

Reuse your FLOPs: Scaling RL on Hard Problems by Conditioning on Very Off-Policy Prefixes

Amrith Setlur^{1,2,*}, Zijian Wang¹, Andrew Cohen^{1,*}, Paria Rashidinejad^{1,3,*,†}, Sang Michael Xie^{1,*,†}

¹FAIR at Meta, ²Carnegie Mellon University, ³University of Southern California

*Work done at Meta, †Authors jointly advised this work.

Typical reinforcement learning (RL) methods for LLM reasoning waste compute on hard problems, where correct *on-policy* traces are rare, policy gradients vanish, and learning stalls. To bootstrap more efficient RL, we consider reusing old sampling FLOPs (from prior inference or RL training) in the form of *off-policy* traces. Standard off-policy methods supervise against off-policy data, causing instabilities during RL optimization. We introduce **PrefixRL**, where we *condition* on the prefix of successful off-policy traces and run on-policy RL to complete them, side-stepping off-policy instabilities. PrefixRL boosts the learning signal on hard problems by modulating the difficulty of the problem through the off-policy prefix length. We prove that the PrefixRL objective is not only consistent with the standard RL objective but also more sample efficient. Empirically, we discover **back-generalization**: training *only* on prefixed problems generalizes to *out-of-distribution* unprefixed performance, with learned strategies often differing from those in the prefix. In our experiments, we source the off-policy traces by rejection sampling with the base model, creating a self-improvement loop. On hard reasoning problems, PrefixRL reaches the same training reward 2× faster than the strongest baseline (SFT on off-policy data then RL), even after accounting for the compute spent on the initial rejection sampling, and increases the final reward by 3×. The gains transfer to held-out benchmarks, and PrefixRL is still effective when off-policy traces are derived from a different model family, validating its flexibility in practical settings.

Date: January 26, 2026

Correspondence: asetlur@cs.cmu.edu, zijianwang@meta.com



1 Introduction

Reinforcement learning (RL) is the de facto method to boost large language model (LLMs) reasoning, especially for math and coding (An et al., 2025; Liu et al., 2025b; Guo et al., 2025). Most successful RL recipes (Ahmadian et al., 2024; Yu et al., 2025) are *on-policy*: sample multiple reasoning traces (rollouts) from the current model and derive updates from correct (and incorrect) traces. This paradigm breaks down on *hard problems* with low pass@ k (e.g., pass@2k \approx 0), where the model rarely samples a correct trace. In this regime, learning stalls as RL spends enormous amounts of sampling FLOPs without receiving any learning signal, and RL rewards plateau.

In practice, we are rarely solving these problems for the first time since earlier RL runs or inference on previous models may have spent compute on the same (or similar) hard problems. The natural question now is how to reuse this ever-growing dataset of *off-policy* traces, which often contains some correct traces even for very hard problems, in order to guide the online RL policy towards higher-rewarding states and accelerate on-policy RL.

A straightforward approach is to treat the off-policy traces as supervision: perform supervised fine-tuning (a.k.a., mid-training or continued pretraining) on the correct off-policy traces followed by standard on-policy RL (Wang et al., 2025d). However, SFT on a small set of correct traces can lead to memorization (Chu et al., 2025) and entropy collapse, which hurts exploration during subsequent RL (Zhang et al., 2025a). Alternatively, we can use off-policy traces directly in RL via importance weighting, but this is often unstable due to high-variance gradient estimates (Liu et al., 2025a; Yan et al., 2025). Both options use off-policy traces as target supervision, and since these off-policy traces are very low probability under the RL policy, this leads to suboptimal RL optimization.

To avoid these pitfalls, we propose **PrefixRL**: run on-policy RL *conditioned* on prefixes of correct off-policy traces instead of directly supervising on them (Figure 1). *First*, we extract and fix a few off-policy prefixes and append them to the original problem to create *prefixed problems*. *Second*, we run on-policy RL on both *no-prefix*

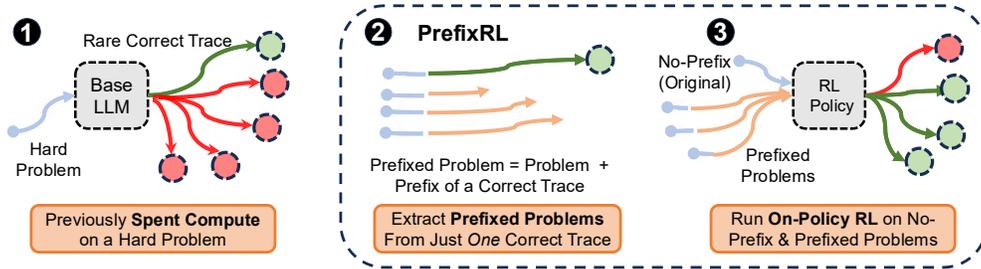


Figure 1 PrefixRL: On-Policy RL Conditioned on Off-Policy Prefixes. We leverage previously spent compute (❶) on hard problems in the form of correct off-policy traces rejection sampled from the base LLM we start RL from. Off-policy traces could also come from other model families or previous RL runs. We append prefixes of a single correct off-policy trace per problem to the original problem, creating prefixed problems (❷). Then, we simply run on-policy RL on prefixed and no-prefix (original) problems (❸). PrefixRL places the RL policy in higher-rewarding states, which boosts the learning signal. Performance transfers from the prefixed to no-prefix problems via a phenomenon we call back-generalization.

(original) problems and prefixed problems, where gradients are masked on the off-policy prefix. The prefixes extracted from correct off-policy traces place the current RL policy in states that are more likely to lead to a correct answer on hard problems, reducing gradient variance and increasing the strength of the learning signal.

PrefixRL is consistent with and more sample-efficient than standard RL. However, it is not immediately clear what the effect on the bias is. In Section 3.1, we prove that when the prefixes are correct and realizable in the model class, (i) maximizers of the PrefixRL objective also maximize performance on the standard RL objective; and (ii) since the prefixes lessen the exploration burden, PrefixRL reduces suboptimality gap with less samples compared to standard RL (by a factor of context length). Overall, PrefixRL changes the on-policy RL objective by using off-policy prefixes solely to guide exploration and unblock training on hard problems.

Back-generalization. Beyond the theory, we empirically find an additional phenomenon behind the gains in PrefixRL we call *back-generalization*, where on-policy RL on *only* prefixed problems substantially boosts test performance on the original no-prefix problems, which were never trained on. Beyond the generalization in the face of train/test mismatch, back-generalization is distinctive for two reasons. First, back-generalization is a type of *generalization via shared parameters* because it alters the next-token distribution on prefixes it was never trained on (impossible in the tabular RL setting). Second, we find that *back-generalization can be even more powerful than standard generalization* in RL (transfer across related problems or environments). We show this in an in-context learning setup, where we run RL on problems prefixed with another problem and reasoning trace in context. We find that training on a problem P1 conditioned on a related problem P2 in context improves generalization from P1 to P2 considerably more than directly running RL on the problem P1 (see Section 4.3).

PrefixRL can discover and learn strategies beyond what is provided in the prefix. Interestingly, the model does not simply back-generalize by imitating the off-policy prefix it is conditioned on. Through controlled experiments, we find that PrefixRL is more compute-efficient than standard RL at amplifying successful strategies and rejecting suboptimal ones, *even* when the suboptimal strategy is explicitly present in the off-policy prefix. As a result of observing non-zero rewards (and advantages) more often, we hypothesize that PrefixRL allows the model to more quickly identify the flaws in the suboptimal strategy and use this insight to find a better strategy (see Section 4.2).

PrefixRL improves both compute efficiency and final performance. In our experiments, we instantiate PrefixRL in a *self-improvement setting* by collecting a dataset of off-policy traces through large-scale rejection sampling on the base model. On hard problems in training, PrefixRL improves compute-efficiency over the strongest mid-training baseline (SFT on the off-policy data followed by on-policy RL) by $2\times$, even when we account for the initial compute spent on collecting the off-policy traces via rejection sampling, and training accuracy by $>45\%$ (over $3\times$ relative) on the original no-prefix problems (Figure 2 (left)). These gains transfer to held-out benchmarks: for example, on AIME '25, PrefixRL improves pass@1 by 12% over the strongest mid-training baseline in a compute-matched comparison. Finally, we find that PrefixRL is still effective when off-policy prefixes are sourced from Qwen3-4B-instruct while the RL policy is a distilled Llama-3.1-8B-instruct. This setting provides similar compute and accuracy gains, demonstrating the flexibility of PrefixRL to the off-policy data source and model size.

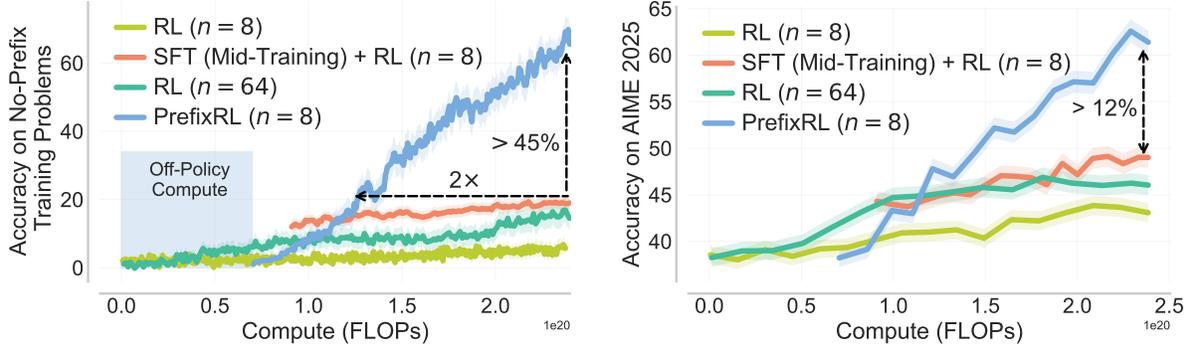


Figure 2 PrefixRL affords a self-improvement pipeline that recycles RL flops on hard problems. We instantiate PrefixRL for self-improvement by collecting a dataset of off-policy traces through large-scale rejection sampling on the base LLM (distilled Llama3.1-8B). In FLOPs-matched training, PrefixRL outperforms the strongest baseline (SFT on rejection-sampled data + RL): 2× higher compute efficiency (including rejection-sampling cost) and >45% (over 3× relative) higher final training accuracy on no-prefix training problems (*left*), with gains transferring to standardized evals such as AIME ’25 (*right*).

2 Preliminaries

Notation. We use \mathbf{x} to denote a problem and $\mathbf{y} = (y_1, \dots, y_H)$ for a response of H tokens sampled auto-regressively from an LLM π in policy class Π , where $\mathbf{y}_{:h}$ refers to the prefix consisting of first h tokens in \mathbf{y} . We use $\pi^0 \in \Pi$ to denote the base (pre-trained) LLM that we want to post-train on a dataset \mathcal{D} of hard problems $\mathcal{D} =: \{\mathbf{x}_i\}_{i=1}^{|\mathcal{D}|}$. We have an outcome reward function $r(\mathbf{x}_i, \mathbf{y})$ which is 1 when \mathbf{y} is correct and 0 when incorrect (*e.g.*, by matching the boxed answer at the end of \mathbf{y} for a math reasoning problem). We define the pass rate @ k (pass@ k) given \mathbf{x} and LLM π as $\mathbb{E}_{\mathbf{y}_1, \dots, \mathbf{y}_k \sim \pi(\cdot|\mathbf{x})} \max(\{r(\mathbf{x}, \mathbf{y}_i)\}_{i=1}^k)$. We say \mathbf{x} is a *hard problem* for π^0 if pass@512 under π^0 is ≈ 0 for \mathbf{x} . We use $J(\pi) := \mathbb{E}_{\mathbf{x} \sim \rho} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot|\mathbf{x})} r(\mathbf{x}, \mathbf{y})$ to denote the performance of π on ρ which is the empirical distribution over \mathcal{D} . Next, we describe the iterative policy gradient RL algorithm that is the backbone of post-training methods used to optimize $J(\pi)$. For a full set of notations, see Appendix A.

Policy gradient RL algorithms. Almost all iterative policy gradient RL algorithms for training LLMs (*e.g.*, GRPO (Guo et al., 2025), PPO-clip (Schulman et al., 2017), REINFORCE (Ahmadian et al., 2024)) start from base LLM π^0 , and in each iteration t they perform a step of gradient ascent on $J(\pi)$, where the gradient $\nabla_{\pi} J(\pi)$ is typically computed on samples drawn from π^t , *i.e.*, they make on-policy updates. *E.g.*, the REINFORCE algorithm we use in our experiments uses the following gradient update in population:

$$\nabla_{\pi} J(\pi)|_{\pi=\pi^t} = \mathbb{E}_{\mathbf{x} \sim \rho} \mathbb{E}_{\mathbf{y} \sim \pi^t(\cdot|\mathbf{x})} [A_{\pi^t}(\mathbf{x}, \mathbf{y}) \cdot \log \pi^t(\mathbf{y}|\mathbf{x})]. \quad (2.1)$$

Following GRPO (Guo et al., 2025), the expectation in (2.1) is approximated by averaging over n independent traces $\mathbf{y}_1, \dots, \mathbf{y}_n \sim \pi^t(\cdot|\mathbf{x})$ and the advantage $A_{\pi^t}(\mathbf{x}, \mathbf{y}) = r(\mathbf{x}, \mathbf{y}) - Q_{\pi^t}(\mathbf{x}, \mathbf{y})$ (where $Q_{\pi^t}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{y}' \sim \pi^t(\cdot|\mathbf{x})} r(\mathbf{x}, \mathbf{y}')$) is also approximated with the empirical average $r(\mathbf{x}, \mathbf{y}_i) - 1/n \sum_{i \in [n]} r(\mathbf{x}, \mathbf{y}_i)$. For a hard problem \mathbf{x} with pass@ $n \approx 0$, it is highly likely that all n traces in the group will fail, making the computed gradient on \mathbf{x} be $\mathbf{0}$ and giving rise to the *stalling regime for RL on hard problems*.

Problem setup. Consider a base model π^0 that we train with RL on a dataset of verifiable problems \mathcal{D} . We aim to maximize $J(\pi)$ over $\pi \in \Pi$ with access to previously spent compute on π^0 (or on models fine-tuned from it), available in the form of correct off-policy traces \mathcal{D}_{off} . This includes any prior sampling or RL training runs where the problem set includes problems in \mathcal{D} , and \mathcal{D}_{off} contains the rare correct traces from the prior compute. We situate our study in the setup of training only on hard (low pass rate) problems since that is where we expect off-policy data to be most useful, though notably, large-scale RL runs would train on a wider mixture.

Source of off-policy traces. To simplify our study, we mainly source the off-policy traces via rejection sampling on the base model π^0 . Concretely, for each problem $\mathbf{x} \in \mathcal{D}$, we collect a single correct trace by sampling from π^0 until we see a correct trace. Therefore, if the pass@1 under π^0 is $p_{\mathbf{x}}$ on problem \mathbf{x} , then in expectation we will sample $1/p_{\mathbf{x}}$ traces for \mathbf{x} to get a correct one. Doing this for every $\mathbf{x} \in \mathcal{D}$ gives us \mathcal{D}_{off} , where $|\mathcal{D}_{\text{off}}| = |\mathcal{D}|$. In theory, we assume that \mathcal{D}_{off} is realizable by the model μ in policy class Π (*e.g.*, μ can be the policy that generates \mathcal{D}_{off} which is the rejection sampling policy over π^0), although in practice \mathcal{D}_{off} can also be curated with sophisticated inference algorithms that interleave sequential and parallel compute, or use some form of oracle feedback, that may not be representable. We empirically show the flexibility of the off-policy source in our experiments in Section 5.1.

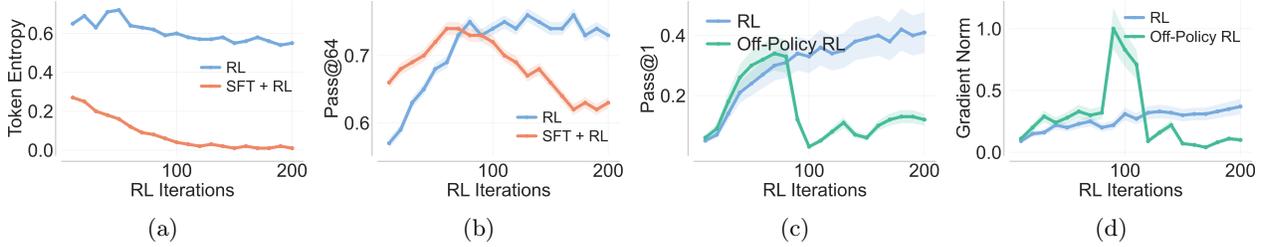


Figure 3 Supervising the policy on \mathcal{D}_{off} can cause diversity collapse or training instabilities during RL: (a, b) Warm-starting the RL run by running SFT on \mathcal{D}_{off} before (mid-training) reduces token entropy (a) and hurts exploration during RL (worse Pass@64 performance in (b)). (c, d) Directly using \mathcal{D}_{off} during online RL by updating the current RL policy with importance-weighted off-policy traces (in addition to on-policy traces) leads to training instabilities. We see the gradient norm (clipped at 1.0) blow up during training (d) and this leads to optimization collapse (c).

3 PrefixRL: On-Policy RL on Very Off-policy Prefixes

In this section, we introduce the **PrefixRL** framework. First, we discuss why typical ways of using off-policy data \mathcal{D}_{off} as supervision targets lead to worse performance. Then, we outline the steps in PrefixRL, which instead *conditions* the RL policy on prefixes from \mathcal{D}_{off} to guide on-policy exploration and increase the probability of success on prefixed problems, which then generalizes to no-prefix (original) problems. Finally, we theoretically prove the correctness of the PrefixRL objective and show how PrefixRL can improve training rewards with less samples spent on exploration compared to standard RL.

SFT on \mathcal{D}_{off} can hurt exploration during RL. A common way to improve RL on hard problems is to first mid-train (SFT) on traces in \mathcal{D}_{off} to warm-start RL training. However, while SFT boosts post-RL pass@1, pass@64 drops after about 100 iterations (Figure 3(a)). This loss of diversity is driven by a sharp entropy collapse after SFT, which decreases further during RL (Figure 3(b)). This suggests that here, RL mainly sharpens the distribution of responses the model can already produce after SFT. In contrast, an early-stopped SFT checkpoint can underfit and encourage random exploration during RL (Wang et al., 2025d). One can partially mitigate this by enlarging the SFT dataset (\mathcal{D}_{off}), but doing so increases the upfront cost of the SFT step.

Off-policy RL using \mathcal{D}_{off} leads to training instabilities. A more direct way to use the off-policy data \mathcal{D}_{off} is to do importance-weighted off-policy RL (Degris et al., 2012) (see Appendix D.1), which accounts for the distribution shift between the current RL and sampling (rejection sampling on base LLM) policies. However, this suffers from large gradient variance or heavily biased gradients due to clipping and token-level weighting rather than sequence-level (Agarwal et al., 2021) (Section 5.2). This can cause training reward collapse and unstable optimization (Figures 3(c),(d)) as we force updates on very unlikely token sequences under the current RL policy, leading to memorization of \mathcal{D}_{off} (Kang et al., 2024a; Setlur et al., 2024).

Creating prefixed problems in PrefixRL. Instead of imitating off-policy traces, PrefixRL runs on-policy RL conditioned on off-policy prefixes (in addition to the original no-prefix problems). Crucially, the gradients are always masked on the off-policy prefix, avoiding the instability of policy gradients on very off-policy tokens. We create a dataset of prefixed problems \mathcal{D}_{pre} by taking prefixes from a correct trace $\mathbf{y}^x \in \mathcal{D}_{\text{off}}$ and appending the first h tokens ($\mathbf{y}_{1:h}^x$) in \mathbf{y}^x to the original problem \mathbf{x} , creating the prefixed problem $\text{concat}(\mathbf{x}, \mathbf{y}_{1:h}^x)$. We create multiple prefixed problems for every original problem by choosing different value of the number of tokens we prefix h . In general, we choose h such that, conditioned on the prefix, the base LLM has a reasonable accuracy under base LLM (see Section 5 for details). We explain this in Section 4, but typically these are states revealing a high-level problem-solving strategy which the base LLM has little probability of sampling on its own (see Figure 4).

