

Supervised Fine-tuning versus Reinforcement Learning: A Comprehensive Survey on Large Language Model Post-training

Anonymous ACL submission

Abstract

Pre-trained Large Language Model (LLM) exhibits broad capabilities, yet, for specific tasks or domains their attainment of higher accuracy and more reliable reasoning generally depends on post-training through Supervised Fine-Tuning (SFT) or Reinforcement Learning (RL). Although often treated as distinct methodologies, recent theoretical and empirical developments demonstrate that SFT and RL are closely connected. This survey presents a comprehensive and unified perspective on LLM post-training with SFT and RL. We first provide an in-depth overview of both techniques, examining their objectives, algorithmic structures, and data requirements. We then systematically analyze their interplay, highlighting frameworks that integrate SFT and RL, hybrid training pipelines, and methods that leverage their complementary strengths. Drawing on a representative set of recent application studies from 2023 to 2025, we identify emerging trends, characterize the rapid shift toward hybrid post-training paradigms, and distill key takeaways that clarify when and why each method is most effective. By synthesizing theoretical insights, practical methodologies, and empirical evidence, this survey establishes a coherent understanding of SFT and RL within a unified framework and outlines promising directions for future research in scalable, efficient, and generalizable LLM post-training.

1 Introduction

Pre-trained LLMs have demonstrated remarkable capabilities across a wide range of tasks, from fact-based question answering (Joshi et al., 2017) to code generation (Jimenez et al., 2024). Despite being trained on corpora containing billions to trillions of tokens, LLMs often require task-specific post-training adaptation to improve accuracy, mitigate erroneous outputs, and handle new tasks. For

instance, fine-tuning can enhance multi-step reasoning by enabling the generation of progressively longer reasoning chains, ultimately leading to more accurate final answers, particularly in tasks that require complex reasoning ability (Guo et al., 2025). It can also enhance practical interaction skills, including performing tasks in household (Song et al., 2024a) or device-control environments (Rawles et al., 2025). Such capabilities rarely emerge from pre-training alone, as the data seldom contain the task-specific patterns or feedback necessary for accurate reasoning or complex interactions, highlighting the importance of post-training.

Current post-training methodologies for LLMs primarily fall into two paradigms: SFT and RL. The objective of SFT (Zhou et al., 2023; Li et al., 2024c; Pang et al., 2025) is to maximize the likelihood of tokens conditioned on context, whereas RL (Ouyang et al., 2022; Guo et al., 2025; Yang et al., 2025a) optimizes a reward signal derived from human or automated preference feedback. Despite their different objectives, in recent years, research efforts have increasingly focused on bridging these two approaches (Wu et al., 2025; Fu et al., 2025b), exploring how their combination (Wu et al., 2025; Qin and Springenberg, 2025; Yan et al., 2025a; Liu et al., 2025a) can enhance performance beyond what either approach achieves in isolation.

This combination is especially pronounced in complex tasks that demand both accuracy and generalization, such as multi-step reasoning. SFT alone can teach models to generate basic Chain-of-Thought (CoT), but may struggle with novel problem structures (Ross and Bagnell, 2010; De Haan et al., 2019). Conversely, RL fine-tuning based on preference feedback can improve step-wise correctness, yet it often requires extensive exploration in the absence of offline demonstrations. By combining SFT and RL, models can leverage the strengths of both approaches, achieving more reliable and robust reasonings. As illustrated in Figure 1, these

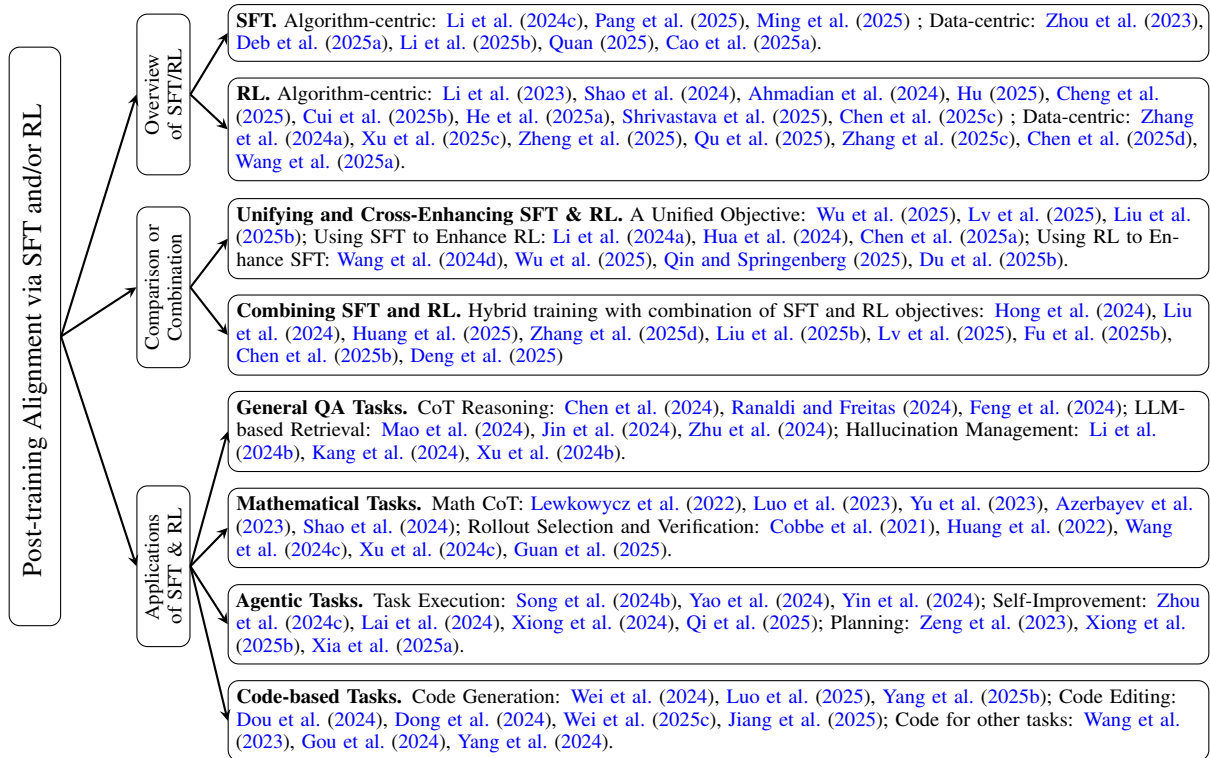


Figure 1: A taxonomy of large language model (LLM) post-training alignment methods via supervised fine-tuning (SFT) and reinforcement learning (RL). We organize prior work along three categories: (1) algorithm-centric versus data-centric approaches within SFT and RL, (2) comparative, unifying, and hybrid frameworks that integrate SFT and RL objectives, and (3) representative downstream application domains, including reasoning, mathematics, agentic behavior, and code-related tasks.

studies highlight that understanding and combining the complementary strengths of SFT and RL is crucial for advancing LLM post-training methods.

While recent surveys offer valuable insights into LLM post-training, the majority of them typically examine SFT or RL separately (Parthasarathy et al., 2024; Tao et al., 2024; Mao et al., 2025; Tie et al., 2025; Zhang et al., 2025b,a), leaving the relationships between these approaches comparatively underexplored. Other works focus on specialized dimensions of post-training, such as vision-centric adaptation (Chu et al., 2025), advances in reasoning (Kumar et al., 2025), agentic behaviors (Du et al., 2025a), or scaling strategies (Lai et al., 2025). In contrast, our survey provides a systematic and integrated perspective on SFT and RL as complementary post-training tools, with particular emphasis on their interplay and practical applications.

Our key **contributions** are threefold:

- We are the first study to systematically summarize and compare SFT and RL in LLM post-training, providing a clear understanding of what SFT and RL are, and how they can be extended from both algorithm-centric and

data-centric perspectives.

- We then establish a unified framework for characterizing SFT and RL, highlighting how they can complement each other or be integrated into hybrid learning approaches.
- From a survey of applications spanning 2023 to 2025, we observe rapid task-domain expansion, growing adoption of integrated SFT–RL training, and a continued shift from API-based labeling to open-weight-generated datasets.

2 Background: SFT and RL

Supervised Fine-Tuning (SFT). SFT is a method for adapting pre-trained LLMs to specific tasks or domains by training on high-quality prompt–response pairs using standard language modeling objectives. This process typically involves collecting or curating a dataset of expert-written demonstrations $D = \{(x, y)\}$, prompt x paired with target responses y , and then fine-tuning the model π_θ :

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [-\log \pi_\theta(y | x)].$$

In essence, SFT trains a model to imitate expert behavior. It is closely related to Behavior Cloning (BC) (Pomerleau, 1991) in traditional reinforcement learning, as both frameworks learn directly from expert demonstrations without relying on explicit reward signals. The goal of SFT is to reproduce expert performance; however, like BC, it suffers from distribution shift, which can lead to compounding errors (Ross and Bagnell, 2010; De Haan et al., 2019), and depends heavily on the quality of the demonstrations.

Reinforcement Learning (RL). RL offers an alternative paradigm in which LLMs are tuned not through direct supervised signals, but by optimizing behaviors according to a reward function. In classical RL, models interact with environments and learn a policy that maximizes cumulative rewards over time. Recent LLMs post-training methods (Ouyang et al., 2022; Guo et al., 2025; Yang et al., 2025a) have adopted similar techniques: an environment is constructed, the model serves as a policy that interacts with this environment, and training proceeds to maximize expected rewards:

$$\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [r(x, y)],$$

where the reward function r is either manually specified or learned from data. Although RL can also refer to offline reinforcement learning without environment interactions, in this work, we use it exclusively to denote online RL, as it plays a predominant role in the current frontier of LLM alignment, including reasoning and agent development.

Here, we summarize the key distinction between SFT and RL in the context of LLM post-training, as commonly described in current literature (Shao et al., 2024; Zhang et al., 2025b): SFT is an **supervised learning** paradigm that trains on expert-annotated prompt–response pairs, Whereas RL is a **reward-driven optimization** paradigm that learns by updating the model from its own generations.

3 SFT and RL: Distinct Methodological Landscapes

In the evolving landscape of LLM post-training, SFT and RL are two primary methodological paradigms. This section presents an overview of these approaches from (1) **Algorithm-centric**: refined training algorithms or loss functions, and (2) **Data-centric**: curated data selection or sophisticated data synthesis.

3.1 SFT

Algorithm-centric SFT. Li et al. (2024c) introduces Entropic Distribution Matching (GEM), which reformulates SFT as a distribution-matching problem with entropy regularization to mitigate overfitting and preserve output diversity. Token Cleaning (Pang et al., 2025) estimates the contribution of each token to model updates and removes uninformative ones, effectively reducing noise in supervision. One-Token Rollout (Ming et al., 2025) is a policy-gradient-inspired variant of SFT that treats each token prediction as a one-step trajectory and leverages the ground-truth token as a reward, introducing on-policy learning signals without the complexity of RL.

Data-centric SFT. Zhou et al. (2023) fine-tunes their model on only 1,000 high-quality and diverse instruction–response pairs, achieving alignment performance comparable to that of much larger models. FisherSFT (Deb et al., 2025b) selects training examples that maximize information gain, achieving efficient learning with limited data. Li et al. (2025b) optimizes data mixing by learning domain-specific weights that minimize validation loss, thereby improving generalization with minimal tuning cost. Quan (2025) proposes to automatically generate context-driven instruction–response pairs to enrich and diversify SFT data without heavy human annotation. Finally, Condor (Cao et al., 2025a) integrates knowledge-guided synthesis and iterative refinement to produce high-fidelity alignment data, further demonstrating the importance of data curation.

3.2 RL

Algorithm-centric RL. Several modified objectives and training procedures have been proposed to extend LLM capabilities beyond what standard RL achieves. A key line of innovation lies in policy optimization algorithms. While Proximal Policy Optimization (PPO) (Schulman et al., 2017), with a learned value critic, has been the dominant approach, recent work favors critic-free methods such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024), which simplifies training by replacing the value network with group-relative normalized advantages. Other methods adopt direct REINFORCE-style updates (Li et al., 2023; Ahmadian et al., 2024; Hu, 2025; Xiong et al., 2025a), foregoing the complex components of PPO. Another active direction focuses on improving train-

ing stability and efficiency through objective regularization. Entropy regularization remains crucial for preventing entropy collapse (Cheng et al., 2025; Cui et al., 2025a). He et al. (2025a); Shrivastava et al. (2025) introduce weighted token entropy to regularize policy updates. In contrast, Cui et al. (2025a) identifies the covariance between an action’s probability and its advantage as the key entropy “driver,” proposing covariance and Kullback–Leibler (KL) clipping to selectively constrain tokens with exceptionally high covariance. Moreover, Cheng et al. (2025) and Chen et al. (2025c) incorporate entropy information into advantage estimation, further stabilizing training dynamics.

Data-centric RL. Several approaches (Zhang et al., 2024a; Xu et al., 2025c) perform rollout selection, demonstrating that selectively training on a subset of informative model rollouts can achieve promising performance, such as high-variance or diverse examples, hence highly data efficient. Other approaches (Zheng et al., 2025; Qu et al., 2025) explore prompt selection (prior to rollout) to reduce computational cost while maintaining performance. Curriculum learning also contributes to this direction. Zhang et al. (2025c); Yao et al. (2025) dynamically select prompts of intermediate difficulty to maximize the learning signal. Moreover, Chen et al. (2025d); Wang et al. (2025a) introduce a distribution-level adaptation approach that prioritizes tasks where the model exhibits the greatest advantage or lowest visitation.

4 Comparison and Combination of SFT and RL

Despite the prevalence of SFT and RL, two commonly adopted post-training techniques for LLMs, they stay relatively disentangled and are usually applied in a sequential manner. For example, Reinforcement Learning from Human Feedback (RLHF) initially applies SFT to inject prior general knowledge into the LLMs, and then conducts RL for promoting the capability in a particular aspect (Ouyang et al., 2022; Bai et al., 2022). However, there is a missing systematical comparison and combination between these two stages, especially from a methodological and theoretical perspective, despite that there have been some pioneering works that try to modify the objectives of each other to complement each other’s strengths. The detailed comparison between different manners to combine SFT and RL together is summarised in Table 1.

