
Distributional Biases in Post-Training: A Markovian Analysis of Reasoning Trajectories

Anonymous Authors¹

Abstract

Foundation models exhibit broad knowledge but limited task-specific reasoning, motivating post-training strategies such as RL with verifiable rewards (RLVR) and test-time scaling (TTS). While recent work highlights the role of *exploration* in improving pass@K, empirical evidence points to a paradox: RLVR and ORM/PRM typically reinforce existing paths rather than expanding the reasoning scope, raising the question of why exploration helps if no new patterns emerge. To reconcile this paradox, we adopt the perspective of Kim et al. (2025), viewing easy (e.g., simplifying a fraction) versus hard (e.g., discovering the some symmetry) reasoning steps as low versus high probability Markov transitions. In this tractable model, pretraining corresponds to tree-graph discovering, while post-training corresponds to CoT reweighting. We provably show that, both RLVR and ORM/PRM would favor heavily to several high-probability paths, and thereby forget rare-but-crucial CoTs. Building on this, we further prove that exploration strategies such as rejecting easy instances and KL regularization help preserve rare CoTs. Empirical simulations corroborate our theoretical results.

1. Introduction

Foundation models provide broad knowledge and versatile capabilities across tasks, yet their task-specific reasoning remains constrained. For reasoning datasets with only 0/1 verifiers, many studies explore post-training strategies, including Reinforcement Learning with Verifiable Reward (RLVR) finetuning (Xin et al., 2024; Shao et al., 2024; Guo et al., 2025a; Yu et al., 2025), as well as inference scaling with Outcome Reward Models (ORM) or Process Reward

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

Models (PRM) (Lightman et al., 2023; Snell et al., 2024), both aiming to obtain task-specific experts.

Recently, a line of work has emphasized maintaining *exploration* and *entropy stability* to prevent entropy collapse, and observed that suitable entropy preservation during post-training yields *systematic performance gains*, such as improved pass@K on math benchmarks (Xiong et al., 2025; Li et al., 2025b; Ren & Sutherland, 2025; Wang et al., 2025b; Cui et al., 2025; Zuo & Zhu, 2025).

However, seemingly *paradoxical* findings emerge: RLVR typically aligns models with target objectives by reinforcing existing reasoning paths, *rather than* expanding their tree-like reasoning scope (Snell et al., 2024; Yue et al., 2025; AI et al., 2025; Gandhi et al., 2025). Similarly, ORM/PRM-guided inference scaling biases models toward pre-existing Chain-of-Thought (CoT) patterns instead of incentivizing genuinely new branches. This raises a natural question:

Why can exploration help, given post-training cannot explore beyond the base model’s tree scope?

Our work takes a first step toward reconciling this tension, motivated the key phenomena below.

Phenomenon 1: Squeezing Effect of RLVR. RLVR implicitly reduces CoT entropy (Li et al., 2025a; Wu et al., 2025a; Deng et al., 2025), shrinking confidence over CoTs with high frequency to be correct, and sometime at the cost of forgetting certain correct CoTs within the base model’s scope (Wu et al., 2025a; Shao et al., 2024; Wen et al., 2025).

Phenomenon 2: Neural Verifier Checks Consistency, Not Accuracy. Inference-scaling with ORM/PRM may avoid the learner’s squeezing effect, yet empirical evidence shows that neural scorers are prone to reward consistency rather than true accuracy (Xu et al., 2025; Guo et al., 2025b). As a result, they would favor CoTs that follow *common, high-frequency* reasoning patterns.

Phenomenon 3: Merits of Rare CoTs. A widely-utilized difficulty measure for reasoning dataset (e.g., GSM8K (Cobbe et al., 2021) and AQuA (Ling et al., 2017)) is the *pass rate* (i.e., the frequency with which a base model correctly solves an instance under parallel attempts) (Tong et al., 2024; Parashar et al., 2025). This implies that hard instances

correspond to rare-but-correct CoTs with low model confidence, whereas common CoTs typically reflect frequent patterns to easier instances.

Taken together, these observations suggest a resolution:

Exploration, even when confined within the existing tree scope, helps prevent the model from entirely forgetting rare CoTs that may be crucial for hard instances, and preserved broad-capability.

Our Contributions. In this work, we rigorously formalize and prove these phenomena within a tractable theoretical framework. Motivated by the view that discrete graphs naturally abstract the sequential structure of complex reasoning (Xu et al., 2019; Sanford et al., 2024; Abbe et al., 2024; Besta et al., 2024), we model each reasoning step as a Markov *state transition* following Kim et al. (2025). Pretraining is framed as a *tree-graph* discovering process over child states across tasks, while post-training CoT generalization is modeled by a Multi-task Tree-structured Markov Chain (TMC). We prove that our toy model captures key **Phenomena 1–3** with population 0/1 reward (expected accuracy), and then provide theoretical justification for exploration techniques such as rejecting easy questions (Yu et al., 2025; Xiong et al., 2025; Zhang et al., 2025a) and KL regularization. Our paper is organized as below.

- Sec. 2 introduces a multi-task Tree-structured Markov chain model to capture diverse CoT reasoning patterns across tasks, explicitly linking instance difficulty with pass rate.
- Sec. 3 analyzes a simple softmax model and shows that RLVRs inherit a *simplicity bias*, over-favoring easier CoTs due to the *advantage-driven squeezing effect*. We further provide theoretical justification for rejecting easy instances and applying KL regularization, which both promote valid hard CoT learning.
- Sec. 4 demonstrates that inference-scaling with ORM/PRM assigns credit to CoTs by their accuracy likelihood, leading to overemphasis on easier CoTs. We further show that PRMs with BoN can be interpreted as special cases of a more general *Doob h’s Transformed-induced PRM* (DPRM). In principle, DPRM is equivalent to soft-BoN (Verdun et al., 2025) asymptotically, enabling adjustable preservation of base-model capabilities and better alignment with hard-to-reason and cross-task patterns.

Discussions of additional related work are in App. A. All proofs are deferred to the appendix.

Humble Remark. While we prove **empirically-observed Phenomena 1–3** and the benefits of those existing techniques in our theory-friendly setting that captures partial

but crucial rationales, we do not overclaim their *direct* applicability to GPT or large-scale models, given the many unmodeled complexities, as discussed in App. B and D.4.

2. Tree-structured Multi-task Reasoning

2.1. Multi-Task CoT as Tree-structured Markov Chains

We propose Tree-structured Markov Chain (TMC) framework to abstractly model the tree-like reasoning capability of base model, following (Kim et al., 2025; Nichani et al., 2024).