PrefixRL training objective. The PrefixRL objective runs on-policy RL both prefixed problems in \mathcal{D}_{pre} and no-prefix ones in \mathcal{D} within a maximum context length of H tokens. Note that the reward $r(\mathbf{x}, \cdot)$ for any prefixed problem $\mathbf{x}_{\text{pre}} \in \mathcal{D}_{\text{pre}}$ is identical to the reward of the corresponding no-prefix counterpart.

$$\text{(PrefixRL)} \quad \max_{\pi} \left(\underbrace{\sum_{\mathbf{x}_{\text{pre}} \in \mathcal{D}_{\text{pre}}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x}_{\text{pre}})} [r(\mathbf{x}_{\text{pre}}, \mathbf{y})]}_{\text{prefixed problems}} + \underbrace{\sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})]}_{\text{Standard RL: No-Prefix Problems}} \right) \quad (3.1)$$

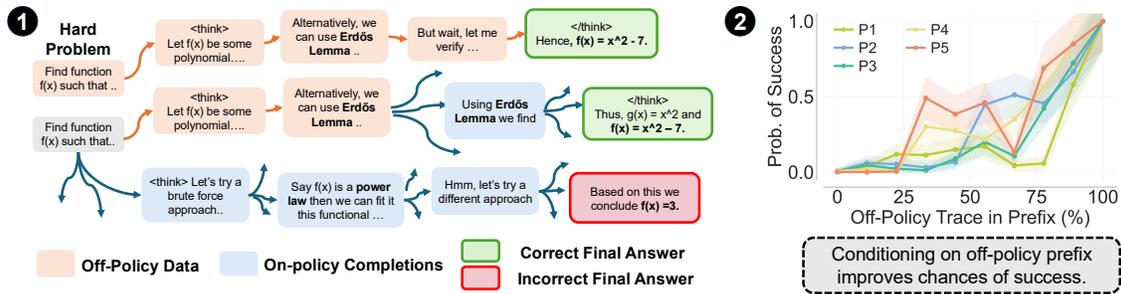


Figure 4 Prefixing on off-policy trace improves probability of future success. When we condition on prefixed problems, we increase the accuracy by placing the policy at key strategy-revealing states (Erdős Lemma in example ①). For five problems (P1, ..., P5) we plot accuracy when conditioning on prefixes of varying lengths, as a proportion of the full off-policy prefix length (②). The accuracy is near zero until a key state is visited, after which it jumps sharply.

Takeaway: PrefixRL conditions on off-policy traces instead of using them as supervision.

- The correct traces in \mathcal{D}_{off} serve as poor supervision targets: using them in the SFT phase hurts exploration during RL and directly updating on them via off-policy RL leads to optimization collapse.
- PrefixRL never updates the RL policy on \mathcal{D}_{off} and instead runs on-policy RL conditioned on prefixed problems (original problems appended with off-policy prefixes in \mathcal{D}_{off}) along with the original problems.

3.1 PrefixRL Objective-Consistency and Sample-Efficiency Guarantees

Prefixing problems with partial correct traces would expectedly place the RL policy in a higher-rewarding state. The central theoretical question, however, is not whether learning the prefixed problems are easier, but instead whether optimizing such *augmented* problem sets provably improves the performance on the *original* RL objective $J(\pi)$, which evaluates policies only on the no-prefix problems. We provide an answer for this next.

In general, training on an altered input distribution could change the objective away from maximizing $J(\pi)$. We show that this is *not* the case for PrefixRL as long as the prefixes come from correct traces generated by a realizable policy. Concretely, we prove: (i) *objective consistency*: every maximizer of the PrefixRL objective is also a maximizer of $J(\pi)$; and (ii) *sample complexity guarantees and improvement over online RL*: for a natural policy gradient variant, PrefixRL achieves a smaller suboptimality bound, which translates to a smaller number of on-policy samples required to reach a given reward $J(\pi)$. In other words, we formally show that **PrefixRL reuses your FLOPs**: it converts information already paid for in logged prefixes into sample-complexity advantages over standard RL.

PrefixRL objective is consistent with standard RL. We make the following assumption that the prefixes are taken from the correct traces and that there exists a realizable policy $\mu \in \Pi$ that can fit the traces completely.

Assumption 3.1 (Realizability and correctness of off-policy traces). Assume that for any $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{off}}$: (i) the trace is correct: $r(\mathbf{x}, \mathbf{y}) = 1$, and (ii) the trace is realizable, i.e., there exists an optimal policy $\mu \in \Pi$ s.t. $\mu(\mathbf{y} | \mathbf{x}) = 1$.

The next theorem states that as long as prefixes are taken from the correct traces generated by a realizable policy, optimizing the PrefixRL objective preserves optimality on $J(\pi)$; see Appendix B.1 for a proof.

Theorem 3.2 (Consistency of the PrefixRL objective). Assume the realizability and correctness of off-policy traces (Assm. 3.1). Then, the maximizer of the PrefixRL objective (3.1) also maximizes standard RL objective $J(\pi)$.

Intuitively, a maximizer of the PrefixRL objective produces correct traces on both no-prefix problems and prefixed problems. Since the prefixes come from correct traces, a good policy should also be able to complete the prefix to get the same reward; thus the two terms in the objective do not conflict with each other. Note that while PrefixRL does not change the global solution, it does not produce the same gradients as the standard RL objective.

PrefixRL is more sample-efficient than standard RL. Having established that PrefixRL does not bias policy optimization, we now quantify the benefits of prefixing in terms of the number of on-policy samples needed to reach a near-optimal policy. We analyze PrefixRL by instantiating the policy update to be natural policy gradient (Kakade, 2001) (NPG) (Algorithm 1 in Appendix B.2). Concretely, each RL iteration alternates between two sub-steps: (i) *policy evaluation* by fitting a critic or Q -function in the class \mathcal{F} using a few completions from the

current policy conditioned on the prefixed problem, and (ii) *policy improvement* using the fitted critic via a KL-regularized mirror-descent update of NPG. In practice, algorithms like PPO, GRPO serve as approximations of the NPG update (Schulman et al., 2017), since its hard to solve the optimization problem for high-dimensional actions.

The next theorem bounds the suboptimality of the policy $\bar{\pi}_T$ returned by PrefixRL (Algorithm 1) in terms of the number of policy updates T , the number of on-policy traces per iteration N , and a single distribution shift quantity between the initial base policy π_0 and policy μ that samples the off-policy data \mathcal{D}_{off} .¹

Theorem 3.3 (Suboptimality gap of PrefixRL). *Under Assumption 3.1, let \mathcal{D}_{off} be sampled by $\mu \in \Pi$. For any $\delta \in [0, 1]$, with probability at least $1 - \delta$, policy $\bar{\pi}_T$ returned by PrefixRL-NPG (Algorithm 1) satisfies:*

$$\max_{\pi \in \Pi} J(\pi) - J(\bar{\pi}_T) \leq \mathcal{O} \left(\sqrt{\frac{\text{KL}(\mu || \pi_0)}{T}} + \sqrt{\frac{1}{N} \cdot \log \left(\frac{T|\mathcal{F}|}{\delta} \right)} \right).$$

Proof² is in Appendix B.2. The bound consists of two terms:

- **Optimization term:** The first term captures PrefixRL’s convergence rate assuming access to an oracle policy evaluator (*i.e.*, $N \rightarrow \infty$ so that the second term vanishes), which is $1/\sqrt{T}$. This rate is impacted by a constant that depends on the initial distribution shift between the base policy π_0 and μ that realizes \mathcal{D}_{off} , and crucially does not accumulate with intermediate policy iterations π_t . Consider the case where the off-policy traces are obtained through rejection sampling π_0 with at most R attempts per problem. In that case, the induced behavior policy μ is π_0 conditioned on success and satisfies $\text{KL}(\mu || \pi_0) = \mathcal{O}(\log R)$. Hence, the optimization term in Theorem 3.3 grows only logarithmically with the rejection budget.
- **Critic approximation term:** The second term captures the statistical error of policy evaluation, which fits $Q_{\pi_t} \in \mathcal{F}$ using N on-policy traces. This term is not impacted by any distribution shift penalty because PrefixRL samples and evaluates traces under the same reset distribution induced by off-policy prefixes from \mathcal{D}_{off} , but standard RL still pays a distribution shift term between π^t and optimal policy π^* .

In the following proposition, we show that there exists a reward function and base LLM such that there is a performance gap between on-policy NPG and PrefixRL-NPG (see Appendix B.3 for the proof).

Proposition 3.4 (Worst-case separation with standard RL). *Let $\bar{\pi}_T^{\text{pre}}$ be the policy returned after T PrefixRL iterations of Algorithm 1 and $\bar{\pi}_T^{\text{std}}$ be the policy after T iterations of standard RL (states in Algorithm 1 $\sim \pi_t$). Then, there exists a reward function r and base LLM π^0 such that $J(\bar{\pi}_T^{\text{pre}}) - J(\bar{\pi}_T^{\text{std}}) \geq 1 - (TN \cdot e^{-H})$ for $TN = o(e^H)$.*

Takeaway: PrefixRL is consistent with standard RL except more sample-efficient.

- Maximizer of the PrefixRL objective also maximizes $J(\pi)$ when traces in \mathcal{D}_{off} are correct and realizable.
- The sample complexity of PrefixRL with natural policy gradient *provably* scales more favorably with longer context windows (horizons) compared to standard RL when \mathcal{D}_{off} is correct and realizable.

4 Back-Generalization Boosts the Learning Signal in PrefixRL

In Section 3.1, we showed that PrefixRL is a consistent objective that leads to the same optimal solution as standard RL while being more sample-efficient. But, different from the algorithm we analyze theoretically, in practice we run PrefixRL conditioned on a handful of off-policy prefixes per problem and yet we see strong improvements. Now, we show that an empirical phenomenon we call **back-generalization** is a strong source of performance improvement behind PrefixRL that is not explained by our theory. Back-generalization is defined as the performance improvement on no-prefix problems when we train *only* on their prefixed counterparts.

Back-generalization from prefixed to no-prefix problems goes beyond stitching. Consider a straightforward *stitching* argument as an explanation for back-generalization (Zhang et al., 2025b; Qu et al., 2025a). According to this, the model simply learns to complete from the off-policy intermediate states better without updating the next-token distributions on no-prefix problems, but still improves performance on no-prefix problems when the model happens to sample the same intermediate states on its own. Note that this argument still holds in the tabular

¹We do not require access to μ and only assume the off-policy traces are realizable under some unknown μ .

²It is straightforward to extend the result showing improvement competing with a suboptimal behavior policy μ .

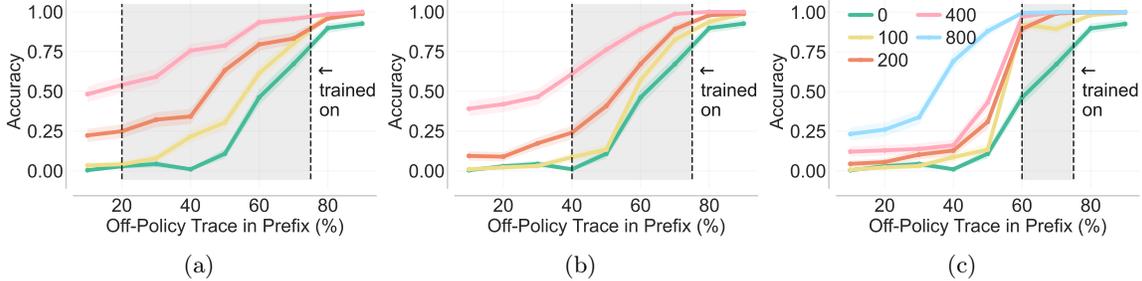


Figure 5 Back-generalization (train-test mismatch): Using Llama3.1-8b-instruct, we run RL only on prefixed problems with prefix length lies in the shaded interval. We evaluate different training step checkpoints across the full range of prefix lengths, including no-prefix problems. Training on longer prefixes improves performance on shorter prefixes and can eventually improve no-prefix, indicating back-generalization (a,b). When training uses only very long prefixes (severe mismatch with no-prefix), back-generalization to no-prefix problems takes more training steps (*e.g.*, 800 iterations) (c).

policy setting. In contrast, we find that back-generalization indeed influences next-token distributions on untrained states (no-prefix problems), which can only occur through favorable function approximation in LLMs.

4.1 PrefixRL Improves No-Prefix Performance Even When Training Only on Prefixed Problems

We run on-policy RL only on prefixed problems where the prefix lengths are distributed uniformly between a fixed band of token-length percentiles of the full off-policy trace, but we evaluate accuracy across the full spectrum of prefix lengths, including the no-prefix endpoint (0% prefixing). In Figure 5, we see generalization to no-prefix problems despite not having trained on them. This transfer from prefixed to no-prefix problems is particularly notable since the prefixes are highly unlikely under the base policy. When the training mixture includes relatively short prefixes, the mismatch is moderate (Figure 5 (a,b)). In this case, performance increases first near the trained band and then progressively improves for shorter prefixes, eventually lifting no-prefix accuracy. When training is restricted to very long prefixes (Figure 5(c)), the train/test mismatch with no-prefix problems is more severe. The transfer is slower in this case, but longer training (*e.g.*, 800 steps) still yields measurable no-prefix gains.

4.2 PrefixRL can Discover New Strategies Beyond What is Present in the Prefixed Problem

Clearly, back-generalization improves performance on unseen shorter prefixes, but the mechanism behind this remains unclear. To this end, we create a simplified setup where we run RL on the prefixed-problems derived from a single off-policy trace (and thus, a single problem) and can track the problem-solving strategies the model uses.

Setup. We run PrefixRL (for 100 iterations) on the prefixes of a single off-policy trace in \mathcal{D}_{off} . Crucially, we don’t run it on the corresponding no-prefix problem. For two hard problems, we use a keyword heuristic to extract a salient “strategy” present in the off-policy trace for each problem. Then we track the prevalence of this keyword in (i) the prefixed-problem (PP) itself, (ii) the model’s response when conditioned on that prefixed-problem (response for PP), and (iii) the early part of the trace on the original problem (untrained states in the model’s response).

Figure 6 probes how strategy usage evolves when training is conducted only on prefixed problems. Because the prefix is sampled from a fixed pool, the keyword frequency in the prefix remains constant by construction

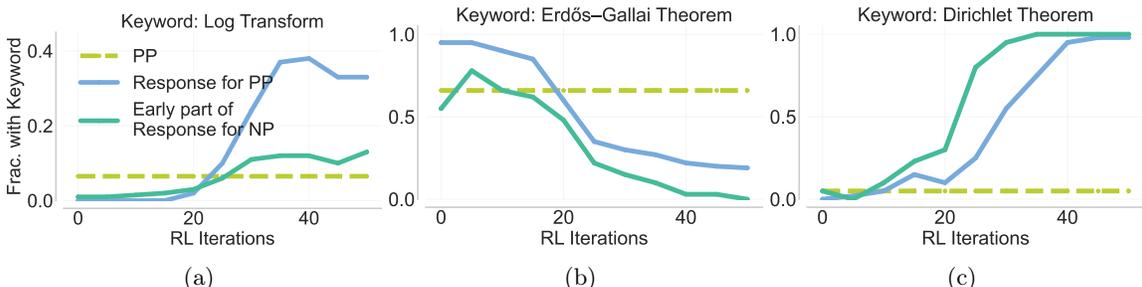


Figure 6 Strong coupling between responses for prefixed and no-prefix problems.: We train only on prefixed-problems (PP) and track the frequency of a strategy-indicating keyword in: (i) PP, (ii) model’s response to PP, and (iii) and the early part (prefixes or states not trained on) of the response to the no-prefix (NP) problem. There is a tight coupling of the strategies present in the responses for PP and NP (throughout RL), yet not purely imitative of the strategy explicitly provided in the prefixed-problem itself: the model can learn new strategies or suppress prefixed ones (*e.g.*, Erdős–Gallai).

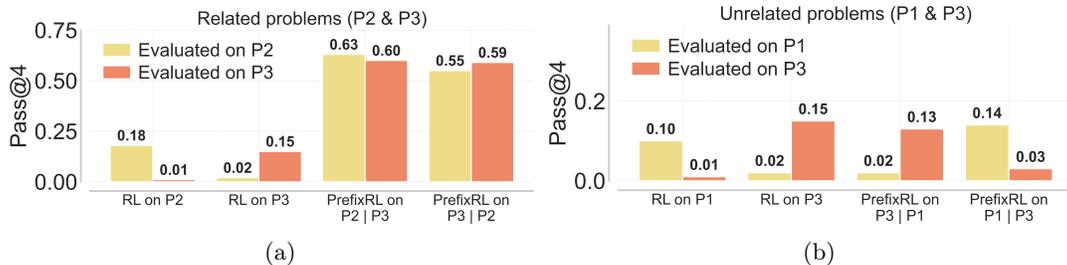


Figure 7 Performance transfer via back-generalization can be stronger than typical generalization in RL: When we prefix on the problem and full solution trace of one (in-context) problem (P2), and run Prefix RL to solve a different but related problem P3 | P2, we are able to improve performance on both P2 and P3 individually, and the performance is much higher compared to running RL either problem individually. The same holds in the opposite direction, when we run PrefixRL on P2 | P3. We do not see these gains when the in-context problem is unrelated in the case of P1 and P3.

(dashed horizontal line). In contrast, the response curves change over training and reveal two consistent patterns:

Strategy usage is strongly correlated between prefixed and no-prefix responses. There is a tight coupling between strategy use on prefixed and no-prefix problems, which is difficult to explain since no-prefix problems are never trained on and many prefixed ones (>90% in (a)) may not even contain the keyword. This instead suggests that PrefixRL updates internal state representations that are shared with or without the prefix.

PrefixRL can unlearn strategies in the prefixed-problem and instead discover new ones. When comparing the strategies in the off-policy data vs. the NP (no-prefix) strategies learned by the model, we see that the model is not simply learning to copy the same strategy that is in the off-policy prefix. Instead, PrefixRL can push the model to adopt strategies that are not present in the prefix, or to actively suppress strategies that the prefix hints at. In our example, the prevalence of samples that use the “Erdős–Gallai theorem” illustrates this. In Figure 6(b), we see that the policy at initialization uses the “Erdős–Gallai theorem” close to 90% of the time on the prefixed-problems since >50% of the prefixed-problems contain references to it. Throughout training, the frequency of traces mentioning “Erdős–Gallai” decreases steadily on responses for the prefixed-problems, indicating that optimization can downweight the off-policy hinted strategy, and instead discover a more useful one: “Dirichlet’s theorem”. In Figure 6(c) we note that despite being conditioned on prefixes that contain references to the suboptimal strategy of using “Erdős–Gallai”, RL optimization upweights the rare (<2%) strategy at initialization (“Dirichlet Theorem”). The trends present in the responses for PP also transfer to the early part of the responses for the NP.

Remark. The above trends support the view that prefixes simply boost the training signal to accelerate training progress and exploration rather than biasing the model towards any particular solution. This supports the theoretical consistency of the PrefixRL objective in Theorem 3.2, since if PrefixRL could only learn strategies present in the off-policy prefixes, then it is unlikely to share the same optimal policy as the standard RL objective.

Takeaway: Back-generalization can transfer strategies different from the one in the prefix.

PrefixRL can discover strategies beyond what is present in the prefixed-problems. Benefitting from function approximation, PrefixRL alters the next-token distribution on unseen states; strategies learned and unlearned on the prefixed-problems are quickly reflected in the responses for the no-prefix ones.

4.3 Which Prefixes Back-Generalize the Most in PrefixRL?: Analysis via In-Context Learning

To study when back-generalization is effective, we analyze it in the *in-context learning* setting, where we run RL on problems prefixed with another problem and reasoning trace in context. This lets us cleanly ablate the relationship between the off-policy prefix and the generated on-policy suffix based on how related the in-context problems are.

Setup. We run RL (for 100 iterations) on a given problem with an entirely different problem (and its solution trace) in its context or prefix. Consider two problem sets: (P1, P3) where P1 and P3 are unrelated sub-problems, and (P2, P3) where P2 and P3 are related and solved with the same high level strategy (see Appendix C for details on P1, P2 and P3). We choose problems that are hard for the base model, with <1% pass@32.