4.1 A Unified Objective

First, we write down the objectives of both SFT and RL respectively in an explicit manner,

$$\begin{aligned}\mathcal{L}_{\text{SFT}} &= \mathbb{E}_{(x,y)\sim\mathcal{D}}[-\log \pi_{\theta}(y|x)], \\ \mathcal{L}_{\text{RL}} &= \mathbb{E}_{x\sim\mathcal{D}_x, y\sim\pi_{\theta}(\cdot|x)}[r(x, y)].\end{aligned}$$

Based on the above formulation, we could derive the gradients for both objectives below,

$$\begin{aligned}\nabla_{\theta}\mathcal{L}_{\text{SFT}} &= \mathbb{E}_{(x,y)\sim\mathcal{D}}[-\nabla_{\theta}\log \pi_{\theta}(y|x)], \\ \nabla_{\theta}\mathcal{L}_{\text{RL}} &= \mathbb{E}_{x\sim\mathcal{D}_x, y\sim\pi_{\theta}(\cdot|x)}[\nabla_{\theta}\log \pi_{\theta}(y|x) \cdot r(x, y)].\end{aligned}$$

As pointed out in recent works (Wu et al., 2025), the objective of SFT could be regarded as a special case of RL, i.e.,

$$\begin{aligned}\nabla_{\theta}\mathcal{L}_{\text{SFT}} &= \mathbb{E}_{(x_i, y_i)\sim\mathcal{D}}[-\nabla_{\theta}\log \pi_{\theta}(y_i|x_i)] \\ &= \mathbb{E}_{x_i\sim\mathcal{D}_x}[\mathbb{E}_{y\sim\pi_{\theta}(\cdot|x)}[-\frac{\mathbb{I}\{y=y_i\}}{\pi_{\theta}(y|x)}\nabla_{\theta}\log \pi_{\theta}(y_i|x_i)]],\end{aligned}$$

where $\mathbb{I}\{y=y_i\}$ is the indicator function, which takes value 1 when the output y sampled from policy $\pi_{\theta}(\cdot|x)$ is exactly the same with the ground-truth response y_i in the SFT training dataset. Therefore, the ratio $-\frac{\mathbb{I}\{y=y_i\}}{\pi_{\theta}(y|x)}$ could be seen as a proxy reward function in SFT.

Based on the above argument, formatting SFT as a special case of RL, the training objective of most post-training could be abstracted into the following formula,

$$\begin{aligned}\max_{\pi} &\mathbb{E}_{x\sim\mathcal{D}_x, y\sim\pi_{\theta}(\cdot|x)}[r(x, y)] \\ &- \beta \text{KL}(\pi_{\theta}(\cdot|x)\|\pi_0(\cdot|x)),\end{aligned}$$

where π_0 is the base (reference) policy model; r is a proxy for the reward, and β is a hyper-parameter. The KL regularization between the policy model and the reference model restricts π_{θ} from deviating too much from a pre-trained checkpoint, primarily for stability reasons. In the following sections, we will continue to discuss the difference and combination of SFT and RL, mainly from the objective perspective.

Therefore, optimizing LLMs with both SFT and RL objectives ultimately collapses to the RL objective, and the tricks that work for one thus have potential to be applied to the other. To mitigate the drop in generalization ability from SFT, one could consider applying importance sampling and online rollout, as in RL. Meanwhile, to help LLMs memorize additional knowledge, one could integrate the

SFT loss into the RL objective. The roles of these two stages should be regarded as a mutually reinforcing and interdependent relationship instead of merely being applied alternatively.

4.2 Leveraging SFT to Enhance RL

Combine offline demonstration and online data

Recent studies have leveraged SFT to enhance RL in LLMs by combining offline demonstrations with online rollouts. [Yan et al. \(2025a\)](#) introduces an off-policy guided framework that augments on-policy updates with reasoning traces, effectively balancing imitation and exploration. SRFT ([Fu et al., 2025b](#)) integrates supervised and reinforcement objectives in a single-stage framework, avoiding inefficiencies of sequential fine-tuning. Similar to LLMs, Vision Language Models (VLMs) are another flourishing area with great potential by extending the single text-based modality of LLMs into vision, like images and videos. [Liu et al. \(2025a\)](#) propose a dynamic memorization–exploration strategy for small VLMs, which adaptively alternates between demonstration imitation and online reward optimization. AMFT ([He et al., 2025b](#)) adopts a meta-learning perspective, dynamically adjusting the imitation–exploration balance to maximize performance. [Lv et al. \(2025\)](#) provides a unified view of post-training, showing that both SFT and RL optimize a common objective, and propose hybrid updates to exploit demonstration and rollout data. BREAD ([Zhang et al., 2025e](#)) bridges SFT and RL through branched rollouts anchored by expert prefixes, reducing reliance on large demonstration sets while improving stability. Similarly, [Huang et al. \(2025\)](#) develops a prefix sampling approach that guides generation into high-quality trajectories before applying RL optimization. Collectively, these methods illustrate complementary strategies, including single-stage integration, adaptive weighting, and prefix-based seeding, that enhance the efficiency and robustness of RL when grounded in SFT demonstrations.

Objective Modification NFT ([Chen et al., 2025a](#)) enables models to learn from both correct and incorrect outputs under supervision while implicitly optimizing a policy to enhance reasoning performance. UFT ([Liu et al., 2025b](#)) unifies supervised and reinforcement fine-tuning into a single process that balances memorization and exploration, achieving better generalization and exponentially improved sample efficiency. ReLIFT ([Zhu](#)

[et al., 2025b](#)) interleaves RL with SFT on the hardest questions, allowing models to acquire capabilities beyond what pure RL can achieve. [Zhu et al. \(2025b\)](#) demonstrates that penalizing incorrect answers alone, through negative sample reinforcement, can substantially improve reasoning performance, often matching or exceeding classical RL methods by suppressing wrong outputs and reallocating probability to plausible alternatives.

4.3 From RL Perspective to Improve SFT

DFT ([Wu et al., 2025](#)) is proposed to rescale each token’s loss by its predicted probability to rectify implicit reward bias in SFT and boost generalization. iw-SFT ([Qin and Springenberg, 2025](#)) interpret curated SFT as optimizing a lower bound on a sparse-reward RL objective and tightening that bound via importance weights. Another work ([Du et al., 2025b](#)) in RLHF proposes a reweighted reward-driven SFT objective through variational inference. In contrast, inspired by the RL, especially Direct Preference Optimization (DPO) ([Rafailov et al., 2023](#)), objective, [Wang et al. \(2024d\)](#) minimizes the distribution difference between the reference model and policy model on an offline dataset.

4.4 Hybrid Training Combining SFT and RL

[Huang et al. \(2025\)](#) proposes to combine SFT and RL together through utilizing the prefix of a ground-truth response to generate new continuation based on the current policy model, and use SFT loss to train the prefix partition while use RL loss to train the newly generated partition. [Lv et al. \(2025\)](#) uses interleaving SFT and RL based on the performance of the policy model during online rollouts. When the performance is above a preset threshold, online RL is preferred for exploration, while correct guidance from SFT is preferred when the performance is bad. [Liu et al. \(2024\)](#) starts from the original DPO ([Rafailov et al., 2023](#)) objective, and directly injects the SFT loss on samples from the base model to mitigate the overoptimization. [Zhang et al. \(2025d\)](#) combines a weighted SFT loss and GRPO ([Shao et al., 2024](#)) loss together to aggregate the off-policy and on-policy training. [Liu et al. \(2025b\)](#) uses expert data in the first several steps, and switches to the new generation in the following steps. SRL ([Deng et al., 2025](#)) proposes to decompose the reasoning process into several intermediate steps, and compare online rollouts with offline expert trajectories, and assign rewards proportional to the matching between them.

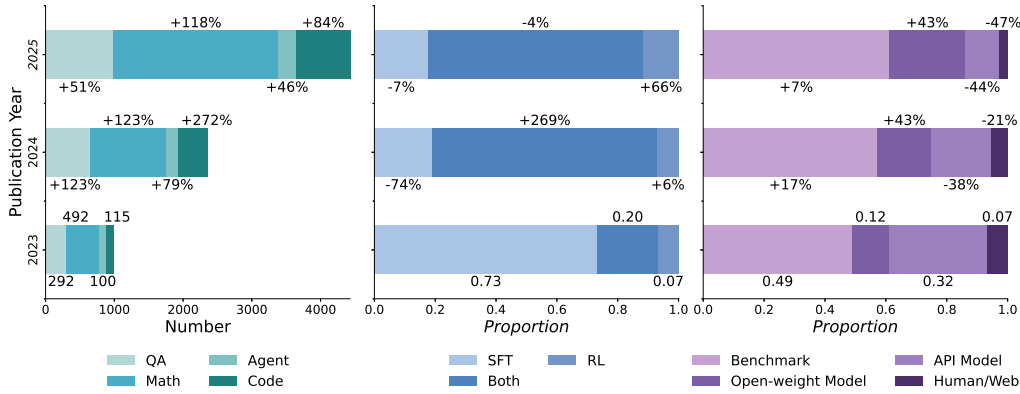


Figure 2: Trends in task focus, training methodologies, and ground-truth data sources from 2023 to 2025. Results show rapid growth across all surveyed domains, with substantial increases in research volume and diversification of application areas; increasing convergence toward hybrid SFT–RL pipelines, supported by more mature training infrastructures, libraries, and preference datasets; and a continued shift from API-based labeling to data generated with increasingly capable open-weight models. Projections for 2025 and all reported proportions are derived from surveyed publications; see Appendix B.4 for further discussion.

5 Applications

We focus on application studies that post-train LLMs using SFT or RL with four domains, where each domain encompasses distinct methodologies and challenges, and these categories capture the breadth of LLM capabilities and support a systematic analysis of performance across tasks. Key trends are summarized in Fig. 2, with additional details provided in Appendix B.

5.1 LLMs for General QA Tasks

This subsection reviews methods for improving LLM Question Answering (QA) performance via reasoning augmentation, e.g., CoT; external knowledge integration, e.g., Retrieval-Augmented Generation (RAG); and hallucination mitigation techniques.

Step-by-Step Reasoning. Smaller Foundation Model (FM) trained only on web data generally lack native CoT reasoning; thus, CoT demonstrations are typically generated using larger models via prompt engineering and subsequently used for SFT to instruction-tune reasoning or query decomposition capabilities. Building upon this approach, some works (Ranaldi and Freitas, 2024; Feng et al., 2024) introduce an additional RL stage, where collected demonstrations are treated as positive examples and newly generated model outputs are treated as negative samples. A key challenge in equipping LLMs with CoT capabilities lies in evaluating the quality of generated reasoning, as providing fine-grained, step-wise supervision in an online setting is often difficult or costly. To mitigate

this, researchers have proposed using surrogate signals, such as answer correctness (Chen et al., 2024; Guo et al., 2025) or reasoning comprehensiveness (Huang et al., 2024), and applying Rejection Fine-Tuning (RFT) or RL on self-generated samples.

LLM-based Retrieval. In RAG pipelines, LLMs are often used for query rewriting prior to retrieval. Prior work shows that FMs can perform this task effectively, with retrieval quality evaluated via downstream responses, and several studies (Mao et al., 2024; Hsu et al., 2024; Yan and Ling, 2025; Xu et al., 2025a; Gao et al., 2025) directly apply RL with task-specific search rewards for model alignment.

Hallucination Management. LLM hallucination refers to the generation of responses that appear plausible but contain incorrect or fabricated information. This phenomenon typically arises when the required factual knowledge is absent from the pre-training corpus, leading the model to infer or construct an answer rather than abstain. To mitigate hallucinations for general question, researchers have leveraged PPO to reduce factual errors by incorporating response corrections (Li et al., 2024b) provided by stronger models. Additionally, to encourage models to explicitly acknowledge uncertainty and suppress hallucinated outputs, prior works (Kang et al., 2024; Xie et al., 2024; Cheng et al., 2024; Xu et al., 2024a) have proposed inference over a collection of sampled responses, often derived from benchmark datasets, and synthesizing “I do not know” outputs for subsequent SFT or RL pipelines. Furthermore, to quantify

and model uncertainty, several studies (Band et al., 2024; Sun et al., 2025) generate multiple candidate responses for each training instance and fine-tune the model using aggregated or summarized outputs as approximations of the underlying uncertainty distribution.

5.2 LLMs for Mathematical Tasks

This subsection reviews methods for enhancing LLMs’ mathematical capabilities, with a focus on joint mathematical reasoning and training using step-wise rollouts.

Mathematical Reasoning. Early efforts for mathematical reasoning abilities can be traced back to Polu and Sutskever (2020) and Cobbe et al. (2021), which studied problems involving algorithmic computation and mathematical induction. Subsequent work further improved performance by SFT FMs on curated mathematical corpora collected from the web (Lewkowycz et al., 2022; Azerbayev et al., 2023), synthetically generated datasets (Yu et al., 2023; Toshniwal et al., 2024a), or data produced by stronger models (Luo et al., 2023; Toshniwal et al., 2024b). More recent studies construct explicit reward models for mathematical evaluation and propose alignment procedures, such as GRPO (Shao et al., 2024) and Critique-DPO (Xu et al., 2024c), to enable downstream RL.

Rollout-Based Training. Similar to standard CoT prompting, developing reliable step-wise reward evaluation for mathematical reasoning remains challenging. Huang et al. (2022) proposes fine-tuning models on their own generated solutions using majority voting. Additionally, many researchers exploited FMs with multiple reasoning paths when solving mathematical problems. For PPO-based methods, Wang et al. (2024c) and Guan et al. (2025) approximate process-level rewards using the probability of producing a correct final answer, while Shen et al. (2025) randomly assigns reflection prompts and incorporates both outcome-based and reflection-based bonus rewards.

5.3 LLMs for Agentic Tasks

This subsection reviews strategies for improving the agentic capabilities of LLMs, with a focus on action selection and long-horizon planning.

Action Selection. Agentic tasks require selecting actions over multiple steps within a given environment. To equip FMs with this capability, pioneering work such as Schick et al. (2023); Yin et al. (2024); Zhang et al. (2024b,c) synthesizes

or unifies action trajectories across diverse environments for SFT. Subsequent studies adopt a sequential paradigm that combines SFT with RL using environment-specific rewards. For example, Xi et al. (2024) incorporates SFT rewards directly into the SFT objective, while hierarchical RL approaches (Zhou et al., 2024c; Hu et al., 2025b) assign rewards at the sentence level while updating policies at the token level. Similarly, DMPO (Shi et al., 2025) introduces step-wise reward discounting, and Qi et al. (2025) employs critic-based curriculum learning to stabilize RL training.

Planning with LLMs. Task planning has emerged as an active research direction, with agentic tasks serving as a natural testbed due to their long action horizons and complex inter-step dependencies. Recent work predominantly leverages outcome-based surrogate rewards and GPT-synthesized sequences to optimize planning behavior. For instance, Xiong et al. (2025b) proposes a hierarchical planning framework with a top-level planner fine-tuned on seed plans and optimized via DPO. Vattikonda et al. (2025) further incorporates RFT to refine high-reward actions, while Wang et al. (2025b) improves planning performance by adjusting exploration strategies and prioritizing critical steps during DPO training.

5.4 LLMs for Coding Tasks

This subsection reviews methods for improving coding capabilities in LLMs, focusing on code generation and code editing.