Definition 2.1 (Tree-structured Markov Chains (TMC)). A process $X = (X_t)_{t \geq 0}$ is defined on a finite state space $S = \bigcup_{l=1}^L S_l$, where $S_l \cap S_{l'} = \emptyset$ for $l \neq l'$, and transitions occur from S_l to S_{l+1} with probability kernel $\mathbb{P}(\cdot | o_l)$ for $o_l \in S_l$. Define $M_0 = |S_1|$ and $M = \max_{l, o_l \in S_l} |C_{o_l}|$, where $C_{o_l} \subset S_{l+1}$ is the high-probability transition subset. The TMC satisfies:

- Root states $o_1 \in S_1$ are sampled with $\mathbb{P}_{\text{TMC}}(o_1) = \Theta(1/M_0)$.
- For $o_l \in S_l$ and $o_{l+1} \in C_{o_l}$, we have $\mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) = \Theta(1/M)$, while if $o'_{l+1} \notin C_{o_l}$, $\mathbb{P}_{\text{TMC}}(o'_{l+1} | o_l) = o(1/M^2)$ (feeble, $\geq c > 0$) or 0.
- The topology ensures that for each $q \in S_1$ there are $n_q = O(1) \geq 1$ high probability CoT traces ($q = o_1, \dots, o_L$), i.e. traces with $o_{l+1} \in C_{o_l}, \forall l \in [L-1]$.

In our TMC (Def. 2.1), *states* represent *logical assertion* (e.g., a sentence or mathematical expression) rather than surface tokens (Kim et al., 2025). Especially, the CoTs with ≥ 1 sparse edge (i.e., edge with feeble transition probability $o(1/M^2)$) are called *hard-to-reason* CoTs, otherwise *easy-to-reason* CoTs. The following definition formalizes the reason we called them “*easy*” or “*hard*” based on their uncertainty, modeling after the widely utilized difficulty measure—namely *pass rate*—for reasoning dataset (e.g., GSM8K (Cobbe et al., 2021) and AQuA (Ling et al., 2017)).

Definition 2.2 (Multi-task Capability in TMC. (Informal)). Let $X = (X_t)_{t \geq 0}$ be a TMC (Def. 2.1), and let \mathcal{T} be a set of tasks. Each task $k \in \mathcal{T}$ is specified by a collection of state tuples (q, a, k) , where all tuples have distinct q and a . Among these CoTs, a nonempty subset is *valid* for (q, a, k) , satisfying:

- (i) all easy-to-reason CoTs are *valid* for (q, a, k) and *invalid* for any $k' \neq k$;
- (ii) every nonzero transition in TMC appears in ≥ 1 *valid* CoT across all tasks;
- (iii) each (q, a, k) induces a QA distribution $\mathcal{D}_a^{q,k}$, and for any sampled instance $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_a^{q,k}$, only a *determined* subset of *valid* CoTs is *correct*, where the probability that any *valid* CoT is *correct* for the (\mathbf{Q}, \mathbf{A}) is **proportional** to its likelihood among *valid* CoTs.

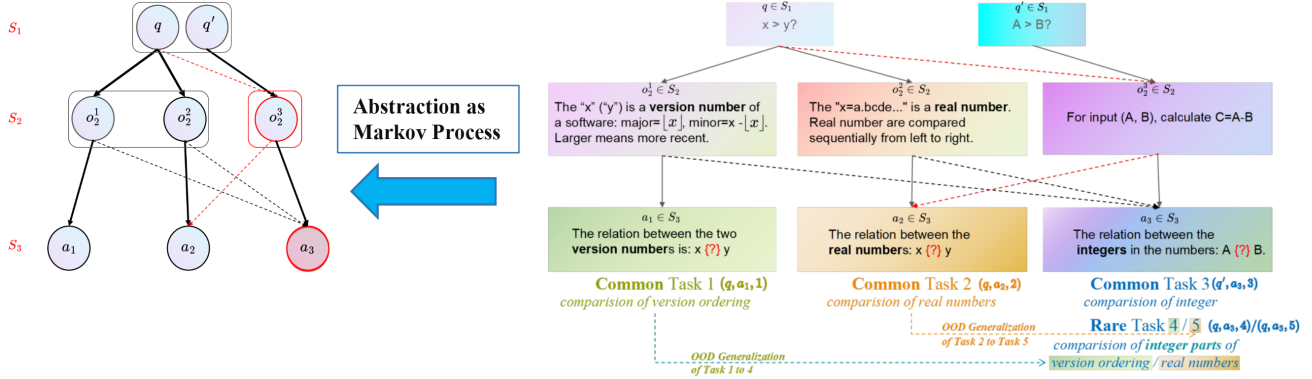


Figure 1. **Left:** abstraction of a 3-layer TMC. Nodes are states grouped into layers S_1 – S_3 ; solid arrows denote high-prob (confident) transitions and dashed arrows denote low-prob (unsure) transitions. A task is specified by $(\mathbf{q}, \mathbf{a}, \mathbf{k})$, where $\mathbf{q} \in \{q, q'\}$ is the question state, $\mathbf{a} \in \{a_1, a_2, a_3\}$ is the answer state, and $\mathbf{k} \in [5]$; **Right:** a concrete illustration of a 5-task, 3-layer Multi-task TMC. x, y in q represent real numbers (with decimals), whereas A, B in q' represent integers. We here use two instances, namely $3.9 > 3.11?$ of q and $3 > 3?$ of q' to describe the five tasks: (1) decimal version ordering ($3.9 < 3.11$); (2) real-number comparison ($3.9 > 3.11$); (3) integer equality ($3 - 3 = 0$); (4) integer-part version ordering (e.g., $3.9 = 3.11$); and (5) integer-part real-number comparison ($3.9 = 3.11$). Tasks 1–3 are common and each admits ≥ 1 easy-to-reason CoT, while Tasks 4–5 are rare and admit *only* hard-to-reason CoTs. For Task 2, there are two *valid* CoTs: $q \rightarrow o_2^2 \rightarrow a_2$ (where o_2^2 merely left-to-right compares number) and $q_2 \rightarrow o_3^2 \rightarrow a_2$ (where o_3^2 performs the arithmetic calculation). The instance $3.9 > 3.11?$ admits both CoTs *correct*. However, for *hard* question instances such as $0.8 + 3.1 > 2.11 + 1.0?$, *only* the hard-to-reason CoT $q_2 \rightarrow o_3^2 \rightarrow a_2$ is *correct*, since it requires explicit calculation—left-to-right token comparison alone doesn’t suffice.

A task is *common* if it admits ≥ 1 *valid* easy-to-reason CoTs, and *rare* otherwise.

Fig. 1 illustrates Def. 2.1 and Def. 2.2; a more detailed version of Def. 2.2 (Def. E.2) is given in Sec. E. Condition (i) avoids major task conflicts, (ii) removes redundancy so every edge contributes, and (iii) links model confidence to *pass rate* per **Phenomenon 3**, enabling error analysis. We distinguish two notions of CoT:

- **Validity**¹: a **task-level** property, indicating whether it solves *any* (\mathbf{Q}, \mathbf{A}) under (q, a, k) .
- **Correctness**: an **instance-level** property, deterministically defined for a specific (\mathbf{Q}, \mathbf{A}) .