Back-generalization occurs when the prefix and suffix are sufficiently related. From Figure 7, when the problems are related (P2 and P3), PrefixRL on P2 given P3 achieves 63% pass@4 on P2 and 60% pass@4 on the untrained in-context problem P3. Running standard RL on P2 alone predictably improves the pass@4 of P2 to 18%

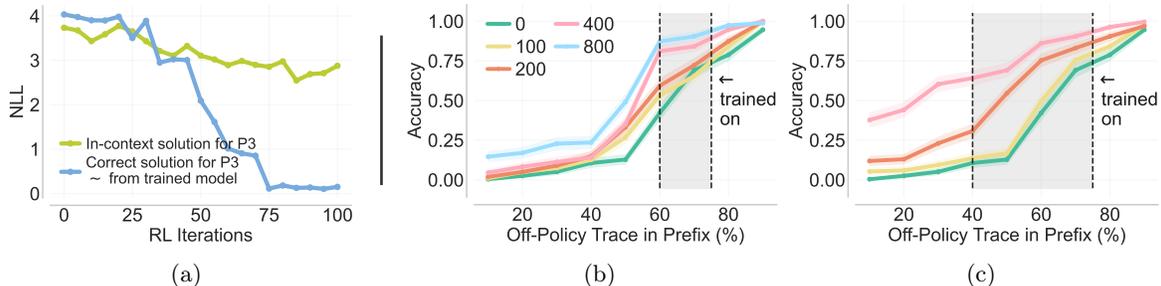


Figure 8 (a) Likelihood of the in-context solution: When asked to solve problem P3, we measure the negative log-likelihood (NLL) of the in-context solution for P3 provided in the context for P2, *i.e.*, when we run PrefixRL on P2 | P3 (setup in Section 4.3). Surprisingly, we find that the likelihood of the in-context trace drops less compared to a correct trace for P3 we sample from the final checkpoint. This indicates that back-generalization does not exactly clone the in-context prefix to improve performance on the no-prefix counterpart. **(b,c): Prefixes sourced from a different model family:** On Llama3.1-8b-instruct we run RL only on prefixes sourced from Qwen3-4b-instruct. The setup is similar to Figure 5 except there the prefixes were sourced from the same base LLM we were training. Different from Figure 5, we find that when the prefix distribution is skewed towards long prefixes back-generalization is weaker despite running RL for 800 iterations.

but the performance transfer to P3 is limited. In contrast, PrefixRL on unrelated problems (P3 | P1 or P1 | P2) performs similarly to doing RL on just P3 and P1 respectively. This suggests that *back-generalization is more effective when the components of the prefixed problem are related*, and in the in-context learning setting, back-generalization can be stronger than standard generalization across the two related problems. So when running RL on P2 | P3, why does performance improve on the related in-context problem (P3) that is also hard?

NLL of the in-context solution changes little. A natural hypothesis is that improvements on the untrained but in-context problem P3 come from memorizing the in-context trace and replaying it at test time. Instead, because this trace is extremely unlikely and never directly trained on, the model does not learn to imitate it (we also saw an example of this in Figure 6(b)). In Figure 8(a), the negative log-likelihood of the in-context P3 solution barely decreases under RL on P2|P3; the final policy instead prefers a different token sequence that still yields the correct answer. Together with the correlation in Figure 6, this suggests strong similarity between prefixed and no-prefix solutions, but weak similarity to the specific off-policy prefix used in the prefixed problem.

Mental model: function approximation and back-generalization. Although the model does not clone the in-context off-policy trace, performance still transfers to the in-context problem. Our speculation is that for long chain-of-thought rollouts, “state” is better viewed as the model’s internal representation induced by the prefix: because the model self-corrects and backtracks, many distinct token sequences can map to similar latent states. Thus, while solving the prefixed problem, the policy can backtrack into representations close to those encountered when solving the original problem directly, but now with positive reward. If the history-conditioned and non-history-conditioned representations are close, rewards observed in the former will shift the next-token distribution in the latter. This explains the overlap between prefixed and no-prefix responses (Figure 6), why transfer is stronger for related pairs like P2|P3 than unrelated ones like P2|P1 (Figure 7), and why the model still may not learn to reproduce the literal off-policy prefix when its rephrased representation is far from the original context (Figure 8(a)).

How strong is back-generalization when prefixes are sourced from a different model family? A natural question is if back-generalization only works when the prefixes themselves are somewhat related to the current policy by sharing the same base model family. In Figure 8, we present the same experiment as in Figure 5 but with off-policy prefixes from Qwen3 for training a Llama model. When we prefix on Qwen3 prefixes that are very long, back-generalization takes more iterations than with off-policy prefixes from the same family (Llama). On a wider range of prefix lengths, back-generalization occurs at a similar rate regardless of the off-policy model family. This suggests that the off-policy model family matters less for back-generalization if the prefix length distribution is wide enough to build a “bridge” to no-prefix problems.

5 Experiments and Results

In this section, we present our main empirical findings from evaluating PrefixRL on math reasoning benchmarks. First, in compute-matched comparisons on hard problems, PrefixRL consistently outperforms standard RL and SFT+RL baselines in training rewards. These improvements transfer to held-out benchmarks such as AIME, HMMT, and IMO-AnswerBench. We also ablate the source of the off-policy data, showing that PrefixRL

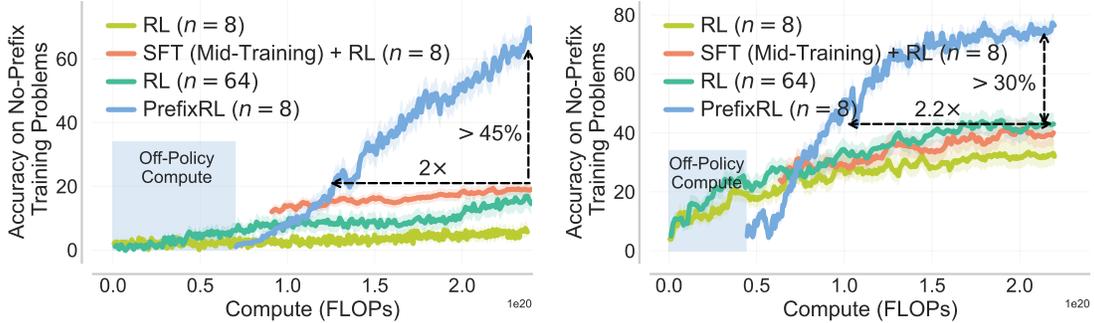


Figure 9 Training compute matched evaluation of PrefixRL. Using prefixed-problems sourced by rejection sampling the base model, we run PrefixRL on the dataset of hard training problems (low pass rate under the base model) and measure accuracy on the no-prefix training problems for **(Left)** Llama-3.1-8b-instruct and **(Right)** Qwen-3-4b-instruct. Even accounting for the initial compute spent on rejection sampling (shaded blue region), PrefixRL improves compute-efficiency by 2 \times over the strongest baseline with and an absolute gain of 45% with Llama and 30% with Qwen.

remains effective when the off-policy prefixes are from a different family. Finally, we analyze training dynamics and compare PrefixRL against other off-policy RL approaches.

Experimental setup. We conduct experiments on two models: **Llama-3.1-8B-instruct** and **Qwen3-4B-instruct**. We focus on “thinking” models, where the model outputs a chain-of-thought followed by a final answer. Note that Llama-3.1-8B-instruct is not a thinking model, so we distill it on OpenThoughtsV3 (Guha et al., 2025) before running all our experiments. For the training problems, we select only hard problems from DAPO (Yu et al., 2025) and OMNI-MATH (levels 6-8) (Gao et al., 2024), where pass@512 of the distilled Llama-3.1-8B model is zero. This results in our fixed training set of 1k problems. We compare against the standard on-policy REINFORCE baseline (Ahmadian et al., 2024), as well as off-policy baselines that use \mathcal{D}_{off} : SFT (mid-training) on \mathcal{D}_{off} followed by standard RL (SFT+RL), importance-weighted off-policy RL (Mahmood et al., 2014), and LUFFY (Yan et al., 2025) which trains with a mixed off-policy and on-policy GRPO objective. For implementation details on PrefixRL and the baselines (*e.g.*, importance-weight computation), please see Appendix D.1.

Off-policy Dataset \mathcal{D}_{off} and prefixed-problems. For each base model, we construct \mathcal{D}_{off} using large-scale rejection sampling: for every training problem, we sample until we obtain one correct trace. We explicitly account for the compute used to curate \mathcal{D}_{off} in the total compute budget allocated to PrefixRL training. Given \mathcal{D}_{off} , we form prefixed problems as follows. For each training problem, we sample three prefixes by truncating its correct off-policy trace at a uniformly random cut point between 40% and 80% of the tokens. The 1k original no-prefix problems and these 3k prefixed-problems together form the training dataset for PrefixRL.

Evaluation. All evaluation results in this section are on the original no-prefix problems. For the plots where we report pass@ k , we estimate it by drawing 256 samples per problem and using the bootstrapped estimate in Chen et al. (2021). Where possible, we include 95% confidence intervals across evaluated problems. Details on FLOPs accounting for our compute-matched plots are in Appendix D.2.

5.1 PrefixRL is More Compute-Efficient on Hard Problems Than Standard RL

With the same training FLOPs, PrefixRL achieves higher training rewards on the no-prefix problems compared to standard RL and mid-training baselines. Crucially, this gain holds even when we account for the initial compute spent on collecting off-policy traces via rejection sampling. This implies that PrefixRL is able to re-allocate the overall sampling and training FLOPs better than standard RL, improving training rewards on hard problems where standard RL stagnates. Importantly, we will show that these gains also transfer to held-out benchmarks.

PrefixRL is 2 \times more compute-efficient and achieves higher training accuracy. Figure 9 shows that in a compute-matched evaluation, PrefixRL achieves higher accuracy on no-prefix problems compared to baselines for both Llama-3.1-8B-instruct (45% greater) and Qwen-3-4B-Instruct (30% greater). After accounting for the initial rejection-sampling cost, PrefixRL improves compute-efficiency by roughly 2 \times over the strongest baseline (SFT+RL). In contrast, standard on-policy RL and SFT+RL only slowly improve the no-prefix accuracy even when the number of samples per problem n is increased from 8 \rightarrow 64. In iteration-matched plots (Figure 11(c)), we see that standard RL and SFT+RL baselines have stable training curves, while higher values of n unsurprisingly

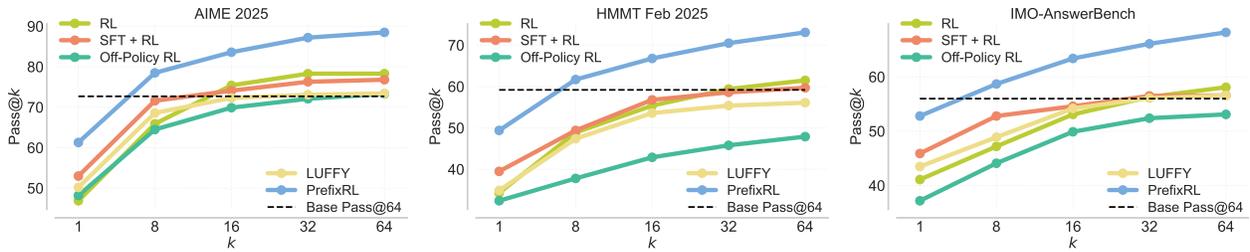


Figure 10 Held-out pass@k of PrefixRL trained Llama models: We measure pass@k (%) on AIME’25, HMMT’25 and 200 problems from IMO-AnswerBench for our distilled Llama-3.1-8B (base model) finetuned with standard on-policy RL, and PrefixRL (with and without off-policy suffix injection). Note that IMO-AnswerBench is sampled from a similar distribution to our training problems, but is still held-out. The horizontal dashed line marks the base model’s pass@64.

attains higher accuracy. So, the compute gains are not explained by unstable or degenerate baseline runs.

PrefixRL improves both pass@1 and pass@k on held-out benchmarks. Figure 10 shows that across AIME’25, HMMT’25, and IMO-AnswerBench, PrefixRL consistently improves pass@k (for $k \leq 64$) over the baselines by over 10% absolute, including off-policy RL methods such as importance-weighting (off-Policy RL) and LUFFY. The gains are visible already at $k=1$ and widen as k increases. This gain is notable since we only train on hard problems, so there is substantial transfer from the better training accuracy on hard problems to solving both easy and hard problems. On AIME’25 with Llama-3.1-8B, PrefixRL improves pass@1 from 38.2 to 61.3, a sizable absolute gain given the benchmark’s difficulty. We observe a similar effect on HMMT’25, where pass@1 increases from 29.2 to 49.4. These small- k improvements are diagnostically important: they suggest that the model is more likely to instantiate the right high-level plan earlier in the trajectory, aligning with our “backward generalization” hypothesis (Section 4). As k grows, the performance gap generally widens (e.g., on AIME’ 25: +18 points at $k=8$ and +28 at $k=64$). This pattern indicates that PrefixRL improves the diversity of the overall search distribution: additional samples explore more promising subspaces rather than repeating low-value trajectories. In other words, PrefixRL enhances both the *mean* performance (seen in $k=1$) and the *tail* (seen as k increases) of the trace distribution.

PrefixRL increases the support of solvable problems over the course of training. Now we ask how much of the gains in Figures 9, 10 come from “sharpening” the model’s output distribution on problems that were already solvable with more samples (pass@k) versus an actual expansion in the support of solvable problems. Figure 11(b) shows that PrefixRL not only improves pass@1 but also steadily improves compute-matched pass@32, while competing baselines largely saturate over the course of training. This suggests that PrefixRL increases the set of problems with non-trivial success probability rather than only converting pass@k into higher pass@1.

Uniform improvement across training problems. The above comparisons on the evolution of pass@32 are also corroborated by the evolution of the pass@1 histogram across training problems in Figure 13(a). Here, we see that RL only improves pass@1 on a narrow band of problems (presumably those that were lucky enough to see a positive sample early in the training run), learning to solve them fully, while making little progress on others. Increasing n (number of traces per problem) allays this to some extent because now we have a higher chance of seeing positive rewards for a greater fraction of the training batch. In fact, prior work (Schaul et al., 2019) in deep RL goes on to show that making non-uniform progress across the training environments are a “source of plateaus” in RL, a phenomenon they call ray interference, where over-optimizing performance on a subset of training problems may severely hurt exploration on the remainder of the training set. Avoiding this, PrefixRL is able to improve pass@1 on a larger fraction of training problems simply by collecting rewards on their prefixed versions and relying on back-generalization to make progress on the no-prefix versions.

PrefixRL is still effective with off-policy prefixes from a different model family. Up to now, we have used PrefixRL with prefixes from the same base LLM that we ran RL with. In many practical settings, it is easier to source off-policy prefixes from another model that is substantially different in training data or architecture; for example, from the open-source community. Motivated by this, we compare PrefixRL with mid-training and standard RL when the off-policy data is sourced from a policy (Qwen3-4b-instruct) that is different from the base LLM (Llama3.1-8B-instruct) in Figure 12. Despite the Qwen prefixes being more off-policy due to having a different training dataset and architecture, they end up being similarly effective on both train and test problems as off-policy prefixes from the Llama base model we run RL on. Note that the Llama base model required more initial compute for rejection sampling to collect the off-policy prefixes, which accounts for most of the difference in the curves. In Appendix D.3,

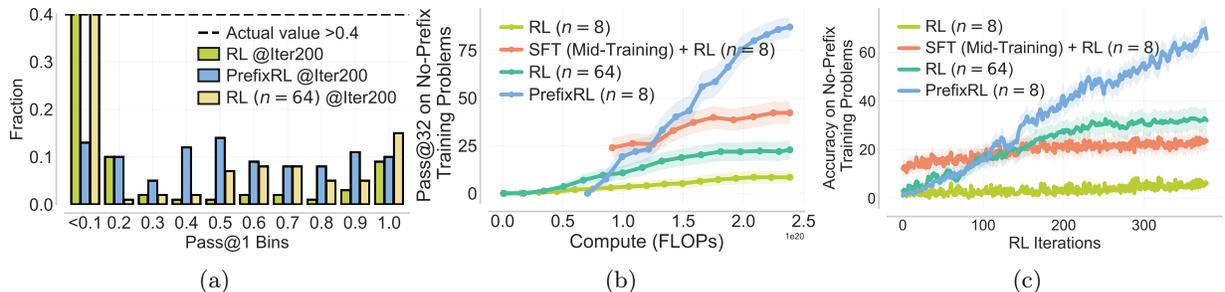


Figure 11 Pass@ k evolution and iteration-matched comparison on training problems. (a) *Uniform pass@1 improvement*: By design, the base LLM places all training problems in the pass@1 bin <math><0.1</math>. After 200 iterations, PrefixRL yields the most uniform gains across problems, while RL concentrates improvements on a small subset with rare successes; increasing n partially mitigates this. (b) *New problems solved*: compute-matched pass@32 plots indicate that PrefixRL steadily expands the set of solvable problems rather than merely converting a fixed pass@ k (for small k) into higher pass@1, whereas the baselines largely saturate on pass@32. (c) *Fair baselines*: iteration-matched reward curves confirm stable training across methods, so the compute-matched gains are not explained by unstable or degenerate baseline runs.

we also find that the reverse direction is effective, where we use Llama to generate off-policy prefixes for Qwen.

5.2 Training Dynamics of PrefixRL are More Stable Compared to Typical Off-Policy Methods

PrefixRL conditions on off-policy data instead of supervising on it akin to importance-weighted off-policy RL, which can often destabilize optimization as we saw in Section 3. As a result, PrefixRL achieves stable RL training while leveraging off-policy data. In this section, we investigate RL training dynamics (like gradient variance) underlying the instabilities of off-policy RL, and the signal-to-noise ratio of policy gradients observed by PrefixRL as a result of off-policy traces placing the current RL policy in states with non-zero advantages.

PrefixRL leverages off-policy data while preserving entropy for RL exploration. It is common practice to mid-train models on reasoning traces from more capable models or its own traces from base model inference or past RL runs (Zelikman et al., 2022). This is typically done to prime the RL initialization and improve coverage over high-reward regions. But this typically comes at the cost of the model’s token-level entropy which impacts its ability to explore during RL. Figure 13 (left) shows the average token-level entropy of the model’s next-token distributions during the RL run. Doing SFT on off-policy data causes the entropy to dramatically decrease and this only continues to drop further to values as low as 0.01 during RL, suggesting that reward maximization during RL is mostly just sharpening the distribution over correct traces the base model can already sample (after running SFT). In contrast, PrefixRL preserves most of the token-level entropy while still making use of off-policy data.

Fewer all-negative batches on hard problems. Figure 13 (middle) plots the fraction of *all-negative* problems (i.e., problems in a training batch for which all n on-policy traces on that problem receive zero reward), measured *only*

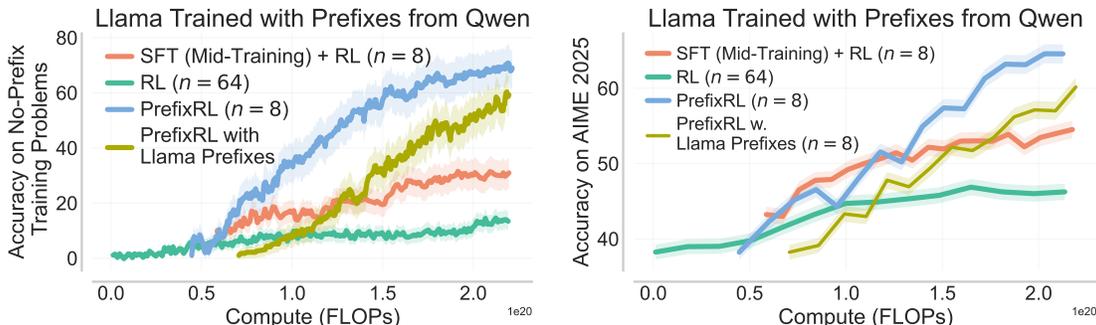


Figure 12 PrefixRL is still effective with off-policy prefixes from a different model family. We train Llama3.1-8b-instruct on prefixed-problems constructed using prefixes of rejection sampled traces from the Qwen-3-4b-instruct (*left*). Interestingly, even though prefixes are more out-of-distribution than those rejection sampled from Llama3.1-8b itself, they are still equally effective in improving on hard problems compared to when the prefixes are sampled from Llama3.1-8b (olive green line). We also plot the performance on AIME when we train Llama with prefixes from Qwen (*right*).