Code Generation. Recent work has explored various approaches to enhance LLM-based code generation. Under SFT, Chaudhary (2023) employs self-instruction to synthesize training data, while Luo et al. (2025) generates more challenging queries by modifying constraints of seed problems. Wei et al. (2024) leverages open-source code snippets to produce diverse and controllable instruction datasets. From the RL perspective, Yang et al. (2025b) introduces a multi-agent framework in which an SFT model generates rollouts for subsequent DPO training.

Code Editing. Code editing aims to refine programs that fail to satisfy intended functionality, whether human-written or LLM-generated. For example, Wei et al. (2025c) leverages human pull requests for GRPO-based training, while Chae et al. (2024) and Jiang et al. (2025) generate iterative refinement rollouts evaluated using unit tests to drive RL optimization.

575 **6 Common Practice and Takeaways**

576 In this section, we summarize common practices
577 and key takeaways from our previous theoretical
578 and application discussions, highlighting the inter-
579 play between SFT and RL, when to choose each
580 method, and training trade-offs.

581 **Relationship Between SFT and RL Objectives.**

582 As discussed in Section 4, SFT can be viewed as
583 a special case of RL under certain formulations,
584 where an indicator function that measures whether
585 the current policy reproduces an offline trajectory
586 acts as a surrogate reward. This shared optimiza-
587 tion structure enables techniques developed for ei-
588 ther SFT or RL to transfer across paradigms. For
589 example, Section 5.3 reviews several studies that
590 jointly optimize weighted SFT and RL objectives,
591 demonstrating consistent improvements over stan-
592 dalone training losses.

593 **Leveraging Expert Data or Policies.**

594 When high-quality expert data are available, SFT is gen-
595 erally preferred as an initial training stage over
596 RL as observed in Section 5, possibly due to its
597 simpler implementation and greater stability. If a
598 strong policy model or prompt is available to gen-
599 erate expert data, incorporating importance sam-
600 pling on failed queries can mitigate distribution
601 shift by treating them as informative positive sam-
602 ples for RFT/RL, bridging the gap between the lim-
603 ited offline dataset and the evolving policy. When
604 a reliable reward model can be trained, the com-
605 mon practice is to first perform SFT and then RL,
606 which typically achieves the highest reported per-
607 formance.

608 **Training Strategies and Trade-offs.**

609 SFT is essential when the policy model is unfamiliar with the
610 downstream task or cannot efficiently generate suf-
611 ficient positive samples. However, SFT may reduce
612 the model’s generalization ability compared with
613 RL. Conversely, directly applying RL can better
614 balance performance gains with exploratory capac-
615 ity, though it may suffer from issues such as entropy
616 collapse (Cui et al., 2025a) or reward hacking (Pan
617 et al., 2024a,b), one possible reason why the ma-
618 jority of pre-training computation is still devoted
619 to SFT.

620 **7 Future Directions and Open Problems**

621 Despite rapid progress, many fundamental ques-
622 tions remain unresolved. We highlight two open
623 challenges and outline promising directions for fu-
624 ture research.

Sample- and compute-efficient methodologies.

625 State-of-the-art SFT and RL pipelines typically re-
626 quire substantial computational resources, large
627 volumes of high-quality data, and extensive roll-
628 out generation, raising both practical and environ-
629 mental concerns. Improving efficiency is there-
630 fore a central challenge. Early progress includes
631 data-efficient SFT based on information-theoretic
632 principles (Deb et al., 2025b) and quantization-
633 aware methods that reduce training cost (Wei et al.,
634 2025b). For RL, recent work explores selective
635 or partial rollouts to reduce computation (Zheng
636 et al., 2025; Xu et al., 2025b). Nevertheless, signif-
637 icantly more research is needed to develop broadly
638 applicable, highly efficient methods.

SFT and RL under sparse or indirect reward

640 **signals.** Many real-world tasks lack well-defined or
641 easily verifiable reward signals. User feedback may
642 be sparse, inconsistent, or costly. In safety-critical
643 settings, it may be infeasible to obtain ground-truth
644 evaluations altogether. Recent work has begun to
645 explore alignment through verbal or otherwise in-
646 direct feedback (Stephan et al., 2024; Wang et al.,
647 2024b). An open challenge is to identify and lever-
648 age additional natural but underexplored sources of
649 implicit supervision, such as self-evaluation signals
650 or user-churn behaviors, to guide both SFT and RL
651 in low-feedback regimes.

8 Conclusion

652 This survey provides a unified view of SFT and
653 RL as post-training strategies for LLMs, examin-
654 ing both their theoretical foundations and practical
655 implementations. SFT offers stable and effi-
656 cient learning from high-quality offline datasets,
657 whereas RL enables reward-driven optimization
658 that enhances generalization and exploration. Re-
659 cent research characterizes SFT as a special case
660 of RL, underscoring the promise of hybrid ap-
661 proaches that balance training stability with adap-
662 tive behavior. Application papers published be-
663 tween 2023 and 2025 indicate a clear shift toward
664 hybrid post-training pipelines that integrate SFT
665 and RL, with scalable workflows centered on open-
666 weight, model-generated rollouts emerging as the
667 dominant paradigm. These developments under-
668 score the need for integrative methodologies that
669 leverage the strengths of both approaches, and this
670 survey lays a foundation by consolidating the cur-
671 rent state of the field, theoretical advances, and
672 versatile LLM applications.

675 Limitations

676 Despite our best efforts to ensure comprehensiveness,
677 this survey may omit some recent advances in
678 SFT/RL research due to the rapid pace of progress
679 and the sheer volume of emerging literature in this
680 field. In the *Applications* section, the use of a paper-
681 filtering strategy may introduce approximation bias,
682 as commonly adopted benchmarks are not fully in-
683 clusive of all real-world scenarios or methodolo-
684 gies. Additional discussion of application approxi-
685 mation limitations is provided in Appendix B.2.

686 References

687 Mohammad Mahdi Abootorabi, Amirhosein Zobeiri,
688 Mahdi Deghani, Mohammadali Mohammadkhani,
689 Bardia Mohammadi, Omid Ghahroodi, Mahdih So-
690 leymani Baghshah, and Ehsaneddin Asgari. 2025.
691 [Ask in any modality: A comprehensive sur-
692 vey on multimodal retrieval-augmented generation.](#)
693 *Preprint*, arXiv:2502.08826.

694 Arash Ahmadian, Chris Cremer, Matthias Gallé,
695 Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,
696 Ahmet Üstün, and Sara Hooker. 2024. Back to ba-
697 sics: Revisiting reinforce style optimization for learn-
698 ing from human feedback in llms. *arXiv preprint*
699 *arXiv:2402.14740*.

700 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster,
701 Marco Dos Santos, Stephen McAleer, Albert Q Jiang,
702 Jia Deng, Stella Biderman, and Sean Welleck. 2023.
703 Llemma: An open language model for mathematics.
704 *arXiv preprint arXiv:2310.10631*.

705 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
706 Askell, Anna Chen, Nova DasSarma, Dawn Drain,
707 Stanislav Fort, Deep Ganguli, Tom Henighan, and 1
708 others. 2022. Training a helpful and harmless assis-
709 tant with reinforcement learning from human feed-
710 back. *arXiv preprint arXiv:2204.05862*.

711 Neil Band, Xuechen Li, Tengyu Ma, and Tatsunori
712 Hashimoto. 2024. [Linguistic calibration of long-
713 form generations.](#) *Preprint*, arXiv:2404.00474.

714 Maosong Cao, Taolin Zhang, Mo Li, Chuyu Zhang,
715 Yunxin Liu, Haodong Duan, Songyang Zhang, and
716 Kai Chen. 2025a. Condor: Enhance llm alignment
717 with knowledge-driven data synthesis and refinement.
718 *arXiv preprint arXiv:2501.12273*.

719 Zouying Cao, Runze Wang, Yifei Yang, Xinbei Ma, Xi-
720 aoyong Zhu, Bo Zheng, and Hai Zhao. 2025b. [Pgp-
721 o: Enhancing agent reasoning via pseudocode-style
722 planning guided preference optimization.](#) *Preprint*,
723 arXiv:2506.01475.

724 Hyunjoo Chae, Taeyoon Kwon, Seungjun Moon,
725 Yongho Song, Dongjin Kang, Kai Tzu iunn Ong,
726 Beong woo Kwak, Seonghyeon Bae, Seung won

Hwang, and Jinyoung Yeo. 2024. [Coffee-gym: An
environment for evaluating and improving natural
language feedback on erroneous code.](#) *Preprint*,
arXiv:2409.19715. 727
728
729
730

Sahil Chaudhary. 2023. Code alpaca: An instruction-
following llama model for code generation. <https://github.com/sahil280114/codealpaca>. 731
732
733

Huayu Chen, Kaiwen Zheng, Qinsheng Zhang, Ganqu
Cui, Yin Cui, Haotian Ye, Tsung-Yi Lin, Ming-Yu
Liu, Jun Zhu, and Haoxiang Wang. 2025a. Bridging
supervised learning and reinforcement learning in
math reasoning. *arXiv preprint arXiv:2505.18116*. 734
735
736
737
738

Jack Chen, Fazhong Liu, Naruto Liu, Yuhan Luo, Erqu
Qin, Harry Zheng, Tian Dong, Haojin Zhu, Yan
Meng, and Xiao Wang. 2025b. Step-wise adaptive
integration of supervised fine-tuning and reinforce-
ment learning for task-specific llms. *arXiv preprint*
arXiv:2505.13026. 739
740
741
742
743
744

Minghan Chen, Guikun Chen, Wenguan Wang, and
Yi Yang. 2025c. Seed-grpo: Semantic entropy en-
hanced grpo for uncertainty-aware policy optimiza-
tion. *arXiv preprint arXiv:2505.12346*. 745
746
747
748

Xiaoyin Chen, Jiarui Lu, Minsu Kim, Dinghuai Zhang,
Jian Tang, Alexandre Piché, Nicolas Gontier, Yoshua
Bengio, and Ehsan Kamalloo. 2025d. Self-evolving
curriculum for llm reasoning. *arXiv preprint*
arXiv:2505.14970. 749
750
751
752
753

Zhaorun Chen, Zhuokai Zhao, Zhihong Zhu, Ruiqi
Zhang, Xiang Li, Bhiksha Raj, and Huaxiu Yao.
2024. [Autoprml: Automating procedural supervision
for multi-step reasoning via controllable question de-
composition.](#) *Preprint*, arXiv:2402.11452. 754
755
756
757
758

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai,
Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei.
2025. Reasoning with exploration: An entropy per-
spective. *arXiv preprint arXiv:2506.14758*. 759
760
761
762

Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wen-
wei Zhang, Zhangyue Yin, Shimin Li, Linyang Li,
Zhengfu He, Kai Chen, and Xipeng Qiu. 2024. [Can
ai assistants know what they don't know?](#) *Preprint*,
arXiv:2401.13275. 763
764
765
766
767

Sanjiban Choudhury and Paloma Sodhi. 2024. [Better
than your teacher: Llm agents that learn from privi-
leged ai feedback.](#) *Preprint*, arXiv:2410.05434. 768
769
770

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Sheng-
bang Tong, Saining Xie, Dale Schuurmans, Quoc V
Le, Sergey Levine, and Yi Ma. 2025. Sft mem-
orizes, rl generalizes: A comparative study of
foundation model post-training. *arXiv preprint*
arXiv:2501.17161. 771
772
773
774
775
776

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, Christopher Hesse, and John Schulman.
2021. [Training verifiers to solve math word prob-
lems.](#) *Preprint*, arXiv:2110.14168. 777
778
779
780
781
782

783	Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, Zhiyuan Liu, Hao Peng, Lei Bai, Wanli Ouyang, Yu Cheng, Bowen Zhou, and Ning Ding. 2025a. The entropy mechanism of reinforcement learning for reasoning language models . <i>Preprint</i> , arXiv:2505.22617.	839
784		840
785		841
786		
787		842
788		843
789		844
790	Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, and 1 others. 2025b. The entropy mechanism of reinforcement learning for reasoning language models. <i>arXiv preprint arXiv:2505.22617</i> .	845
791		
792		846
793		847
794		848
795		849
796	Wenqian Cui, Dianzhi Yu, Xiaoqi Jiao, Ziqiao Meng, Guangyan Zhang, Qichao Wang, Yiwen Guo, and Irwin King. 2025c. Recent advances in speech language models: A survey . <i>Preprint</i> , arXiv:2410.03751.	850
797		851
798		852
799		853
800		
801	Pim De Haan, Dinesh Jayaraman, and Sergey Levine. 2019. Causal confusion in imitation learning. <i>Advances in neural information processing systems</i> , 32.	854
802		855
803		856
804	Rohan Deb, Kiran Thekumparampil, Kousha Kalantari, Gaurush Hiranandani, Shoham Sabach, and Branislav Kveton. 2025a. Fishersft: Data-efficient supervised fine-tuning of language models using information gain. <i>arXiv preprint arXiv:2505.14826</i> . ICML Poster (2025).	857
805		858
806		
807		859
808		860
809		861
810	Rohan Deb, Kiran Thekumparampil, Kousha Kalantari, Gaurush Hiranandani, Shoham Sabach, and Branislav Kveton. 2025b. Fishersft: Data-efficient supervised fine-tuning of language models using information gain . <i>Preprint</i> , arXiv:2505.14826. ICML Poster (2025).	862
811		863
812		864
813		865
814		866
815		867
816	Yihe Deng, I Hsu, Jun Yan, Zifeng Wang, Rujun Han, Gufeng Zhang, Yanfei Chen, Wei Wang, Tomas Pfister, Chen-Yu Lee, and 1 others. 2025. Supervised reinforcement learning: From expert trajectories to step-wise reasoning. <i>arXiv preprint arXiv:2510.25992</i> .	868
817		
818		869
819		870
820		871
821		872
822	Zhirui Deng, Zhicheng Dou, Yutao Zhu, Ji-Rong Wen, Ruibin Xiong, Mang Wang, and Weipeng Chen. 2024. From novice to expert: Llm agent policy optimization via step-wise reinforcement learning . <i>Preprint</i> , arXiv:2411.03817.	873
823		
824		874
825		875
826		876
827	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms . <i>Preprint</i> , arXiv:2305.14314.	877
828		878
829		
830	Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. Self-play with execution feedback: Improving instruction-following capabilities of large language models . <i>Preprint</i> , arXiv:2406.13542.	879
831		880
832		881
833		882
834		883
835	Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Wei Shen, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Xuanjing Huang, and Tao Gui. 2024. Stepcoder: Improve code generation with reinforcement learning from compiler feedback . <i>Preprint</i> , arXiv:2402.01391.	884
836		885
837		886
838		887
		888
	Shangheng Du, Jiabao Zhao, Jinxin Shi, Zhentao Xie, Xin Jiang, Yanhong Bai, and Liang He. 2025a. A survey on the optimization of large language model-based agents. <i>arXiv preprint arXiv:2503.12434</i> .	889
		890
		891
		892
	Yuhao Du, Zhuo Li, Pengyu Cheng, Zhihong Chen, Yuejiao Xie, Xiang Wan, and Anningzhe Gao. 2025b. Simplify rlhf as reward-weighted sft: A variational method. <i>arXiv preprint arXiv:2502.11026</i> .	
		850
	Peiyuan Feng, Yichen He, Guanhua Huang, Yuan Lin, Hanchong Zhang, Yuchen Zhang, and Hang Li. 2024. Agile: A novel reinforcement learning framework of llm agents . <i>Preprint</i> , arXiv:2405.14751.	851
		852
		853
	Dayuan Fu, Keqing He, Yejie Wang, Wentao Hong, Zhuoma Gongque, Weihao Zeng, Wei Wang, Jingang Wang, Xunliang Cai, and Weiran Xu. 2025a. Agentrefine: Enhancing agent generalization through refinement tuning . <i>Preprint</i> , arXiv:2501.01702.	854
		855
		856
		857
		858
	Yuqian Fu, Tinghong Chen, Jiajun Chai, Xihuai Wang, Songjun Tu, Guojun Yin, Wei Lin, Qichao Zhang, Yuanheng Zhu, and Dongbin Zhao. 2025b. Srf: A single-stage method with supervised and reinforcement fine-tuning for reasoning. <i>arXiv preprint arXiv:2506.19767</i> .	859
		860
		861
		862
		863
		864
	Jingsheng Gao, Linxu Li, Weiyuan Li, Yuzhuo Fu, and Bin Dai. 2025. Smartrag: Jointly learn rag-related tasks from the environment feedback . <i>Preprint</i> , arXiv:2410.18141.	865
		866
		867
		868
	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yuju Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. Tora: A tool-integrated reasoning agent for mathematical problem solving . <i>Preprint</i> , arXiv:2309.17452.	869
		870
		871
		872
		873
	Xinyu Guan, Li Lina Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. 2025. rstar-math: Small llms can master math reasoning with self-evolved deep thinking . <i>Preprint</i> , arXiv:2501.04519.	874
		875
		876
		877
		878
	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, and 1 others. 2025. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. <i>Nature</i> , 645(8081):633–638.	879
		880
		881
		882
		883
	Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang, Jiacheng Xu, Wei Shen, and 1 others. 2025a. Skywork open reasoner 1 technical report. <i>arXiv preprint arXiv:2505.22312</i> .	884
		885
		886
		887
		888
	Lixuan He, Jie Feng, and Yong Li. 2025b. Amft: Aligning llm reasoners by meta-learning the optimal imitation-exploration balance. <i>arXiv preprint arXiv:2508.06944</i> .	889
		890
		891
		892