Why this definition? First, **Phenomenon 1** directly motivates us to formally bridge uncertainty and *pass rate* across tasks. Yet, not all rare outputs are useful—some may be not *correct* for all instances of the current task, surfacing only due to *shared* reasoning states across tasks. This motivates our definition of *validity* to distinguish useful task-specific CoTs. Second, the key property inherent from the pass rate is that, easy-to-reason CoTs cover most instances in $\mathcal{D}_a^{q,k}$, but some instances can still *only* be *correctly* solved by hard-to-reason CoTs. Our Multi-task TMC framework highlights the importance of such rare reasoning paths, consistent with large-scale evidence that many errors on GSM8K, AQUA, and MATH arise from *misapplied common patterns* (e.g., assuming overlapping events are independent) per observed

¹Our *validity* also speaks that no all hard-to-reason CoTs useful, see App. B a discussion.

in Sun et al. (2025); more relevant empirical examples appear in Rmk. E.1. We can then define outcome signal to verify *correctness* as follow.

Outcome Reward Signal. Let $\mathbf{o}_l \in \mathbb{R}^{|S_l|}$ be the one-hot encoding of observation o_l , and $\mathbf{o} = (\mathbf{o}_1, \dots, \mathbf{o}_L)^\top$ the full trajectory. In mathematical reasoning tasks, the *correctness* of a CoT trace is deterministic and verifiable (Yue et al., 2025; Xiong et al., 2025; Setlur et al., 2025a; 2024; 2025b), typically via formal systems such as *Lean4* (Yang et al., 2023). Hence, for any QA pair $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$ in task $k \in \mathcal{T}$, we define the $R_{(\mathbf{Q}, \mathbf{A})}^k(\cdot) : \mathbb{R}^{L \times |S|} \rightarrow \{0, 1\}$ as

$$R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) = \mathbb{1}(\mathbf{o} \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}), \quad (1)$$

where $\mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}$ is the deterministic set of *correct* CoTs for \mathbf{Q} within the *valid* CoT collection for task tuple (q, a, k) —in practice, *Lean4* only certifies the overall correctness of a CoT trace—providing an outcome reward signal—without verifying individual process reasoning steps.

2.2. Pretrained Base Model

Base Model. For simplicity, following (Kim et al., 2025), we model the LLM base model using a straightforward linear softmax predictor:

$$\hat{p}_\theta(\cdot|x) = \text{softmax}(\langle \theta, x \rangle), \quad (2)$$

where $\theta \in \mathbb{R}^{|S| \times |S|}$ and $x \in \{0, 1\}^{|S|}$ is a one-hot vector. This formulation is theoretically tractable and plausible, as noted by Li et al. (2025b); Ren & Sutherland (2025);

Chen et al. (2025), which highlight that the LLM’s final layer employs logits $h_\theta(\cdot, \mathbf{x})$, encoded in the last token, to generate a softmax distribution over the vocabulary as the predictive probability for the next token. Following (Kim et al., 2025), we train the base model through entropy loss as below, akin to the next-token prediction process despite in the Markov chain setting.

Theorem 2.3 (Informal Version of Thm. F.1). *Let $X_0 \sim \text{Unif}(S \setminus S_L)$ and $X_1 \sim \mathbb{P}(\cdot|X_0)$ be random samples from the TMC X in Def. 2.1, the softmax predictor trained by cross-entropy $L_{CE} = \mathbb{E}_{X_0, X_1}[\log \hat{p}_{\theta^{(t-1)}}(X_1|X_0)]$ via Alg. 1 achieves the following: (1) After $T_1 = \tilde{O}(M^2)$ iterations, the uniform convergence error of the predictor is $\tilde{O}(\sqrt{M/T})$. (2) After thresholding, the predictor converges linearly to the true probabilities with error decaying as $\tilde{O}(e^{-\Omega(T)})$.*

Similar to the treatment in (Kim et al., 2025), in the subsequent sections, we suggest that the pretrained θ^* achieve the exact transition probability as the TMC model $\hat{p}_{\theta^*} = \mathbb{P}$ after pretraining. This is plausible given the longer timescales of pretraining relative to finetuning and inference.

3. Simplicity Bias of RLVR Finetuning: Provable Challenge and Solution

In this section, we first analyze the inherent simplicity biases of the standard RLVR finetunings, and then provide theoretical justifications for certain strategies that can alleviate this issue. Throughout, the expectation $\mathbb{E}[\cdot]$ is operated on $\mathbf{o}_1^i \sim P^k(\mathcal{Q}^k)$, $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$, $\{\mathbf{o}^i\}_{i=2}^G \sim \hat{p}_{\theta}^k(\mathcal{O}|\mathbf{o}_1^i)$.

REINFORCE & RAFT. In terms of mathematics dataset, the standard REINFORCE objective maximizes the correctness of the sampled CoTs (Xiong et al., 2025; Setlur et al., 2025a) (i.e., $R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) = 1$ in our case for task $k \in \mathcal{T}$). Separately, RAFT (Rejection Sampling Finetuning) (Dong et al., 2023; Touvron et al., 2023; Yuan et al., 2023) maximize cross-entropy on successful CoT sampled from current policy. Their objectives in our TMC case are

$$\mathcal{J}_{\text{REINFORCE}}(\theta) = \mathbb{E} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right], \quad (3)$$

$$\mathcal{J}_{\text{RAFT}}(\theta) = \mathbb{E} \left[\sum_{l=1}^{L-1} \log \hat{p}_{\theta}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right]. \quad (4)$$

PPO & GRPO. Proximal Policy Optimization (PPO) (Schulman et al., 2017; OpenAI, 2018) and Group Relative Policy Optimization (GRPO) (Shao et al., 2024) both optimize clipped surrogate objectives with temperature $\beta > 0$:

$$\mathcal{J}_{\text{PPO}}(\theta) = \mathbb{E} \left[\frac{1}{L} \sum_{l=1}^{L-1} \min \left(\frac{\hat{p}_{\theta}(\mathbf{o}_{l+1}|\mathbf{o}_l^i)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}|\mathbf{o}_l^i)} A_{l+1}^{\hat{p}_{\theta}, k}, \text{clip} \left(\frac{\hat{p}_{\theta}(\mathbf{o}_{l+1}|\mathbf{o}_l^i)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}|\mathbf{o}_l^i)}, 1 - \epsilon, 1 + \epsilon \right) A_{l+1}^{\hat{p}_{\theta}, k} \right) \right] \quad (5)$$

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} \left[\frac{1}{GL} \sum_{i=1, l=1}^{G, L-1} \min \left(\frac{\hat{p}_{\theta}(\mathbf{o}_{l+1}^i|\mathbf{o}_l^i)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}^i|\mathbf{o}_l^i)} \hat{A}_{i, l+1}^k, \text{clip} \left(\frac{\hat{p}_{\theta}(\mathbf{o}_{l+1}^i|\mathbf{o}_l^i)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}^i|\mathbf{o}_l^i)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i, l+1}^k \right) \right] \quad (6)$$

Here, the RL advantage (Levine, 2018) at reasoning step l for task k and predictor \hat{p}_{θ} is

$$\begin{aligned} A_{l+1}^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) &= Q^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) - V^{\hat{p}_{\theta}, k}(\mathbf{o}_l), \\ V^{\hat{p}_{\theta}, k}(\mathbf{o}_l) &:= \mathbb{E} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \middle| \mathbf{o}_l \right], \\ Q^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) &:= \mathbb{E} \left[V^{\hat{p}_{\theta}, k}(\mathbf{o}_l) \middle| \mathbf{o}_{l+1} \right] \end{aligned} \quad (7)$$

PPO (5) typically estimates $A_{l+1}^{\hat{p}_{\theta}, k}$ via GAE with an additional critic model, while GRPO (6) employs group-normalized advantages $\hat{A}_{i, l+1}^k = (R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}^i) - \mu) / \sigma$ computed across sampled CoTs, with μ, σ the group mean and std across G sampled CoTs.

