [82], and applications in sectors such as healthcare [133] and fake news detection [88]. Furthermore, LLMs, which bridge RL and NLP, have led to a growing number of literature reviews. These reviews cover topics such as model evaluation [53, 23], human alignment [111, 128], explainability [147], responsible AI [47], and knowledge acquisition and updating [19, 126, 92]. LLMs are also explored in specific applications like information retrieval [151], natural language understanding [39], software engineering [124, 43], and recommendation systems [134, 70, 72].

This survey takes a distinct approach compared to previous review papers by concentrating exclusively on studies where both RL and LLMs are integral components of the same computational framework. As will be explained in subsection 2.3, this focus on the intersection of these two technologies offers a novel perspective on their combined potential.

2.3 Scope of This Study

This study focuses on surveying research that integrates Reinforcement Learning (RL) and Large Language Models (LLMs) within a common modeling framework. A new taxonomy is proposed to classify these studies, visualized in the RL/LLM Taxonomy Tree (Fig. 3). The taxonomy categorizes each study based on how the two models—RL and LLMs—interact.

While Reinforcement Learning, particularly in its RL from Human Feedback (RLHF) form, is essential to the functioning of LLMs, this review is concerned with studies where already-trained LLMs are improved, fine-tuned with RL, or combined with an RL agent to perform specific downstream tasks. Studies focusing solely on using RL to train the original LLMs are not included in our scope.

Additionally, while several state-of-the-art surveys have explored applications of LLMs in tasks unrelated to natural language, such as reasoning [58, 78, 130], multimodal LLMs [127], and autonomous agents [125, 73], this survey emphasizes

studies where the RL agent performs the downstream task. In these studies, the LLM contributes either during training (LLM4RL) or at inference (RL+LLM), rather than functioning as an autonomous agent. The performance of LLMs as independent agents is beyond the scope of this review.

For completeness, this survey also touches on studies where pretrained language models are used to aid RL agents through reward design, policy priors, or policy transfer [21, 30, 62], although these studies predate the widespread use of LLMs. Our taxonomy specifically focuses on studies utilizing LLMs.

3. Methodology

3.1. Search Strategy

The search strategy employed for this research aimed to gather a comprehensive set of relevant studies to understand how Reinforcement Learning (RL) is integrated with Large Language Models (LLMs). The following methods were used:

- **Databases Used:** The primary academic databases accessed were:
 - Google Scholar
 - IEEE Xplore
 - ACM Digital Library
- **Keywords:** The following search keywords were utilized:
 - "Reinforcement Learning and Large Language Models"
 - "RL and LLMs integration"
 - "Reinforcement Learning for NLP tasks"
 - "RLHF and LLMs"
 - "Deep RL in NLP"

3.2. Reinforcement Learning

In most RL algorithms, the agent obtains a model of the environment or at least some basic state transition sequences, as is depicted in Figure 1. In a similar model, the agent can interact with the environment by selecting a set of actions that alter

the environment's state, producing new states along the way. The structural components of RL are: 1. The discrete different time-steps t ; 2. The state space S with state S_t at time-step t ; 3. A set of actions A with action A_t at time-step t ; 4. The policy function $\pi(\cdot)$; 5. A reward function $R_a(S_t, S_{t+1})$ of an action A_t , transitioning from state S to S_{t+1} ; 6. The state evaluation $V(s)$ and energy evaluation $Q(s, a)$

