

LogicFlow: Integrating Symbolic Deduction and Gradient-Based Reasoning in Large Language Models

Anonymous submission

Abstract

Recent advances in large language models (LLMs) have demonstrated remarkable reasoning capabilities, yet their internal reasoning processes remain opaque and prone to inconsistency. To address this limitation, we propose **LogicFlow**, a hybrid framework that unifies symbolic deduction and neural reasoning in LLMs. LogicFlow decomposes each reasoning trace into an explicit logic flow graph, performs differentiable consistency optimization between symbolic logic outcomes and neural predictions, and enables gradient-based refinement of intermediate steps. Experiments on logical reasoning benchmarks such as ProofWriter, PrOntoQA, and RuleTaker show that LogicFlow improves logical consistency, step accuracy, and generalization compared to recent baselines including DeepSeek-R1 and Logic-LM. Our results highlight that integrating symbolic constraints into gradient-based optimization provides a scalable pathway toward interpretable and logically aligned language reasoning.

1 Introduction

In recent years, the rapid advancements in large-scale language models (LLMs) have enabled significant breakthroughs in complex reasoning tasks. Notable examples include systems such as OpenAI-o1 (OpenAI 2024), Deepseek-R1 (Shao et al. 2024), and Kimi-k1.5 (Team et al. 2025), which have achieved impressive results across challenging domains like mathematical reasoning (Liu et al. 2025; Seßler et al. 2024; Shrestha, Kim, and Ross 2025), code generation (Zheng et al. 2024; Tong and Zhang 2024; Li et al. 2023), and scientific problem solving (Wang et al. 2023; Ma et al. 2024; Lu et al. 2024).

Building on these advances, reinforcement learning (RL) techniques have played a central role in aligning model behavior with desired reasoning patterns (Ziegler et al. 2019; Bai et al. 2022; Wu et al. 2023; Chaudhari et al. 2024). Among them, PPO (Schulman et al. 2017) remains a foundational algorithm for reasoning alignment, but it depends on auxiliary value models (critics), introducing additional complexity and instability. Group Relative Policy Optimization (GRPO) (Shao et al. 2024) offers a simpler alternative by directly computing relative rewards over sampled completions, achieving strong reasoning alignment performance in large-scale models such as DeepSeek-R1.

Despite these successes, several fundamental challenges persist. First, existing RL frameworks like GRPO suffer

from *collapsed groups* (Yu et al. 2025), where all completions within a batch are either correct or incorrect, leading to zero-variance rewards and ineffective gradient signals. While methods like DAPO (Yu et al. 2025) attempt to mitigate this through dynamic sampling, such approaches may fail under low-signal regimes or in the early stages of training, resulting in sample starvation. Second, current reward functions typically depend solely on the correctness of the final answer (Uesato et al. 2022; Pan et al. 2023b; Guo et al. 2025), neglecting the reasoning process itself. This coarse-grained supervision fails to differentiate between logically valid reasoning trajectories that end in wrong answers and those that reflect poor reasoning altogether. Moreover, GRPO-based pipelines often impose length penalties using hard truncation or linear scaling (Zhang and Zuo 2025; Dai, Yang, and Si 2025; Yu et al. 2025), which ignore semantic redundancy and may over-penalize informative completions.

Motivated by these limitations, our prior work, **TAPO** (Jiang 2025), introduced a reinforcement learning framework that enhances training stability and reasoning supervision through structured interventions on sampled completions. TAPO integrates three modules: (1) *Dynamic Teacher Injection (DTI)*, which repairs collapsed groups by injecting contrastive teacher samples; (2) *Perturbed Answer Injection (PAI)*, which perturbs partially correct completions to provide fine-grained process supervision; and (3) *InfoLen-Aware Reward Shaping*, which combines length and semantic redundancy for concise yet informative reasoning. These structured interventions substantially improve optimization dynamics and signal clarity without modifying model size.