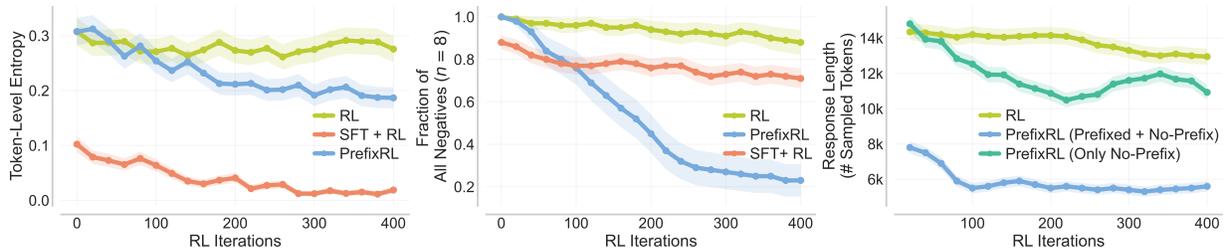


Figure 13 RL Training Dynamics. (Left) Compared to RL and PrefixRL, supervised learning on off-policy data (SFT+RL) leads to much lower token-level entropy during RL, which may hurt exploration. (Middle) PrefixRL has much lower “all negative ratio”, or the number of prompts with all zero rewards (zero advantage) during training. Note that the ratio is measured only on problems without any prefixes. (Right) When we look at the average response length on only no-prefix problems (green), we see PrefixRL generates shorter responses, mostly because correct solutions are usually shorter and standard RL has worse accuracy. Blue shows the average length across all problems, which is lower due to prefixing. This shortening is one way that PrefixRL reduces gradient variance and gains on compute efficiency compared to standard RL.

over no-prefix problems. PrefixRL consistently reduces this ratio relative to on-policy RL. This reduction reveals an underlying shift in the unconditional policy: as training on strategy-revealing states proceeds, the model becomes more likely to *enter* regions where non-zero advantages are attainable (either due to the prefix revealing “useful” strategies that are further reinforced with positive rewards, or revealing likely but incorrect strategies that are quickly down-weighted and unlearned (Section 4.2)), thereby breaking the stalling regime (Section 2).

PrefixRL achieves better accuracies with fewer sampled tokens. Figure 13 (right) tracks the average sampled tokens per trace across all prompts (with and without prefixes). On the no-prefix problems, PrefixRL eventually maintains a lower sampled token budget per trace while achieving higher reward rates, implying better *iteration efficiency*. This is perhaps expected since correct traces are biased to be shorter and qualitatively, once the model internalizes the strategy, it reaches decisive steps earlier, which reduces “unproductive wandering” later in the horizon. Moreover, since PrefixRL trains on 3:1 mixture of prefixed to no-prefix problems, the average number of tokens sampled (across all problem types) per batch is less than $1/2$ of the RL run. As a result of short length RL, the gradient variance for PrefixRL is much lower compared to standard RL, despite the PrefixRL not being biased, *i.e.*, it shares the same set of optimal policies as standard RL (Section 3.1).

PrefixRL has higher signal-to-noise ratio during RL training. Figure 14 shows the gradient norm and gradient standard deviation of three methods: PrefixRL, importance weighted off-policy RL, and standard RL. For off-policy RL, in addition to the on-policy samples, we also compute the batch gradient using the off-policy traces for the problems in the batch. To correct for the distribution mismatch, one can use the unbiased importance-weighted update that would apply the correction at the sequence level and without clipping the importance weights (Tan et al., 2025). On the other hand, such an approach suffers from very high variance and following prior work (Liu et al., 2025a) we choose to lower the variance with a biased token-level correction

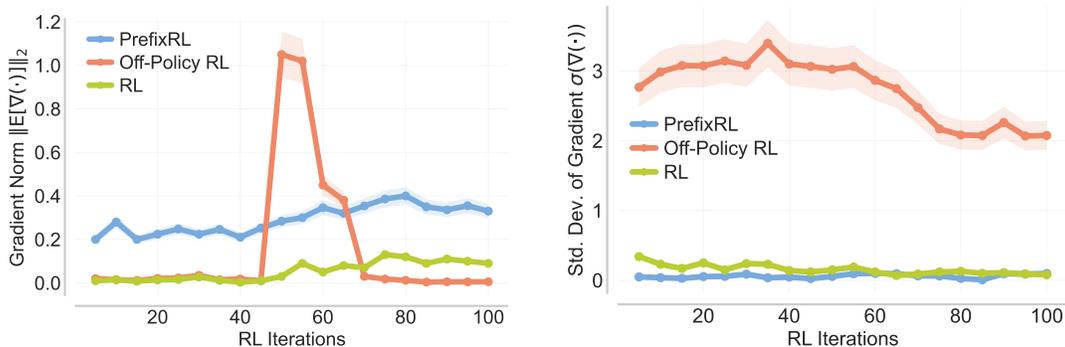


Figure 14 PrefixRL signal-to-noise ratio and safely training on off-policy suffixes. PrefixRL simultaneously has higher gradient norm (left) and lower gradient variance (right) than both standard RL and importance-weighted off-policy RL. For off-policy RL in particular, the importance weighting causes the gradient variance to be much higher and causes a gradient norm spike, which is a sign of training instability. PrefixRL avoids these issues while still leveraging off-policy data.

and weight-clipping (see Appendix D.1 for details). For each approach, we compute the norm of the expected gradient norm and the standard deviation of the sampled gradient by keeping a moving average of the first and second moments of the gradient (coordinate-wise), then summing across the coordinates. For the variance / standard deviation, this is equivalent to estimating the trace of the covariance matrix (see Appendix D.4).

Off-policy RL suffers from very noisy and biased (due to token-level correction) gradients and even suffers a gradient norm spike, a sign of training instability. Since the off-policy traces are all rejection sampled, they are very unlikely under the base policy as well subsequent RL iterates. As a result, the token-level importance weighting term multiplies the gradient on the off-policy tokens with a very low value (as low as 0.001), which makes the norm of the expected gradient also small. Consequently, despite using off-policy data, the gradient norm for the off-policy data is similar in magnitude to that of on-policy RL, at least until the spike (Figure 14 (left)). Recall that in Figure 13 (right), we showed that PrefixRL samples fewer tokens during training, which also contributes to the lower gradient variance as we see in Figure 14 (right). At the same time, since the off-policy prefixes place the current RL policy in states which are likely to see a non-zero advantage (Figure 13 (middle)), the expected gradient norm for PrefixRL is also higher. To conclude, PrefixRL simultaneously achieves a higher gradient norm and lower gradient variance, implying that it benefits from a higher signal-to-noise ratio (norm of the expected gradient over batch gradient’s standard deviation).

6 Related Work and Discussion

Learning from off-policy LLM rollouts. When on-policy search stalls due to over-sharpening or “over-thinking,” a common approach is to supervise on human or oracle-provided traces (Lightman et al., 2023; Corrado et al., 2024), but teacher-driven methods inherit the teacher’s capacity limit (Agarwal et al., 2024) and often require reward shaping (Yan et al., 2025), entropy control (Wang et al., 2025a), and heavy hyperparameter tuning (Zhang et al., 2025a); moreover, for hard problems, long model-compatible chains of thought are scarce and mismatches can collapse response diversity (Kang et al., 2024b). When off-policy data come from “close enough” (in KL divergence) policies as in Async RL, reuse becomes more efficient (Fu et al., 2025; Khatri et al., 2025), yet large importance weights and high gradient variance pose instability risks (Agarwal et al., 2021), so practical systems cap behavior-policy staleness to only a few RL iterations (Sheng et al., 2024). These constraints motivate approaches that do not treat off-policy trajectories as direct supervision targets; related directions condition on subgoals or plans (Hong et al., 2025), higher-level abstractions (Qu et al., 2025b), or partial solutions (Amani et al., 2025; Chen et al., 2025b; Li et al., 2025). Different from the above, PrefixRL conditions on off-policy prefixes from long-thinking traces, as opposed to training on them before or during RL. Instead of suffering from instability due to supervising on off-policy data (Sections 3 and 5.2), PrefixRL benefits from them via back-generalization.

Conditioning on hints to improve on-policy RL. A related line of work augments prompts with hints or partially revealed human solutions to “guide” on-policy RL (Chen et al., 2025b; Li et al., 2025; Qu et al., 2025a). AdaBack (Amani et al., 2025) adaptively searches for the minimal hint that improves performance over human-written solutions, but is hard to scale to long-context “thinking models” and large datasets. Similarly, QuestA (Li et al., 2025) uses answer-hinted prompts derived from human solutions. In general, these methods are only feasible when we have access to solution traces written by a human or a more capable teacher model. In contrast, PrefixRL enables a self-improvement loop by not relying on external sources and instead reusing compute from prior models. Moreover, our work also analyzes the back-generalization phenomenon that may be shared across these methods, showing that it cannot be explained by some of the “stitching” arguments made in prior works (Zhang et al., 2025b).

Resetting to off-policy states in RL. The idea of “resetting” current RL policy to off-policy states is not new in RL (Kakade, 2003; Bagnell et al., 2003; Nair et al., 2018; Salimans and Chen, 2018; Yin et al., 2022; Uchendu et al., 2023; Silver et al., 2016a,b; Agarwal et al., 2019; Daumé III and Marcu, 2005; Daumé III et al., 2009). Chang et al. (2024) also applied the resetting idea for post-training LLMs with human feedback. While similar in principle, their work does not study computational gains accounting for the initial compute spent on collecting off-policy traces. In fact, in their case the data is human labeled whereas our work lies more in a self-improvement setting. Our contribution is to instantiate this perspective for RL of reasoning LLMs. We show that a relatively small dataset of *correct* off-policy traces is sufficient to enable effective resets that make hard, low-pass-rate problems trainable even when on-policy rollouts almost never succeed. We also show that PrefixRL yields a strictly better allocation of compute, even after accounting for the inference cost of collecting the off-policy traces.

Improving exploration on hard problems in LLM reasoning. Small models fine-tuned with RL can outperform much larger base models (Liu et al., 2025b; Luo et al., 2025), largely by reinforcing long chain-of-thought behaviors like self-correction (Qu et al., 2024) and reflection (Gandhi et al., 2025). Yet, without careful controls, RL often under-explores and leaves hard instances underprobed; empirically this appears as a drop in $\text{pass}@k$ versus the base model (Yue et al., 2025; Zhao et al., 2025). One response is to regularize training to curb over-sharpening via intrinsic-motivation bonuses (Gao et al., 2025), entropy (Wang et al., 2025b), count-based signals (Song et al., 2025), or objectives that directly optimize $\text{pass}@n$ (Chow et al., 2024; Balashankar et al., 2025), but these still inherit sparse-reward limits and depend on easy problems for signal (He et al., 2024). A complementary thread (Setlur et al., 2025b) exploits base-model asymmetries, *e.g.*, the verification–generation gap (Setlur et al., 2025a; Song et al., 2024), and can combine with negative-gradient dynamics to chain such asymmetries across updates (Zhu et al., 2025); nevertheless, models often “under-think” (Wang et al., 2025c), persisting with wrong high-level plans despite more rollouts. In contrast, PrefixRL avoids carefully tuned auxiliary exploration objectives by reshaping the start-state distribution directly. Empirically, we do not observe the $\text{pass}@k$ regressions often induced by over-sharpening or over exploration with token-level entropy regularizers. In the worst case, uninformative prefixes recover standard on-policy RL (Section 4.3).

7 Conclusion

We showed that PrefixRL can recycle FLOPs in the form of correct off-policy traces from prior sampling or training, which may be rare and have taken a lot of compute to obtain. PrefixRL achieves this by running on-policy RL conditioned on off-policy prefixes that place the current RL policy in higher-rewarding states and boost the learning signal. As a result, it makes on-policy RL more efficient on hard problems. We see PrefixRL as a step beyond the typical off-policy paradigm of supervising directly on off-policy data, instead relying on the powerful back-generalization mechanism to incorporate off-policy data as conditioning context while doing on-policy updates.

Acknowledgements. We would like to thank Sean Bell, Ankur Pai, Aviral Kumar, Risabh Agarwal, Saurabh Garg, Wen Sun, Sergey Levine, Yuxiao Qu, Ian Wu, Rohan Maheshwari, and Yuandong Tian for helpful discussions and thoughtful feedback on our work.

References

- Alekh Agarwal, Nan Jiang, Sham M. Kakade, and Wen Sun. Reinforcement learning: Theory and algorithms. Technical report / book draft, 2019. URL <https://rltheorybook.github.io/>. Online manuscript; frequently updated.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes, 2024. URL <https://arxiv.org/abs/2306.13649>.
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms, 2024. URL <https://arxiv.org/abs/2402.14740>.
- Mohammad Hossein Amani, Aryo Lotfi, Nicolas Mario Baldwin, Samy Bengio, Mehrdad Farajtabar, Emmanuel Abbe, and Robert West. RL for reasoning by adaptively revealing rationales, 2025. URL <https://arxiv.org/abs/2506.18110>.
- Chenxin An, Zhihui Xie, Xiaonan Li, Lei Li, Jun Zhang, Shansan Gong, Ming Zhong, Jingjing Xu, Xipeng Qiu, Mingxuan Wang, and Lingpeng Kong. Polaris: A post-training recipe for scaling reinforcement learning on advanced reasoning models, 2025. URL <https://hkunlp.github.io/blog/2025/Polaris>.
- J. Andrew Bagnell, Sham Kakade, Andrew Y. Ng, and Jeff Schneider. Policy search by dynamic programming. In *Advances in Neural Information Processing Systems*, 2003. Often cited as NIPS 2003 / proceedings volume published in 2004.
- Yu Bai, Chi Jin, and Tiancheng Yu. Near-optimal reinforcement learning with self-play. *Advances in neural information processing systems*, 33:2159–2170, 2020.

- Ananth Balashankar, Ziteng Sun, Jonathan Berant, Jacob Eisenstein, Michael Collins, Adrian Hutter, Jong Lee, Chirag Nagpal, Flavien Prost, Aradhana Sinha, Ananda Theertha Suresh, and Ahmad Beirami. Inalign: Inference-aware language model alignment, 2025. URL <https://arxiv.org/abs/2412.19792>.
- Jonathan D Chang, Wenhao Shan, Owen Oertell, Kianté Brantley, Dipendra Misra, Jason D Lee, and Wen Sun. Dataset reset policy optimization for rlhf. *arXiv preprint arXiv:2404.08495*, 2024.
- Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu, Chao Wang, Cheng Zhu, et al. Minimax-m1: Scaling test-time compute efficiently with lightning attention. *arXiv preprint arXiv:2506.13585*, 2025a.
- Justin Chih-Yao Chen, Becky Xiangyu Peng, Prafulla Kumar Choubey, Kung-Hsiang Huang, Jiaxin Zhang, Mohit Bansal, and Chien-Sheng Wu. Nudging the boundaries of llm reasoning, 2025b. URL <https://arxiv.org/abs/2509.25666>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Yinlam Chow, Guy Tennenholtz, Izzeddin Gur, Vincent Zhuang, Bo Dai, Sridhar Thiagarajan, Craig Boutilier, Rishabh Agarwal, Aviral Kumar, and Aleksandra Faust. Inference-aware fine-tuning for best-of-n sampling in large language models. *arXiv preprint arXiv:2412.15287*, 2024.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025.
- Nicholas E. Corrado, Yuxiao Qu, John U. Balis, Adam Labiosa, and Josiah P. Hanna. Guided data augmentation for offline reinforcement learning and imitation learning, 2024. URL <https://arxiv.org/abs/2310.18247>.
- Hal Daumé III and Daniel Marcu. Learning as search optimization: Approximate large margin methods for structured prediction. In *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, 2005.
- Hal Daumé III, John Langford, and Daniel Marcu. Search-based structured prediction. *Machine Learning*, 75:297–325, 2009. doi: 10.1007/s10994-009-5106-x.
- Thomas Degris, Martha White, and Richard S Sutton. Off-policy actor-critic. *arXiv preprint arXiv:1205.4839*, 2012.
- Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun Mei, Jiashu Wang, Tongkai Yang, Binhang Yuan, and Yi Wu. Areal: A large-scale asynchronous reinforcement learning system for language reasoning, 2025. URL <https://arxiv.org/abs/2505.24298>.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. *arXiv preprint arXiv:1812.02900*, 2018.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D. Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars, 2025. URL <https://arxiv.org/abs/2503.01307>.
- Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. Omni-math: A universal olympiad level mathematic benchmark for large language models, 2024. URL <https://arxiv.org/abs/2410.07985>.
- Jingtong Gao, Ling Pan, Yejing Wang, Rui Zhong, Chi Lu, Qingpeng Cai, Peng Jiang, and Xiangyu Zhao. Navigate the unknown: Enhancing llm reasoning with intrinsic motivation guided exploration, 2025. URL <https://arxiv.org/abs/2505.17621>.
- Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su, Wanxia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave, Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy, Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for reasoning models, 2025. URL <https://arxiv.org/abs/2506.04178>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.

- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems, 2024. URL <https://arxiv.org/abs/2402.14008>.
- Josef Hofbauer and Sylvain Sorin. Best response dynamics for continuous zero-sum games. *Discrete and Continuous Dynamical Systems Series B*, 6(1):215, 2006.
- Joey Hong, Anca Dragan, and Sergey Levine. Planning without search: Refining frontier llms with offline goal-conditioned rl. *arXiv preprint arXiv:2505.18098*, 2025.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the nineteenth international conference on machine learning*, pages 267–274, 2002.
- Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Sham Machandranath Kakade. *On the Sample Complexity of Reinforcement Learning*. PhD thesis, University College London, 2003.
- Katie Kang, Amrith Setlur, Dibya Ghosh, Jacob Steinhardt, Claire Tomlin, Sergey Levine, and Aviral Kumar. What do learning dynamics reveal about generalization in llm reasoning? *arXiv preprint arXiv:2411.07681*, 2024a.
- Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. Unfamiliar finetuning examples control how language models hallucinate, 2024b.
- Devvrit Khatry, Lovish Madaan, Rishabh Tiwari, Rachit Bansal, Sai Surya Duvvuri, Manzil Zaheer, Inderjit S. Dhillon, David Brandfonbrener, and Rishabh Agarwal. The art of scaling reinforcement learning compute for llms, 2025. URL <https://arxiv.org/abs/2510.13786>.
- Vijaymohan Konda and John N. Tsitsiklis. *Actor-Critic Algorithms*. PhD thesis, USA, 2002. AAI0804543.
- Jiazheng Li, Hong Lu, Kaiyue Wen, Zaiwen Yang, Jiakuan Gao, Hongzhou Lin, Yi Wu, and Jingzhao Zhang. Questa: Expanding reasoning capacity in llms via question augmentation. *arXiv preprint arXiv:2507.13266*, 2025.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023.
- Jiacai Liu, Yingru Li, Yuqian Fu, Jiawei Wang, Qian Liu, and Yu Shen. When speed kills stability: Demystifying rl collapse from the inference-training mismatch, 2025a. URL <https://yingru.notion.site/>.
- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models, 2025b. URL <https://arxiv.org/abs/2505.24864>.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. URL <https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>. Notion Blog.
- A Rupam Mahmood, Hado P Van Hasselt, and Richard S Sutton. Weighted importance sampling for off-policy learning with linear function approximation. *Advances in neural information processing systems*, 27, 2014.
- Ashvin Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, Shane Legg, and Volodymyr Mnih. Overcoming exploration in reinforcement learning with demonstrations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018. doi: 10.1109/ICRA.2018.8463167. arXiv:1709.10089.
- Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching language model agents how to self-improve. *arXiv preprint arXiv:2407.18219*, 2024.
- Yuxiao Qu, Amrith Setlur, Virginia Smith, Ruslan Salakhutdinov, and Aviral Kumar. How to explore to scale rl training of llms on hard problems? [urlhttps://blog.ml.cmu.edu/2025/11/26/how-to-explore-to-scale-rl-training-of-llms-on-hard-problems](https://blog.ml.cmu.edu/2025/11/26/how-to-explore-to-scale-rl-training-of-llms-on-hard-problems), 2025a. CMU MLD Blog.
- Yuxiao Qu, Anikait Singh, Yoonho Lee, Amrith Setlur, Ruslan Salakhutdinov, Chelsea Finn, and Aviral Kumar. Rlad: Training llms to discover abstractions for solving reasoning problems, 2025b. URL <https://arxiv.org/abs/2510.02263>.
- Tim Salimans and Richard Chen. Learning Montezuma’s Revenge from a single demonstration. *arXiv preprint*, 2018. arXiv:1812.03381.