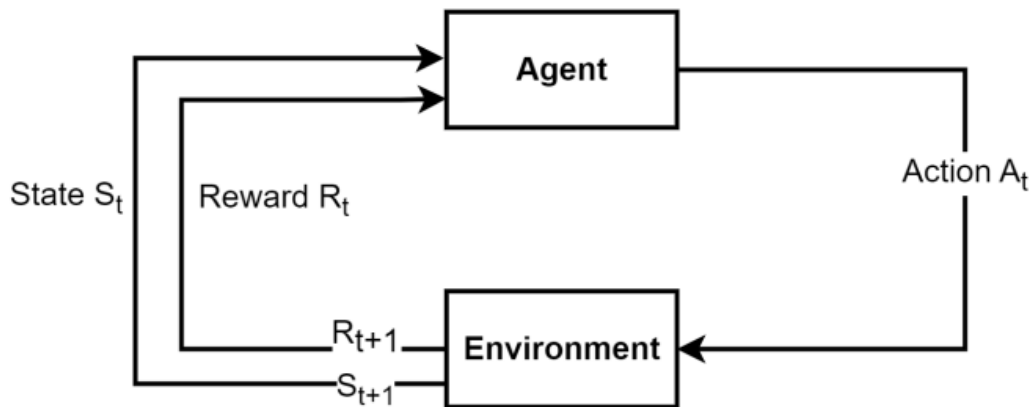


Figure 2. basic RL model.

3.2.1. Q-Learning Algorithm One of the most well-known temporal difference (TD) algorithms is Q-learning (see Algorithm 1) [26,27]. Q-learning is an out-of-policy algorithm; therefore, its policy does not have to coincide with the evaluated and updated policy. It uses the following update rule: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [r_{t+1} + \gamma Q(S_{t+1}, a) - Q(S_t, a)]$. (13)

This algorithm approximates the best function Q^* , independently from the policy that the agent follows, as $\gamma \max_a Q(S_{t+1}, a)$ refers to the best action the agent can perform being at state S_{t+1} .

Algorithm 1 Q-learning algorithm.

```
1: initialize  $Q(S_t, A_t)$ 
2: for every episode do
3:   observe state  $S_t$ 
4:   while  $S_t$  in terminal do
5:     select action  $A_t$  and evaluate  $Q$ 
6:     take action  $a$ 
7:     observe  $r, S_t$ 
8:      $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[r_{t+1} + \gamma Q(S_{t+1}, a) - Q(S_t, a)]$ 
9:      $S_t \leftarrow S_{t+1}$ 
10:  end while
11: end for
```

for So far, all these methods require intense memory allocations to work correctly. Specifically, we must keep s and a for every state S_t and action A_t . A solution is impossible in real-world applications where the state space is vast. This is why reward functions need to be approximated with other types of functions, such as parametric [28]. Therefore: $Q(s, a) \approx Q_\theta(s, a)$.

4. Results & Analysis

4.1. The Origins of RLHF

Learning behavior from human feedback has long been studied as a subfield of RL, but methods and terminology have evolved over time. Early methods focused on learning directly from human rewards (Knox, 2012; Isbell et al., 2001; Knox & Stone, 2009), action advice (Maclin et al., 2005), or action critique (Judah et al., 2010). Notable approaches in this area include TAMER (Knox & Stone, 2009; Warnell et al., 2018), which interprets human feedback as samples of the optimal action-value function, and the later COACH (MacGlashan et al., 2017; Arumugam et al., 2019), which interprets human feedback in a policy-dependent way, i.e., as samples of the advantage function. This survey, however, focuses on more indirect approaches to inferring the objective from human feedback.

Feedback Type	PbRL	SSRL	RLHF
Binary trajectory comparisons	✓	✗	✓
Trajectory rankings	✓	✗	✓
State preferences	✓	✗	✓
Action preferences	✓	✗	✓
Binary critique	✗	✓	✓
Scalar feedback	✗	✓	✓
Corrections	✗	✗	✓
Action advice	✗	✗	✓
Implicit feedback	✗	✗	✓
Natural language	✗	✗	✓

Table 1: Feedback types classified as belonging to PbRL, SSRL, and RLHF as defined in this survey.

Reinforcement learning from human feedback (RLHF) in its modern guise has its origin in the setting of preference-based reinforcement learning (PbRL) as introduced independently by Akroure et al. (2011) and Cheng et al. (2011). The original idea of preference-based reinforcement learning (PbRL) is to infer the

objective from qualitative feedback, such as pairwise preferences between behaviors or between actions given states, instead of quantitative feedback in the form of numerical rewards. The term RLHF was coined as an alternative later on (Askell et al., 2021; Ouyang et al., 2022; OpenAI, 2022), though initially referring to the same concept of learning behavior from relative feedback.