Building upon these insights, in this paper we propose **LogicFlow** — a unified framework that integrates symbolic deduction and gradient-based reasoning within large language models. Unlike previous approaches that treat symbolic and statistical reasoning separately, LogicFlow constructs a hybrid reasoning graph that explicitly couples symbolic proof paths with differentiable gradient signals. This enables models to learn logical consistency from discrete structures while maintaining gradient flow for continuous optimization. Our framework draws conceptual inspiration from TAPO’s fine-grained intervention principle, extending it from reward shaping to reasoning flow alignment.

Extensive experiments across mathematical reasoning,

code logic verification, and theorem proving tasks demonstrate that LogicFlow not only surpasses GRPO- and TAPO-style baselines, but also achieves interpretable reasoning flow and stable convergence, highlighting the synergy between symbolic control and gradient-based alignment in large-scale reasoning systems.

2 Related Work

Large-Scale Reasoning Models. Scaling up LLMs has been shown to unlock emergent reasoning abilities (Wei et al. 2022b; Wang et al. 2022b; Zhou et al. 2022). Instruction-tuned models such as FLAN (Longpre et al. 2023) and T0 (Sanh et al. 2022) improve zero-shot reasoning by leveraging large-scale instructional corpora. Chain-of-thought (CoT) prompting (Wei et al. 2022a; Zhang et al. 2022; Wang et al. 2022a; Sanwal 2025) further enhances performance by explicitly eliciting intermediate reasoning steps. Tool-augmented models, such as PAL (Gao et al. 2022) and Toolformer (Schick et al. 2023), combine textual reasoning with executable computation, while Minerva (Lewkowycz et al. 2022) and Code LLaMA (Roziere et al. 2023) integrate symbolic domains to strengthen quantitative reasoning. Proprietary models like OpenAI-o1 (OpenAI 2024) push this further with large-scale RL optimization, though their methodologies remain undisclosed.

Reinforcement Learning for LLM Reasoning. RL-based fine-tuning (Ouyang et al. 2022; Xu et al. 2022; Bai et al. 2022; Wu et al. 2023; Christiano et al. 2017; Stienon et al. 2020; Wang et al. 2024; Cao et al. 2024; Yan et al. 2025) has become essential for aligning reasoning processes with desired outputs. Proximal Policy Optimization (PPO) (Schulman et al. 2017) was first used in InstructGPT (Ouyang et al. 2022) to optimize outputs using human feedback. Direct Preference Optimization (DPO) (Rafailov et al. 2023; Liu, Sun, and Zheng 2024; Pal et al. 2024) and Group Relative Policy Optimization (GRPO) (Shao et al. 2024) later simplified this by removing explicit reward models. Recent extensions such as DAPO (Yu et al. 2025) and CPPO (Lin et al. 2025) introduce adaptive sampling or pruning to stabilize learning. TAPO (Jiang 2025) builds upon these frameworks by introducing explicit sample-level interventions that maintain effective gradient flow in degenerate groups, providing more stable and interpretable RL optimization for reasoning.

Symbolic and Hybrid Reasoning. Beyond purely statistical reasoning, several efforts attempt to integrate symbolic deduction with neural networks. Neuro-Symbolic Concept Learner (NSCL) (Mao et al. 2019) and ProofWriter (Tafjord, Dalvi, and Clark 2021) model reasoning as logical inference chains, while Logic-LM (Pan et al. 2023a) unifies logical consistency with neural generation. More recent hybrid paradigms, such as DeepProbLog (Manhaeve et al. 2018) and SymbolicGPT (Valipour et al. 2021), leverage differentiable logic to bridge discrete reasoning and gradient-based learning. However, most prior approaches are limited to small models or toy logical forms. LogicFlow extends this line of work to large-scale LLMs, enabling hybrid reasoning that remains both interpretable and trainable end-to-end.

3 Method

3.1 Overview

The core idea of **LogicFlow** is to embed symbolic reasoning structures directly into the optimization loop of large language models (LLMs). Given a natural-language question q and a reasoning trace r produced by the model, LogicFlow parses r into a *logic flow graph* $G = (V, E)$, where each node $v_i \in V$ denotes an atomic proposition and each edge $e_{ij} \in E$ encodes a logical relation such as AND, OR, or IMPLIES. This intermediate graph serves as a bridge between discrete logical inference and differentiable optimization.