- Tom Schaul, Diana Borsa, Joseph Modayil, and Razvan Pascanu. Ray interference: a source of plateaus in deep reinforcement learning. *CoRR*, abs/1904.11455, 2019. URL <http://arxiv.org/abs/1904.11455>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. Rl on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold. *Advances in Neural Information Processing Systems*, 37:43000–43031, 2024.
- Amrith Setlur, Nived Rajaraman, Sergey Levine, and Aviral Kumar. Scaling test-time compute without verification or rl is suboptimal, 2025a. URL <https://arxiv.org/abs/2502.12118>.
- Amrith Setlur, Matthew Y. R. Yang, Charlie Snell, Jeremy Greer, Ian Wu, Virginia Smith, Max Simchowitz, and Aviral Kumar. e3: Learning to explore enables extrapolation of test-time compute for llms, 2025b. URL <https://arxiv.org/abs/2506.09026>.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016a. doi: 10.1038/nature16961.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016b.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.
- Yuda Song, Yifei Zhou, Ayush Sekhari, J Andrew Bagnell, Akshay Krishnamurthy, and Wen Sun. Hybrid rl: Using both offline and online data can make rl efficient. *arXiv preprint arXiv:2210.06718*, 2022.
- Yuda Song, Hanlin Zhang, Carson Eisenach, Sham Kakade, Dean Foster, and Udaya Ghai. Mind the gap: Examining the self-improvement capabilities of large language models. *arXiv preprint arXiv:2412.02674*, 2024.
- Yuda Song, Julia Kempe, and Remi Munos. Outcome-based exploration for llm reasoning, 2025. URL <https://arxiv.org/abs/2509.06941>.
- Hongze Tan, Jianfei Pan, Jinghao Lin, Tao Chen, Zhihang Zheng, Zhihao Tang, and Haihua Yang. Gtpo and grpo-s: Token and sequence-level reward shaping with policy entropy. *arXiv preprint arXiv:2508.04349*, 2025.
- Ikechukwu Uchendu, Yujia Li, Shibin Yuan, Yuke Zhu, Sergey Levine, Karol Hausman, and Chelsea Finn. Jump-start reinforcement learning. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, volume 202 of *Proceedings of Machine Learning Research*, 2023.
- Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025a.
- Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, Weizhu Chen, Shuohang Wang, Simon Shaolei Du, and Yelong Shen. Reinforcement learning for reasoning in large language models with one training example, 2025b. URL <https://arxiv.org/abs/2504.20571>.
- Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, et al. Thoughts are all over the place: On the underthinking of o1-like llms. *arXiv preprint arXiv:2501.18585*, 2025c.
- Zengzhi Wang, Fan Zhou, Xuefeng Li, and Pengfei Liu. Octothinker: Mid-training incentivizes reinforcement learning scaling, 2025d. URL <https://arxiv.org/abs/2506.20512>.
- Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. Learning to reason under off-policy guidance, 2025. URL <https://arxiv.org/abs/2504.14945>.

- Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity. In *Proceedings of the 18th Annual Symposium on Foundations of Computer Science*. IEEE Computer Society, 1977. doi: 10.1109/SFCS.1977.24.
- Dong Yin, Brendan Hao, Yasin Abbasi-Yadkori, Nevena Lazic, and Csaba Szepesvári. Efficient local planning with linear function approximation. In *Proceedings of The 33rd International Conference on Algorithmic Learning Theory (ALT)*, volume 167 of *Proceedings of Machine Learning Research*, 2022.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025. URL <https://arxiv.org/abs/2503.14476>.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?, 2025. URL <https://arxiv.org/abs/2504.13837>.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Wenhao Zhang, Yuexiang Xie, Yuchang Sun, Yanxi Chen, Guoyin Wang, Yaliang Li, Bolin Ding, and Jingren Zhou. On-policy rl meets off-policy experts: Harmonizing supervised fine-tuning and reinforcement learning via dynamic weighting, 2025a. URL <https://arxiv.org/abs/2508.11408>.
- Xuechen Zhang, Zijian Huang, Yingcong Li, Chenshun Ni, Jiasi Chen, and Samet Oymak. Bread: Branched rollouts from expert anchors bridge sft & rl for reasoning. *arXiv preprint arXiv:2506.17211*, 2025b.
- Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: RL post-training amplifies behaviors learned in pretraining, 2025. URL <https://arxiv.org/abs/2504.07912>.
- Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. The surprising effectiveness of negative reinforcement in llm reasoning, 2025. URL <https://arxiv.org/abs/2506.01347>.

Appendix

A Additional Notation

Markov decision process. We use \mathbf{x} to denote an input problem and $\mathbf{y} = (y_1, \dots, y_H)$ for a response of H tokens, and if $\mathbf{y} \sim \pi(\cdot | \mathbf{x})$, then \mathbf{y} is sampled auto-regressively from the LLM π fed with input \mathbf{x} . Each token in this response \mathbf{y} belongs to a set of tokens or actions \mathcal{A} . The state \mathbf{s}_h at time step h is given by $(\mathbf{x}, y_1, y_2, \dots, y_h)$, where the initial state \mathbf{s}_0 is just the problem \mathbf{x} . The set of states across all time steps is denoted by the class \mathcal{S} . We use d_h^π to denote the distribution over states \mathbf{s}_h at time step h by rolling out the policy π auto-regressively for h time steps. For compactness, we write the trajectory-level log-likelihood $\log \pi(\mathbf{y} | \mathbf{x}) = \sum_{h=1}^{|\mathbf{y}|} \log \pi(y_h | \mathbf{x}, \mathbf{y}_{<h})$. For each problem we have access to outcome reward function $r(\mathbf{x}_i, \mathbf{y})$ to check whether the final answer in response \mathbf{y} is correct/incorrect (1/0) for the problem \mathbf{x}_i (e.g., by matching the boxed answer in the end of \mathbf{y} for math problems).

Dataset of hard problems and off-policy traces. We use \mathcal{D} to denote a dataset of N hard problems $\mathcal{D} =: \{\mathbf{x}_i\}_{i=1}^N$. We use π^0 to denote the base pre-trained LLM that we initialize the RL algorithm, π^t as the policy after t RL iterations and \mathcal{D}_{off} as the dataset of off-policy traces. Finally, we define the pass rate @ k for problem \mathbf{x} and LLM π as $\mathbb{E}_{\mathbf{y}_1, \dots, \mathbf{y}_k \sim \pi(\cdot | \mathbf{x})} \max(\{r(\mathbf{x}, \mathbf{y}_i)\}_{i=1}^k)$. In the main paper, we define the set of hard problems as those with pass@ $k \approx 0$ under the base LLM π^0 . See Section 5 for how we select these low pass rate hard problems for training.

B Omitted Proofs

In this section, we present proofs for our theoretical results in Section 3.1. We begin with the proof for Theorem 3.2 which implies that the PrefixRL objective is consistent with standard RL, and any solution for our PrefixRL objective, is also a maximizer of the standard RL objective which just maximizes $J(\pi)$. Following this, we show the proof for Theorem 3.3 which bounds the suboptimality gap of an algorithm using natural policy gradient (NPG) to optimize the PrefixRL objective. Note that this is slightly different from the policy gradients we use in practice, but is nevertheless insightful in informing a formal mental model for the gains behind PrefixRL. We also provide a proof for Proposition 3.4 that lower bounds the performance gap between PrefixRL and standard RL in the worst case. In the end we list auxiliary lemmas useful for analysis.

B.1 Proof of Theorem 3.2

Let the standard (no-prefix) RL objective be

$$J(\pi) = \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})], \quad J^* = \max_{\pi \in \Pi} J(\pi). \quad (\text{B.1})$$

For each $\mathbf{x} \in \mathcal{D}$, Assumption 3.1 gives a single correct trace \mathbf{y}^x such that $(\mathbf{x}, \mathbf{y}^x) \in \mathcal{D}_{\text{off}}$. Given any cut index h , define the prefixed problem

$$\mathbf{x}_{\text{pre}} = \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h}), \quad (\text{B.2})$$

and define its reward by evaluating the full transcript:

$$r(\mathbf{x}_{\text{pre}}, \mathbf{z}) := r(\mathbf{x}, (\mathbf{y}^x)_{1:h} \circ \mathbf{z}).$$

The PrefixRL objective is

$$J_{\text{pre}}(\pi) = \sum_{\mathbf{x}_{\text{pre}} \in \mathcal{D}_{\text{pre}}} \mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x}_{\text{pre}})} [r(\mathbf{x}_{\text{pre}}, \mathbf{z})] + \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})]. \quad (\text{B.3})$$

A uniform upper bound. Fix any $\mathbf{x} \in \mathcal{D}$ and any prefix $(\mathbf{y}^x)_{1:h}$ used to form a prefixed problem. For any policy $\pi \in \Pi$ defined on such prefixed problems, construct a policy $\tilde{\pi} \in \Pi$ for the no-prefix problem \mathbf{x} that deterministically emits the prefix $(\mathbf{y}^x)_{1:h}$ and then samples the suffix from $\pi(\cdot | \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h}))$. By the reward definition,

$$\mathbb{E}_{\mathbf{y} \sim \tilde{\pi}(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})] = \mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h}))} [r(\text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h}), \mathbf{z})].$$

Therefore, for every π and every such prefixed problem $\mathbf{x}_{\text{pre}} = \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h})$,

$$\mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x}_{\text{pre}})} [r(\mathbf{x}_{\text{pre}}, \mathbf{z})] \leq \max_{\pi' \in \Pi} \mathbb{E}_{\mathbf{y} \sim \pi'(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})]. \quad (\text{B.4})$$

Let $m(\mathbf{x})$ be the number of prefixed problems in \mathcal{D}_{pre} derived from \mathbf{x} . Summing the above inequality over all prefixed problems and grouping by their originating \mathbf{x} gives

$$\sum_{\mathbf{x}_{\text{pre}} \in \mathcal{D}_{\text{pre}}} \mathbb{E}[r(\mathbf{x}_{\text{pre}}, \cdot)] \leq \sum_{\mathbf{x} \in \mathcal{D}} m(\mathbf{x}) \max_{\pi' \in \Pi} \mathbb{E}_{\mathbf{y} \sim \pi'(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})]. \quad (\text{B.5})$$

Define the constant

$$C := \sum_{\mathbf{x} \in \mathcal{D}} m(\mathbf{x}) \max_{\pi' \in \Pi} \mathbb{E}_{\mathbf{y} \sim \pi'(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})], \quad (\text{B.6})$$

which is independent of π . Then for every $\pi \in \Pi$ we have the uniform upper bound

$$J_{\text{pre}}(\pi) \leq C + J(\pi). \quad (\text{B.7})$$

Tightness using Assumption 3.1. By realizability in Assumption 3.1, there exists a policy $\mu \in \Pi$ such that

$$\mu(\mathbf{y}^x | \mathbf{x}) = 1, \quad \forall \mathbf{x} \in \mathcal{D}. \quad (\text{B.8})$$

In particular, for any cut h , conditioning μ on the prefix $(\mathbf{y}^x)_{1:h}$ yields the deterministic continuation:

$$\mu(\mathbf{z} | \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h})) = \mathbf{1}\{\mathbf{z} = (\mathbf{y}^x)_{h+1:|\mathbf{y}^x|}\}. \quad (\text{B.9})$$

Therefore, on every prefixed problem $\mathbf{x}_{\text{pre}} = \text{concat}(\mathbf{x}, (\mathbf{y}^x)_{1:h})$,

$$\mathbb{E}_{\mathbf{z} \sim \mu(\cdot | \mathbf{x}_{\text{pre}})} [r(\mathbf{x}_{\text{pre}}, \mathbf{z})] = r(\mathbf{x}, \mathbf{y}^x) = 1, \quad (\text{B.10})$$

where the last equality uses correctness of \mathcal{D}_{off} in Assumption 3.1. Hence,

$$\sum_{\mathbf{x}_{\text{pre}} \in \mathcal{D}_{\text{pre}}} \mathbb{E}_{\mathbf{z} \sim \mu(\cdot | \mathbf{x}_{\text{pre}})} [r(\mathbf{x}_{\text{pre}}, \mathbf{z})] = \sum_{\mathbf{x} \in \mathcal{D}} m(\mathbf{x}). \quad (\text{B.11})$$

Moreover, by (B.8) and correctness,

$$J(\mu) = \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{\mathbf{y} \sim \mu(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})] = \sum_{\mathbf{x} \in \mathcal{D}} r(\mathbf{x}, \mathbf{y}^x) = |\mathcal{D}|. \quad (\text{B.12})$$

Therefore $J^* = |\mathcal{D}|$ and $J(\mu) = J^*$. Since rewards are in $[0, 1]$, for each \mathbf{x} we have $\max_{\pi' \in \Pi} \mathbb{E}_{\mathbf{y} \sim \pi'(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})] = 1$, and thus the constant simplifies to $C = \sum_{\mathbf{x} \in \mathcal{D}} m(\mathbf{x})$. Combining with (B.11) yields

$$J_{\text{pre}}(\mu) = C + J(\mu) = C + J^*. \quad (\text{B.13})$$

Step 3: Concluding the consistency of PrefixRL objective. Let $\hat{\pi} \in \text{argmax}_{\pi \in \Pi} J_{\text{pre}}(\pi)$ be any maximizer of the PrefixRL objective. By optimality of $\hat{\pi}$, (B.13), and the upper bound (B.7),

$$C + J(\hat{\pi}) \geq J_{\text{pre}}(\hat{\pi}) \geq J_{\text{pre}}(\mu) = C + J^*. \quad (\text{B.14})$$

Cancelling C yields $J(\hat{\pi}) \geq J^*$, hence $J(\hat{\pi}) = J^*$. Therefore $\hat{\pi} \in \text{argmax}_{\pi \in \Pi} J(\pi)$, proving that any maximizer of the PrefixRL objective also maximizes the standard RL objective.

Algorithm 1 PrefixRL with Natural Policy Gradients

Require: Base policy π^0 , off-policy data \mathcal{D}_{off} , horizon H , iterations T , step size η , Q function class \mathcal{F} .

- 1: Initialize the iterative algorithm with base policy: $\pi^1 \leftarrow \pi^0$.
- 2: **for** $t = 1, \dots, T$ **do**
- 3: Initialize dataset $\mathcal{D}_t \leftarrow \{\}$.
- 4: **for** $i = 1 \dots n$ **do**
- 5: Sample (s_h, a_h^{off}) uniformly across state-action pairs in \mathcal{D}_{off} . ▷ sample prefixed problem
- 6: $a_h \leftarrow a_h^{\text{off}}$ with probability $1/2$ and $\sim \pi^t(\cdot | s_h)$ otherwise.
- 7: Execute $\pi^t(\cdot | s_h, a_h)$ from step $h+1$ through H to obtain the full trace with reward r .
- 8: $\mathcal{D}_t \leftarrow \mathcal{D}_t \cup (s_h, a_h, r)$.
- 9: **Critic fit (regression oracle):**
- 10: $\hat{Q}^t \leftarrow \operatorname{argmin}_{f \in \mathcal{F}} \sum_{(s,a,r) \in \mathcal{D}_t} (f(s,a) - r)^2$.
- 11: **Natural policy update (mirror ascent):** ▷ performed state-wise
- 12: $\pi^{t+1}(\cdot | s) \leftarrow \operatorname{argmin}_p \langle -\hat{Q}^t(s, \cdot), p \rangle + \frac{1}{\eta} \text{KL}(p \| \pi^t(\cdot | s))$.
- 13: **end for**
- 14: **end for**
- 15: **return** $\bar{\pi}_T \leftarrow \frac{1}{T} \sum_{t=1}^T \pi^t$. ▷ return mixture policy

B.2 Proof of Theorem 3.3

In this section we present our proof for Theorem 3.3 which bounds the performance suboptimality for PrefixRL. In particular, we bound this gap for an algorithm that conforms to the PrefixRL workflow (Algorithm 1) but uses natural policy gradient (NPG) (Kakade, 2001) to update the policy iteratively (starting from the base LLM π^0). In our practical implementation we use REINFORCE to compute on-policy gradients. Next we introduce the setup and describe the key steps in Algorithm 1, some differences with practice and the full proof.

Setup. We use π^0 to denote the base LLM we start RL training with, and μ as the policy used to IID sample the dataset of off-policy traces \mathcal{D}_{off} , one trace for each problem in $\mathbf{x} \in \mathcal{D}$. We assume that \mathcal{D}_{off} is realizable (Assumption 3.1) which implies that it only consists of correct off-policy traces. Let \mathcal{F} denote the class of Q -functions induced by all policies in the policy class Π , and H be the maximum context length or horizon H of the auto-regressive Markov decision process (MDP) induced by the policies in Π and reward function r .

Description of PrefixRL with NPG (Algorithm 1). In each iteration of Algorithm 1 we *first* collect a dataset of state, action and reward triplets \mathcal{D}_t . *Second*, we fit a critic or Q function \hat{Q}^t on this dataset of size N (step 9). *Third*, we use the fitted Q function to perform a state-wise mirror ascent or natural policy update (step 11) in order to get the subsequent RL iterate. We collect the N traces in \mathcal{D}_t by uniformly sampling an off-policy state s_h (prefixed problem) from \mathcal{D}_{off} . Then, we rollout the current RL policy π^t conditioned on state s_h to sample a single action or token a_h (step 6). To estimate the Q function under the current RL policy at this state-action pair we now complete the rollout till time step H and collect a reward (step 7).

Difference with practice: Algorithm 1 uses Q functions instead of direct rewards. The update in NPG is similar to REINFORCE except that we use N on-policy samples to first estimate a Q function (step 10) in the Q function class \mathcal{F} for the current RL iterate π^t and then use the estimated Q function to update the policy and get π_{t+1} using mirror ascent (step 12). This is a bit different from REINFORCE where we compute the policy gradient using only the rewards attained by the N on-policy traces and perform a step of gradient.

Difference with practice: Algorithm 1 samples new prefixed problems from \mathcal{D}_{off} . In practice we construct a prefixed problems from a fixed dataset of off-policy traces \mathcal{D}_{off} and also use the set of prefixed problems in \mathcal{D}_{pre} are fixed throughout RL training. In contrast, Algorithm 1 samples off-policy states (prefixed problems) from the dataset of off-policy traces \mathcal{D}_{off} . This difference is pretty minor but perhaps underscores the performance improvements driven by back-generalization in being able to improve performance on the original no-prefix problems despite PrefixRL only using a small fraction of all possible off-policy prefixes in \mathcal{D}_{off} .

Comparison with Chang et al. (2024). Our proof technique follows Chang et al. (2024), adapting to the setting of verifiable rewards with our different “reset” policy (which we refer to as prefix policy μ), and removing the requirement of KL divergence between the current and the reset policy. Since our off-policy dataset consists

of only realizable correct traces we will need much weaker assumptions. Following are some key differences compared to Algorithm 3 in [Chang et al. \(2024\)](#) that allows us to prove the suboptimality gap with weaker assumptions. *First*, we sample the prefix from the comparator policy (in other words the prefix generating policy is realizable and lies in the class of optimal policies). This ensures sufficient coverage for the distribution of Q-function regression (ensuring small error in fitting the critic) over states visited by a “good” policy even though the current RL iterate is far from it. *Second*, we output the mixture policy (standard in self-play literature ([Bai et al., 2020](#); [Hofbauer and Sorin, 2006](#))). *Finally*, unlike [Chang et al. \(2024\)](#), we don’t require a bound on the KL divergence against the SFT policy or the policy trained on the off-policy data.

Assumptions needed for Theorem 3.3. Now, we list the assumptions we make in our analysis of the suboptimality gap of PrefixRL. In general, they are milder than the assumptions in [Chang et al. \(2024\)](#).

- Assumption [B.1](#) is pretty standard in the analysis of actor-critic methods ([Konda and Tsitsiklis, 2002](#)) and only requires that our critic function class is expressive enough to realize the Q function induced by any policy in Π . Note that since rewards are binary and terminal the Q -value at any state $\in [0,1]$.
- Assumption [B.2](#) is a milder form of the typical assumption on the coverage over states visited by the optimal policy. Here, we only assume that there is an optimal policy that can fit the dataset \mathcal{D}_{off} we collected. Typically the coverage assumption places a uniform bound on the likelihood ratio over the state distributions of the optimal policy and the current RL policy (d^{π^*}/d^π) as in [Chang et al. \(2024\)](#).
- Assumption [B.3](#) is necessary to ensure that the KL between the prefix generating policy (empirical distribution over \mathcal{D}_{off}) and the base policy π^0 is finite. If the size of the dataset $\mathcal{D}_{\text{off}} \rightarrow \infty$ and the samples in \mathcal{D}_{off} are all drawn IID from a policy $\nu \in \Pi$, then this assumption requires that cross-entropy between ν and π^0 is finite.