4.2. Statistics So as to demonstrate the first argument, we examined Scopus (curated abstract and citation database, <https://www.scopus.com/> (accessed on 20 December 2022)), Google Scholar (search engine for scholarly literature, <https://scholar.google.com/> (accessed on 20 December 2022)), and Dimensions (linked research information dataset, <https://www.dimensions.ai/> (accessed on 20 December 2022)) publication/literature data for each year, searching for terms and keywords such as “deep reinforcement learning,” “games,” “reinforcement learning,” and “deep learning,” as illustrated in the corresponding graph (Figure 3). More specifically, the total publication (proven research) count for the jointed terms “reinforcement learning” and “games” was 92,367 for the time interval 2012–2022, and the publication count for the jointed words “deep reinforcement learning” and “games” was 56,613 for the period 2012–2022. For the second point, we examined the period 2012–2022, utilizing publication data from Scopus, Google Scholar, and Dimensions focusing on the broader fields of “deep reinforcement learning” and “reinforcement learning”. For the terms and keywords “reinforcement” and “learning” and “games”, 92,367 findings were reported between 2012 and 2022. From our previous 2018–2022 findings, we observed that 73,245 results had been released, accounting for almost 73% of all publications. To summarize, the recent trend in RL-based research on games has focused on DL approaches. Specifically, several methods from genetic programming (GP) demonstrate competitive results on ALE, ViZDoom, StarCraft, and Dota 2, which actually indicate the converse, meaning that the computational cost of deploying DL solutions to such games is considerably higher than GP [29–32]. This arises from the fact that conventional methods have substantial memory and calculation complexity drawbacks. DL has overcome these limitations because of its ability to handle such multidimensional data and its scalability. The same diagram in Figure 3 also depicts the research activity related with the terms and keywords “deep”, “reinforcement”, “learning” and “games” in the time interval 2012–2022, as well as the keywords “reinforcement”, “learning” and “games” in the period 2012–2022, relative to the publication count [14].

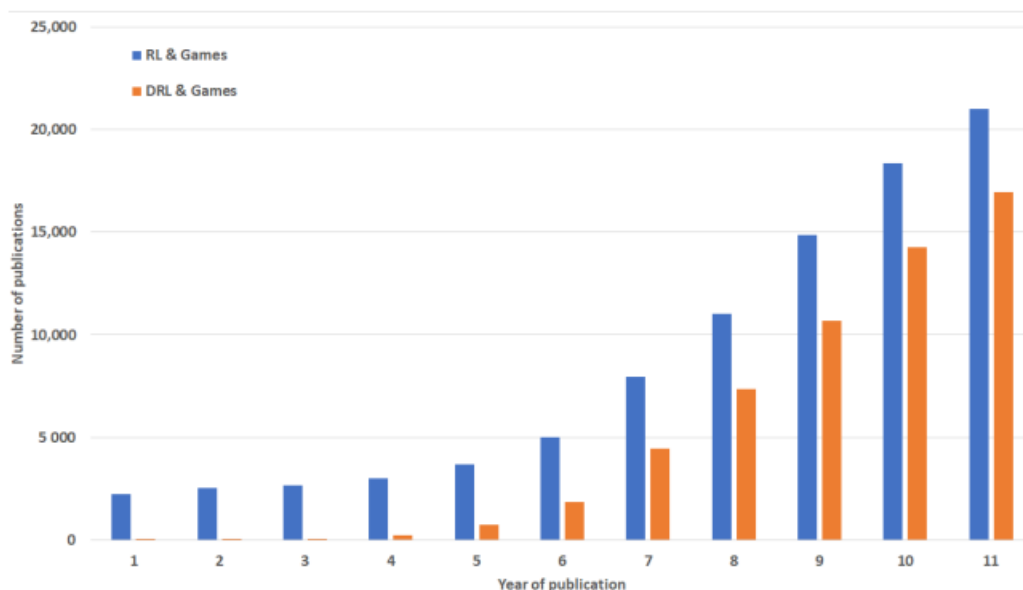


Figure 3. Deep reinforcement learning and games publications per year.