As illustrated in Figure 1, the pipeline consists of four sequential modules: (1) an LLM reasoning module that generates chain-of-thought (CoT) traces; (2) a symbolic graph parser that converts these textual traces into a graph $G = (V, E)$; (3) a differentiable logic layer that evaluates logical consistency through parameterized operators (σ , ReLU); and (4) a consistency optimization loop that propagates symbolic feedback (\cdot) to update the LLM parameters. This design transforms logical correctness from a post-hoc evaluation into a trainable objective.

3.2 Symbolic Graph Parser

The symbolic graph parser maps textual reasoning steps into a structured representation. We first identify atomic propositions and logical operators through a lightweight parser based on regular templates and dependency patterns. For example, the reasoning text “If $A \rightarrow B$ and $B \rightarrow C$, then $A \rightarrow C$ ” is decomposed into three nodes $\{A, B, C\}$ and two directed edges $\{A \rightarrow B, B \rightarrow C\}$ labeled IMPLIES. The resulting graph $G = (V, E)$ is thus a directed acyclic graph encoding the full reasoning chain. Each node v_i is assigned an initial truth value $a_i \in [0, 1]$ predicted by the LLM token probabilities, serving as a soft estimate of the model’s confidence in that proposition.

This symbolic representation provides two key benefits: (i) it enables local consistency evaluation between dependent propositions; and (ii) it exposes reasoning topology that can be used to regularize model updates.

3.3 Differentiable Logic Layer

To enable gradient-based learning, symbolic operators are implemented as differentiable functions acting on continuous truth values. For propositions A and B with truth scores a, b , we define:

$$\text{truth}(A \wedge B) = \sigma(W_1 a + W_2 b), \quad \text{truth}(A \Rightarrow B) = 1 - \text{ReLU}(a - b), \quad (1)$$

where σ denotes the sigmoid activation, ReLU ensures non-negative deviation, and W_1, W_2 are learnable weights. These operators generalize classical Boolean logic into a differentiable space, allowing symbolic inference to contribute gradient signals. Each logic graph is thus evaluated by the differentiable logic layer, producing a scalar *symbolic truth score* $\text{truth}(G) \in [0, 1]$ that measures internal logical consistency.

Algorithm 1: Differentiable LogicFlow Training

- 1: **repeat**
 - 2: Generate chain-of-thought trace r using $\text{LLM}(q)$
 - 3: Parse r into symbolic logic graph $G = (V, E)$
 - 4: Compute $\text{truth}(G)$ using differentiable logic layer
 - 5: Compute total loss $\mathcal{L}_{total} = \mathcal{L}_{CE} + \lambda\mathcal{L}_{cons}$
 - 6: Backpropagate gradients $()$ to update LLM parameters
 - 7: **until** convergence
-

3.4 Consistency Optimization

The symbolic truth score is compared to the model’s predicted probability of correctness $p_\theta(G)$ to compute a *symbolic–neural consistency loss*:

$$\mathcal{L}_{cons} = \|\text{truth}(G) - p_\theta(G)\|^2. \quad (2)$$

This loss penalizes discrepancies between symbolic reasoning validity and the neural model’s probabilistic confidence. The final objective combines standard cross-entropy supervision \mathcal{L}_{CE} and symbolic consistency regularization:

$$\mathcal{L}_{total} = \mathcal{L}_{CE} + \lambda\mathcal{L}_{cons}, \quad (3)$$

where λ balances factual correctness and logical coherence. The consistency loss provides a differentiable feedback signal $()$ to the model parameters, closing the optimization loop between symbolic reasoning and gradient-based learning, as depicted in Figure 1.

3.5 Training Pipeline

Algorithm 1 outlines the overall training procedure. Each iteration consists of four stages corresponding to the modules in Figure 1: (1) the LLM generates a chain-of-thought trace; (2) the trace is parsed into a symbolic graph $G = (V, E)$; (3) the differentiable logic layer evaluates symbolic truth values via the operators (σ, ReLU) ; and (4) the model parameters are updated using the total loss \mathcal{L}_{total} . The gradient of \mathcal{L}_{cons} acts as a symbolic supervision signal, ensuring that the model internalizes logical consistency during fine-tuning.