The following theorem shows that the above methods are inherently biased toward easy-to-reason CoTs per **Phenomenon 1** (Wu et al., 2025a; Deng et al., 2025), resulting failure over hard instances.

Theorem 3.1 (Squeezing Effect of RL-finetuning). *Consider a base model θ^* defined in Sec. 2.2 and a targeted task $k \in \mathcal{T}$ with $\Theta(M)$ valid hard-to-reason CoTs. Suppose we apply one of the following finetuning algorithms: REINFORCE, RAFT, PPO, or GRPO (without KL regularization) with access to the expected gradient oracle. For PPO/GRPO, assume the advantage is estimated accurately and the clipping threshold are functioning. Then, for any $\epsilon > 0$, there exists $t \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ such that for any valid hard to reason CoT \mathbf{o}^{hard} for task k , we have*

$$\Pr(\mathbf{o}_{2:L}^{\text{hard}} \sim \hat{p}_{\theta^{k, (t)}}^k(\cdot|\mathbf{o}_1^{\text{hard}})) \leq \epsilon.$$

Therefore, for any (\mathbf{Q}, \mathbf{A}) of task k , if all correct CoTs solving (\mathbf{Q}, \mathbf{A}) are hard-to-reason, then the finetuned model $\hat{p}_{\theta^{k, (t)}}$ satisfies

$$\mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^{k, (t)}}^k(\cdot|\mathbf{o}_1)} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right] \leq \epsilon.$$

Sketch of Proof. The key observation is the following proposition, showing that along any easy-to-reason CoT for a task, the hard-to-learn CoT deviate from it would have smaller advantage.

Proposition 3.2 (Advantage Gap between Easy and Hard CoT). *Let X be a Multi-task TMC as in Def. 2.1 and 2.2, fix a common task state tuple (q, a, k) . Then, for the shared states $\mathbf{o}_l, l \in [L-1]$ of any valid easy-to-reason CoT \mathbf{o}^{easy} and hard-to-learn CoT \mathbf{o}^{hard} , then there exists $c > 0$, such that $A_{l+1}^{\hat{p}_{\theta^*}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{easy}}) \geq c > A_{l+1}^{\hat{p}_{\theta^*}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{hard}}), \forall l \in [L-1]$.*

We then denote PO as the algorithm of the PPO/GRPO in Thm. 3.1, through standard policy gradient derivation with notation $\nabla := \nabla_{\theta}$; $\hat{p}_k := \hat{p}_{\theta}$, it holds that

$$\begin{aligned} \nabla \mathcal{J}_{\text{REINFORCE}} &= \sum_{l=1}^{L-1} \mathbb{E}[\nabla \log \hat{p}_k(\mathbf{o}_{l+1} | \mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})], \\ \nabla \mathcal{J}_{\text{RAFT}} &= \sum_{l=1}^{L-1} \mathbb{E}[(1 + \log \hat{p}_k(\mathbf{o}_{l+1} | \mathbf{o}_l)) \nabla \log \hat{p}_k(\mathbf{o}_{l+1} | \mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})], \\ \nabla \mathcal{J}_{\text{PO}} &= \sum_{l=1}^{L-1} \mathbb{E}[(1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_k, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon) \\ &\quad \cdot A_{l+1}^{\hat{p}_k, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \nabla \log \hat{p}_k(\mathbf{o}_{l+1} | \mathbf{o}_l)] \end{aligned} \quad (8)$$

where $\nabla \log \hat{p}_k(\mathbf{o}_{l+1} | \mathbf{o}_l) = \mathbf{o}_{l+1} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_k(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbf{o}'_{l+1}$ holds in linear case. For valid hard CoTs $\mathbf{o}_{l+1}^{\text{hard}, i}$ and easy CoTs $\mathbf{o}_{l+1}^{\text{easy}, j}$ sharing the same \mathbf{o}_l but a different $\mathbf{o}_{l+1}^{\text{hard}, i}, \mathbf{o}_{l+1}^{\text{easy}, j}$ at $l \in [L-1]$, where i, j are index of hard and easy valid CoTs, the logits update difference under $\mathcal{J}_{\text{REINFORCE}}(\theta)$ is

$$\begin{aligned} &\eta[\nabla_{\theta_{\mathbf{o}_{l+1}^{\text{hard}, i}, \mathbf{o}_l}} \mathcal{J}_{\text{REINFORCE}}(\theta) - \nabla_{\theta_{\mathbf{o}_{l+1}^{\text{easy}, j}, \mathbf{o}_l}} \mathcal{J}_{\text{REINFORCE}}(\theta)] \\ &= \eta[(A_{l+1}^{\hat{p}_k, k}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{hard}, i}) - A_{l+1}^{\hat{p}_k, k}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{easy}, j})) \\ &\quad + V_{\hat{p}_k, k}(\mathbf{o}_l)(\hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{hard}, i} | \mathbf{o}_l) - \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}, j} | \mathbf{o}_l))] < 0. \end{aligned}$$

Here, the inequality follows from Prop. 3.2 together with $\hat{p}_{\theta}(\mathbf{o}_{l+1}^{\text{hard}, i} | \mathbf{o}_l) \leq \hat{p}_{\theta}(\mathbf{o}_{l+1}^{\text{easy}, j} | \mathbf{o}_l)$. As a result, the ratio $\hat{p}_{\theta}(\mathbf{o}_{l+1}^{\text{hard}, i} | \mathbf{o}_l) / \hat{p}_{\theta}(\mathbf{o}_{l+1}^{\text{easy}, j} | \mathbf{o}_l)$ strictly decreases after each gradient update. From Eq.(8), RAFT's gradient further amplifies this gap through the $(1 + \log(p))$ factor, while PPO's update similarly magnifies it via the term $(1 + (2\mathbb{1}(A \geq 0) - 1)\epsilon)A$. By induction, the disparity between easy and hard CoTs compounds over iterations, and the convergence proof then follows directly.