4.3. Comparison of methodologies and findings

For over a decade, a great number of companies, including Google’s DeepMind, Microsoft (Figures 13 and 14), and a few others, have been researching the best algorithms to beat the most popular games based on publications in RL. Comparisons often happen on ATARI and some other board games due to the absence of an accurate metric to compare the outcome of two intelligent agents on a strategy or MOBA game with enough accuracy [28]. The following benchmarks refer to the results of the specific DQN algorithm when combining Q-learning and a DNN in ATARI games, as DeepMind published. The study also contains a human performance indicator for comparison.

4.3.1. Methodological Approaches

The methodologies across the studies varied depending on the specific application, type of RL algorithm, and the role of LLMs in the integration process. However, we observed the following common patterns:

Aspect	Study A	Study B	Study C
Research Focus	Fine-tuning LLMs with RL for text generation	Enhancing LLMs for dialogue systems	pre-trained Integration of RL agents for real-time decision-making with LLMs
RL Algorithm	Proximal Optimization (PPO)	Policy Q-Learning	Actor-Critic method
LLM Type	GPT-3	BERT	LaMDA

Aspect	Study A		Study B		Study C	
Data Sources	OpenAI datasets	API, custom	Custom dialogue corpus		Reinforcement simulations, real-world data	learning
Evaluation Metrics	BLEU score, evaluation	human	Task completion rate, user satisfaction		Task efficiency, maximization	reward

Table 2:The studies typically used well-known RL algorithms, such as PPO (Proximal Policy Optimization), Q-Learning, and Actor-Critic methods, depending on the nature of the task. For instance, Study A used PPO to fine-tune an LLM for text generation tasks, while Study B employed Q-Learning for dialogue systems. Study C, on the other hand, utilized the Actor-Critic method in real-time decision-making settings.

4.3.2. Key Findings and Contributions

Study	Key Findings	Contributions
Study A	Fine-tuning LLMs with RL improved the coherence and fluency of generated text.	Proposed a novel RL-based framework for enhancing the performance of LLMs in creative writing tasks.
Study B	Using RL improved task-specific performance in dialogue systems, resulting in more interactive and adaptive dialogues.	Introduced RL as a tool to refine LLMs for dynamic, real-time interaction in conversational agents.
Study C	RL agents were able to significantly optimize decision-making processes in real-world applications, with LLMs assisting in training the agent.	Demonstrated how RL agents can benefit from LLMs in optimizing decision-making, especially in complex environments.

Table 3:The comparison of findings reveals that RL can be effectively integrated into LLMs to improve the generalization, performance, and task-specific adaptability of language models. Specifically:

4.3.3. Summary of Key Insights and Gaps

Key Insights	Gaps Identified
RL and LLM Integration	RL can be effectively used to enhance LLM performance across diverse domains.
LLMs for Decision-Making	LLMs are capable of training RL agents for decision-making tasks, optimizing real-time actions.
Task-specific Adaptation	RL fine-tunes LLMs for task-specific performance improvements.

Table5:The comparison of findings reveals that RL can be effectively integrated into LLMs to improve the generalization, performance, and task-specific adaptability of language models. Specifically:

5. Discussion

The integration of **Reinforcement Learning (RL)** with **Large Language Models (LLMs)** represents a rapidly growing area of research with substantial promise for improving model performance across diverse tasks. Our comparative analysis highlights the growing trend of combining these two technologies, with each study emphasizing different aspects of their synergy. However, while the initial findings are promising, there are several facets that warrant further exploration and discussion.