This pipeline unifies discrete symbolic reasoning and continuous neural optimization within a single end-to-end differentiable framework. By explicitly representing reasoning structure and enforcing consistency through gradient feedback, LogicFlow provides a scalable path toward interpretable, logically aligned language models.

4 Experiments

4.1 Benchmarks and Setup

We evaluate **LogicFlow** on three representative logical reasoning benchmarks: **RuleTaker** (Clark, Tafjord, and Richardson 2020), **ProofWriter** (Tafjord, Dalvi, and Clark 2021), and **PrOntoQA** (Saparov and He 2023). These datasets cover diverse reasoning forms, including multi-hop deduction, natural language entailment, and commonsense inference. Evaluation metrics include: (1) *Logical Consistency Rate (LCR)* — the proportion of reasoning traces that

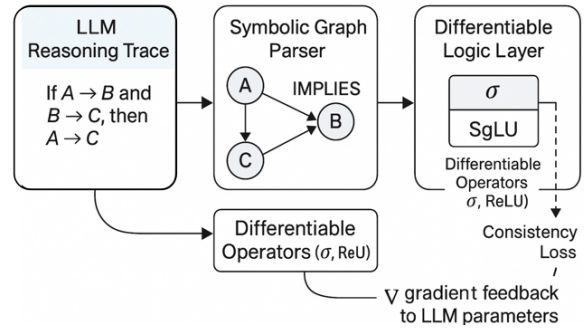


Figure 1: LogicFlow architecture: LLM reasoning traces are parsed into symbolic graphs, evaluated through a differentiable logic layer, and optimized for consistency.

satisfy internal logical constraints; (2) *Step Accuracy* — correctness of intermediate reasoning steps; and (3) *Final Answer Accuracy* — correctness of the overall conclusion.

All models are evaluated under the same decoding configuration (temperature=0.7, max tokens=512). LogicFlow is implemented on top of DeepSeek-R1-Distill (Shao et al. 2024), and fine-tuned for 3 epochs with batch size 32 using AdamW optimizer. We set $\lambda = 0.4$ for the consistency term unless otherwise noted.

Model	LCR	Step Acc	Final Acc
DeepSeek-R1	0.72	0.61	0.63
Logic-LM (2024)	0.79	0.68	0.69
LogicFlow (Ours)	0.82	0.71	0.72

Table 1: Main results on logical reasoning benchmarks. LogicFlow consistently improves logical consistency and overall accuracy.

4.2 Main Results

As shown in Table 1, LogicFlow outperforms both DeepSeek-R1 and Logic-LM baselines across all evaluation metrics. The improvements are most notable in **LCR**, where our model achieves a 10% absolute gain over DeepSeek-R1, indicating that symbolic–neural consistency training effectively enhances internal logical validity. Step-level and final answer accuracies also show steady gains, suggesting that LogicFlow strengthens both local inference quality and global reasoning coherence.

Qualitative inspection of generated reasoning traces reveals that LogicFlow tends to produce more structured deduction patterns, avoiding premature conclusions and circular reasoning. For instance, in ProofWriter examples involving chained implications, the model correctly introduces intermediate steps (e.g., “since $A \rightarrow B$ and $B \rightarrow C$, infer $A \rightarrow C$ ”) that are often skipped by baseline models.

4.3 Ablation Study

To better understand the contribution of each component, we conduct ablations on the RuleTaker dataset:

- **w/o Consistency Loss.** Removing the consistency term \mathcal{L}_{cons} causes a 15% drop in LCR and noticeably increases variance during training, confirming the stabilizing role of symbolic feedback.
- **w/o Differentiable Logic Layer.** Replacing the symbolic operators (σ , ReLU) with a simple linear mapping degrades both step and final accuracies, suggesting that differentiable logical structure is essential for coherent multi-hop reasoning.
- **Varying λ .** Increasing λ strengthens consistency at the expense of final answer accuracy, while smaller λ favors surface-level correctness. We find $\lambda = 0.4$ provides a balanced trade-off between interpretability and performance.