Solution 1: Rejection of Easy Questions. Recent studies (Yu et al., 2025; Xiong et al., 2025; Zhang et al., 2025a) show that rejecting instances, where *all* parallelly sampled CoTs are correct, improves performance. In our setting, such instances correspond to those solvable by already well-learned easy CoTs. By discarding them and retaining only hard CoT correct-only instances, the model gradually shifts its focus toward harder reasoning paths. Formally, we define `RL-rej` as any algorithm in Thm. 3.1 augmented with rejection: whenever a sampled CoT has probability mass above $M^{-1}(1 - \epsilon)$ by the current model, it is discarded. This ensures training emphasizes harder CoTs gradually in the small learning rate regime, prevents collapse into easy ones, and in the end secure all valid CoTs with probability at least $M^{-1}(1 - \epsilon)$. We summarize this finding per below.

Corollary 3.3 (`RL-rej` Enables Hard-CoT Learning). *Under the identical setting and assumptions of Thm. 3.1, consider applying `RL-rej`. Then, for any $\epsilon \geq 0$, there exists $t \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ such that for any valid hard to reason CoT \mathbf{o}^{hard} for task $k \in \mathcal{T}$, we have*

$$\Pr(\mathbf{o}_{2:L}^{\text{hard}} \sim \hat{p}_{\theta^{(t)}}^k(\cdot | \mathbf{o}_1^{\text{hard}})) \geq \frac{1 - \epsilon}{M}.$$

Therefore, for any (\mathbf{Q}, \mathbf{A}) of task k with ≥ 1 correct CoTs, the finetuned model $\hat{p}_{\theta^{(t)}}$ satisfies

$$\mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^{(t)}}^k(\cdot | \mathbf{o}_1)} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})] \geq \frac{1 - \epsilon}{M}.$$

That is, with $K \geq \frac{M}{1 - \epsilon} (\log(\epsilon^{-1}))$, we have `pass@K` performance no worse than $1 - \epsilon$.

Notably, after sufficient iterations, the algorithms in Thm. 3.1 and Cor. 3.3 concentrate probability mass on valid CoTs of the targeted task up to $\Theta(1 - \epsilon)$ from its start state. Consequently, the generation probability of CoTs for other tasks sharing some state would be less than $o(\epsilon)$, eroding cross-task capability. In what follows, we discuss an alternative exploration approach, which in design can preserve such meta-capabilities.

Solution 2: KL-regularization. It is also worth noting that GRPO typically is equipped with a KL regularization term, as in Eq.(6). The formulation of KL-regularized Reinforcement Learning has been noticed as a distribution optimization (Fan et al., 2023; Black et al., 2024; Clark et al., 2024; Uehara et al., 2024; Marion et al., 2024; Kawata et al., 2025). In theory, the solution is a tilted (or Gibbs) distribution (Csiszár, 1975), as characterized below.

Lemma 3.4 (Optimal Sampling of GRPO Variants). *For each task $k \in \mathcal{T}$, let θ^* denote the pretrained Foundation Model. Then the GRPO induces an optimal step-wise sampling distribution:*

$$\hat{p}_{\theta^*}^{\text{PO}}(\mathbf{o}_{l+1} | \mathbf{o}_l) \propto \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \cdot \exp\left(\hat{r} \cdot \frac{A_{l+1}^{\hat{p}_{\theta^*}, k}(\mathbf{o}_{l+1})}{\beta}\right), \quad (9)$$

where $\hat{r} \leq \Theta(M)$ and $A_{l+1}^{\hat{p}_{\theta^*}, k}(\mathbf{o}_{l+1})$ is defined in Eq. (7).

Notably, the induced Gibbs distribution is governed by the KL-regularization temperature $\beta > 0$: a larger β reduces the gap between CoTs with high and low advantage. The following corollary formalizes this intuition, showing that $\hat{p}_{\theta^*}^{\text{PO}}$ can, in principle, preserve the broad capability.

Corollary 3.5 (KL-regularization Enables Hard-CoT learning and Maintain Cross-task Capability). *Consider a base model θ^* defined in Sec. 2.2, a targeted task $k \in \mathcal{T}$ and a different task $k' \neq k$, denote $\hat{p}_{\theta^*}^{\text{PO}}$ as the learner in Eq.(9). For any start state $\mathbf{o}_1 = q$ of task k , suggest the number of CoTs starting from \mathbf{o}_1 is $N_{\mathbf{o}_1}$. Then for any ϵ' satisfying $1/N_{\mathbf{o}_1} > \epsilon' \geq \epsilon > 0$, denote $\hat{p}_{\theta^{k, (t)}}^k$ as the PPO/GRPO in Thm. 3.1 with ϵ , then there exists $\beta = \Omega(ML/\log(\epsilon'^{-1}))$, such that*

1. **Capable of Hard CoTs:** For instance (\mathbf{Q}, \mathbf{A}) with only some hard-to-reason CoTs correct:

$$\begin{aligned} \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^*}^{\text{PO}}(\cdot | \mathbf{o}_1)} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})] &\geq \epsilon' \geq \epsilon \\ &\geq \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^{k, (t)}}^k(\cdot | \mathbf{o}_1)} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})]. \end{aligned}$$

2. **Preserve Multi-task:** For instance $(\mathcal{Q}, \mathbf{A})$ belonging to untargeted task $k' \neq k$:

$$\begin{aligned} \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^{\text{PO}}(\cdot | \mathbf{o}_1)} [R_{(\mathbf{Q}, \mathbf{A})}^{k'}(\mathbf{o})] &\geq \epsilon' \\ &\geq \epsilon \geq \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(\cdot | \mathbf{o}_1)} [R_{(\mathbf{Q}, \mathbf{A})}^{k'}(\mathbf{o})] \end{aligned}$$

The pass@K performance of any task could be adjusted by temperature β given K and ϵ' .

Other Solutions. We also provided discussions of the benefit of *Evolution Strategy (ES) finetuning* (Qiu et al., 2025) and *representation-based exploration finetuning* (Tuyls et al., 2025) in App. D.4.

4. Simplicity Bias of Population Reward Inference-Scaling: Challenge and Solution

ORM Mode. During inference, the *outcome reward model* (ORM) evaluates entire paths via an outcome-level reward $R_{\text{out}}^k(\mathbf{o})$ (e.g., a neural scorer), guiding solution generation through **Best-of-N (BoN) sampling** (Lightman et al., 2023). We define the natural ORM as $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}[R^k(\mathbf{o})]$ (i.e., the expectation of instance-level rewards). Statistically, R_{out} is the Bayes-optimal L^2 predictor (and the MLE under Gaussian noise). A neural scorer $R_{\theta}^k(\cdot)$ is then trained to approximate $R_{\text{out}}^k(\cdot)$ by $\arg \min_{\theta} \mathbb{E}[(R_{\theta}^k(\mathbf{o}) - R_{\text{out}}^k(\mathbf{o}))^2]$ which under standard conditions converges to R_{out}^k .