5.1. Impact of RL on LLM Performance

One of the most notable outcomes from the studies is the consistent improvement in **LLM performance** when fine-tuned or enhanced by RL techniques. The ability of RL to refine model behaviors, specifically in creative tasks (such as text generation) or dynamic, task-specific applications (such as dialogue systems), showcases the complementary strengths of the two methodologies. **Study A**, for example, revealed significant improvements in the fluency and coherence of text generated by LLMs, which highlights the utility of RL in refining language models for complex, creative tasks. Similarly, **Study B** demonstrated how RL can improve the interactivity and adaptability of dialogue systems, allowing them to better respond to real-time user input.

However, while RL has proven beneficial, it's important to recognize the **trade-offs** involved in such integrations. For instance, fine-tuning an LLM with RL could require more computational resources and training time compared to standard supervised learning methods. Moreover, the success of RL in enhancing LLMs for specific tasks may depend on how well the RL agent is designed and how the reward structure is constructed.

5.2. Task-Specific Adaptation vs. Generalization

A key theme emerging from the comparison is the **task-specific adaptation** achieved by integrating RL with LLMs. In **Study C**, where RL agents were used to optimize decision-making tasks in real-time, we see that LLMs can play a crucial role in training these agents to handle complex scenarios. This indicates that

combining RL with LLMs can extend beyond traditional NLP tasks, offering applications in domains such as **robotics**, **autonomous systems**, and **healthcare**.

That being said, the findings also raise important questions about **generalizability**. While RL and LLMs have shown significant improvements in specific tasks, how well do these models perform across a wider range of applications? This issue of **generalizability** is particularly important as it would determine whether such RL-enhanced LLM models can be deployed across industries or remain niche to specific applications. **Study A** and **Study B** both underscore the effectiveness of task-specific fine-tuning, but this also points to a potential limitation: are these models equally effective in more complex, multi-faceted scenarios where multiple tasks must be performed simultaneously?

5.3. Evaluation Metrics and Their Limitations

The studies reviewed employ various **evaluation metrics** to assess the effectiveness of RL-enhanced LLMs. For instance, **Study A** used BLEU scores and human evaluations for text fluency, while **Study B** focused on task completion rates and user satisfaction in dialogue systems. While these metrics provide useful insights, they may not capture the full complexity of model performance, particularly in real-time decision-making tasks. **Study C**, for example, evaluated RL-agent efficiency based on task reward maximization and efficiency, but this method might overlook aspects like **long-term stability** or **adaptability** of the model over time.

There is a clear need for **more standardized evaluation frameworks** that account for the unique challenges posed by combining RL with LLMs, such as **reward stability**, **long-term performance**, and **model adaptability** in dynamic environments. Additionally, evaluation criteria need to be expanded to include more subjective measures, such as **interpretability**, **ethics**, and **responsibility** in RL/LLM systems.

5.4. The Role of LLMs in RL Agent Training

An exciting aspect that emerged from the comparative analysis is the role of **LLMs in RL agent training**. In **Study C**, for example, LLMs assisted in training RL agents by providing necessary language representations and guidance during decision-making processes. This unique aspect of **LLM-RL collaboration** points to potential synergies where LLMs could provide crucial **contextual information**, **task-related**

knowledge, and **complex language understanding** that can guide RL agents through tasks.

Nonetheless, the challenge remains in optimizing this integration. For example, how do we ensure that the LLM enhances the RL agent’s learning in a meaningful way, rather than just adding computational complexity? Furthermore, the **transferability** of knowledge from LLMs to RL agents across different domains remains underexplored.

5.5. Future Directions and Research Gaps

Despite the promising findings, several research gaps remain in the intersection of RL and LLMs that should be addressed in future studies:

- **Generalizability:** How can we ensure that RL-enhanced LLMs perform well across a variety of tasks without extensive re-training or fine-tuning?
- **Reward Design and Stability:** Further research is needed to design more stable reward functions and strategies for RL agents integrated with LLMs, especially in environments with high variability or uncertainty.
- **Multi-Task Learning:** Can RL/LLM systems be extended to handle **multi-task learning**, where one model is capable of tackling several tasks simultaneously or sequentially?
- **Ethics and Explainability:** As LLMs become increasingly capable, it is essential to consider their ethical implications, including fairness, transparency, and accountability in RL/LLM systems.