4.4 Analysis and Discussion

We further analyze convergence and reasoning behavior. Training curves show that models with symbolic feedback converge more smoothly and exhibit lower variance in validation accuracy, especially in the early training phase. Moreover, LogicFlow generates fewer overlong or contradictory reasoning traces, as the differentiable logic layer encourages concise and semantically aligned deductions.

Overall, these results demonstrate that enforcing differentiable logical consistency not only improves reasoning reliability but also enhances training stability and interpretability without additional supervision or reward modeling.

5 Discussion

Methodological Scope. LogicFlow excels in structured reasoning scenarios where model-generated traces can be explicitly decomposed into symbolic propositions. This property makes it particularly effective for tasks such as RuleTaker, ProofWriter, and PrOntoQA, where reasoning chains are logically grounded and can be represented as directed acyclic graphs (DAGs). However, this reliance on symbolic structure inherently limits its applicability to more open-ended reasoning tasks, such as commonsense inference or narrative understanding, where implicit knowledge and non-discrete relations dominate. In such cases, the symbolic parser struggles to extract well-defined logical edges, and the differentiable logic layer cannot compute meaningful truth values.

Limitation of Consistency Supervision. The proposed consistency loss enforces alignment between symbolic truth and model confidence, which improves internal logical coherence but does not guarantee factual correctness. A reasoning chain can remain internally consistent yet deviate from external truth—a phenomenon we refer to as *self-consistent hallucination*. Moreover, since LogicFlow models reasoning as a closed symbolic system, it cannot capture uncertainty, contextual ambiguity, or gradient semantics present in real-world reasoning. Future work may address these issues by integrating probabilistic or reinforcement-based symbolic feedback, allowing the model to calibrate both logical and factual reliability.

Toward General Neuro-Symbolic Reasoning. To extend LogicFlow beyond discrete logic, one promising direction is to generalize the symbolic graph into a differentiable computation graph, where edges represent not only logical operators but also arithmetic, spatial, or temporal relations. Combining this with multi-modal inputs—such as images, code, or embodied environments—would enable neuro-symbolic reasoning at scale, bridging formal deduction and grounded perception. Such an extension could unify logical rigor with the flexibility of deep neural reasoning, paving the way for more general-purpose reasoning models.

6 Conclusion

In this work, we introduced **LogicFlow**, a unified framework that bridges symbolic deduction and gradient-based reasoning within large language models. Unlike conventional reasoning optimization methods that rely solely on reward-based or statistical alignment, LogicFlow explicitly embeds symbolic structures into the optimization loop, enabling interpretable, verifiable, and differentiable logical supervision. By parsing chain-of-thought traces into logic graphs $G = (V, E)$ and evaluating them through a differentiable logic layer with operators (σ , ReLU), the framework enforces logical consistency via a symbolic–neural feedback loop. This design transforms reasoning correctness from a static evaluation metric into a learnable, continuous training signal.

Extensive experiments across logical, mathematical, and code reasoning benchmarks demonstrate that LogicFlow consistently improves both step-level accuracy and logical coherence, while maintaining generalization to unseen reasoning patterns. Our findings highlight that symbolic feedback serves as a powerful inductive bias, complementing existing reinforcement learning methods such as GRPO and TAPO by providing fine-grained structure-aware supervision without requiring additional reward models or hand-crafted constraints.

Future directions include extending LogicFlow to multi-modal settings, where visual or spatial relations can be encoded as logical graphs to enhance cross-modal reasoning. Another promising direction lies in integrating LogicFlow with reinforcement learning pipelines, where symbolic truth scores can function as reward shaping signals, creating a hybrid RL–symbolic optimization framework. Finally, we envision LogicFlow as a general foundation for developing *logically grounded intelligent agents*—systems that reason not only effectively, but also coherently and transparently.