PRM Mode. Instead of outcome-level scoring, the *process reward model* (PRM) provides intermediate rewards along the reasoning trajectory: $R_{\text{pro}}^k(\mathbf{o}_l) = g(\mathbf{o}_1, \dots, \mathbf{o}_l)$, $l \in \{1, \dots, L\}$, where $g(\cdot)$ estimates step-wise utility (Shao et al., 2024; Snell et al., 2024; Wang et al., 2024; Li et al., 2023). PRM can be integrated into structured decoding, e.g., **BoN** (Lightman et al., 2023) (selecting top PRM-scoring step) or **Beam Search (BS)** (Snell et al., 2024) (augmenting beam scores). Since process-level annotations are costly, most approaches design PRMs heuristically via **likelihood-based estimates**, which predict the expected final correctness given the current prefix:

$$R_{\text{likelihood}}^k(\mathbf{o}_l) = V^{\hat{p}_{\theta^*}}(\mathbf{o}_l) = \mathbb{E}[R^k(\mathbf{o}) | \mathbf{o}_l], \quad (10)$$

for all $\mathbf{o}_l \in S_l$. The expectation, which is operated on $\mathbf{o}_1 \sim P^k(\mathcal{Q}^k)$, $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a,q}^{q,k}$, $\mathbf{o} \sim \hat{p}_{\theta^*}^k(O | \mathbf{o}_1)$, is typically approximated by Monte Carlo rollouts or by training a neural scorer R_{θ} with squared loss, i.e.

$$\arg \min_{\theta} \mathbb{E}[(R_{\theta}(\mathbf{o}_l) - R_{\text{likelihood}}^k(\mathbf{o}_l))^2].$$

Indeed, our following theorem shows that the above two ‘‘population rewards’’ (i.e., expectation-based ORM/PRM) check *consistency* instead of *correctness*, per **Phenomenon 2** (Xu et al., 2025).

Theorem 4.1 (Failure of Inference-Scaling with ORM/PRM). *Under the setting of Thm. 3.1, consider the ORM $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}[R^k(\mathbf{o})]$, the PRM $R_{\text{likelihood}}^k(\mathbf{o}_l)$ and inference methods: (i) ORM + BoN, (ii) PRM + BoN (step-wise), or (iii) PRM + BS with width N and beam size $B \geq 1$. For any instance (\mathbf{Q}, \mathbf{A}) of task (\mathbf{o}_1, a, k) , suppose all correct CoTs are hard-to-reason and their ≥ 1 sparse edges diverge from shared states with some valid easy-to-reason CoT. Then, for any $\epsilon > 0$:*

- If $N \geq \Omega(\log(\epsilon)/\log(\frac{M^L - M}{M^L}))$, method (i) fails with probability at least $1 - \epsilon$.
- If $N \geq \Omega(\log(\epsilon)/\log(\frac{M-1}{M}))$, methods (ii) and (iii) fail with probability at least $1 - \epsilon$.

Sketch of Proof. Our key observation is the following Prop. 4.2, which reveals that population rewards systematically favor easy CoTs, assigning higher scores to \mathbf{o}^{easy} than \mathbf{o}^{hard} .

Proposition 4.2 (Population Rewards Favor Easy CoTs). *Under the same settings as Thm. 4.1, let \mathbf{o}^{easy} be any valid easy-to-reason CoT and \mathbf{o}^{hard} any valid hard-to-reason CoT under (q, a, k) . Then for $\forall \mathbf{o}_{l-1}^{\text{easy}} = \mathbf{o}_{l-1}^{\text{hard}}, \mathbf{o}_l^{\text{easy}} \neq \mathbf{o}_l^{\text{hard}}$:*

$$\begin{aligned} R_{\text{out}}^k(\mathbf{o}^{\text{easy}}) &> R_{\text{out}}^k(\mathbf{o}^{\text{hard}}), \\ R_{\text{likelihood}}^k(\mathbf{o}_l^{\text{easy}}) &> R_{\text{likelihood}}^k(\mathbf{o}_l^{\text{hard}}). \end{aligned}$$

Sketch of Proof. The first inequality follows from Def. 2.2(iii): an easy-to-reason CoT has a larger probability of being correct over the distribution, whereas a hard-to-reason CoT, carries higher uncertainty and thus a smaller population-level chance of correctness. The second inequality follows from Prop. 3.2 by noting that

$$R_{\text{likelihood}}^k(\mathbf{o}_l) = A_{l+1}^{\hat{p}_{\theta^*}, k}(\mathbf{o}_{l-1}, \mathbf{o}_l) + V^{\hat{p}_{\theta^*}, k}(\mathbf{o}_{l-1}).$$

Given Prop. 4.2, the remaining proofs for Thm. 4.1 follow by choosing N sufficiently large so that \mathbf{o}^{easy} (or $\mathbf{o}_l^{\text{easy}}$) is sampled at least once across the N parallel trials.

Solution: Gibbs Sampling. Soft Best-of-N sampling (Soft-BoN) (Verdun et al., 2025) is designed to approximate the gibbs distribution $P_{\text{Gibbs}}^k(\mathbf{o})$ with $O(N^{-1})$ error, which is defined as

$$P_{\text{Gibbs}}^k(\mathbf{o}) \propto (\hat{p}_{\theta^*}(\mathbf{o}) \exp(\lambda R_{\text{out}}^k(\mathbf{o}))), \quad (11)$$

for $\mathbf{o}_1 = q \sim P^k(\mathcal{Q}_k)$. Akin to Eq.(9), the distribution $P_{\text{Gibbs}}^k(\mathbf{o})$ also can control the trade-off between reward maximization and the divergence from the base model’s predictive power.

Corollary 4.3. (Csiszár, 1975) *Consider a base model θ^* defined in Sec. 2.2, a targeted task $k \in \mathcal{T}$ and ORM $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}[R^k(\mathbf{o})]$. For $\lambda > 0$, Eq.(11) solves:*

$$\max_{P_{\text{new}}^k} \mathbb{E}_{P_{\text{new}}^k} [R_{\text{out}}^k(\mathbf{o})] - \frac{1}{\lambda} D_{\text{KL}}(P_{\text{new}}^k \| \hat{p}_{\theta^*}). \quad (12)$$

Table 1. Theoretical comparison between RLVR and inference-scaling under our TMC setting. The first column indicates pass@K performance. The second and third columns assess whether a method assigns highest credit to easy-to-reason CoTs and whether it can also sample hard-to-reason CoTs with suitable temperature. The fourth column evaluates whether the method preserves the base model’s multi-task capability. The results suggest that post-training methods tend to favor easy-to-reason CoTs, and that only methods capable of sampling hard-to-reason CoTs can achieve satisfactory pass@K.

Methods	Succeed w. pass@K	Prefer Easy CoT	Capable of Hard CoT	Preserve Multi-task
REINFORCE (Eq.(3)) / RAFT (Eq.(4))	✗	✓	✗	✗
PPO (Eq.(5))	✗	✓	✗	✗
RL-rej (Sec 3)	✓	✓	✓	✗
KL-regularized PO (Eq.(6))	✓	✓	✓	✓
ORM/PRM-BoN/BS (Sec 4)	✗	✓	✗	✗
Soft-BoN/DPRM-AS (Sec 4)	✓	✓	✓	✓

Indeed, through the statistical merit of Doob’s h-transform techniques (Uehara et al., 2024; Kawata et al., 2025; Rogers & Williams, 2000; Chopin et al., 2023; Heng et al., 2024), we provably show that there is a principled framework to design process reward, which could mathematically generate the same CoT distribution as Eq.(11) per below.