893	Tao He, Hao Li, Jingchang Chen, Runxuan Liu, Yixin Cao, Lizi Liao, Zihao Zheng, Zheng Chu, Jiafeng Liang, Ming Liu, and Bing Qin. 2025c. Breaking the reasoning barrier a survey on LLM complex reasoning through the lens of self-evolution . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 7377–7417, Vienna, Austria. Association for Computational Linguistics.	
894		
895		
896		
897		
898		
899		
900		
901	Joey Hong, Anca Dragan, and Sergey Levine. 2024. Q-sft: Q-learning for language models via supervised fine-tuning. <i>arXiv preprint arXiv:2411.05193</i> .	
902		
903		
904	Sheryl Hsu, Omar Khattab, Chelsea Finn, and Archit Sharma. 2024. Grounding by trying: Llms with reinforcement learning-enhanced retrieval . <i>Preprint</i> , arXiv:2410.23214.	
905		
906		
907		
908	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models . <i>Preprint</i> , arXiv:2106.09685.	
909		
910		
911		
912	Jian Hu. 2025. Reinforce++: A simple and efficient approach for aligning large language models. <i>arXiv preprint arXiv:2501.03262</i> .	
913		
914		
915	Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruixuan Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao, Xiangxin Zhou, Ziyu Zhao, Yuhuai Li, Shengze Xu, Shenzhi Wang, Xinchun Xu, Shuofei Qiao, Zhaokai Wang, Kun Kuang, Tiejong Zeng, Liang Wang, and 10 others. 2025a. OS agents: A survey on MLLM-based agents for computer, phone and browser use . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7436–7465, Vienna, Austria. Association for Computational Linguistics.	
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926	Zican Hu, Wei Liu, Xiaoye Qu, Xiangyu Yue, Chunlin Chen, Zhi Wang, and Yu Cheng. 2025b. Divide and conquer: Grounding llms as efficient decision-making agents via offline hierarchical reinforcement learning . <i>Preprint</i> , arXiv:2505.19761.	
927		
928		
929		
930		
931	Ermo Hua, Biqing Qi, Kaiyan Zhang, Kai Tian, Xingtai Lv, Ning Ding, and Bowen Zhou. 2024. Intuitive fine-tuning: Towards simplifying alignment into a single process. <i>arXiv preprint arXiv:2405.11870</i> .	
932		
933		
934		
935	Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. <i>arXiv preprint arXiv:2210.11610</i> .	
936		
937		
938		
939	Lei Huang, Xiaocheng Feng, Weitao Ma, Liang Zhao, Yuchun Fan, Weihong Zhong, Dongliang Xu, Qing Yang, Hongtao Liu, and Bing Qin. 2024. Advancing large language model attribution through self-improving . <i>Preprint</i> , arXiv:2410.13298.	
940		
941		
942		
943		
944	Zeyu Huang, Tianhao Cheng, Zihan Qiu, Zili Wang, Yinghui Xu, Edoardo M Ponti, and Ivan Titov. 2025. Blending supervised and reinforcement fine-tuning with prefix sampling. <i>arXiv preprint arXiv:2507.01679</i> .	
945		
946		
947		
948		
	Nan Jiang, Xiaopeng Li, Shiqi Wang, Qiang Zhou, Soneya Binta Hossain, Baishakhi Ray, Varun Kumar, Xiaofei Ma, and Anoop Deoras. 2025. Ledex: Training llms to better self-debug and explain code . <i>Preprint</i> , arXiv:2405.18649.	949
		950
		951
		952
		953
	Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. Swe-bench: Can language models resolve real-world github issues? <i>Preprint</i> , arXiv:2310.06770.	954
		955
		956
		957
		958
	Jiajie Jin, Yutao Zhu, Yujia Zhou, and Zhicheng Dou. 2024. Bider: Bridging knowledge inconsistency for efficient retrieval-augmented llms via key supporting evidence . <i>Preprint</i> , arXiv:2402.12174.	959
		960
		961
		962
	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension . <i>Preprint</i> , arXiv:1705.03551.	963
		964
		965
		966
	Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. 2024. Unfamiliar finetuning examples control how language models hallucinate . <i>Preprint</i> , arXiv:2403.05612.	967
		968
		969
		970
	Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, Caiming Xiong, and Shafiq Joty. 2025. A survey of frontiers in llm reasoning: Inference scaling, learning to reason, and agentic systems . <i>Preprint</i> , arXiv:2504.09037.	971
		972
		973
		974
		975
		976
	Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip HS Torr, Fahad Shahbaz Khan, and Salman Khan. 2025. Llm post-training: A deep dive into reasoning large language models. <i>arXiv preprint arXiv:2502.21321</i> .	977
		978
		979
		980
		981
		982
	Hanyu Lai, Xiao Liu, Junjie Gao, Jiale Cheng, Zehan Qi, Yifan Xu, Shuntian Yao, Dan Zhang, Jinhua Du, Zhenyu Hou, and 1 others. 2025. A survey of post-training scaling in large language models. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2771–2791.	983
		984
		985
		986
		987
		988
		989
	Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, and Jie Tang. 2024. Autowebglm: A large language model-based web navigating agent . <i>Preprint</i> , arXiv:2404.03648.	990
		991
		992
		993
		994
	Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In <i>International Conference on Machine Learning</i> , pages 18319–18345. PMLR.	995
		996
		997
		998
		999
		1000
	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo	1001
		1002
		1003

1004	Gutman-Solo, and 1 others. 2022. Solving quantitative reasoning problems with language models. <i>Advances in neural information processing systems</i> , 35:3843–3857.	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2025. Wizardcoder: Empowering code large language models with evol-instruct . <i>Preprint</i> , arXiv:2306.08568.	1058
1005			1059
1006			1060
1007			1061
1008	Jiaxiang Li, Siliang Zeng, Hoi-To Wai, Chenliang Li, Alfredo Garcia, and Mingyi Hong. 2024a. Getting more juice out of the sft data: Reward learning from human demonstration improves sft for llm alignment. <i>Advances in Neural Information Processing Systems</i> , 37:124292–124318.	Xingtai Lv, Yuxin Zuo, Youbang Sun, Hongyi Liu, Yuntian Wei, Zhekai Chen, Lixuan He, Xuekai Zhu, Kaiyan Zhang, Bingning Wang, and 1 others. 2025. Towards a unified view of large language model post-training. <i>arXiv preprint arXiv:2509.04419</i> .	1062
1009			1063
1010			1064
1011			1065
1012			1066
1013			1067
1014	Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2024b. The dawn after the dark: An empirical study on factuality hallucination in large language models . <i>Preprint</i> , arXiv:2401.03205.	Zexiong Ma, Chao Peng, Pengfei Gao, Xiangxin Meng, Yanzhen Zou, and Bing Xie. 2025. Sorft: Issue resolving with subtask-oriented reinforced fine-tuning . <i>Preprint</i> , arXiv:2502.20127.	1068
1015			1069
1016			1070
1017			1071
1018			1072
1019	Moxin Li, Yong Zhao, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, Tat-Seng Chua, and Yang Deng. 2025a. Knowledge boundary of large language models: A survey . <i>Preprint</i> , arXiv:2412.12472.	Shengyu Mao, Yong Jiang, Boli Chen, Xiao Li, Peng Wang, Xinyu Wang, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. 2024. Rafe: Ranking feedback improves query rewriting for rag . <i>Preprint</i> , arXiv:2405.14431.	1073
1020			1074
1021			1075
1022			1076
1023			1077
1024	Yuan Li, Zhengzhong Liu, and Eric Xing. 2025b. Data mixing optimization for supervised fine-tuning of large language models. <i>arXiv preprint arXiv:2508.11953</i> .	Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. 2025. A survey on lora of large language models. <i>Frontiers of Computer Science</i> , 19(7):197605.	1078
1025			1079
1026			1080
1027			1081
1028	Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Ruoyu Sun, and Zhi-Quan Luo. 2024c. Entropic distribution matching in supervised fine-tuning of llms: Less overfitting and better diversity. <i>arXiv preprint arXiv:2408.16673</i> . Version v1; latest version v2 (2025).	Rui Ming, Haoyuan Wu, Shoubo Hu, Zhuolun He, and Bei Yu. 2025. One-token rollout: Guiding supervised fine-tuning of llms with policy gradient. <i>arXiv preprint arXiv:2509.26313</i> .	1082
1029			1083
1030			1084
1031			1085
1032			1086
1033			1087
1034	Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. 2023. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. <i>arXiv preprint arXiv:2310.10505</i> .	Shikhar Murty, Hao Zhu, Dzmitry Bahdanau, and Christopher D. Manning. 2025. Nnetnav: Unsupervised learning of browser agents through environment interaction in the wild . <i>Preprint</i> , arXiv:2410.02907.	1088
1035			1089
1036			1090
1037			1091
1038			1092
1039	Jiazhen Liu, Yuchuan Deng, and Long Chen. 2025a. Empowering small vlms to think with dynamic memorization and exploration. <i>arXiv preprint arXiv:2506.23061</i> .	Dang Nguyen, Jian Chen, Yu Wang, Gang Wu, Namyong Park, Zhengmian Hu, Hanjia Lyu, Junda Wu, Ryan Aponte, Yu Xia, Xintong Li, Jing Shi, Hongjie Chen, Viet Dac Lai, Zhouhang Xie, Sungchul Kim, Ruiyi Zhang, Tong Yu, Mehrab Tanjim, and 11 others. 2025. GUI agents: A survey . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 22522–22538, Vienna, Austria. Association for Computational Linguistics.	1093
1040			1094
1041			1095
1042			1096
1043	Mingyang Liu, Gabriele Farina, and Asuman Ozdaglar. 2025b. Uft: Unifying supervised and reinforcement fine-tuning. <i>arXiv preprint arXiv:2505.16984</i> .	Fan Nie, Lan Feng, Haotian Ye, Weixin Liang, Pan Lu, Huaxiu Yao, Alexandre Alahi, and James Zou. 2025. Weak-for-strong: Training weak meta-agent to harness strong executors . <i>Preprint</i> , arXiv:2504.04785.	1097
1044			1098
1045			1099
1046	Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. 2024. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer. <i>Advances in Neural Information Processing Systems</i> , 37:138663–138697.	Tianyue Ou, Frank F. Xu, Aman Madaan, Jiarui Liu, Robert Lo, Abishek Sridhar, Sudipta Sengupta, Dan Roth, Graham Neubig, and Shuyan Zhou. 2024. Synatra: Turning indirect knowledge into direct demonstrations for digital agents at scale . <i>Preprint</i> , arXiv:2409.15637.	1100
1047			1101
1048			1102
1049			1103
1050			1104
1051			1105
1052	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. <i>arXiv preprint arXiv:2308.09583</i> .	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	1106
1053			1107
1054			1108
1055			1109
1056			1110
1057			1111
			1112
			1113