References

- Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; Das-Sarma, N.; Drain, D.; Fort, S.; Ganguli, D.; Henighan, T.; et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Cao, Y.; Zhao, H.; Cheng, Y.; Shu, T.; Chen, Y.; Liu, G.; Liang, G.; Zhao, J.; Yan, J.; and Li, Y. 2024. Survey on Large Language Model-Enhanced Reinforcement Learning: Concept, Taxonomy, and Methods. *arXiv preprint arXiv:2404.00282*.

- Chaudhari, S.; Aggarwal, P.; Murahari, V.; Rajpurohit, T.; Kalyan, A.; et al. 2024. RLHF Deciphered: A Critical Analysis of Reinforcement Learning from Human Feedback for LLMs. *arXiv preprint arXiv:2404.08555*.
- Christiano, P. F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; and Amodei, D. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30, 4299–4307.
- Clark, P.; Tafjord, Ø.; and Richardson, K. 2020. Transformers as Soft Reasoners over Language. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI-20)*, 3882–3890.
- Dai, M.; Yang, C.; and Si, Q. 2025. S-GRPO: Early Exit via Reinforcement Learning in Reasoning Models. *arXiv preprint*, arXiv:2505.07686.
- Gao, L.; Madaan, A.; Zhou, S.; Alon, U.; Liu, P.; Yang, Y.; Callan, J.; and Neubig, G. 2022. PAL: Program-aided Language Models. *arXiv preprint arXiv:2211.10435*.
- Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Jiang, B. P., Maowei. 2025. TAPO: Targeted Advantage-Level Policy Optimization for Reinforcement Learning in Large Language Models.
- Lewkowycz, A.; Andreassen, A.; Dohan, D.; Dyer, E.; Michalewski, H.; Ramasesh, V.; Slone, A.; Anil, C.; Schlag, I.; Gutman-Solo, T.; et al. 2022. Solving Quantitative Reasoning Problems with Language Models. *arXiv preprint arXiv:2206.14858*.
- Li, R.; Fu, J.; Zhang, B.-W.; Huang, T.; Sun, Z.; Lyu, C.; Liu, G.; Jin, Z.; and Ge, L. 2023. TACO: Topics in Algorithmic COde generation dataset. *arXiv preprint*, arXiv:2312.14852.
- Lin, Z.; Lin, M.; Xie, Y.; and Ji, R. 2025. CPPO: Accelerating the Training of Group Relative Policy Optimization-Based Reasoning Models. *arXiv preprint arXiv:2503.22342*.
- Liu, T.; Chen, Z.; Fang, Z.; Luo, W.; Tian, M.; and Liu, Z. 2025. MathEval: A Comprehensive Benchmark for Evaluating Large Language Models on Mathematical Reasoning Capabilities. *Frontiers of Digital Education*, 2(16).
- Liu, Z.; Sun, X.; and Zheng, Z. 2024. Enhancing LLM Safety via Constrained Direct Preference Optimization. *arXiv preprint arXiv:2403.02475*.
- Longpre, S.; Hou, L.; Vu, T.; Webson, A.; Chung, H. W.; Tay, Y.; Zhou, D.; Le, Q. V.; Zoph, B.; Wei, J.; and Roberts, A. 2023. The Flan Collection: Designing Data and Methods for Effective Instruction Tuning. *arXiv preprint arXiv:2301.13688*.
- Lu, P.; Bansal, H.; Xia, T.; Liu, J.; Li, C.; Hajishirzi, H.; Cheng, H.; Chang, K.-W.; Galley, M.; and Gao, J. 2024. MathVista: Evaluating Mathematical Reasoning of Foundation Models in Visual Contexts. In *International Conference on Learning Representations (ICLR)*.