Assumption B.1 (Realizability of Q -function class). *There is a finite Q -function class $\mathcal{F} \subseteq \{f: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]\}$, and that Q -function induced by any policy is realized in this class, i.e., $Q^\pi \in \mathcal{F} \forall \pi \in \Pi$.*

Assumption B.2 (Correctness and realizability of \mathcal{D}_{off}). *We say that \mathcal{D}_{off} is correct if it contains a single correct trace \mathbf{y} for every $\mathbf{x} \in \mathcal{D}$ and realizable if \exists some policy $\mu \in \Pi$ such that $\mu(\mathbf{y} | \mathbf{x}) = 1, \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{off}}$.*

Assumption B.3 (Bounded likelihood of \mathcal{D}_{off} under π^0). *The KL divergence between base LLM π^0 and the policy $\mu \in \Pi$ that perfectly fits the data is $\text{KL}(\mu | \pi^0) < \infty$. In other words, this assumes that the samples in \mathcal{D}_{off} have a bounded likelihood under the base LLM π^0 , i.e., $\text{KL}(\mu | \pi^0) = \frac{1}{|\mathcal{D}_{\text{off}}|} \sum_{(\mathbf{x}, \mathbf{y}^x) \in \mathcal{D}_{\text{off}}} -\log \pi^0(\mathbf{y}^x | \mathbf{x}) < \infty$.*

Proof. We prove the guarantee against the comparator policy μ from Assumption [B.2](#). Since \mathcal{D}_{off} is correct and realizable by μ (i.e., for each $\mathbf{x} \in \mathcal{D}$ the unique correct trace \mathbf{y}^x satisfies $\mu(\mathbf{y}^x | \mathbf{x}) = 1$), we have $J(\mu) = J^*$. Thus it suffices to upper bound:

$$J(\mu) - J(\bar{\pi}_T). \tag{B.15}$$

State distribution induced by \mathcal{D}_{off} . Let $\mathbf{s} \in \mathcal{S}$ denote an autoregressive prefix-state. Algorithm 1 samples prefix-states by drawing (\mathbf{x}, \mathbf{y}) uniformly from \mathcal{D}_{off} and then sampling a prefix of \mathbf{y} (according to the algorithm’s prefix-selection rule). Because μ deterministically generates the same trace \mathbf{y} in \mathcal{D}_{off} for each \mathbf{x} , this state distribution coincides with the state-visitation distribution of μ ; we denote it by $d_{\mathbf{s}}^\mu$.

$$d_{\mathbf{s}}^\mu \equiv d_{\mathbf{s}}^{\text{off}}. \tag{B.16}$$

Let π^0 be the base LLM and $\{\pi^t\}_{t=1}^T$ be the iterates produced by NPG / mirror descent with stepsize η and critic \hat{Q}^t , and define the averaged policy:

$$\bar{\pi}_T := \frac{1}{T} \sum_{t=1}^T \pi^t. \tag{B.17}$$

Performance difference lemma under $d_{\mathbf{s}}^\mu$. Applying performance difference Lemma [B.4](#) with $(\pi, \pi') = (\mu, \pi^t)$ yields:

$$J(\mu) - J(\pi^t) = \mathbb{E}_{\mathbf{s} \sim d_{\mathbf{s}}^\mu} \mathbb{E}_{\mathbf{a} \sim \mu(\cdot | \mathbf{s})} [A^{\pi^t}(\mathbf{s}, \mathbf{a})]. \tag{B.18}$$

Using $A^{\pi^t}(\mathbf{s}, \mathbf{a}) = Q^{\pi^t}(\mathbf{s}, \mathbf{a}) - V^{\pi^t}(\mathbf{s})$ and the identity $\mathbb{E}_{\mathbf{a} \sim \pi^t(\cdot | \mathbf{s})}[A^{\pi^t}(\mathbf{s}, \mathbf{a})] = 0$, (B.18) can be rewritten as

$$J(\mu) - J(\pi^t) = \mathbb{E}_{\mathbf{s} \sim d_{\mathbf{s}}^{\mu}} \left[\langle Q^{\pi^t}(\mathbf{s}, \cdot), \mu(\cdot | \mathbf{s}) - \pi^t(\cdot | \mathbf{s}) \rangle \right]. \quad (\text{B.19})$$

B.2.1 Critic Estimation Error.

Fix an iteration t . The critic is fit by least squares over a finite class \mathcal{F} (Assumption B.1) using N i.i.d. samples $(\mathbf{s}_k, \mathbf{a}_k, z_k)$ where $\mathbf{s}_k \sim d_{\mathbf{s}}^{\mu}$, $\mathbf{a}_k \sim \rho^t(\cdot | \mathbf{s}_k)$ (see discussion below on ρ^t), and z_k is an unbiased target for $Q^{\pi^t}(\mathbf{s}_k, \mathbf{a}_k)$. Because rewards are *terminal and binary* in $\{0, 1\}$, we have

$$0 \leq Q^{\pi^t}(\mathbf{s}, \mathbf{a}) \leq 1, \quad 0 \leq \widehat{Q}^t(\mathbf{s}, \mathbf{a}) \leq 1, \quad (\text{B.20})$$

so we may take $R=1$ in Lemma B.6. Therefore, setting $\delta_t := \delta/(2T)$ and applying Lemma B.6 with $\mathcal{H} = \mathcal{F}$, with probability at least $1 - \delta_t$,

$$\mathbb{E}_{\mathbf{s} \sim d_{\mathbf{s}}^{\mu}, \mathbf{a} \sim \rho^t(\cdot | \mathbf{s})} \left[(\widehat{Q}^t(\mathbf{s}, \mathbf{a}) - Q^{\pi^t}(\mathbf{s}, \mathbf{a}))^2 \right] \leq \frac{256}{N} \log \left(\frac{2|\mathcal{F}|}{\delta_t} \right) = \frac{256}{N} \log \left(\frac{4T|\mathcal{F}|}{\delta} \right). \quad (\text{B.21})$$

Behavior distribution ρ^t and pointwise domination. At iteration t , Algorithm 1 forms critic data by first sampling $(\mathbf{s}_h, a_h^{\text{off}})$ uniformly from \mathcal{D}_{off} and then sampling

$$a_h = \begin{cases} a_h^{\text{off}} & \text{w.p. } \frac{1}{2}, \\ a_h \sim \pi^t(\cdot | \mathbf{s}_h) & \text{w.p. } \frac{1}{2}. \end{cases} \quad (\text{B.22})$$

Let $\mu(\cdot | \mathbf{s}_h)$ denote the (deterministic) conditional action distribution induced by \mathcal{D}_{off} , i.e.,

$$\mu(a | \mathbf{s}_h) := \mathbf{1}\{a = a_h^{\text{off}}\}. \quad (\text{B.23})$$

Then the induced action-sampling (behavior) distribution used for critic fitting is the mixture

$$\rho^t(\cdot | \mathbf{s}_h) := \frac{1}{2} \mu(\cdot | \mathbf{s}_h) + \frac{1}{2} \pi^t(\cdot | \mathbf{s}_h). \quad (\text{B.24})$$

Consequently, for every state \mathbf{s} and action a , we have the pointwise lower bounds

$$\rho^t(a | \mathbf{s}) \geq \frac{1}{2} \mu(a | \mathbf{s}), \quad \rho^t(a | \mathbf{s}) \geq \frac{1}{2} \pi^t(a | \mathbf{s}), \quad (\text{B.25})$$

and hence the pointwise domination inequalities

$$\mu(a | \mathbf{s}) \leq 2\rho^t(a | \mathbf{s}), \quad \pi^t(a | \mathbf{s}) \leq 2\rho^t(a | \mathbf{s}). \quad (\text{B.26})$$

In particular, (B.25) also implies absolute continuity: if $\rho^t(a | \mathbf{s}) = 0$ then $\mu(a | \mathbf{s}) = \pi^t(a | \mathbf{s}) = 0$, so $\mu(\cdot | \mathbf{s}) \ll \rho^t(\cdot | \mathbf{s})$ and $\pi^t(\cdot | \mathbf{s}) \ll \rho^t(\cdot | \mathbf{s})$.

Then Cauchy–Schwarz and Jensen applied to (B.21) yield, for $\pi \in \{\mu, \pi^t\}$,

$$\left| \mathbb{E}_{\mathbf{s} \sim d_{\mathbf{s}}^{\mu}, \mathbf{a} \sim \pi(\cdot | \mathbf{s})} [\widehat{Q}^t(\mathbf{s}, \mathbf{a}) - Q^{\pi^t}(\mathbf{s}, \mathbf{a})] \right| \leq \epsilon_{\text{crt}}, \quad \epsilon_{\text{crt}} := 16\sqrt{2} \sqrt{\frac{1}{N} \log \left(\frac{4T|\mathcal{F}|}{\delta} \right)}. \quad (\text{B.27})$$

Taking a union bound over $t \in [T]$, with probability at least $1 - \delta/2$, (B.27) holds for all t .

B.2.2 Mirror Ascent and NPG Optimization Error.

The mirror-descent update at state \mathbf{s} is the KL-regularized maximization

$$\pi^{t+1}(\cdot | \mathbf{s}) = \arg \max_{p(\cdot | \mathbf{s})} \left\{ \eta \langle \widehat{Q}^t(\mathbf{s}, \cdot), p(\cdot | \mathbf{s}) \rangle - \text{KL}(p(\cdot | \mathbf{s}) \| \pi^t(\cdot | \mathbf{s})) \right\}.$$

The first-order optimality condition implies that for any $p(\cdot | \mathbf{s})$,

$$\left\langle -\eta \widehat{Q}^t(\mathbf{s}, \cdot) + \nabla_r \text{KL}(r(\cdot | \mathbf{s}) \| \pi^t(\cdot | \mathbf{s})) \Big|_{r=\pi^{t+1}}, p - \pi^{t+1} \right\rangle \geq 0. \quad (\text{B.28})$$

Set $p = \mu(\cdot | \mathbf{s})$ and rearrange (B.28) to obtain

$$\eta \langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu(\cdot | \mathbf{s}) - \pi^{t+1}(\cdot | \mathbf{s}) \rangle \leq \left\langle \nabla_r \text{KL}(r \| \pi^t) \Big|_{r=\pi^{t+1}}, \mu - \pi^{t+1} \right\rangle. \quad (\text{B.29})$$

Apply the KL three-point identity (Lemma (B.62)) with $p = \mu(\cdot | \mathbf{s})$, $r = \pi^{t+1}(\cdot | \mathbf{s})$, and $q = \pi^t(\cdot | \mathbf{s})$ to rewrite the right-hand side:

$$\left\langle \nabla_r \text{KL}(r \| \pi^t) \Big|_{r=\pi^{t+1}}, \mu - \pi^{t+1} \right\rangle = \text{KL}(\mu \| \pi^t) - \text{KL}(\mu \| \pi^{t+1}) - \text{KL}(\pi^{t+1} \| \pi^t). \quad (\text{B.30})$$

Combining (B.29) and (B.30) gives

$$\eta \langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^{t+1} \rangle \leq \text{KL}(\mu \| \pi^t) - \text{KL}(\mu \| \pi^{t+1}) - \text{KL}(\pi^{t+1} \| \pi^t). \quad (\text{B.31})$$

Using (B.20) and Pinsker's inequality, we bound the shift term

$$\eta \langle \widehat{Q}^t(\mathbf{s}, \cdot), \pi^{t+1} - \pi^t \rangle \leq \eta \|\widehat{Q}^t(\mathbf{s}, \cdot)\|_\infty \|\pi^{t+1} - \pi^t\|_1 \leq \frac{\eta^2}{2} + \text{KL}(\pi^{t+1} \| \pi^t), \quad (\text{B.32})$$

where the last inequality uses $\|\widehat{Q}^t\|_\infty \leq 1$ and $\|\pi^{t+1} - \pi^t\|_1^2 \leq 2\text{KL}(\pi^{t+1} \| \pi^t)$. Adding (B.31) and (B.32) cancels $\text{KL}(\pi^{t+1} \| \pi^t)$ and yields

$$\eta \langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^t \rangle \leq \text{KL}(\mu \| \pi^t) - \text{KL}(\mu \| \pi^{t+1}) + \frac{\eta^2}{2}. \quad (\text{B.33})$$

Taking expectation over $\mathbf{s} \sim d_s^\mu$ and summing over $t = 1, \dots, T$ gives

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\mathbf{s} \sim d_s^\mu} \left[\langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^t \rangle \right] \leq \frac{D_1}{\eta T} + \frac{\eta}{2}, \quad (\text{B.34})$$

where

$$D_1 := \mathbb{E}_{\mathbf{s} \sim d_s^\mu} \left[\text{KL}(\mu(\cdot | \mathbf{s}) \| \pi^1(\cdot | \mathbf{s})) \right]. \quad (\text{B.35})$$

Using $\pi^1 = \pi^0$ and the definition of $\text{KL}(\mu \| \pi^0)$ under the \mathcal{D}_{off} -induced state distribution (Assumption B.3), we identify

$$D_1 = \text{KL}(\mu \| \pi^0). \quad (\text{B.36})$$

B.2.3 Combining Critic Error and Optimization Error.

Starting from (B.19), add and subtract \widehat{Q}^t :

$$\begin{aligned} J(\mu) - J(\pi^t) &= \mathbb{E}_{\mathbf{s} \sim d_s^\mu} \left[\langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^t \rangle \right] \\ &\quad + \mathbb{E}_{\mathbf{s} \sim d_s^\mu} \left[\langle Q^{\pi^t}(\mathbf{s}, \cdot) - \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^t \rangle \right]. \end{aligned} \quad (\text{B.37})$$

On the high-probability event where (B.27) holds for both $\pi = \mu$ and $\pi = \pi^t$, the critic-error term is bounded by

$$\mathbb{E}_{\mathbf{s} \sim d_s^\mu} \left[\langle Q^{\pi^t} - \widehat{Q}^t, \mu - \pi^t \rangle \right] \leq \left| \mathbb{E}_{\mathbf{s} \sim d_s^\mu, \mathbf{a} \sim \mu} [\widehat{Q}^t - Q^{\pi^t}] \right| + \left| \mathbb{E}_{\mathbf{s} \sim d_s^\mu, \mathbf{a} \sim \pi^t} [\widehat{Q}^t - Q^{\pi^t}] \right| \leq 2\epsilon_{\text{crt}}, \quad (\text{B.38})$$

which introduces no extra factor of H .

Averaging (B.37) over $t=1,\dots,T$ and using (B.17), we obtain

$$J(\mu) - J(\bar{\pi}_T) \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\mathbf{s} \sim d_{\mathbf{s}}^{\mu}} \left[\langle \widehat{Q}^t(\mathbf{s}, \cdot), \mu - \pi^t \rangle \right] + 2\epsilon_{\text{crt}}. \quad (\text{B.39})$$

Applying (B.34) to the first term in (B.39) yields

$$J(\mu) - J(\bar{\pi}_T) \leq \frac{D_1}{\eta T} + \frac{\eta}{2} + 2\epsilon_{\text{crt}}. \quad (\text{B.40})$$

Choose $\eta := \sqrt{\frac{2D_1}{T}}$ to balance the first two terms in (B.40), giving

$$\frac{D_1}{\eta T} + \frac{\eta}{2} = \sqrt{\frac{2D_1}{T}}. \quad (\text{B.41})$$

Combining (B.36), (B.27), and (B.40), and recalling $J(\mu) = J^*$, we conclude that with probability at least $1 - \delta$,

$$J(\pi^*) - J(\bar{\pi}_T) = J^* - J(\bar{\pi}_T) \leq \mathcal{O} \left(\sqrt{\frac{\text{KL}(\mu \parallel \pi^0)}{T}} + \sqrt{\frac{1}{N} \log \left(\frac{T|\mathcal{F}|}{\delta} \right)} \right). \quad (\text{B.42})$$

□

B.3 Proof of Proposition 3.4

Now, we prove our separation result in Proposition 3.4 that lower bounds the performance gap between standard RL and PrefixRL. Here, standard RL runs Algorithm 1 but now without any access to the off-policy dataset \mathcal{D}_{off} . In each iteration, the critic training dataset \mathcal{D}_t is now populated with (\mathbf{s}, a, r) tuples where both \mathbf{s} and a are sampled from the current policy π^t . So, unlike PrefixRL we never sample the state or prefix from \mathcal{D}_{off} . In this simple worst-case instance we present below there is a single trajectory in \mathcal{D}_{off} that is also what that the optimal policy samples with probability 1 and attains performance of $J(\pi^*) = 1$.

Proof. We present (i) the MDP instance together with a choice of *base policy* that generates the off-policy trace, and then (ii) an exponential lower bound for standard on-policy RL without \mathcal{D}_{off} , and (iii) a horizon-independent (non-exponential) upper bound for PrefixRL with \mathcal{D}_{off} .

MDP instance (hidden rewarding binary string) and base policy. Fix a horizon H and an unknown binary string $\mathbf{b} = (b_1, \dots, b_H) \in \{0, 1\}^H$. Let the state space be

$$\mathcal{S} = \{s_0, s_1, \dots, s_H\}, \quad (\text{B.43})$$

where s_{h-1} encodes the first $h-1$ actions taken so far (s_0 is the start state). The action space is $\mathcal{A} = \{0, 1\}$. Transitions are deterministic: from s_{h-1} , taking action $a_h \in \{0, 1\}$ moves to s_h . The episode ends at s_H with terminal reward

$$r = \mathbf{1}\{(a_1, \dots, a_H) = (b_1, \dots, b_H)\}. \quad (\text{B.44})$$

Thus, exactly one length- H action sequence earns reward 1.

Let π^* be the deterministic policy that selects b_h at s_{h-1} for each $h \in [H]$. Then $J(\pi^*) = 1$. For the PrefixRL part, we also choose a *base policy* μ and an off-policy dataset \mathcal{D}_{off} : we set $\mu := \pi^*$ and let \mathcal{D}_{off} contain the unique successful trajectory of μ , equivalently the state-action pairs

$$\mathcal{D}_{\text{off}} = \{(s_{h-1}, b_h) : h \in [H]\}. \quad (\text{B.45})$$

B.3.1 Exponential lower bound for standard on-policy RL (Algorithm 1 without \mathcal{D}_{off}).

We analyze an on-policy variant of Algorithm 1 in which *there is no off-policy dataset*: at each iteration t , the algorithm samples N full episodes only from its current policy π^t , observes only terminal rewards $r^{(t,i)} \in \{0,1\}$, fits a critic, and updates the policy. Let $\hat{\pi}_T$ denote the (possibly randomized) policy output after T iterations (so the total number of full episodes is TN). We prove that for any such algorithm, there exists an instance \mathbf{b} for which the expected suboptimality gap is at least $1 - (TN+1)2^{-H}$.

Yao's minimax setup. By Yao's minimax principle Yao (1977), it suffices to fix an arbitrary adaptive algorithm and analyze its expected performance when the instance is random:

$$\mathbf{b} \sim \text{Unif}(\{0,1\}^H), \quad (\text{B.46})$$

and we write $\mathbb{P}_{\mathbf{b}}, \mathbb{E}_{\mathbf{b}}$ for probability/expectation over this draw (and over the algorithm's internal randomness).