1114	Alexander Pan, Erik Jones, Meena Jagadeesan, and Jacob Steinhardt. 2024a. Feedback loops with language models drive in-context reward hacking . <i>Preprint</i> , arXiv:2402.06627.	1167
1115		1168
1116		1169
1117		1170
1118	Jane Pan, He He, Samuel R. Bowman, and Shi Feng. 2024b. Spontaneous reward hacking in iterative self-refinement . <i>Preprint</i> , arXiv:2407.04549.	1171
1119		1172
1120		
1121	Jinlong Pang, Na Di, Zhaowei Zhu, Jiaheng Wei, Hao Cheng, Chen Qian, and Yang Liu. 2025. Token cleaning: Fine-grained data selection for llm supervised fine-tuning. <i>arXiv preprint arXiv:2502.01968</i> .	
1122		
1123		
1124		
1125	Venkatesh Balavadhani Parthasarathy, Ahtsham Zafar, Aafaq Khan, and Arsalan Shahid. 2024. The ultimate guide to fine-tuning llms from basics to breakthroughs: An exhaustive review of technologies, research, best practices, applied research challenges and opportunities. <i>arXiv preprint arXiv:2408.13296</i> .	
1126		
1127		
1128		
1129		
1130		
1131	Stanislas Polu and Ilya Sutskever. 2020. Generative language modeling for automated theorem proving. <i>arXiv preprint arXiv:2009.03393</i> .	
1132		
1133		
1134	Dean A Pomerleau. 1991. Efficient training of artificial neural networks for autonomous navigation. <i>Neural computation</i> , 3(1):88–97.	
1135		
1136		
1137	Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Wenyi Zhao, Yu Yang, Xinyue Yang, Jiadai Sun, Shuntian Yao, Tianjie Zhang, Wei Xu, Jie Tang, and Yuxiao Dong. 2025. Webrl: Training llm web agents via self-evolving online curriculum reinforcement learning . <i>Preprint</i> , arXiv:2411.02337.	
1138		
1139		
1140		
1141		
1142		
1143	Chongli Qin and Jost Tobias Springenberg. 2025. Supervised fine tuning on curated data is reinforcement learning (and can be improved). <i>arXiv preprint arXiv:2507.12856</i> .	
1144		
1145		
1146		
1147	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, and 1 others. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. <i>arXiv preprint arXiv:2307.16789</i> .	
1148		
1149		
1150		
1151		
1152	Yun Qu, Qi Wang, Yixiu Mao, Vincent Tao Hu, Björn Ommer, and Xiangyang Ji. 2025. Can prompt difficulty be online predicted for accelerating rl finetuning of reasoning models? <i>arXiv preprint arXiv:2507.04632</i> .	
1153		
1154		
1155		
1156		
1157	Shanghaoran Quan. 2025. Automatically generating numerous context-driven sft data for llms across diverse granularity. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 25074–25082.	
1158		
1159		
1160		
1161		
1162	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.	
1163		
1164		
1165		
1166		
	Leonardo Ranaldi and Andre Freitas. 2024. Self-refine instruction-tuning for aligning reasoning in language models . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , page 2325–2347. Association for Computational Linguistics.	1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
	Christopher Rawles, Sarah Clinckemaulle, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyi Campbell-Ajala, Daniel Toyama, Robert Berry, Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. 2025. Androidworld: A dynamic benchmarking environment for autonomous agents . <i>Preprint</i> , arXiv:2405.14573.	1181
		1182
		1183
		1184
		1185
	Mathieu Rita, Florian Strub, Rahma Chaabouni, Paul Michel, Emmanuel Dupoux, and Olivier Pietquin. 2024. Countering reward over-optimization in llm with demonstration-guided reinforcement learning . <i>Preprint</i> , arXiv:2404.19409.	1186
		1187
		1188
		1189
		1190
	Stéphane Ross and Drew Bagnell. 2010. Efficient reductions for imitation learning. In <i>Proceedings of the thirteenth international conference on artificial intelligence and statistics</i> , pages 661–668. JMLR Workshop and Conference Proceedings.	1191
		1192
		1193
		1194
		1195
		1196
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. <i>Advances in Neural Information Processing Systems</i> , 36:68539–68551.	1197
		1198
		1199
		1200
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	1201
		1202
		1203
		1204
		1205
		1206
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	1207
		1208
		1209
		1210
		1211
		1212
	Maohao Shen, Guangtao Zeng, Zhenting Qi, Zhang-Wei Hong, Zhenfang Chen, Wei Lu, Gregory Wornell, Subhro Das, David Cox, and Chuang Gan. 2025. Satori: Reinforcement learning with chain-of-action-thought enhances llm reasoning via autoregressive search . <i>Preprint</i> , arXiv:2502.02508.	1213
		1214
		1215
		1216
	Wentao Shi, Mengqi Yuan, Junkang Wu, Qifan Wang, and Fuli Feng. 2025. Direct multi-turn preference optimization for language agents . <i>Preprint</i> , arXiv:2406.14868.	1217
		1218
		1219
		1220
		1221
	Vaishnavi Shrivastava, Ahmed Awadallah, Vidhisha Balachandran, Shivam Garg, Harkirat Behl, and Dimitris Papailiopoulos. 2025. Sample more to think less: Group filtered policy optimization for concise reasoning. <i>arXiv preprint arXiv:2508.09726</i> .	

1222	Yifan Song, Weimin Xiong, Xiutian Zhao, Dawei Zhu, Wenhao Wu, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024a. Agentbank: Towards generalized llm agents via fine-tuning on 50000+ interaction trajectories . <i>Preprint</i> , arXiv:2410.07706.	1279
1223		1280
1224		1281
1225		1282
1226		1283
1227	Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024b. Trial and error: Exploration-based trajectory optimization for llm agents . <i>Preprint</i> , arXiv:2403.02502.	1284
1228		1285
1229		1286
1230		1287
1231	Moritz Stephan, Alexander Khazatsky, Eric Mitchell, Annie S Chen, Sheryl Hsu, Archit Sharma, and Chelsea Finn. 2024. Rlvf: Learning from verbal feedback without overgeneralization . <i>Preprint</i> , arXiv:2402.10893.	1288
1232		1289
1233		1290
1234		1291
1235		1292
1236	Xin Sun, Jianan Xie, Zhongqi Chen, Qiang Liu, Shu Wu, Yuehe Chen, Bowen Song, Weiqiang Wang, Zilei Wang, and Liang Wang. 2025. Divide-then-align: Honest alignment based on the knowledge boundary of rag . <i>Preprint</i> , arXiv:2505.20871.	1293
1237		1294
1238		1295
1239		1296
1240		1296
1241	Quanwei Tang, Sophia Yat Mei Lee, Junshuang Wu, Dong Zhang, Shoushan Li, Erik Cambria, and Guodong Zhou. 2025. A comprehensive graph framework for question answering with mode-seeking preference alignment . <i>Preprint</i> , arXiv:2506.17951.	1297
1242		1298
1243		1299
1244		1300
1245		1301
1246	Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. A survey on self-evolution of large language models. <i>arXiv preprint arXiv:2404.14387</i> .	1302
1247		1303
1248		1304
1249		1305
1250		1306
1251	Guiyao Tie, Zeli Zhao, Dingjie Song, Fuyang Wei, Rong Zhou, Yurou Dai, Wen Yin, Zhejiang Yang, Jianguye Yan, Yao Su, and 1 others. 2025. A survey on post-training of large language models. <i>arXiv e-prints</i> , pages arXiv–2503.	1307
1252		1308
1253		1309
1254		1310
1255		1310
1256	Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisanin, Alexan Ayrapetyan, and Igor Gitman. 2024a. Openmathinstruct-2: Accelerating ai for math with massive open-source instruction data . <i>Preprint</i> , arXiv:2410.01560.	1311
1257		1312
1258		1313
1259		1314
1260		1315
1261	Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. 2024b. Openmathinstruct-1: A 1.8 million math instruction tuning dataset . <i>Advances in Neural Information Processing Systems</i> , 37:34737–34774.	1316
1262		1317
1263		1318
1264		1319
1265		1320
1266	Dheeraj Vattikonda, Santhoshi Ravichandran, Emiliano Penalzoza, Hadi Nekoei, Megh Thakkar, Thibault Le Sellier de Chezelles, Nicolas Gontier, Miguel Muñoz-Mármol, Sahar Omidi Shayegan, Stefania Raimondo, Xue Liu, Alexandre Drouin, Laurent Charlin, Alexandre Piché, Alexandre Lacoste, and Massimo Caccia. 2025. How to train your llm web agent: A statistical diagnosis . <i>Preprint</i> , arXiv:2507.04103.	1321
1267		1321
1268		1322
1269		1323
1270		1324
1271		1325
1272		1325
1273		1322
1274		1323
1275	Huaijie Wang, Shibo Hao, Hanze Dong, Shenao Zhang, Yilin Bao, Ziran Yang, and Yi Wu. 2024a. Offline reinforcement learning for llm multi-step reasoning . <i>Preprint</i> , arXiv:2412.16145.	1324
1276		1325
1277		1322
1278		1323
	Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Mathcoder: Seamless code integration in llms for enhanced mathematical reasoning . <i>Preprint</i> , arXiv:2310.03731.	1324
		1325
	Kuan Wang, Yadong Lu, Michael Santacroce, Yeyun Gong, Chao Zhang, and Yelong Shen. 2024b. Adapting llm agents with universal feedback in communication . <i>Preprint</i> , arXiv:2310.01444.	1284
		1285
		1286
		1287
	Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. 2024c. Math-shepherd: Verify and reinforce llms step-by-step without human annotations . <i>Preprint</i> , arXiv:2312.08935.	1288
		1289
		1290
		1291
		1292
	Zhenting Wang, Guofeng Cui, Yu-Jhe Li, Kun Wan, and Wentian Zhao. 2025a. Dump: Automated distribution-level curriculum learning for rl-based llm post-training . <i>arXiv preprint arXiv:2504.09710</i> .	1293
		1294
		1295
		1296
	Zhichao Wang, Bin Bi, Zixu Zhu, Xiangbo Mao, Jun Wang, and Shiyu Wang. 2024d. Uft: Unifying fine-tuning of sft and rlhf/dpo/una through a generalized implicit reward function . <i>arXiv preprint arXiv:2410.21438</i> .	1297
		1298
		1299
		1300
		1301
	Zilong Wang, Jingfeng Yang, Sreyashi Nag, Samarth Varshney, Xianfeng Tang, Haoming Jiang, Jingbo Shang, and Sheikh Muhammad Sarwar. 2025b. Rro: Llm agent optimization through rising reward trajectories . <i>Preprint</i> , arXiv:2505.20737.	1302
		1303
		1304
		1305
		1306
	Hui Wei, Zihao Zhang, Shenghua He, Tian Xia, Shijia Pan, and Fei Liu. 2025a. Plangenllms: A modern survey of llm planning capabilities . <i>Preprint</i> , arXiv:2502.11221.	1307
		1308
		1309
		1310
	Quan Wei, Chung-Yiu Yau, Hoi-To Wai, Yang Katie Zhao, Dongyeop Kang, Youngsuk Park, and Mingyi Hong. 2025b. Roste: An efficient quantization-aware supervised fine-tuning approach for large language models . <i>Preprint</i> , arXiv:2502.09003.	1311
		1312
		1313
		1314
		1315
	Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I. Wang. 2025c. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution . <i>Preprint</i> , arXiv:2502.18449.	1316
		1317
		1318
		1319
		1320
		1321
	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2024. Magicoder: Empowering code generation with oss-instruct . <i>Preprint</i> , arXiv:2312.02120.	1322
		1323
		1324
		1325
	Yongliang Wu, Yizhou Zhou, Zhou Ziheng, Yingzhe Peng, Xinyu Ye, Xinting Hu, Wenbo Zhu, Lu Qi, Ming-Hsuan Yang, and Xu Yang. 2025. On the generalization of sft: A reinforcement learning perspective with reward rectification . <i>arXiv preprint arXiv:2508.05629</i> .	1326
		1327
		1328
		1329
		1330
		1331

1332	Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. 2024. Agentgym: Evolving large language model-based agents across diverse environments . <i>Preprint</i> , arXiv:2406.04151.	1385
1333		1386
1334		1387
1335		1388
1336		1389
1337		1390
1338		
1339		
1340	Yu Xia, Jingru Fan, Weize Chen, Siyu Yan, Xin Cong, Zhong Zhang, Yaxi Lu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2025a. Agentrm: Enhancing agent generalization with reward modeling . <i>Preprint</i> , arXiv:2502.18407.	1391
1341		1392
1342		1393
1343		1394
1344		
1345	Yu Xia, Yiran Shen, Junda Wu, Tong Yu, Sungchul Kim, Ryan A. Rossi, Lina Yao, and Julian McAuley. 2025b. Sand: Boosting llm agents with self-taught action deliberation . <i>Preprint</i> , arXiv:2507.07441.	1395
1346		1396
1347		1397
1348		1398
1349	Qiming Xie, Zengzhi Wang, Yi Feng, and Rui Xia. 2024. Ask again, then fail: Large language models' vacillations in judgment . <i>Preprint</i> , arXiv:2310.02174.	1399
1350		1400
1351		1401
1352	Zhihui Xie, Jie chen, Liyu Chen, Weichao Mao, Jingjing Xu, and Lingpeng Kong. 2025. Teaching language models to critique via reinforcement learning . <i>Preprint</i> , arXiv:2502.03492.	1402
1353		1403
1354		1404
1355		1405
1356	Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, and 1 others. 2025a. A minimalist approach to llm reasoning: from rejection sampling to reinforce. <i>arXiv preprint arXiv:2504.11343</i> .	1406
1357		1407
1358		1408
1359		1409
1360		1410
1361	Weimin Xiong, Yifan Song, Qingxiu Dong, Bingchan Zhao, Feifan Song, Xun Wang, and Sujian Li. 2025b. Mpo: Boosting llm agents with meta plan optimization . <i>Preprint</i> , arXiv:2503.02682.	1411
1362		1412
1363		1413
1364		1414
1365	Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch every step! llm agent learning via iterative step-level process refinement . <i>Preprint</i> , arXiv:2406.11176.	1415
1366		1416
1367		1417
1368		1418
1369		1419
1370	Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai Fan, Lu Chen, and Kai Yu. 2024a. Rejection improves reliability: Training llms to refuse unknown questions using rl from knowledge feedback . <i>Preprint</i> , arXiv:2403.18349.	1420
1371		1421
1372		1422
1373		1423
1374		1424
1375	Ran Xu, Wenqi Shi, Yuchen Zhuang, Yue Yu, Joyce C. Ho, Haoyu Wang, and Carl Yang. 2025a. Collab-rag: Boosting retrieval-augmented generation for complex question answering via white-box and black-box llm collaboration . <i>Preprint</i> , arXiv:2504.04915.	1425
1376		1426
1377		1427
1378		1428
1379		1429
1380	Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. 2024b. Sayself: Teaching llms to express confidence with self-reflective rationales . <i>Preprint</i> , arXiv:2405.20974.	1430
1381		1431
1382		1432
1383		1433
1384		1434
	Yifan Xu, Xiao Liu, Xinghan Liu, Zhenyu Hou, Yueyan Li, Xiaohan Zhang, Zihan Wang, Aohan Zeng, Zhengxiao Du, Wenyi Zhao, Jie Tang, and Yuxiao Dong. 2024c. Chatglm-math: Improving math problem-solving in large language models with a self-critique pipeline . <i>Preprint</i> , arXiv:2404.02893.	1435
	Yixuan Even Xu, Yash Savani, Fei Fang, and J. Zico Kolter. 2025b. Not all rollouts are useful: Down-sampling rollouts in llm reinforcement learning . <i>Preprint</i> , arXiv:2504.13818.	1436
	Yixuan Even Xu, Yash Savani, Fei Fang, and Zico Kolter. 2025c. Not all rollouts are useful: Down-sampling rollouts in llm reinforcement learning. <i>arXiv preprint arXiv:2504.13818</i> .	1437
	Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. 2025a. Learning to reason under off-policy guidance. <i>arXiv preprint arXiv:2504.14945</i> .	1438
	Shi-Qi Yan and Zhen-Hua Ling. 2025. Rpo: Retrieval preference optimization for robust retrieval-augmented generation . <i>Preprint</i> , arXiv:2501.13726.	1439
	Yibo Yan, Jiamin Su, Jianxiang He, Fangteng Fu, Xu Zheng, Yuanhuiyi Lyu, Kun Wang, Shen Wang, Qingsong Wen, and Xuming Hu. 2025b. A survey of mathematical reasoning in the era of multimodal large language model: Benchmark, method & challenges . <i>Preprint</i> , arXiv:2412.11936.	1440
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	1441
	Guang Yang, Yu Zhou, Xiang Chen, Xiangyu Zhang, Terry Yue Zhuo, and Taolue Chen. 2024. Chain-of-thought in neural code generation: From and for lightweight language models . <i>Preprint</i> , arXiv:2312.05562.	1442
	Jian Yang, Wei Zhang, Jiayi Yang, Yibo Miao, Shanghaoran Quan, Zhenhe Wu, Qiyao Peng, Liqun Yang, Tianyu Liu, Zeyu Cui, Binyuan Hui, and Junyang Lin. 2025b. Multi-agent collaboration for multilingual code instruction tuning . <i>Preprint</i> , arXiv:2502.07487.	1443
	Zhengdong Yang, Shuichiro Shimizu, Yahan Yu, and Chenhui Chu. 2025c. When large language models meet speech: A survey on integration approaches . <i>Preprint</i> , arXiv:2502.19548.	1444
	Jiarui Yao, Yifan Hao, Hanning Zhang, Hanze Dong, Wei Xiong, Nan Jiang, and Tong Zhang. 2025. Optimizing chain-of-thought reasoners via gradient variance minimization in rejection sampling and rl. <i>arXiv preprint arXiv:2505.02391</i> .	1445
	Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran	1446