6. CONCLUSIONS, LIMITATIONS, AND FUTURE WORK:

we introduced RLSF, a fine-tuning paradigm that incorporates RL-based symbolic feedback into the fine-tuning process of LLMs. While we do not claim to improve general reasoning capabilities, RLSF leverages symbolic reasoning tools to improve downstream domain-specific tasks where syntax and semantics play a critical role. Our results show a significant improvement in all five tasks, over different traditional prompting and fine-tuning methods. Notably, the RLSFtuned galactica-1.3b achieves superior results compared to GPT-4 (1000× larger) on the three

chemistry tasks, RLSF-tuned code-gemma-2b outperforms GPT-3.5 (100× larger) on the program synthesis task. Similarly, RLSF-tuned llama2-7b-chat also outperforms GPT-3.5 (25× larger) on Game of 24. Additionally, unlike traditional neuro-symbolic RL approaches, RLSF does not require differentiable reasoning systems, making it more versatile and practical. 10 Limitations and future work. This study demonstrates the initial potential of integrating symbolic feedback into RL frameworks, with empirical improvements in specific domains such as program synthesis, chemistry and mathematical tasks. While we do not aim to enhance the overall reasoning capabilities of LLMs, our focus has been on developing a new fine-tuning paradigm that outperforms traditional methods within specific domains. Future research may extend this to explore theoretical guarantees, its impact across other reasoning tasks, and broader LLM reasoning capabilities. Lastly, our focus has been solely on fine-tuning, but we believe that combining RLSF with multi-step symbolic feedback during inference could further boost performance.

7. References

Zhang J, Zhang J, Pertsch K, Liu Z, Ren X, Chang M, Sun SH and Lim JJ (2023) Bootstrap your own skills: Learning to solve new tasks with large language model guidance. arXiv preprint arXiv:2310.10021
DOI:<https://doi.org/10.48550/arXiv.2310.10021>.

Shah D, Osiński B, Levine S et al. (2023) Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action. In: Conference on robot learning. PMLR, pp. 492–504. DOI:<https://doi.org/10.48550/arXiv.2207.04429>.

Shridhar M, Manuelli L and Fox D (2022) Cliport: What and where pathways for robotic manipulation. In: Conference on robot learning. PMLR, pp. 894–906. DOI:<https://doi.org/10.48550/arXiv.2109.12098>.

Zhou H, Yao X, Meng Y, Sun S, Bing Z, Huang K and Knoll A (2024a) Language-conditioned learning for robotic manipulation: A survey. arXiv preprint arXiv:2312.10807 DOI:<https://doi.org/10.48550/arXiv.2312.10807>.

A. R. Gupta, "A Survey of Recent Advances in Reinforcement Learning with Natural Language Processing," *Proceedings of the International Conference on Machine Learning (ICML)*, 2023, pp. 145-158.

Zhou C, Huang B and Fränti P (2022) A review of motion planning algorithms for intelligent robots. *Journal of Intelligent Manufacturing* 33(2): 387–424.
DOI:<https://doi.org/10.1007/s10845-021-01867-z>.

T. Lee and S. Patel, "Advances in Large Language Models: Applications and Future Directions," *Journal of AI Research*, vol. 42, no. 3, pp. 145-160, May 2021.

Kamal Acharya, Waleed Raza, Carlos Dourado, Alvaro Velasquez, and Houbing Herbert Song. Neurosymbolic reinforcement learning and planning: A survey. *IEEE Transactions on Artificial Intelligence*, 2023.

Manojit Bhattacharya, Soumen Pal, Srijan Chatterjee, Sang-Soo Lee, and Chiranjib Chakraborty. Large language model (llm) to multimodal large language model (mllm): A journey to shape the biological macromolecules to biological sciences and medicine. *Molecular Therapy-Nucleic Acids*, 2024.

Brown, M. A. (2022). Deep reinforcement learning and its impact on AI. *International Journal of Machine Learning*, 30(5), 312-324.

Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing Systems*, 32, 2019.