Definition 4.4. Doob’s h-Transform-induced Process Reward Model (DPRM). Consider a base model θ^* defined in Sec. 2.2, a targeted task $k \in \mathcal{T}$ and ORM $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}[R^k(\mathbf{o})]$. The DPRM Adjusted Sampling (DPRM-AS) defines the process reward at step l via harmonic function

$$h_k(\mathbf{o}_l) = \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o})) \mid \mathbf{o}_l],$$

and sample according to step-wise distribution adjustment:

$$\begin{aligned} R_{\text{DPRM}}^k(\mathbf{o}_l) &= \frac{1}{\lambda} \log h_k(\mathbf{o}_l), \quad \forall \lambda > 0 \\ \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) &= \hat{p}_{\theta^*}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) \cdot \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)} \\ &\propto \hat{p}_{\theta}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1})), \end{aligned} \quad (13)$$

where the first-step is initialized as

$$\hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_1 \mid \mathbf{o}_0) := \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}_k)} [\exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_1))].$$

Then, the induced distribution

$$P_{\text{DPRM}}^k(\mathbf{o}) := \prod_{l=0}^{L-1} \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} \mid \mathbf{o}_l)$$

satisfies $P_{\text{DPRM}}^k(\mathbf{o}) = P_{\text{Gibbs}}^k(\mathbf{o})$.

Computational Cost. By $P_{\text{DPRM}}^k(\mathbf{o}) = P_{\text{Gibbs}}^k(\mathbf{o})$ as shown above, the Soft-BoN method is a realization of the induced distribution, with a convergence rate $O(N^{-1})$ (Verdun et al., 2025). Notably, $R_{\text{DPRM}}^k(\cdot)$ and $R_{\text{likelihood}}^k(\cdot)$ are both conditioned expectations (i.e., $\log \mathbb{E}[\exp(\cdot) \mid \cdot]$ and $\mathbb{E}[\cdot \mid \cdot]$), and thus $R_{\text{DPRM}}^k(\cdot)$ does not induce computational overhead, despite extra but negligible evaluations of exp and log. Inherently, there is an asymptotic equivalence between them when using certain sampling strategies, which we formalized as below.

Corollary 4.5. Under the same settings as Def. 4.4, for $0 < \lambda < \infty$, it holds that **BoN/BS** with $R_{\text{DPRM}}^k(\mathbf{o}_l)$ is equivalent to **BoN/BS** with $R_{\text{likelihood}}^k(\mathbf{o}_l)$.

Sketch of Proof. The key observation is by the monotonicity of $\exp(\cdot)$ and $\log(\cdot)$, it holds that

$$\arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} R_{\text{DPRM}}^k(\mathbf{o}_l) = \arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} R_{\text{likelihood}}^k(\mathbf{o}_l),$$

where $S_l^{\text{BoN}} = \{\mathbf{o}_l^1, \dots, \mathbf{o}_l^N\}$ is the set of **BoN** candidates.

Through the similar techniques in Cor. 3.5, we then show that the broad capability is preserved by $P_{\text{DPRM}}^k(\mathbf{o}) = P_{\text{Gibbs}}^k(\mathbf{o})$ as below.

Corollary 4.6 (Gibbs Distribution Preserves Meta-Capability). Under the same settings in Cor. 3.5, for any ϵ' satisfying $1/N_{\mathbf{o}_1} > \epsilon' \geq \epsilon > 0$. Then there exists $\lambda = O(\log(\epsilon'^{-1})/ML)$, denote $\hat{p}_{\text{IS}}^k(\cdot)$ as any of the inference predictors (i)-(iii) in Cor. 3.5 with $N \geq \Omega(\log(\epsilon)/\log(\frac{M^L-M}{ML}))$, it holds that

- Capable of Hard CoTs:** $\mathbb{E}_{\mathbf{o} \sim P_{\text{Gibbs}}^k(\mathbf{o})} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})] \geq \epsilon' \geq \epsilon \geq \mathbb{E}_{\mathbf{o} \sim \hat{p}_{\text{IS}}^k(\mathbf{o})} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})]$.
- Preserve Multi-task:** $\mathbb{E}_{\mathbf{o} \sim P_{\text{Gibbs}}^k(\mathbf{o})} [R_{(\mathbf{Q}, \mathbf{A})}^{k'}(\mathbf{o})] \geq \epsilon' \geq \epsilon \geq \mathbb{E}_{\mathbf{o} \sim \hat{p}_{\text{IS}}^k(\mathbf{o})} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})]$.

The pass@K of Soft-BoN ($\frac{1}{N}$) error to gibbs sampling (Verdun et al., 2025) for any task could then be adjusted by temperature β given K and ϵ' (Wu et al., 2025b).

5. Empirical Simulations

To validate our theoretical findings, we run simulations on an abstract Tree-structured Markov Chain (TMC) with two tasks (**TASK 1** is the target), as shown in Tab. 2 below.

The TMC has $L = 4$ layers with two nodes each ($|S_l| = 2$). In layers 1–3, each state has one high-probability outgoing edge. Pretraining runs for $T_1 = 2000$ and $T_2 = 500$ steps (error < 0.001); fine-tuning for $T = 1000$ steps with learning rate 0.05. Estimation of the Rewards/advantages use

Table 2. Task definition and CoTs characteristics in our Multi-Task TMC simulation.

Task	Path Index	State Transition	Type	Probability	Expected Correctness over $\mathcal{D}_a^{q,k}$
TASK1	0	$S_1[0] \rightarrow S_2[0] \rightarrow S_3[0] \rightarrow S_4[0]$	EASY-To-REASON	0.413223	0.727995
	1	$S_1[0] \rightarrow S_2[0] \rightarrow S_3[1] \rightarrow S_4[0]$	HARD-To-REASON	0.075131	0.132363
	2	$S_1[0] \rightarrow S_2[1] \rightarrow S_3[0] \rightarrow S_4[0]$	HARD-To-REASON	0.004132	0.007280
	3	$S_1[0] \rightarrow S_2[1] \rightarrow S_3[1] \rightarrow S_4[0]$	HARD-To-REASON	0.075131	0.132363
TASK2	0	$S_1[1] \rightarrow S_2[0] \rightarrow S_3[0] \rightarrow S_4[1]$	EASY-To-REASON	0.413223	0.955691
	1	$S_1[1] \rightarrow S_2[0] \rightarrow S_3[1] \rightarrow S_4[1]$	HARD-To-REASON	0.007513	0.017376
	2	$S_1[1] \rightarrow S_2[1] \rightarrow S_3[0] \rightarrow S_4[1]$	HARD-To-REASON	0.004132	0.009557
	3	$S_1[1] \rightarrow S_2[1] \rightarrow S_3[1] \rightarrow S_4[1]$	HARD-To-REASON	0.007513	0.017376

Table 3. CoT generation statistics for different strategies on TASK1 and TASK2. Values represent percentages of valid easy CoTs, valid hard CoTs, and invalid CoTs generated by each method.