1439	Xu, Phil Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. 2024. Retroformer: Retrospective large language agents with policy gradient optimization . <i>Preprint</i> , arXiv:2308.02151.	1495
1440		1496
1441		1497
1442		1498
1443	Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2024. Agent lumos: Unified and modular training for open-source language agents . <i>Preprint</i> , arXiv:2311.05657.	1499
1444		1500
1445		1501
1446		1502
1447		1503
1448	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguang Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models . <i>arXiv preprint arXiv:2309.12284</i> .	1504
1449		1505
1450		1506
1451		1507
1452		1508
1453		
1454	Siyu Yuan, Zehui Chen, Zhiheng Xi, Junjie Ye, Zhengyin Du, and Jiecao Chen. 2025. Agent-r: Training language model agents to reflect via iterative self-training . <i>Preprint</i> , arXiv:2501.11425.	1509
1455		1510
1456		1511
1457		1512
1458		1513
1459	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning: Enabling generalized agent abilities for llms . <i>Preprint</i> , arXiv:2310.12823.	1514
1460		1515
1461		1516
1462		1517
1463	Chuheng Zhang, Wei Shen, Li Zhao, Xuyun Zhang, Lianyong Qi, Wanchun Dou, and Jiang Bian. 2024a. Policy filtration in rlhf to fine-tune llm for code generation .	1518
1464		1519
1465		1520
1466		1521
1467		1522
1468	Fan Zhang, Hao Chen, Zhihong Zhu, Ziheng Zhang, Zhenxi Lin, Ziyue Qiao, Yefeng Zheng, and Xian Wu. 2025a. A survey on foundation language models for single-cell biology . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 528–549, Vienna, Austria. Association for Computational Linguistics.	1523
1469		1524
1470		1525
1471		1526
1472		1527
1473		
1474	Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Ming Zhu, Juntao Tan, Thai Hoang, Zuxin Liu, Liangwei Yang, Yihao Feng, Shirley Kokane, Tulika Awalgaonkar, Juan Carlos Niebles, Silvio Savarese, Shelby Heinecke, Huan Wang, and Caiming Xiong. 2024b. Agentohana: Design unified data and training pipeline for effective agent learning . <i>Preprint</i> , arXiv:2402.15506.	1528
1475		1529
1476		1530
1477		
1478		1531
1479		1532
1480		1533
1481		1534
1482	Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu, Juan Carlos Niebles, Shelby Heinecke, and 3 others. 2024c. xlam: A family of large action models to empower ai agent systems . <i>Preprint</i> , arXiv:2409.03215.	1535
1483		1536
1484		1537
1485		1538
1486		
1487		1539
1488		1540
1489		1541
1490		1542
1491		1543
1492		1544
1493	Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, and 1 others. 2025b. A survey of reinforcement learning for large reasoning models . <i>arXiv preprint arXiv:2509.08827</i> .	1545
1494		1546
		1547
	Shijie Zhang, Guohao Sun, Kevin Zhang, Xiang Guo, and Rujun Guo. 2025c. Clpo: Curriculum learning meets policy optimization for llm reasoning . <i>arXiv preprint arXiv:2509.25004</i> .	1548
		1549
		1500
		1501
		1502
		1503
		1504
		1505
		1506
		1507
		1508
		1509
		1510
		1511
		1512
		1513
		1514
		1515
		1516
		1517
		1518
		1519
		1520
		1521
		1522
		1523
		1524
		1525
		1526
		1527
		1528
		1529
		1530
		1531
		1532
		1533
		1534
		1535
		1536
		1537
		1538
		1539
		1540
		1541
		1542
		1543
		1544
		1545
		1546
		1547

1548 Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen,
1549 Danqi Chen, and Yu Meng. 2025b. The surprising
1550 effectiveness of negative reinforcement in llm reason-
1551 ing. *arXiv preprint arXiv:2506.01347*.

1552 Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng,
1553 Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, Hua-
1554 jun Chen, and Ningyu Zhang. 2025c. [Knowa-
1555 gent: Knowledge-augmented planning for llm-based
1556 agents](#). *Preprint*, arXiv:2403.03101.

1557 Zhihong Zhu, Yunyan Zhang, Xianwei Zhuang, Fan
1558 Zhang, Zhongwei Wan, Yuyan Chen, Qingqing Long,
1559 Yefeng Zheng, and Xian Wu. 2025d. [Can we trust AI
1560 doctors? a survey of medical hallucination in large
1561 language and large vision-language models](#). In *Find-
1562 ings of the Association for Computational Linguistics:
1563 ACL 2025*, pages 6748–6769, Vienna, Austria.
1564 Association for Computational Linguistics.

Table 1: Objective comparison across training paradigms. Here for integration, (SFT \rightarrow) RL means modifying RL objective based on SFT objective, and vice versa.

Method	Paper	Objective / Gradient Expression	Data	Integration
DFT	Wu et al. (2025)	$\mathbb{E}_{(x_i, y_i) \sim \mathcal{D}} [-\text{sg}(\pi_\theta(y_i x_i)) \nabla_\theta \log \pi_\theta(y_i x_i)]$	Offline dataset	(RL \rightarrow) SFT
Prefix Sampling	Huang et al. (2025)	$-\sum_{y_n^1, \dots, y_n^{n-1}, y_n^L} \alpha \nabla_\theta \log \pi_\theta(y_n^t x, y_n^{<t})$ $-\sum_{y_n^L} \beta \nabla_\theta \log \pi_\theta(y_n^t x, y_n^{<t})$	Offline dataset + online rollout	RL + SFT
CHORD	Zhang et al. (2025d)	$\mathcal{L}_{\text{CHORD}}(\theta) = (1 - \mu) \mathcal{L}_{\text{GRPO}}(\theta) + \mu \mathcal{L}_{\text{SFT}-\phi}(\theta)$, where $\mathcal{L}_{\text{SFT}-\phi} = -\mathbb{E}_{(x, y^*) \sim \mathcal{D}_{\text{SFT}}} [\sum_{t=1}^{ y^* } \phi(y_t^*; \pi_\theta) \cdot \log \pi_\theta(y_t^* x, y_{<t}^*)]$; and $\phi(y_t^*; \pi_\theta) = \pi_\theta(y_t^* x, y_{<t}^*) (1 - \pi_\theta(y_t^* x, y_{<t}^*))$ $\mathcal{J}_{\text{UFT}} = \mathbb{E}_{\substack{t, x, l = x_t^* \\ (x_h, y_h)_{h=1}^{H-1} \sim \pi}} [\mathcal{J}_{\text{value}}(x_h, y_h)_{h=l}^{H-1}]$ $-\beta \sum_{h=l}^{H-1} \text{KL}(\pi(\cdot x_h) \ \pi^{\text{ref}}(\cdot x_h)) + \beta \sum_{h=0}^{l-1} \log \pi(x_h^* x_h^*)$	Offline dataset + online rollout	RL + SFT
UFT	Liu et al. (2025b)		Offline dataset + online rollout	RL + SFT
Proximal SFT	Zhu et al. (2025a)	$\mathcal{L}^{\text{PSFT}}(\theta) = \mathbb{E}_{(x_t, y_t) \sim \mathcal{D}} \left[\min \left(\frac{\pi_\theta(y_t x_t)}{\pi_{\theta_{\text{old}}}(y_t x_t)}, \text{clip} \left(\frac{\pi_\theta(y_t x_t)}{\pi_{\theta_{\text{old}}}(y_t x_t)}, 1 - \epsilon, 1 + \epsilon \right) \right) \right]$	Offline dataset	(RL \rightarrow) SFT
UFT-SFT	Wang et al. (2024d)	$\mathcal{L}_{\text{UFT-SFT}}(\pi_\theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} \left\{ \left[\sigma \left(\beta \log \frac{\pi_\theta(y x)}{\pi_{\text{ref}}(y x)} \right) - 1 \right]^2 \right\}$	Offline dataset	(RL \rightarrow) SFT
iw-SFT	Qin and Springenberg (2025)	$\mathcal{L}_{\text{iw-SFT}} = \mathbb{E}_{\tau \in \mathcal{D}} \left[-\frac{g(\tau)}{\pi_{\text{ref}}(\tau)} \log p(\tau; \theta) \right]$	Offline dataset	(RL \rightarrow) SFT
VAR	Du et al. (2025b)	$\mathcal{L} = -\mathbb{E} \left[\frac{\pi_{\text{ref}}(y x) \exp(\frac{1}{2} r(x, y))}{Z(x)} \log \pi_\theta(y x) \right]$	Offline dataset	SFT
HPT	Lv et al. (2025)	$\mathcal{L} = \alpha \mathcal{L}_{\text{RL}} + \beta \mathcal{L}_{\text{SFT}}$	Offline dataset + online rollout	RL + SFT
RPO	Liu et al. (2024)	$\mathcal{L}_{\text{RPO}}(\theta) = \beta \beta \cdot \mathbb{E}_{\tilde{x} \sim d_0, y^0 \sim \pi_{\text{base}}(\cdot x)} [-\log \pi_\theta(y^0 x)] + \mathcal{L}_{\mathcal{D}}(\beta \cdot \log \frac{\pi_\theta}{\pi_{\text{ref}}})$, where $\mathcal{L}_{\mathcal{D}}(r) = -\mathbb{E}_{\mathcal{D}} [r \log(\sigma(r(x_i, y_i^0)) - r(x_i, y_i^0)) + (1 - p_i) \log(\sigma(r(x_i, y_i^0)) - r(x_i, y_i^0))]$	Offline (preference) dataset	SFT + DPO
SRFT	Fu et al. (2025b)	$\mathcal{L}_{\text{SRFT}}(\theta) = \mathcal{L}_{\text{SFT}}^{\text{demo}}(\theta) + \mathcal{L}_{\text{RL}}^{\text{demo}}(\theta) + \mathcal{L}_{\text{RL}}^{\text{self-rollout}}(\theta)$, where $\mathcal{L}_{\text{RL}}^{\text{self-rollout}}$ stands for SFT loss on new rollouts from policy model, and GRPO advantage estimation is used in $\mathcal{L}_{\text{RL}}^{\text{demo}}$ with concatenated group from offline dataset and online rollout	Offline dataset + online rollout	SFT + RL
Q-SFT	Hong et al. (2024)	$\mathcal{L}_{\text{CE}}(\phi) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [\log \pi_\phi(y x)]$ $\mathcal{L}_{\text{TD}}(\theta) = \mathbb{E}_{(x, y, r, x') \sim \mathcal{D}} \left[(r + \gamma \max_{y'} Q_\theta(x', y') - Q_\theta(x, y))^2 \right]$	Offline dataset	SFT + Q learning
IRL-RFT	Li et al. (2024a)	$\frac{1}{\beta} \nabla_\theta r(x_t, k; y_t, k; \theta_t, k) - \frac{1}{\beta} \nabla_\theta r(x_t, k; \tilde{y}_t, k; \theta_t, k)$, where $y \sim \pi_{\text{expert}}(\cdot x_t, k)$, $\tilde{y} \sim \pi^t(\cdot x_t, k)$, $\pi^t(y x) \propto \exp(r(x, y; \theta_t, \kappa))$	Online rollout	RL
SASR	Chen et al. (2025b)	$\mathcal{L}(\theta) = \frac{1}{2} \sum_{s=1}^S [(1 - I(t)) \cdot \mathcal{L}_{\text{SFT}}(\theta) + I(t) \cdot \mathcal{L}_{\text{GRPO}}(\theta)]$	Offline dataset + online rollout	SFT + RL
IFT	Hua et al. (2024)	$\mathcal{L}(\hat{T}_\theta(s_n^*, \rho_0)) = -\sum_{t=n}^N \log \mathcal{T}_\theta(a_t^*, \delta_\theta(s_t^*))$	Offline dataset	RL
NFT	Chen et al. (2025a)	$\mathcal{L}_{\text{D}}^{\text{NFT}}(\theta) = -\sum_{x, y, r} \omega(x) \sum_t \left[r \log R_\theta^t(x, y) + (1 - r) \log \max_v \left(\frac{1 - \hat{r}_x R_\theta^t(x, y)}{1 - \hat{r}_x} \right) \right]$, where $R_\theta^t(x, y) = \frac{\pi(y_t x, y_{<t})}{\pi(y_t x, y_{<t})}$, $\hat{r}_x = \frac{1}{K} \sum_{y x} r(x, y)$	Online rollout	(SFT \rightarrow) RL
SRL	Deng et al. (2025)	$R = \frac{2 \sum_{(i, j, n) \in \text{MatchingBlocks}^n}{r(\mathbf{y}'_{\text{step } p_k}, \mathbf{y}_{\text{step } p_k})}}{ \mathcal{S}_1 + \mathcal{S}_2 }$, $r(\mathbf{y}'_{\text{step } p_k}, \mathbf{y}_{\text{step } p_k}) = \begin{cases} R(\mathbf{y}'_{\text{step } p_k}, \mathbf{y}_{\text{step } p_k}) & \text{if } \mathbf{y}' \text{ follows format,} \\ -1 & \text{otherwise.} \end{cases}$	Offline dataset + online rollout	SFT + RL

A More on Comparing and Combining SFT and RL

We summarize the main notation used throughout this section and Table 1. For ease of reference, commonly used acronyms are listed in the acronym table, and a consolidated summary of symbols is provided in Table 2.