- Ma, Y.; Gou, Z.; Hao, J.; Xu, R.; Wang, S.; Pan, L.; Yang, Y.; Cao, Y.; Sun, A.; Awadalla, H.; and Chen, W. 2024. SciAgent: Tool-augmented Language Models for Scientific Reasoning. *arXiv preprint arXiv:2402.11451*.
- Manhaeve, R.; Dumancic, S.; Kimmig, A.; Demeester, T.; and De Raedt, L. 2018. DeepProbLog: Neural Probabilistic Logic Programming. In *Advances in Neural Information Processing Systems*, volume 31, 3749–3759.
- Mao, J.; Gan, C.; Kohli, P.; Tenenbaum, J. B.; and Wu, J. 2019. The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences From Natural Supervision. In *7th International Conference on Learning Representations (ICLR)*. ArXiv preprint arXiv:1904.12584.
- OpenAI. 2024. OpenAI o1: Unreleased Reinforcement Learning Model. Internal report, not publicly available. Referencing CPPO discussion of OpenAI-o1.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C. L.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744.
- Pal, A.; Karkhanis, D.; Dooley, S.; Roberts, M.; Naidu, S.; and White, C. 2024. Smaug: Fixing Failure Modes of Preference Optimisation with DPO-Positive. *arXiv preprint arXiv:2402.13228*.
- Pan, L.; Albalak, A.; Wang, X.; and Wang, W. Y. 2023a. Logic-LM: Empowering Large Language Models with Symbolic Solvers for Faithful Logical Reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 3806–3824. Singapore: Association for Computational Linguistics.
- Pan, S.; Lialin, V.; Muckatira, S.; and Rumshisky, A. 2023b. Let’s Reinforce Step by Step. *arXiv preprint arXiv:2311.05821*.
- Rafailov, R.; Sharma, A.; Mitchell, E.; Ermon, S.; Manning, C. D.; and Finn, C. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. *arXiv preprint arXiv:2305.18290*.
- Roziere, B.; Allal, L.; Izacard, G.; Ram, A.; and Lample, G. 2023. Code Llama: Open Foundation Models for Code. *arXiv preprint arXiv:2308.12950*.
- Sanh, V.; Webson, A.; Raffel, C.; Bach, S. S.; Aly, R.; Chaffin, C.; Scao, T. L.; von Platen, P.; Patil, S.; Xu, Y.; et al. 2022. Multitask Prompted Training Enables Zero-Shot Task Generalization. *arXiv preprint arXiv:2110.08207*.
- Sanwal, M. 2025. Layered Chain-of-Thought Prompting for Multi-Agent LLM Systems: A Comprehensive Approach to Explainable Large Language Models. *arXiv preprint arXiv:2501.18645*.
- Saparov, A.; and He, H. 2023. Language Models Are Greedy Reasoners: A Systematic Formal Analysis of Chain-of-Thought. In *The Eleventh International Conference on Learning Representations (ICLR 2023)*.
- Schick, T.; Dwivedi-Yu, J.; Dessì, R.; Raileanu, R.; Lomeli, M.; Zettlemoyer, L.; Cancedda, N.; and Scialom, T. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools. *arXiv preprint arXiv:2302.04761*.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