Per-rollout success probability is 2^{-H} . Fix any rollout index (t,i) . Condition on the full interaction history up to this rollout and on the algorithm's internal randomness. Under this conditioning, the action string $\mathbf{a}^{(t,i)} \in \{0,1\}^H$ is some random element (with an arbitrary distribution induced by the algorithm), while \mathbf{b} remains uniform and independent. Therefore,

$$\mathbb{P}_{\mathbf{b}}[r^{(t,i)} = 1 \mid \text{history}] = \mathbb{P}_{\mathbf{b}}[\mathbf{a}^{(t,i)} = \mathbf{b} \mid \text{history}] = \sum_{\mathbf{a} \in \{0,1\}^H} \mathbb{P}[\mathbf{a}^{(t,i)} = \mathbf{a} \mid \text{history}] \cdot \mathbb{P}_{\mathbf{b}}[\mathbf{b} = \mathbf{a}] = 2^{-H}. \quad (\text{B.47})$$

Taking expectation over the history yields the unconditional version:

$$\mathbb{P}_{\mathbf{b}}[r^{(t,i)} = 1] = 2^{-H}. \quad (\text{B.48})$$

Probability of ever seeing a reward-1 rollout. There are exactly TN rollouts in total. By a union bound and (B.48) we get the following upper bound on the probability of seeing a reward 1 rollout across all T steps,

$$\mathbb{P}_{\mathbf{b}}[\exists t \leq T, i \leq N : r^{(t,i)} = 1] \leq \sum_{t=1}^T \sum_{i=1}^N \mathbb{P}_{\mathbf{b}}[r^{(t,i)} = 1] = TN \cdot 2^{-H}. \quad (\text{B.49})$$

Bound the expected value of the returned policy. Let $\mathbf{a} \sim \hat{\pi}_T$ denote the length- H string generated by rolling out $\hat{\pi}_T$ from s_0 . On instance \mathbf{b} , $J(\hat{\pi}_T) = \mathbb{P}[\mathbf{a} = \mathbf{b}]$. We can decompose this reward on two events: whether the training interaction ever produced a reward-1 rollout or not:

$$\mathbb{E}_{\mathbf{b}}[J(\hat{\pi}_T)] \leq \mathbb{P}_{\mathbf{b}}[\exists t, i : r^{(t,i)} = 1] \cdot 1 + \mathbb{P}_{\mathbf{b}}[\forall t, i : r^{(t,i)} = 0] \cdot \sup_{\hat{\pi}_T} \mathbb{E}_{\mathbf{b}}[J(\hat{\pi}_T) \mid \forall t, i : r^{(t,i)} = 0]. \quad (\text{B.50})$$

On the event $\{\forall t, i : r^{(t,i)} = 0\}$, the algorithm has never observed the rewarding string. Each zero-reward rollout rules out *at most one* candidate string, namely the realized action string $\mathbf{a}^{(t,i)}$. Hence, conditioning on $\{\forall t, i : r^{(t,i)} = 0\}$ only implies that \mathbf{b} is not in a set of at most TN excluded strings. Under the prior $\mathbf{b} \sim \text{Unif}(\{0,1\}^H)$, the posterior is uniform over the remaining candidates, so

$$\max_{\mathbf{a} \in \{0,1\}^H} \mathbb{P}_{\mathbf{b}}[\mathbf{b} = \mathbf{a} \mid \forall t, i : r^{(t,i)} = 0] \leq \frac{1}{2^H - TN}, \quad (\text{B.51})$$

and therefore

$$\sup_{\hat{\pi}_T} \mathbb{E}_{\mathbf{b}}[J(\hat{\pi}_T) \mid \forall t, i : r^{(t,i)} = 0] = \sup_{\hat{\pi}_T} \mathbb{E}_{\mathbf{b}}[\mathbb{P}(\mathbf{a} = \mathbf{b} \mid \hat{\pi}_T, \forall t, i : r^{(t,i)} = 0)] \leq \frac{1}{2^H - TN}. \quad (\text{B.52})$$

Substituting (B.49) and (B.52) into (B.50) yields

$$\mathbb{E}_{\mathbf{b}}[J(\hat{\pi}_T)] \leq \mathbb{P}_{\mathbf{b}}[\exists t, i : r^{(t,i)} = 1] \cdot 1 + \mathbb{P}_{\mathbf{b}}[\forall t, i : r^{(t,i)} = 0] \cdot \frac{1}{2^H - TN} \leq TN \cdot 2^{-H} + \frac{1}{2^H - TN}. \quad (\text{B.53})$$

In particular, whenever $TN \leq 2^{H-1}$ we have $\frac{1}{2^{H-TN}} \leq 2^{-(H-1)}$, and thus

$$\mathbb{E}_{\mathbf{b}}[J(\hat{\pi}_T)] \leq TN \cdot 2^{-H} + 2^{-(H-1)} \leq (TN+2)2^{-(H-1)}. \quad (\text{B.54})$$

Suboptimality lower bound and fix an instance. Since $J(\pi^*) = 1$, we obtain from (B.53) the gap bound

$$\mathbb{E}_{\mathbf{b}}[J(\pi^*) - J(\hat{\pi}_T)] \geq 1 - (TN+2)2^{-(H-1)}. \quad (\text{B.55})$$

By Yao’s minimax principle (Yao, 1977), there exists a fixed hidden string \mathbf{b} for which the same bound holds for the algorithm on that instance. This proves the exponential lower bound for standard on-policy RL without access to \mathcal{D}_{off} .

B.3.2 Horizon-independent (non-exponential) upper bound for PrefixRL.

We now analyze Algorithm 1 on the same instance *with* \mathcal{D}_{off} . As discussed, the separation mechanism is that Algorithm 1 explicitly samples states from \mathcal{D}_{off} , which are precisely the prefix states along the π^* trajectory, thereby forcing visitation of optimal-trajectory states.

Choosing π^0 . Let π^0 denote the initialization policy of Algorithm 1, *i.e.* $\pi^1 = \pi^0$. We choose π^0 to be the uniform policy on $\{0,1\}$ at every state (*crucially, π^0 is not instance dependent*):

$$\pi_0(0|s) = \pi_0(1|s) = \frac{1}{2} \quad \text{for all } s \in \{s_0, \dots, s_{H-1}\}. \quad (\text{B.56})$$

Since $\mu = \pi^*$ is deterministic on the states in \mathcal{D}_{off} , the KL term in Theorem 3.1 is not exponential in H . Under the state-averaged convention used in Theorem 3.1,

$$\text{KL}(\mu || \pi_0) := \mathbb{E}_{s \sim d_s^\mu} \left[\text{KL}(\mu(\cdot|s) || \pi_0(\cdot|s)) \right] = \log 2, \quad (\text{B.57})$$

while under the summed convention it is $H \log 2$; in either case it is not exponential in H .

Invoking PrefixRL guarantee in Theorem 3.3. All assumptions required by Theorem 3.3 hold on this instance with $\mu = \pi^*$ and the corresponding $\mathcal{D}_{\text{off}} = \{\mathbf{b}\}$. Therefore, applying Theorem 3.3 to Algorithm 1 yields (with probability at least $1 - \delta$) the bound

$$J(\mu) - J(\bar{\pi}_T) \leq \sqrt{\frac{2\text{KL}(\mu || \pi_0)}{T}} + \mathcal{O}\left(\sqrt{\frac{1}{N} \log\left(\frac{T|\mathcal{F}|}{\delta}\right)}\right), \quad (\text{B.58})$$

where $\bar{\pi}_T$ is the iterate-averaged policy output by Algorithm 1. Since $\mu = \pi^*$ on this instance, $J(\mu) = J(\pi^*) = 1$, and (B.58) implies a non-exponential suboptimality bound for PrefixRL. In particular, with the choice (B.56) we have $\text{KL}(\mu || \pi_0) = \log 2$ (or $H \log 2$ under the summed convention), so the bound has no 2^{-H} term.

Separation mechanism. The on-policy exponential lower bound arises because, without \mathcal{D}_{off} , an algorithm only observes nonzero reward if it guesses the entire length- H string correctly in a single episode. In contrast, Algorithm 1 with \mathcal{D}_{off} repeatedly samples *states along the π^* trajectory* via $\sim \mathcal{D}_{\text{off}}$ and trains at those states using the mixture distribution of Algorithm 1. Overall, this yields a non-exponential sample complexity, establishing a worst-case separation.

B.3.3 Worst-Case Separation Result Between Standard RL and PrefixRL.

The analysis in the above subsection can be viewed either for a fixed hidden string \mathbf{b} , or under a randomized instance distribution. In particular, let $\mathbf{b} \sim \text{Unif}(\{0,1\}^H)$, and note that under this randomization the off-policy dataset $\mathcal{D}_{\text{off}} = \mathcal{D}_{\text{off}}(\mathbf{b})$ also changes with \mathbf{b} since it contains the unique successful trajectory $(s_{h-1}, b_h)_{h=1}^H$.

Lower bound for standard on-policy RL under random \mathbf{b} . For any on-policy algorithm that runs for T iterations with N full episodes per iteration (so TN total episodes) and does not have access to \mathcal{D}_{off} , the expected value of the suboptimality gap satisfies

$$\mathbb{E}_{\mathbf{b}}[J(\pi_{\mathbf{b}}^*) - J(\hat{\pi}_T)] \geq 1 - (TN+2)2^{-(H-1)}, \quad (\text{B.59})$$

where $\pi_{\mathbf{b}}^*$ denotes the optimal policy for instance \mathbf{b} and $TN \leq 2^{H-1}$.

Upper bound for PrefixRL under random \mathbf{b} . Now consider Algorithm 1 (PrefixRL) run on the same randomized instance, where $\mathcal{D}_{\text{off}} = \mathcal{D}_{\text{off}}(\mathbf{b})$ is provided to the algorithm. Choose the initialization π_0 to be the uniform policy (so that $\text{KL}(\mu_{\mathbf{b}} \parallel \pi_0)$ is not exponential in H), with $\mu_{\mathbf{b}} := \pi_{\mathbf{b}}^*$ as the base policy that generates $\mathcal{D}_{\text{off}}(\mathbf{b})$. Invoking the previously proved Theorem 3.1 yields (with probability at least $1 - \delta$ over the algorithm's sampling)

$$\mathbb{E}_{\mathbf{b}}[J(\pi_{\mathbf{b}}^*) - J(\bar{\pi}_T)] \leq \sqrt{\frac{2\mathbb{E}_{\mathbf{b}}[\text{KL}(\mu_{\mathbf{b}} \parallel \pi_0)]}{T}} + \mathcal{O}\left(\sqrt{\frac{1}{N} \log\left(\frac{T|\mathcal{F}|}{\delta}\right)}\right), \quad (\text{B.60})$$

and for the uniform initialization π_0 we have $\mathbb{E}_{\mathbf{b}}[\text{KL}(\mu_{\mathbf{b}} \parallel \pi_0)] = \log 2$ under the state-averaged convention (or $H \log 2$ under the summed convention), which is not exponential in H .

Applying Yao's minimax principle. Equations (B.59) and (B.60) establish an *average-case* separation under the uniform distribution over instances \mathbf{b} (with \mathcal{D}_{off} coupled to \mathbf{b} in the PrefixRL case). By Yao's minimax principle, this implies that there exists a fixed instance \mathbf{b} (and the above choice of initialization π_0 , e.g. uniform) for which the same separation holds for any randomized standard on-policy algorithm without \mathcal{D}_{off} versus Algorithm 1 with \mathcal{D}_{off} . \square

B.4 Auxiliary Lemmas

Lemma B.4 (Performance difference lemma; (Kakade and Langford, 2002)). *For all policies π, π' and initial state distribution ρ ,*

$$\mathbb{E}_{\mathbf{s}_0 \sim \rho} [V^\pi(\mathbf{s}_0) - V^{\pi'}(\mathbf{s}_0)] = \mathbb{E}_{\mathbf{s}_h \sim d_s^\pi} \mathbb{E}_{a_h \sim \pi(\cdot | \mathbf{s}_h)} [A^{\pi'}(\mathbf{s}_h, a_h)]. \quad (\text{B.61})$$

Proof. See proof of Lemma 6.1 in Kakade and Langford (2002). \square

Lemma B.5 (Three-point identity for KL). *Let p, q, r be distributions on a common measurable space such that $p \ll r \ll q$ and all quantities below are finite. Then*

$$\text{KL}(p \parallel q) = \text{KL}(p \parallel r) + \text{KL}(r \parallel q) + \langle p - r, \nabla_r \text{KL}(r \parallel q) \rangle, \quad (\text{B.62})$$

. For discrete distributions p, q, r we have:

$$\nabla_r \text{KL}(r \parallel q)(x) = \log \frac{r(x)}{q(x)} + 1. \quad (\text{B.63})$$

Equivalently we can state this the three-point identity for discrete distributions as:

$$\text{KL}(p \parallel q) = \text{KL}(p \parallel r) + \text{KL}(r \parallel q) + \left\langle p - r, \log \frac{r}{q} \right\rangle, \quad (\text{B.64})$$

since $\langle p - r, 1 \rangle = \int (p - r) = 0$ (or $\sum_i (p_i - r_i) = 0$).

Proof. We work in the continuous case; the discrete case is identical with integrals replaced by sums. Recall

$$\text{KL}(a \parallel b) = \int a(x) \log \frac{a(x)}{b(x)} dx.$$

Consider the difference between the left-hand side and the first two KL terms:

$$\begin{aligned}
\text{KL}(p\|q) - \text{KL}(p\|r) - \text{KL}(r\|q) &= \int p \log \frac{p}{q} dx - \int p \log \frac{p}{r} dx - \int r \log \frac{r}{q} dx \\
&= \int p \left(\log \frac{p}{q} - \log \frac{p}{r} \right) dx - \int r \log \frac{r}{q} dx \\
&= \int p \log \frac{r}{q} dx - \int r \log \frac{r}{q} dx \\
&= \int (p-r) \log \frac{r}{q} dx \\
&= \left\langle p-r, \log \frac{r}{q} \right\rangle.
\end{aligned}$$

This proves (B.64). To obtain the gradient form (B.62), note that for $\text{KL}(r\|q) = \int r \log(r/q) dx$ the pointwise functional derivative with respect to r is

$$\nabla_r \text{KL}(r\|q)(x) = \log \frac{r(x)}{q(x)} + 1,$$

so

$$\langle p-r, \nabla_r \text{KL}(r\|q) \rangle = \left\langle p-r, \log \frac{r}{q} \right\rangle + \langle p-r, 1 \rangle.$$

Finally, $\langle p-r, 1 \rangle = \int (p-r) dx = 1-1=0$, yielding (B.62). \square

Lemma B.6 (Lemma 15 in Song et al. (2022)). *Fix any $R > 0$, $\delta \in (0, 1)$, and assume we have a class of real-valued functions $\mathcal{H}: \mathcal{X} \rightarrow [-R, R]$. Suppose we have K i.i.d. samples $\{(x_k, y_k)\}_{k=1}^K$ where $x_k \sim \rho$ and y_k is sampled via the conditional probability $p(\cdot | x_k)$:*

$$y_k \sim p(\cdot | x_k) := h^*(x_k) + \epsilon_k,$$

where $h^* \in \mathcal{H}$ and $\{\epsilon_k\}_{k=1}^K$ are independent random variables such that $\mathbb{E}[y_k | x_k] = h^*(x_k)$. Additionally, suppose that $\max_k |y_k| \leq R$ and $\max_x |h^*(x)| \leq R$. Then the least squares solution

$$\hat{h} \leftarrow \operatorname{argmin}_{h \in \mathcal{H}} \sum_{k=1}^K (h(x_k) - y_k)^2$$

satisfies, with probability at least $1 - \delta$,

$$\mathbb{E}_{x \sim \rho} \left[(\hat{h}(x) - h^*(x))^2 \right] \leq \frac{256R^2 \log(2|\mathcal{H}|/\delta)}{K}.$$

The proof is the same as in Song et al. (2022) and thus is omitted here.

C Additional Experiment Details and Results on Back-Generalization

This appendix provides (i) the concrete hard problems used in the in-context back-generalization experiment in Section 4, and (ii) an additional train test mismatch study where the off-policy prefixes are sourced from a different model family (Figure 8).

C.1 Hard problems used in the in-context back-generalization experiment

Section 4 studies a meta-learning style setting where the policy is trained on a *prefixed problem* consisting of an in-context example (a full solved hard problem) followed by a target hard problem. We compare transfer when

the in-context hard problem is structurally similar to the target problem versus when it is unrelated (Figure 7). We use the following three problems:

Hard Problem P1 (Pass@16 for base model is 0.119 and is different from P3)

A sphere tangent to the xy -plane has center having z -coordinate > 0 . If it is projected from $P = (0, b, a)$ to the xy -plane, it gives the conic section $y = x^2$. If $a = p/q$ for integers p, q what is $p+q$? Answer: 3.

Hard Problem P2 (Pass@16 for base model is 0.074 and is similar to P3)

League has 30 teams (East 16, West 14). Inside each division everyone has played others once. If we add interleague games, what is smallest k for which every team gets exactly k games? Answer: 29.

Hard Problem P3 (Pass@16 for base model is 0.063)

Amongst 300 people, no one has more than $k-1$ friends. What is the smallest k for which it might be impossible to create some new friendships so that everyone has exactly k friends? Answer: 151.

Relatedness criterion (P2 is similar to P3; P1 is different from P3). P2 and P3 are both naturally expressed as graph feasibility problems with degree constraints, and their solutions rely on reasoning about global constraints induced by local degree requirements (regularity, parity, and obstruction arguments). In contrast, P1 is a geometry and conic projection problem whose solution structure does not share this graph-theoretic scaffold. Figure 7 uses this controlled notion of similarity to separate two effects: when the prefix and suffix share a compatible representation (P2→P3), training on prefixed problems yields substantially stronger transfer than when they do not (P1→P3).

C.2 Back-generalization under model-family mismatch

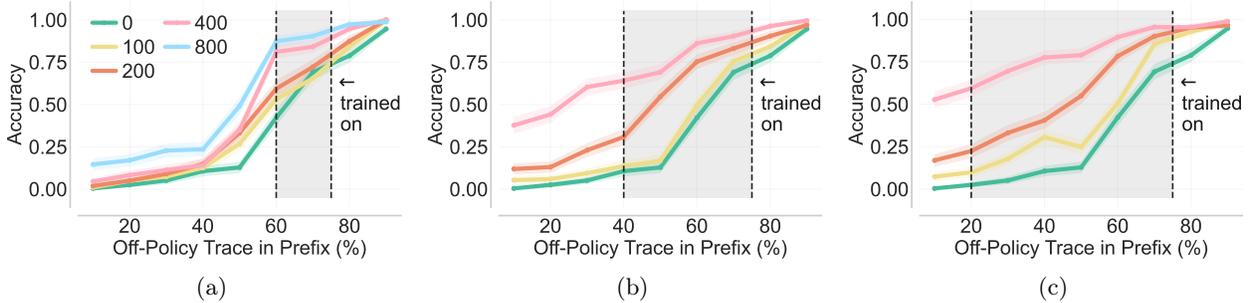


Figure 15 Back-generalization (train-test mismatch): On Llama3.1-8b-instruct we run RL only on prefixed problems sourced from Qwen3-4b-instruct whose prefix length (percent of the off-policy trace) lies in the shaded interval. At evaluation, we test across the full range of prefix lengths, including no-prefix problems (0% prefix). The performance at different RL training iterations (0, 100, 200, 400 and 800) is represented with different colors. Similar to Figure 5, where the prefixes are also sourced from Llama-3.1-8b When the mismatch is moderate, training on longer prefixes improves performance on shorter prefixes and can eventually improve no-prefix, indicating back-generalization (b, c). But different from Figure 5, we find that when the prefix distribution is skewed towards long prefixes back-generalization is weaker despite running RL for 800 iterations.

This section reports an additional train test mismatch experiment analogous to Figure 5, but where the off-policy prefixes used to form prefixed problems are sourced from a different model family. Concretely, we train Llama3.1-8B-Instruct with PrefixRL while constructing prefixed problems from correct off-policy traces generated by Qwen3-4B-Instruct.

Setup. We first collect a pool of correct off-policy traces from Qwen3-4B-Instruct. Each trace induces a family of *prefix states* by truncating the trace to a specified prefix length (reported as a percentage of the full trace). A prefixed problem is formed by conditioning the training policy on such a prefix state and then asking it to complete the solution for the same underlying problem. During training, we restrict prefix lengths to lie in a band (the shaded interval in Figure 15), and we run on-policy RL *only* on these prefixed problems. At evaluation, we sweep prefix lengths across the full range, including the no-prefix endpoint (0% prefix), and report accuracy at multiple RL training iterations.

Cross-family prefixes can still induce back-generalization. Figure 15 shows that cross-family prefixes can still induce back-generalization, but the effect is less robust than in the same-family setting of Figure 5. When the mismatch is moderate (training includes prefixes that are not concentrated exclusively at the very end of traces), improvements appear near the trained prefix region and then propagate to shorter prefixes, eventually improving no-prefix performance (Figure 15b,c). In contrast, when the training prefix distribution is skewed toward long prefixes, back-generalization is weaker, and the transfer to shorter prefixes and to no-prefix remains limited even after long training (800 iterations) (Figure 15a).

Choosing the right prefix distribution. A key difficulty in this setting is that the training prefixes are not sampled from the current policy. They are injected from an external model (Qwen), and therefore correspond to states that can be very unlikely under the evolving Llama policy. When training concentrates on very late prefixes, the policy can improve primarily on a narrow slice of the prefixed state distribution without sufficiently shaping behavior on earlier states that dominate no-prefix rollouts. This makes cross-family prefixing most effective when the training mixture covers a sufficiently broad range of prefix lengths, rather than concentrating only on long, heavily conditioned prefixes.