A.1 Notation

The offline training dataset used in SFT and related settings is denoted by $\mathcal{D} = \{(x_i, y_i)\}$, where x_i represents the input prompt and y_i denotes the corresponding model-generated response. We use π_θ to denote a policy (language) model parameterized by θ , and $\pi_\theta(y_i | x_i)$ to denote the conditional probability of generating response y_i given prompt x_i . The reference policy is denoted by π^{ref} , which is typically a frozen base model or an SFT checkpoint.

In the RL setting, $r(x_i, y_i)$ denotes the reward assigned to a prompt-response pair (x_i, y_i) , and $Q_\theta(x, y)$ represents the action-value (Q) function under policy π_θ , with state x (prompt) and action y (response). We use \mathcal{J} to denote cumulative reward-based objectives.

The sigmoid function is defined as $\sigma(x) = 1/(1 + e^{-x})$. The operator $\text{sg}(\cdot)$ denotes the stop-gradient operation, which preserves the forward-pass value while preventing gradient propagation during backpropagation. Additional symbols, operators, and hyperparameters appearing in Table 1 are summarized in Table 2 for completeness.

A.2 Preliminaries

We consider the optimization of Pretrained LLMs, also referred to as FMs, which are large-scale pre-trained models adapted to downstream tasks such as QA, reasoning with CoT, and database augmentation tasks such as RAG. Training typically proceeds in multiple stages, starting with SFT on curated prompt-response pairs, which can be viewed as a form of BC from expert demonstrations.

Beyond SFT, model alignment and performance are often further improved using RL techniques, particularly RLHF, where a learned or human-provided reward signal guides optimization. Popular policy-gradient-based algorithms include PPO and its variants such as GRPO, which stabilize updates through clipping or relative normalization. Alternative optimization frameworks include preference-based methods such as DPO

and rejection-based methods such as RFT, which bypass explicit reward modeling. Information-theoretic objectives, such as GEM, and regularization via the KL divergence are commonly used to control deviation from a reference policy.

The RL process itself can be formalized as a MDP, defined by the tuple $(\mathcal{S}, \mathcal{A}, P, r, H)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, P denotes the transition dynamics, r is the reward function, and H is the horizon length, which may be infinite under idealized analysis. The policy $\pi_\theta(\cdot | x)$ represents an LLM that induces a distribution over responses conditioned on an input prompt x . This formulation naturally extends to multi-modal settings, including VLM, and other structured decision-making scenarios.

Acronyms from the Main Paper

BC Behavior Cloning.	3	1631
CoT Chain-of-Thought.	1, 6, 7	1632
DPO Direct Preference Optimization.	5, 7	1633
FM Foundation Model.	6, 7	1634
GEM Entropic Distribution Matching.	3	1635
GRPO Group Relative Policy Optimization.	3, 5, 7	1636
KL Kullback–Leibler.	4	1638
LLM Large Language Model.	1–8	1639
PPO Proximal Policy Optimization.	3, 6, 7	1640
QA Question Answering.	6	1641
RAG Retrieval-Augmented Generation.	6	1642
RFT Rejection Fine-Tuning.	6–8	1643
RL Reinforcement Learning.	1–8	1644
RLHF Reinforcement Learning from Human Feedback.	4, 5	1645
SFT Supervised Fine-Tuning.	1–8	1647
VLMs Vision Language Models.	5	1648

Symbol	Description
<i>Models and Policies</i>	
π_θ, π_ϕ	Current policy/model parameterized by θ or ϕ
$\pi_{\text{ref}}, \pi_{\text{base}}$	Reference or base policy (typically frozen)
π_{old}	Policy from a previous optimization step
π_{expert}	Expert policy used for imitation or IRL
<i>Data and Sequences</i>	
x, y	Input prompt and generated response
$y_t, y_{<t}$	Token at time t and the prefix sequence
$\mathcal{D}, \mathcal{D}^+$	Dataset (offline, SFT, or preference-based)
τ	Trajectory or full sequence (x, y)
<i>Functions and Objectives</i>	
\mathcal{L}, \mathcal{J}	Loss or objective function
$r(x, y), R$	Reward function or scalar reward value
$Q(x, y)$	Action-value function (Q-function)
$Z(x)$	Partition function (normalization constant)
$\sigma(\cdot)$	Sigmoid activation function
<i>Operators and Hyperparameters</i>	
\mathbb{E}	Mathematical expectation
∇_θ	Gradient with respect to parameters θ
$\text{sg}(\cdot)$	Stop-gradient operator
$\text{KL}(\cdot \parallel \cdot)$	Kullback-Leibler divergence
α, β, μ, η	Weighting coefficients or temperature parameters
ϵ	Clipping or threshold constant

Table 2: Summary of Notation and Symbols

B Application Survey Details

This appendix section provides detailed information on our survey of LLM applications across professional domains, including the datasets, evaluation criteria, and publication trends that underpin our analysis. The goal is to contextualize the results presented in the main text by describing (i) how different task domains vary in input complexity, reasoning demands, and output characteristics, (ii) our methodology for categorizing research publications by domain using benchmark-oriented searches, and (iii) observed trends in publication volume, methodological approaches, and model us-

age over time. By presenting these details, we aim to offer a transparent and reproducible account of the survey process, enabling deeper insight into the evolving landscape of LLM applications.

B.1 Applications Across Professional Domains

To characterize the difficulty of each domain, we consider three qualitative factors: input complexity, intermediate reasoning requirements (chains of thought, CoT), and output complexity. Table 3 summarizes these characteristics across focompur representative domains. Note that Section 5 focuses on papers addressing text-only tasks. We make this choice because text-only tasks enable clearer side-by-side comparison, whereas multimodal settings introduce additional sources of uncertainty in model selection and fusion techniques, which lie outside the scope of this work. Nevertheless, we acknowledge that numerous impactful works closely related to the surveyed area extend beyond text-only modalities. For further references, we direct readers to Section D, which highlights some recently published surveys as a starting point.

General QA assesses fact-checking and language understanding with variable-length inputs and moderate reasoning, prioritizing accuracy over verbosity. Mathematical tasks demand precise answers supported by extended logical reasoning. Agentic tasks require rich environmental context to enable multi-step planning before producing compact outputs. Code generation relies on context-heavy inputs and moderate reasoning to infer intent and produce syntactically correct, semantically coherent code.

Across these domains, the overarching question is how to optimize performance given varying input scales, reasoning depths, and output complexities. The relative weight of these factors differs by task, reflecting the diverse requirements and shifting priorities of professional LLM applications.

B.2 Benchmark-Oriented Paper Search

To estimate publication counts for each domain, we conducted a benchmark-oriented keyword search over arXiv preprints in the Computer Science categories from 1 January 2023 to 30 June 2025. The arXiv metadata dataset (Cornell-University/arxiv) is provided under the Creative Commons CC0 1.0 Universal Public Domain Dedication, which applies to the dataset’s metadata as hosted on Kaggle. Original arXiv content remains subject to the licenses specified by individual authors.

Table 3: Summary of LLM application domains and their characteristics.

Domain	Input	CoT / Reasoning	Output
General QA	Moderate	Moderate	Moderate
Mathematical	Moderate	Long	Moderate
Agentic	Long	Long(Multi-Step)	Moderate
Code Gen.	Long	Moderate	Long

We accessed scientific articles via Google Cloud Public Datasets ([gs://arxiv-dataset](https://arxiv-dataset)). The initial dataset comprises approximately 195K papers, with 64K, 78K, and 52K papers from 2023, 2024, and 2025, respectively. Our approach is motivated by the observation that recent LLM post-training research consistently relies on a set of widely used, standardized datasets that facilitate reproducible and comparable evaluations across methods. As this practice has become increasingly common, dataset mentions serve as a reasonable proxy for assigning papers to topical domains.

Following Table 7 of [Du et al. \(2025a\)](#), we use 26 datasets spanning four domains as domain-specific query keys. A paper is assigned to a domain if it contains at least five mentions of any dataset associated with that domain. Mutual exclusivity is not enforced: although relatively rare, papers that evaluate across multiple domains are counted in every applicable category.

The datasets used as search keys are listed in Table 4, and the resulting paper counts for each category are shown in Table 5. We adopt a threshold of five keyword occurrences because a paper that genuinely employs a dataset typically introduces it, reports results on it, and provides comparative or analytical discussion, making five mentions a conservative filter against purely methodological contributions. To project the full-year counts for 2025, we note that the first six months of 2023 and 2024 account for 47.38% and 49.81% of their respective annual publication totals. Accordingly, we estimate the 2025 full-year count by doubling the observed number from the first half of 2025.

To contextualize our estimation procedure, we acknowledge both its strengths and limitations. Our keyword-based counting method provides a transparent, scalable, and reproducible way to approximate the volume of research activity across domains, with higher keyword-mention thresholds generally yielding more reliable domain assignments. As shown in Table 5, the number of papers

decreases consistently as the threshold increases, reflecting increasing specificity but reduced recall. We report results primarily using the threshold of keyword matches greater than five, though counts under alternative thresholds are also available. However, this approach also inherits several limitations: low thresholds may inflate counts by including papers that mention a dataset only in passing, while high thresholds may exclude legitimate work that uses a benchmark but references it sparsely. Additionally, cross-domain papers contribute to multiple categories, and naming variations may lead to undercounting unless normalized. Despite these caveats, the method remains a practical heuristic for capturing broad research trends in benchmark-driven LLM evaluation.

B.3 Detailed Paper Application Survey

B.3.1 LLMs for General QA Tasks

LLMs demonstrate strong performance in general question answering by enhancing reasoning capabilities, mitigating hallucinations, and effectively leveraging external knowledge. Current research emphasizes sequential sub-question reasoning, robust answer generation under uncertainty, and retrieval-augmented methods to improve information access. The following sections present representative approaches in each of these areas, with a focus on model-tuning strategies. For studies that do not involve model tuning, we refer readers to [Ke et al. \(2025\)](#).

Step-by-Step Reasoning. Most works decompose a question into sub-questions and train LLMs to answer them sequentially to produce more accurate CoTs. For instance, [Chen et al. \(2024\)](#) proposes generating sub-questions with GPT for each sentence in the ground-truth answer, while [Ranaldi and Freitas \(2024\)](#); [Huang et al. \(2024\)](#) use model rollouts to create preference data based on the quality of the final answers. Furthermore, [Feng et al. \(2024\)](#); [Zhou et al. \(2024d\)](#) align LLM behavior in deciding when to reflect, answer, or ask follow-up

Table 4: Datasets used as search keys for benchmark-oriented paper classification.

Agentic
webshop, webarena, wind2web, miniwob++, scienceworld, alfordworld, tdw-mat, c-wah, alfred, rlcad
QA
hotpotqa, strategyqa, triviaqa, pubmedqa, musique, 2wikimultihopqa, qasper
Math
gsm8k, asdiv, svamp, aime
Code
swe-bench, humaneval, livecodebench, bird, intercodesql

questions.

Hallucination Management. Several works investigate how LLMs behave under uncertainty, particularly regarding hallucination. For instance, [Li et al. \(2024b\)](#) treats GPT generations as positive samples for RLHF, whereas [Kang et al. \(2024\)](#) suggests using conservative RL (assigning low scores when the model is unfamiliar) to mitigate hallucination during RLHF. Other approaches ([Xie et al., 2024](#); [Cheng et al., 2024](#); [Xu et al., 2024a](#); [Sun et al., 2025](#)) propose alternative data pipelines for the RL stage, yielding more robust answers to challenging follow-ups and enabling refusal in cases of unknown scenarios. Moreover, [Band et al. \(2024\)](#); [Xu et al. \(2024b\)](#) equip the LLM with confidence estimation capabilities, typically synthesized from multiple rollouts and GPT-generated root causes.

LLM-based Retrieval. [Mao et al. \(2024\)](#); [Hsu et al. \(2024\)](#); [Yan and Ling \(2025\)](#); [Xu et al. \(2025a\)](#); [Gao et al. \(2025\)](#) post-train LLMs’ query generation with retrieval-based rewards, leveraging them as query rewriters to improve retrieval effectiveness. Other studies ([Jin et al., 2024](#); [Zhu et al., 2024](#)) propose alternative evaluation criteria for extracting key information from retrieved documents. Additionally, [Tang et al. \(2025\)](#) converts long documents into hierarchical graphs and applies RL to the policy model for RAG using mode-seeking DPO algorithms.

B.3.2 LLMs for Mathematical Tasks

Recent efforts have enhanced LLMs for mathematical reasoning to support multi-step logic, symbolic manipulation, and rigorous verification.

Math-Based Capabilities. [Polu and Sutskever \(2020\)](#); [Lewkowycz et al. \(2022\)](#); [Azerbayev et al. \(2023\)](#) fine-tuned LLMs on mathematics-specific

corporato enhance mathematical reasoning and proof generation. Building on this, [Luo et al. \(2023\)](#) applied the Evol-Instruct framework to math question–answer datasets. Similarly, [Yu et al. \(2023\)](#); [Toshniwal et al. \(2024a\)](#) improved Q&A pairs via rephrasing, backward reasoning, and other data augmentation strategies. In parallel, [Toshniwal et al. \(2024b\)](#) distilled reasoning abilities from proprietary models by collecting successful solution trajectories. Additionally, [Shao et al. \(2024\)](#) used FastText-based filtering in Common Crawl and employed GRPO to further enhance problem-solving performance.

Rollout Selection and Verification. [Cobbe et al. \(2021\)](#) proposed training separate verifier models to assess the correctness of model-generated solutions, selecting answers with the highest verification scores. [Huang et al. \(2022\)](#) fine-tuned LLMs on high-confidence data, while [Xu et al. \(2024c\)](#) filtered successful self-generated solutions for iterative training. [Wang et al. \(2024c\)](#) introduced MCTS to generate rollouts and annotate intermediate reasoning steps with probabilities of reaching correct answers. Similarly, [Guan et al. \(2025\)](#) learned a value function from MCTS probabilities to guide exploration and DPO data selection, with [Xu et al. \(2024c\)](#) further refining DPO via a math-specific critic model. Finally, [Shen et al. \(2025\)](#) fine-tuned models to select actions during mathematical reasoning.

B.3.3 LLMs for Agentic Tasks

Agentic tasks challenge LLMs to act as autonomous decision-makers in dynamic environments, requiring continuous state perception, contextual integration, and multi-step planning. Success hinges on long-term reasoning, history man-

Table 5: Paper counts under different keyword-mention thresholds across four domains.