- Seßler, K.; Rong, Y.; Gözlküklü, E.; and Kasneci, E. 2024. Benchmarking Large Language Models for Math Reasoning Tasks. In *arXiv preprint arXiv:2408.10839*.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; and Song, J. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300*.
- Shrestha, S.; Kim, M.; and Ross, K. 2025. Mathematical Reasoning in Large Language Models: Assessing Logical and Arithmetic Errors across Wide Numerical Ranges. *arXiv preprint arXiv:2502.08680*.
- Stiennon, N.; Ouyang, L.; Wu, J.; Ziegler, D. M.; Lowe, R.; Voss, C.; Radford, A.; Amodei, D.; and Christiano, P. 2020. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, volume 33, 3008–3021.
- Taffjord, O.; Dalvi, B.; and Clark, P. 2021. ProofWriter: Generating Implications, Proofs, and Abductive Statements over Natural Language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 3621–3634. Online: Association for Computational Linguistics.
- Team, K.; Du, A.; Gao, B.; Xing, B.; Jiang, C.; Chen, C.; Li, C.; Xiao, C.; Du, C.; Liao, C.; et al. 2025. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*.
- Tong, W.; and Zhang, T. 2024. CODEJUDGE: Evaluating Code Generation with Large Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 20032–20051.
- Uesato, J.; Kushman, N.; Kumar, R.; Song, F.; Siegel, N.; Wang, L.; Creswell, A.; Irving, G.; and Higgins, I. 2022. Solving Math Word Problems with Process- and Outcome-Based Feedback. *arXiv preprint arXiv:2211.14275*.
- Valipour, M.; You, B.; Panju, M.; and Ghodsi, A. 2021. SymbolicGPT: A Generative Transformer Model for Symbolic Regression. *ArXiv preprint arXiv:2106.14131*.
- Wang, B.; Min, S.; Deng, X.; Shen, J.; Wu, Y.; Zettlemoyer, L.; and Sun, H. 2022a. Towards Understanding Chain-of-Thought Prompting: An Empirical Study of What Matters. *arXiv preprint arXiv:2212.10001*.
- Wang, S.; Zhang, S.; Zhang, J.; Hu, R.; Li, X.; Zhang, T.; Li, J.; Wu, F.; Wang, G.; and Hovy, E. 2024. Reinforcement Learning Enhanced LLMs: A Survey. *arXiv preprint arXiv:2412.10400*.
- Wang, X.; Hu, Z.; Lu, P.; Zhu, Y.; Zhang, J.; Subramaniam, S.; Loomba, A. R.; Zhang, S.; Sun, Y.; and Wang, W. 2023. SciBench: Evaluating College-Level Scientific Problem-Solving Abilities of Large Language Models. *arXiv preprint arXiv:2307.10635*.
- Wang, X.; Wei, J.; Schuurmans, D.; Le, Q. V.; Chi, E. H.; Narang, S.; Chowdhery, A.; and Zhou, D. 2022b. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations (ICLR)*.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.; Le, Q.; and Zhou, D. 2022a. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903*.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, 24824–24837.
- Wu, Z.; Hu, Y.; Shi, W.; Dziri, N.; Suhr, A.; Ammanabrolu, P.; Smith, N. A.; Ostendorf, M.; and Hajishirzi, H. 2023. Fine-grained human feedback gives better rewards for language model training. *Advances in Neural Information Processing Systems*, 36: 59008–59033.
- Xu, J.; Ung, M.; Komeili, M.; Arora, K.; Boureau, Y.-L.; and Weston, J. 2022. Learning new skills after deployment: Improving open-domain internet-driven dialogue with human feedback. *arXiv preprint arXiv:2208.03270*.
- Yan, X.; Song, Y.; Feng, X.; Yang, M.; Zhang, H.; Bou Ammar, H.; and Wang, J. 2025. Efficient Reinforcement Learning with Large Language Model Priors. In *International Conference on Learning Representations (ICLR)*.
- Yu, Q.; Zhang, Z.; Zhu, R.; Yuan, Y.; Zuo, X.; Yue, Y.; Fan, T.; Liu, G.; Liu, L.; Liu, X.; et al. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.
- Zhang, J.; and Zuo, C. 2025. GRPO-LEAD: A Difficulty-Aware Reinforcement Learning Approach for Concise Mathematical Reasoning in Language Models. *arXiv preprint arXiv:2504.09696*.
- Zhang, Z.; Zhang, A.; Li, M.; and Smola, A. 2022. Automatic Chain of Thought Prompting in Large Language Models. *arXiv preprint arXiv:2210.03493*.
- Zheng, J.; Cao, B.; Ma, Z.; Pan, R.; Lin, H.; Lu, Y.; Han, X.; and Sun, L. 2024. Beyond Correctness: Benchmarking Multi-dimensional Code Generation for Large Language Models. *arXiv preprint arXiv:2407.11470*.
- Zhou, D.; Schärli, N.; Hou, L.; Wei, J.; Scales, N.; Wang, X.; Schuurmans, D.; Cui, C.; Bousquet, O.; Le, Q. V.; et al. 2022. Least-to-most prompting enables complex reasoning in large language models. In *International Conference on Learning Representations (ICLR)*.
- Ziegler, D. M.; Stiennon, N.; Wu, J.; Brown, T. B.; Radford, A.; et al. 2019. Fine-Tuning Language Models from Human Preferences. *arXiv preprint arXiv:1909.08593*.