D Additional Experiments and Details for Results in Section 5

D.1 Implementation details for PrefixRL and Baselines

D.1.1 Constructing Prefixed Problems: An Example

Original Problem (No-Prefix)

```
<|im_start|>user
Consider the cube whose vertices are the eight points  $(x,y,z)$  for which each of  $x,y$ , and  $z$  is either 0 or 1 . How many ways are there to color its vertices black or white such that, for any vertex, if all of its neighbors are the same color then it is also that color? Two vertices are neighbors if they are the two endpoints of some edge of the cube. Please reason step by step, and put your final answer within boxed.
<|im_end|>
```

Prefixed Problem

```
<|im_start|>user
Consider the cube whose vertices are the eight points  $(x,y,z)$  for which each of  $x,y$ , and  $z$  is either 0 or 1 . How many ways are there to color its vertices black or white such that, for any vertex, if all of its neighbors are the same color then it is also that color? Two vertices are neighbors if they are the two endpoints of some edge of the cube. Please reason step by step, and put your final answer within boxed. <|im_end|>
<|im_start|>assistant
<think> This is a complex or challenging question, and it is difficult to provide a direct and correct answer. I need to think about it. Well, so we have a cube with vertices at all 0-1 triples, so 8 vertices total. Each vertex has degree 3, right? Because in a cube, each vertex is connected to three others (differing in one coordinate). The problem says: color vertices black or white such that for any vertex, if all its neighbors are the same color, then it's also that color. So let's parse that condition.
Let's rephrase the condition as a logical implication to make it clearer. For a vertex  $v$ , let  $N(v)$  be its neighbors (3 neighbors). Let  $c(v)$  be color, say 0=white, 1=black for simplicity (though labels might not matter, just two colors). The condition is: if  $c(u) = c(w)$  for all  $u, w$  in  $N(v)$ , then  $c(v) = c(u)$ . So in other words, it's not allowed for a vertex to be a different color from all its neighbors when all neighbors are the same color. Wait, actually, "if all neighbors are same color, then it's that color" – so the only way this fails is if there's a vertex where all neighbors are, say, black, but the vertex is white. Similarly, all neighbors white but vertex black is bad.
So forbidden configurations: a vertex with 3 black neighbors being white, or 3 white neighbors being black. All other configurations are allowed? Let's check: if a vertex has neighbors of different colors, then the condition is vacuously true? Wait, no: the condition is a conditional. The antecedent is "all neighbors same color", so if the antecedent is false (i.e., neighbors are not all same color), then the implication is true regardless of the vertex color. Yes! That's key. So the only forbidden things are vertices where all three neighbors are white and
```

the vertex is black, or all three neighbors are black and the vertex is white. For any other coloring (where for every vertex, either it has mixed neighbors or if it has monochromatic neighbors, it's the same color as those neighbors), it's okay.

So maybe another way: for each vertex, when we look at its neighborhood (the three adjacent vertices), if the neighborhood is monochromatic (all one color), then the vertex must be that color. If the neighborhood is not monochromatic (has both colors), no restriction on the vertex color (well, it's already colored, but the condition doesn't impose anything here).

So our goal is to count all 256 colorings satisfying this condition for every vertex.

Maybe it's helpful to model this as a graph. The cube graph Q_3 , 8 vertices, each degree 3. We need colorings $f: V(Q_3) \rightarrow \{0,1\}$ such that for all v in V , if $f(N(v))$ is constant, then $f(v) = f(N(v))$. So equivalently, for all v , if $f(N(v))$ not in $f(\text{neighbor})$ for some neighbor (i.e., $N(v)$ not monochromatic), then no problem; if $N(v)$ monochromatic, then $f(v)$ must equal that color.

Let me try to list all possible colorings or find a pattern. Maybe start with colorings that are "monochromatic": all black or all white. Let's check if these work.

All black: take any vertex, all neighbors are black, so condition says it must be black – which it is. Good. Similarly all white: same thing, all neighbors white, so vertex must be white – works. So 2 colorings here.

Now what about colorings with some black, some white. Let's think about smaller cubes? Wait, maybe not. Let's consider linear extensions or something, but maybe better to think about the structure of the cube.

Another idea: in graph theory, such conditions can sometimes be related to "fixed points" under some operation, but maybe not. Alternatively, let's consider the possible colorings by their "stability" or by what constraints they impose.

Wait, let's take an example. Suppose we have a single black vertex, say the origin $(0,0,0)$. Then its neighbors are $(1,0,0)$, $(0,1,0)$, $(0,0,1)$, all white (by assumption). Now check $w = (1,0,0)$, which is white. What are its neighbors? $(0,0,0)$ [black], $(1,1,0)$, $(1,0,1)$. So w has one black neighbor, two unknowns so far. Let's say we color $(1,1,0)$ and $(1,0,1)$ white. Then w has neighbors: black, white, white – so not monochromatic, so condition is satisfied for w . Now check the other endpoints: $(1,1,0)$ is white, neighbors: $(1,0,0)$ [white], $(0,1,0)$ [white], $(1,1,1)$ [unknown]. If we color $(1,1,1)$ black, then $(1,1,0)$ has neighbors: white, white, black – not monochromatic, okay. $(1,0,1)$ similar: neighbors white, white, black (at $(1,1,1)$), so okay. Now $(1,1,1)$ is black, neighbors: $(1,1,0)$ [white], $(1,0,1)$ [white], $(0,1,1)$ [unknown]. If we color $(0,1,1)$ white, then $(1,1,1)$ has neighbors white, white, white – all white, but it's black! Uh-oh, contradiction!

D.1.2 Off-policy RL Baselines

In our work, we use CISPO (Chen et al., 2025a) to compute the policy gradient following Khatri et al. (2025) which found it to work best on long RL runs. Moreover, CISPO can also handle off-policy updates, *i.e.*, update π^t on a trajectory τ sampled from $\mu \neq \pi^t$ with an importance weighting term common in off-policy RL (Fujimoto et al., 2018). For each \mathbf{x} in batch \mathcal{B} , CISPO samples k responses $\{\mathbf{y}_i^x\}_{i=1}^k$ where $\mathbf{y}_i^x \sim \mu(\cdot | \mathbf{x})$, then the CISPO policy gradient is given by:

$$\text{(CISPO)} \quad \frac{1}{\text{token-sum}} \sum_{\mathbf{x} \in \mathcal{B}} \sum_{i=1}^k \sum_{h=1}^k \left(\text{stop-grad}(\max(w(\mathbf{x}, \mathbf{y}_{i,h}^x), \varepsilon_{\text{high}})) \cdot A(\mathbf{x}, \mathbf{y}_i^x) \cdot \nabla_{\pi} \log \pi^t(\mathbf{y}_{i,h}^x | \mathbf{x}, \mathbf{y}_{i,h}^x) \right),$$

$$\text{where } w(\mathbf{x}, \mathbf{y}_{i,h}^x) = \frac{\pi^t(\mathbf{y}_{i,h}^x | \mathbf{x}, \mathbf{y}_{i,<h}^x)}{\mu(\mathbf{y}_{i,h}^x | \mathbf{x}, \mathbf{y}_{i,<h}^x)} \quad \text{and} \quad \text{token-sum} = \sum_{\mathbf{x} \in \mathcal{B}} \sum_{i \in [k]} |\mathbf{y}_i^x|. \quad (\text{D.1})$$

In (D.1), the advantage $A(\mathbf{x}, \mathbf{y}_i^x)$ is computed by subtracting the baseline $\bar{r}(x) = \frac{1}{k} \sum_{i=1}^k r(\mathbf{x}, \mathbf{y}_i^x)$ from $r(\mathbf{x}, \mathbf{y}_i^x)$. The per-token importance weight $w(\mathbf{x}, \mathbf{y}_{i,h}^x)$ accounts for the distribution shift between the current policy π^t and the sampler μ and is 1 for on-policy traces where $\pi^t = \mu$. To reduce gradient variance from importance weights, it is clipped at $\varepsilon_{\text{high}}$ and in practice we set it to a value of 0.01.

In our setup, the off-policy dataset \mathcal{D}_{off} is constructed by *rejection sampling* the base policy π^0 : for each prompt \mathbf{x} , we repeatedly sample $\mathbf{y} \sim \pi^0(\cdot | \mathbf{x})$ until we obtain one *correct* trajectory (according to the verifier), and store that successful trajectory in \mathcal{D}_{off} . This procedure induces an *accepted* (conditional) behavior distribution

$$\mu_{\text{acc}}(\mathbf{y} | \mathbf{x}) = \pi^0(\mathbf{y} | \mathbf{x}, r(\mathbf{x}, \mathbf{y}) = 1) = \frac{\pi^0(\mathbf{y} | \mathbf{x}) \mathbf{1}\{r(\mathbf{x}, \mathbf{y}) = 1\}}{a(\mathbf{x})}, \quad a(\mathbf{x}) := \Pr_{\mathbf{y} \sim \pi^0(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y}) = 1]. \quad (\text{D.2})$$

Thus, when we treat accepted trajectories as "samples from μ " in (D.1), the correct sequence-level importance

ratio for an accepted trajectory \mathbf{y} is

$$\frac{\pi^t(\mathbf{y}|\mathbf{x})}{\mu_{\text{acc}}(\mathbf{y}|\mathbf{x})} = a(\mathbf{x}) \cdot \frac{\pi^t(\mathbf{y}|\mathbf{x})}{\pi^0(\mathbf{y}|\mathbf{x})} \quad (\text{since } \mathbf{1}\{r(\mathbf{x},\mathbf{y})=1\}=1 \text{ for accepted } \mathbf{y}). \quad (\text{D.3})$$

The key point is that the acceptance-probability correction $a(\mathbf{x})$ is a *trajectory-level* factor: it appears once per accepted sequence, not once per token. In practice, we estimate $a(\mathbf{x})$ directly from the rejection-sampling effort. Let $R(\mathbf{x})$ be the number of rollout attempts required to obtain one correct trace for \mathbf{x} during dataset construction; then $\hat{a}(\mathbf{x}) \approx 1/R(\mathbf{x})$. CISPO, however, uses per-token importance weights $w(\mathbf{x}, y_{i,h}^{\mathbf{x}})$ (Eq. (D.1)) and aggregates gradients across tokens. A simple way to incorporate the acceptance correction is to multiply *each token* in an accepted trajectory by $\hat{a}(\mathbf{x})$:

$$\tilde{w}(\mathbf{x}, y_{i,h}^{\mathbf{x}}) := \hat{a}(\mathbf{x}) \cdot \frac{\pi^t(y_{i,h}^{\mathbf{x}}|\mathbf{x}, \mathbf{y}_{i,<h}^{\mathbf{x}})}{\pi^0(y_{i,h}^{\mathbf{x}}|\mathbf{x}, \mathbf{y}_{i,<h}^{\mathbf{x}})}. \quad (\text{D.4})$$

This heuristic is slightly biased relative to the sequence-level ratio in (D.3): multiplying every token effectively scales an accepted trajectory’s total contribution by approximately $\hat{a}(\mathbf{x}) \cdot |\mathbf{y}|$ (modulo the global token normalization), whereas the exact correction would apply $\hat{a}(\mathbf{x})$ once at the trajectory level. When accepted trajectories have heterogeneous lengths, this introduces a mild length-dependent bias. Empirically, we found this approximation to be stable, and it preserves the intended qualitative effect: prompts that are harder under π^0 (larger $R(\mathbf{x})$, smaller $\hat{a}(\mathbf{x})$) receive smaller off-policy gradient mass, reflecting the fact that an accepted sample from μ_{acc} is “more selective” for those prompts.

LUFFY baseline: mixed-policy GRPO with policy shaping. We also compare against LUFFY (Yan et al., 2025), which incorporates off-policy reasoning traces by mixing them with on-policy rollouts inside a GRPO-style objective. LUFFY computes advantages using *group computation* over a mixed set of rollouts: for each prompt it combines N_{on} on-policy samples with N_{off} off-policy traces, and normalizes rewards using the mean and standard deviation over the union of the two groups. In our reproduction, we follow LUFFY’s “fair” setting by using 8 total samples per prompt with a 1-off-policy / 7-on-policy split, rollout batch size 128, update batch size 64, and rollout temperature 1.0. For optimization, LUFFY uses a constant learning rate of 10^{-6} and trains for 500 steps. We also remove KL regularization (setting $\beta=0$) and uses an entropy-loss coefficient of 0.01. Finally, LUFFY introduces *policy shaping* via a regularized importance-sampling transformation controlled by a parameter γ , and for this we use $\gamma=0.1$, chosen after sweeioing over $\{0.01, 0.1, 0.2\}$.

D.1.3 Hyperparameter Details for PrefixRL

We use the REINFORCE (Ahmadian et al., 2024) on-policy algorithm for PrefixRL and standard RL. For this, we set the training batch size of 128 with, a constant learning rate of 1×10^{-6} . We turn off any KL regularization and also disable entropy regularization (entropy coefficient 0). We also use a gradient clipping of 1.0. We set the sampling temperature for training at 1.0 for Llama3.1-8b-instruct and 0.7 for Qwen-3-4b-instruct. At test-time we sample with a temperature of 0.7 for both models, including the inference to collect data for rejection sampling. For all our PrefixRL runs we use $n=8$ rollouts per prompt in the batch. We use the same for standard RL, and off-policy RL, except specified otherwise. In all our runs in Section 5 and Section 4 we run RL training for 400 iterations, except for our RL runs in Section 4.3 and Section 4.2 where run the training for 100 iterations.

Before we run RL, we finetune the Llama3.1-8b-instruct model on OpenThoughtsV3 (Guha et al., 2025). For this, we first filter the dataset to only retain responses of token length < 24192 . Then, we run SFT for 5 epochs on this data at peak learning rate of $8e-5$. We use a batch size of 512 traces per batch and a cosine learning rate schedule that has a linear warm up (for 10% of the total iterations) followed by a cosine decay to a learning rate of $8e-6$. We use a hold out validation set to measure the negative log-likelihood loss during training, and pick the earliest checkpoint with the least validation loss as the final distilled model.

D.2 FLOPs Accounting for our Compute-Matched Performance Plots

We compute FLOPs using the standard Transformer approximation: processing D tokens with a model of N trainable parameters costs $\approx 2ND$ FLOPs for a *forward-only* pass (sampling/inference) and $\approx 6ND$ FLOPs for a *training update* (forward + backward + gradient computation) (Snell et al., 2024).

Definitions. We use N to denote the number of trainable parameters of the model whose compute is being measured and D to denote the total number of tokens processed by that model in the relevant stage, summed over all sequences (prompt + generated tokens).

Per-iteration compute. At RL iteration t , let $D_{\text{samp}}^{(t)}$ be the total number of tokens generated/evaluated during rollout sampling, and let $D_{\text{upd}}^{(t)}$ be the total number of tokens used in gradient-based optimization. We estimate FLOPs as

$$\text{FLOPs}^{(t)} = 2ND_{\text{samp}}^{(t)} + 6ND_{\text{upd}}^{(t)}. \quad (\text{D.5})$$

Cumulative compute. The x-axis in our compute-matched plots reports cumulative FLOPs after T iterations:

$$\text{FLOPs}_{\leq T} = \sum_{t=1}^T \text{FLOPs}^{(t)}. \quad (\text{D.6})$$

If a method includes up-front rejection sampling to construct \mathcal{D}_{off} , we add that forward-only cost (also using $2ND$) to (D.6). So, if we need to sample R times before we accept a correct trace for a problem \mathbf{x} , then the total upfront compute spent on the problem is $2RND$.

D.3 Qwen with Llama Prefixes

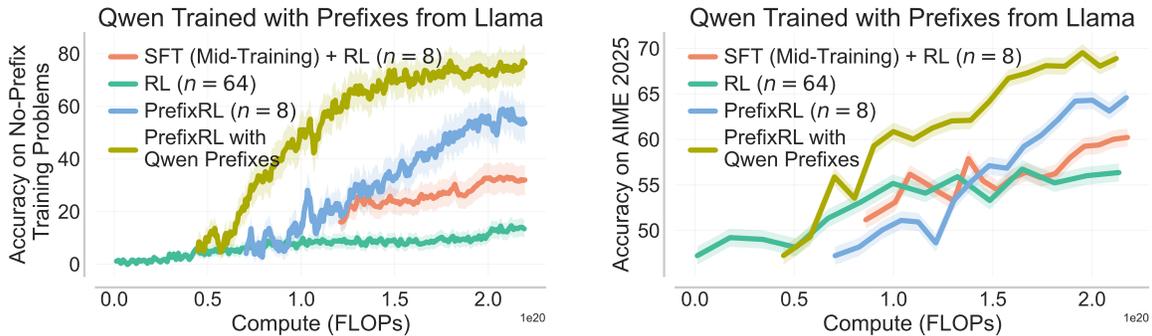


Figure 16 Training Qwen model with prefixes from Llama. Opposite of our experiment in Figure 12, here we train Qwen3-4b-instruct base LLM using off-policy traces sourced by rejection sampling Llama-3.1-8b-instruct model. We find that while off-policy prefixes from Llama are not as effective at improving Qwen as the other way around. Moreover, we also end up spending more compute upfront on rejection sampling the traces with the less capable Llama base model. We show accuracy on no-prefix training problems (*left*) and AIME 2025 (*right*).

Figure 16 complements the cross-family back-generalization results in Figure 8 by reversing the direction of prefix transfer: instead of training Llama using Qwen-sourced off-policy prefixes, we train Qwen using prefixes sourced from Llama. The key takeaway is asymmetric transfer. While Qwen prefixes can still drive back-generalization when optimizing Llama (Section C.2), Llama prefixes are noticeably less effective for improving Qwen in both training accuracy (no-prefix training problems) and standardized evaluation (AIME 2025).

This asymmetry is consistent with the back-generalization discussion: PrefixRL relies on prefix states that are injected from an external distribution, and the degree to which learning transfers to no-prefix behavior depends on how informative and “compatible” those prefix states are with the target model’s internal representations and solution strategies. When prefixes are sourced from a more capable model family (here, Qwen), they tend to encode higher-quality intermediate reasoning states, and RL can more readily propagate improvements from prefixed states to earlier states and ultimately to the no-prefix setting. In contrast, prefixes sourced from a less capable model (here, Llama) are both (i) harder to obtain via rejection sampling and (ii) less likely to contain strategy-revealing intermediate states that meaningfully constrain the continuation policy, resulting in weaker transfer when training Qwen.

A second practical implication is compute efficiency. Since PrefixRL amortizes training over a fixed pool of off-policy traces, the total compute depends not only on the on-policy RL phase but also on the *upfront* cost

of harvesting correct traces. Rejection sampling from the weaker base model increases this upfront cost, and Figure 16 shows that even after paying that cost, the resulting prefixes yield smaller downstream gains for Qwen. Together with Figure 8, these results suggest that cross-family prefix sourcing is most attractive when (a) the source model is strong enough to produce diverse correct traces at reasonable cost, and (b) the injected prefix states align with the target model sufficiently well to allow back-generalization to propagate to the no-prefix distribution.

D.4 Computing the gradient norm and standard deviation metrics

To quantify training signal-to-noise, we track (i) the *norm of the expected gradient* and (ii) the *standard deviation of the sampled gradient* throughout RL training, using the same procedure described at the end of Section 5.

Let $g_t \in \mathbb{R}^N$ denote the (flattened) stochastic policy gradient computed at iteration t from the current minibatch (including any on-policy and/or off-policy contributions, depending on the method). We maintain exponential moving averages (EMA) of the first and second moments *coordinate-wise*:

$$m_t = \beta m_{t-1} + (1 - \beta) g_t, \tag{D.7}$$

$$v_t = \beta v_{t-1} + (1 - \beta) (g_t \odot g_t), \tag{D.8}$$

where \odot denotes elementwise multiplication and $\beta \in (0, 1)$ is a fixed smoothing constant.

Gradient norm. We report the norm of the mean gradient estimate as

$$\text{GRADNORM}_t = \|m_t\|_2. \tag{D.9}$$

Gradient standard deviation. We estimate the (coordinate-wise) variance as $s_t = v_t - m_t \odot m_t$ and report

$$\text{GRADSTD}_t = \sqrt{\sum_{i=1}^N \max\{(s_t)_i, 0\}} = \sqrt{\text{tr}(\widehat{\text{Cov}}(g_t))}. \tag{D.10}$$

Equivalently, this standard-deviation metric corresponds to estimating the trace of the gradient covariance matrix via first/second moments.