Keyword count	QA			Math			Agentic			Code		
	23	24	25	23	24	25	23	24	25	23	24	25
>0	771	1644	1308	12673	20405	13177	1310	1826	1373	334	1091	1302
>1	580	1321	1088	3870	7143	4889	330	541	437	209	767	921
>2	453	1022	858	1673	3295	2684	162	310	260	178	621	717
>3	384	861	730	932	1957	1953	125	234	207	148	549	604
>4	338	746	621	655	1391	1587	108	196	167	130	481	513
>5	292	652	558	492	1098	1366	100	174	145	115	428	458
>10	165	342	326	221	538	794	68	122	92	69	254	252
>20	56	114	90	97	195	302	34	71	46	37	126	93

agement, and adaptive feedback, as each action influences future outcomes.

Autonomous Task Execution. Schick et al. (2023); Qin et al. (2023); Zeng et al. (2023); Yin et al. (2024); Zhang et al. (2024b,c); Song et al. (2024a); Ou et al. (2024) developed custom pipelines to behavior-clone full trajectories with open-weight models, supporting task-specific capabilities, whereas Yao et al. (2024); Murty et al. (2025) fine-tuned instruction-following models using actor-reflector outputs from proprietary APIs. To improve execution reliability, Zhong et al. (2024) used LLM judges to filter low-quality actions, Fu et al. (2025a) excluded erroneous steps during loss computation, and Xia et al. (2025b) added reflective annotations. Extending this, Zhu et al. (2025c) introduced knowledge-based self-learning with dynamic memory. In reinforcement learning, Choudhury and Sodhi (2024) curated preference pairs by synthesizing reasoning traces for ground-truth actions, while others generated preference data iteratively: ETO (Song et al., 2024b) used successful and failed trajectories for DPO, and Yuan et al. (2025) applied MCTS with corrections at the first error step.

Self-Improvement and Feedback Integration. Some studies use sequential SFT followed by RL. Zhou et al. (2024b) leveraged LLM API reflections, while Lai et al. (2024) showed early RL is effective after initial SFT using a 70B teacher model. Hierarchical RL approaches assign rewards at the sentence or utterance level while updating policies at the token level (Zhou et al., 2024c; Hu et al., 2025b), with DMPO (Shi et al., 2025) adding step-wise discounting and Qi et al. (2025) using critic-based curriculum learning. OREO (Wang et al., 2024a) enables multi-step reasoning via maximum-entropy learning, complemented by IRL and PPO tuning (Deng et al., 2024) and reward

calibration (Rita et al., 2024). Some works jointly train SFT and RL losses with variants including masked or weighted objectives, outcome- and step-based DPO, and pseudocode planning (Wang et al., 2024b; Xi et al., 2024; Xiong et al., 2024; Cao et al., 2025b).

Planning and Long-Horizon Control. Xiong et al. (2025b) proposed a hierarchical planning framework with a top-level planner fine-tuned on seed plans and optimized via DPO. Vattikonda et al. (2025) added Rejected Sampling Fine-Tuning (RFT) to further refine high-reward actions, while Wang et al. (2025b) adjusted exploration and prioritized critical steps for DPO. For test-time scaling, Zhou et al. (2024a) integrated general and agentic data (Zeng et al., 2023) with multi-path reasoning, and Xia et al. (2025a) combined policy fine-tuning with reward-model training, showing that beam search with explicit reward modeling improved performance.

B.3.4 LLMs for Code Generation Tasks

LLMs have shown remarkable potential in code generation, editing, and task-specific reasoning, which are often improved by structured code data, executable feedback, or multi-agent frameworks for more robust and verifiable outputs.

Code Generation Capabilities. Recent studies have explored methods to improve code generation with LLMs. Chaudhary (2023) employ self-instruction to synthesize training data, while Luo et al. (2025) generate harder queries by modifying constraints of seed questions. Wei et al. (2024) leverage open-source snippets to produce diverse and controllable instruction data. Yang et al. (2024) show that preceding CoT reasoning with code explanations enhances accuracy. More recently, Yang et al. (2025b) introduce a multi-agent framework in which an SFT model generates rollouts for DPO training.

Code Editing Capabilities. Code editing focuses on refining code that fails to meet its intended functionality, whether human-written or LLM-generated. [Wei et al. \(2025c\)](#) use human pull requests to guide GRPO-based specialization of LLMs. [Chae et al. \(2024\)](#); [Jiang et al. \(2025\)](#) generate refinement rollouts evaluated against unit tests. [Dou et al. \(2024\)](#) decompose code generation into verifiable subtasks for fine-grained alignment, while [Ma et al. \(2025\)](#) propose a four-stage pipeline combining rejection sampling with rule-based RL rewards. [Xie et al. \(2025\)](#) fine-tune an auxiliary LLM to produce verbal critiques that assist the primary model in correcting errors.

Code Reuse for Other Tasks. Code can also be leveraged to enhance a variety of tasks by serving as a verifiable and interpretable reward signal. For instance, [Dong et al. \(2024\)](#) use code execution to assess the instruction-following ability of LLMs. [Nie et al. \(2025\)](#) train an SLM to leverage LLMs via environment feedback by simulating workflows as executable Python code. [Wang et al. \(2023\)](#); [Lai et al. \(2023\)](#); [Gou et al. \(2024\)](#) perform supervised instruction tuning on code-augmented math or data science tasks.

B.4 Trend Analysis Details

This subsection covers a detailed analysis of the procedure-wise trend from surveyed papers, summarized in Figure 2.

Rapid Growth Across Domains. Between 2023 and 2025, research activity accelerates sharply across all major task domains, though the magnitude of growth varies substantially. QA studies more than double from 292 in 2023 to 652 in 2024 (+123%), with projections indicating continued expansion to 983 in 2025 (+118%). Math-related research grows even faster in absolute scale, rising from 492 to 1,098 (+123%) and then reaching 2,399 papers in 2025, which is a near 5× increase from 2023. Agent-focused work, while smaller in volume, shows steady expansion from 100 to 179 (+79%) and further to 261 (+46%). The most dramatic surge occurs in code-related studies, which climb from 115 in 2023 to 428 in 2024 (+272%) and nearly double again to 786 in 2025 (+84%). This trajectory reflects the rapid maturation of code-centric benchmarks, tools, and evaluation pipelines, and highlights coding as one of the most rapidly diversifying application domains for LLMs.

Convergence Toward Hybrid Training. The landscape of training methodologies undergoes a

marked consolidation toward hybrid approaches that combine SFT with RL or other post-training strategies. In 2023, SFT dominates at 73.3%, with hybrid (“Both”) methods representing only 20.0%. By 2024, hybrid training expands by 269%, becoming the most common approach at 73.8%, superseding pure SFT (19.1%) and growing further to 70.6% in 2025. This shift illustrates a collective movement toward more sophisticated multi-stage training pipelines that exploit complementary strengths of SFT and RL. Meanwhile, RL-only methods, which starts initially 6.7% of studies, gain modest traction, reaching 7.1% by 2024 and 11.8% by 2025, reflecting ongoing improvements in open-source RLHF frameworks and increased accessibility of preference data.

Shift from Proprietary to Open Models. A pronounced realignment in model choice accompanies these methodological changes. Between 2023 and 2025, reliance on proprietary API-based models declines sharply. The trend starts at 32.2% in 2023 to 19.9% in 2024, and down to 11.1% by 2025, where researchers pivot toward open-weight models and standardized benchmarks. Open-weight usage more than doubles, growing from 12.2% in 2023 to 17.5% in 2024 and reaching 25.0% in 2025. Benchmarks remain the most commonly used resource category but also increase steadily (48.9% to 61.1% over three years), underscoring the field’s movement toward reproducible, transparent experimentation. In parallel, the use of human-curated or web-scraped datasets continues to shrink (6.7% to 2.8%), reflecting both improved benchmark coverage and growing concerns around legality, license compliance, and data provenance. Overall, these trends indicate a strong and persistent shift toward open, standardized, and reproducible workflows as community norms solidify.

C Hardware Requirements for SFT/RL

For researchers intending to perform practical experimentation with post-training alignment, consulting community-reported hardware guidance can be valuable. It is important to note, however, that GPU requirements vary substantially depending on model size, training approaches, e.g., full fine-tuning versus parameter-efficient methods such as LoRA ([Hu et al., 2021](#)) or QLoRA ([Dettmers et al., 2023](#)), choice of optimizer, and software frameworks. Moreover, these requirements continue to evolve alongside advances in

2046 hardware and software tools.

2047 For SFT, a commonly cited heuristic indicates
2048 that full fine-tuning of a transformer-based large
2049 language model requires approximately 16GB of
2050 VRAM per billion parameters when using half-
2051 precision (FP16). This high memory demand arises
2052 from the need to store gradients and optimizer
2053 states, and is substantially greater than the roughly
2054 2GB per billion parameters typically needed for
2055 inference. Parameter-efficient approaches, such as
2056 LoRA and QLoRA, can reduce this requirement
2057 significantly, often into the range of 5–20GB on
2058 contemporary GPUs, depending on model size and
2059 precision settings¹.

2060 For RL training, frameworks such as *verl* provide
2061 community-reported guidance on hardware
2062 requirements across different model scales. Practical
2063 experimentation indicates that RL generally demands
2064 substantially more memory than SFT due to
2065 additional overhead from rollout generation, policy
2066 gradients, and reward models. Heuristic estimates
2067 suggest that medium-sized models (1–3B parameters)
2068 can be trained on a single 80GB GPU using
2069 parameter-efficient methods like LoRA, while
2070 larger models (7–14B) typically require 2–4 such
2071 GPUs. Very large models (32B+) often necessitate
2072 4–8 high-memory GPUs, with full-parameter
2073 training potentially exceeding 600GB of aggregate
2074 VRAM. Roughly, LoRA-based RL training consumes
2075 15–25GB VRAM per billion parameters, compared
2076 with 16GB/B for full SFT fine-tuning, illustrating
2077 that RL generally incurs 1.5–3× higher memory
2078 demands per parameter. These numbers should be
2079 treated as approximate guidelines, as actual
2080 requirements depend on batch size, rollout length,
2081 optimizer choice, and memory-saving strategies
2082 such as gradient checkpointing or quantization².

2084 It is essential to treat these hardware guidelines
2085 as approximate starting points rather than strict
2086 requirements. Actual resource needs will depend
2087 on factors including model architecture, task
2088 complexity, batch size, rollout length, and ongoing
2089 improvements in software and hardware technologies.

2090 D More References

2091 This section provides additional references that
2092 extend beyond the text-only scope emphasized in

¹<https://modal.com/blog/how-much-vram-need-fine-tuning>

²https://verl.readthedocs.io/en/latest/perf/device_tuning.html

2093 Section 5. These works highlight broader develop-
2094 ments in multimodal learning, retrieval-augmented
2095 methods, hallucination mitigation, and domain-
2096 specific datasets and capabilities. Together, they
2097 offer a more comprehensive view of the rapidly
2098 evolving ecosystem surrounding LLM reasoning
2099 and application domains.

2100 One of the most practically oriented directions is
2101 multimodal LLM research, which integrates modal-
2102 ities such as speech and vision beyond pure text.
2103 This line of work serves as a natural extension to
2104 the text-only tasks surveyed earlier, offering in-
2105 sights into how additional sensory channels intro-
2106 duce new modeling challenges and opportunities.
2107 For instance, [Cui et al. \(2025c\)](#) and [Yang et al. \(2025c\)](#)
2108 summarize methods for leveraging and
2109 aligning LLMs with speech-augmented multimodal
2110 data, enabling richer understanding and generation
2111 capabilities.

2112 Beyond multimodality, another cluster of re-
2113 lated literature addresses system-level improve-
2114 ments that refine how LLMs access information
2115 and maintain reliability. [Abootorabi et al. \(2025\)](#)
2116 survey advances in retrieval-augmented genera-
2117 tion (RAG) across various modalities, providing a
2118 complementary perspective to reasoning-centric ap-
2119 proaches. Meanwhile, [Li et al. \(2025a\)](#) investigates
2120 hallucination through the lens of knowledge bound-
2121 aries, and [Zhu et al. \(2025d\)](#) explores hallucination
2122 in conjunction with multimodal trustworthiness. In
2123 parallel, [He et al. \(2025c\)](#) review methodologies for
2124 constructing SFT and RL datasets specifically tai-
2125 lored for reasoning, bridging data collection prac-
2126 tices with model behavior.

2127 Finally, several specialized surveys focus on
2128 domain-specific datasets and complex task settings,
2129 which enrich the understanding of how LLM rea-
2130 soning varies across fields. For example, [Yan et al. \(2025b\)](#)
2131 discusses mathematical reasoning datasets
2132 and benchmarks that support structured and veri-
2133 fiable inference. Complementing this, [Wei et al. \(2025a\)](#)
2134 reviews LLM-based planning techniques
2135 applicable to both question answering and agentic
2136 tasks. Further extending the discussion to envi-
2137 ronments requiring perception–action loops, [Hu et al. \(2025a\)](#);
2138 [Nguyen et al. \(2025\)](#) examine mul-
2139 timodal–agent interactions, including GUI-based
2140 interfaces that illustrate how reasoning, percep-
2141 tion, and action co-evolve. Together, these surveys
2142 broaden the contextual understanding of LLM ca-
2143 pabilities beyond the domains highlighted in the
2144 main text.

E Use of AI Assistants

We employed LLM-based tools exclusively for language polishing during manuscript preparation and for assisting with full-text filtering of SFT/RL-related papers for human review in the Applications section (Section 5). All outputs generated by LLMs were carefully reviewed by the authors, who take full responsibility for any potential errors or hallucinations.

The prompt used with gpt-oss-120b to identify and classify SFT/RL-related papers is provided in the following verbatim. For all prompt requests to this model, we use FP4 precision, a context length of 128k, a temperature of 0.8, and set top-k and top-p to 40 and 0.95, respectively.

Human readers subsequently verified all candidate papers to confirm their publication status and the relevance of their training paradigms (i.e., whether they involve SFT, RL, or both).

You are a research paper analyzer.

From the given scientific text, extract and classify the following:

1. The **proposed method** (ignore ablations).
2. The **comparison methods/baselines** the paper evaluates against.
3. The **datasets/benchmarks** used for evaluation.
4. The **training type** of the proposed method, exactly one of:
 - Reinforcement Learning (RL)
 - Supervised Fine-Tuning (SFT)
 - Both RL and SFT (SFT+RL)
 - Prompt Optimization
 - none-text Modality (vision, speech, multimodal, etc)
 - Other Methods
 - Survey (no new method, only summarization of others)

Output format:

<thought>

YOUR THOUGHT

</thought>

<answer>

{

 "method": "...",

 "baselines": ["..."],

 "benchmarks": ["..."],

 "training_type": one of ["SFT", "RL", "SFT+RL", "Prompt", "Modality", "Other", "Survey"]

}

</answer>