Strategy	TASK1 Valid Easy CoTs (%)	TASK1 Valid Hard CoTs (%)	TASK1 Invalid CoTs (%)	TASK2 Valid Easy CoTs (%)	TASK2 Valid Hard CoTs (%)	TASK2 Invalid CoTs (%)
Base Model	21.67%	8.07%	70.27%	20.03%	1.10%	78.87%
Finetuned Methods						
REINFORCE	94.33%	3.43%	2.23%	1.52%	0.87%	97.62%
RAFT	95.22%	2.33%	2.45%	2.30%	0.92%	96.78%
PPO (Eq.11)	91.82%	5.40%	2.78%	2.23%	1.03%	96.73%
RL-rej (Sec.3.1)	49.62%	17.42%	32.97%	30.63%	2.27%	67.10%
GRPO-KL (Eq.13)	46.47%	16.27%	37.27%	54.18%	1.68%	44.13%
Inference Scaling Methods						
Soft-BoN	8.98%	19.30%	71.72%	7.00%	17.27%	75.73%
ORM-BoN w. $R_{\text{out}}^k(\cdot)$	21.00%	7.30%	71.70%	20.23%	0.97%	78.80%
PRM-BoN w. $R_{\text{likelihood}}^k(\cdot)$	99.13%	0.87%	0.00%	13.42%	36.77%	49.82%
DPRM-BoN	99.52%	0.48%	0.00%	13.40%	37.02%	49.58%
DPRM-AS (implemented by step-wise Soft-BoN)	17.23%	36.10%	46.67%	12.02%	38.02%	49.97%

1000/200 Monte Carlo samples; temperature $\lambda = 0.5$; BoN uses $N = 15$. Training and testing each use 200 question instances sampled per Def. 2.2; BoN and Gibbs-style methods are fully enumerated. **TASK 1** requires reaching $S_4[0]$ from $S_1[0]$; **TASK 2** requires $S_4[1]$ from $S_1[1]$. A CoT is valid here if it connects the start and end states; among valid paths, only the one via $S_2[0] \rightarrow S_3[0]$ is *easy*, all others are *hard*. We report the proportions of easy, hard, and invalid CoTs from $S_1[0]$ (TASK 1) and $S_1[1]$ (TASK 2), as well as the expected correctness over the population per Def. 2.2.

Findings in Tab. 3. REINFORCE, RAFT, and PPO heavily favor easy-to-reason CoTs in TASK 1, suppressing hard-to-reason CoTs in TASK 1 and valid CoTs in TASK 2, showing clear *simplicity bias* and *forgetting*. In contrast, diversity-promoting methods (RL-rej, GRPO-KL, Soft-BoN, DPRM-AS) balance easy/hard-to-reason CoTs in TASK 1 and preserve TASK 2 CoT’s generation capability, thanks to shared sparse edges in the TMC (two nodes per layer, see Table 2). ORM/PRM-BoN, relying on population rewards $R_{\text{out}}^k(\cdot)$, also overfavor easy-to-reason CoTs; PRM-BoN and DPRM-BoN behave similarly, as predicted. Further details are available in App. C.

6. Conclusion, Limitations, and Future Work

We introduced a Tree-structured Markov toy framework to model foundation model’s diverse multi-task reasoning pat-

terns, and theoretically validates **Phenomenon 1-3**: both RLVR and inference-scaling exhibit a *simplicity bias*, favoring easier, common reasoning paths (consistency) rather than true correctness. Building on this, we demonstrated the benefit of various exploration strategies—mitigating this bias and preserving rare but crucial CoTs—as summarized in Tab. 1. Our analysis further highlights a sharp contrast with traditional RL (e.g., AlphaGo (Silver et al., 2016)): whereas RL advantages promote effective state-space exploration in standard RL, in post-training they instead push models to overemphasize easy (high-pass-rate) paths within the model’s scope (Yue et al., 2025). This negative insight may also explain why Setlur et al. (2025a) employ independent models that reinterpret RL advantage differently for finetuning and PRM scoring.

The central clue of the bias lies in the expectation (population)-based reward estimators, namely $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}[R^k(\mathbf{o})]$ and $R_{\text{likelihood}}^k(\mathbf{o}_l) = V^{\hat{p}\theta^*}(\mathbf{o}_l) = \mathbb{E}[R^k(\mathbf{o}) | \mathbf{o}_l]$. While these estimators are Bayes-optimal in the L^2 sense, they inherently favor frequent patterns, thereby downweighting rare-but-valuable CoTs. This bias highlights the necessity of more reliable reward designs, as also discovered by Xu et al. (2025). Also, our current TMC framework is deliberately abstract and restrictive (see App. B and D.4), and could be generalized to more realistic models. Another promising direction is to apply TMC analysis to reflective behavior and *aha* moments (Yu et al., 2025).

References

- 440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
- Abbe, E., Adsera, E. B., and Misiakiewicz, T. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 2552–2623. PMLR, 2023.
- Abbe, E., Bengio, S., Lotfi, A., Sandon, C., and Saremi, O. How far can transformers reason? The locality barrier and inductive scratchpad. In *Advances in Neural Information Processing Systems*, 2024.
- AI, E., ; Shah, D. J., Rushton, P., Singla, S., Parmar, M., Smith, K., Vanjani, Y., Vaswani, A., Chaluvaraju, A., Hojel, A., Ma, A., Thomas, A., Polloreno, A., Tanwer, A., Sibai, B. D., Mansingka, D. S., Shivaprasad, D., Shah, I., Stratos, K., Nguyen, K., Callahan, M., Pust, M., Iyer, M., Monk, P., Mazarakis, P., Kapila, R., Srivastava, S., and Romanski, T. Rethinking reflection in pre-training. *arXiv preprint arXiv:2504.04022*, 2025.
- Allen-Zhu, Z. and Li, Y. Backward feature correction: How deep learning performs deep (hierarchical) learning. *arXiv preprint arXiv:2001.04413*, 2023.
- Arpit, D., Jastrzkebski, S., Ballas, N., Krueger, D., Bengio, E., Kanwal, M. S., Maharaj, T., Fischer, A., Courville, A., Bengio, Y., et al. A closer look at memorization in deep networks. In *International conference on machine learning*, pp. 233–242. PMLR, 2017.
- Bae, S., Hong, J., Lee, M. Y., Kim, H., Nam, J., and Kwak, D. Online difficulty filtering for reasoning oriented reinforcement learning. *arXiv preprint arXiv:2504.03380*, 2025.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: solving elaborate problems with large language models. In *Proceedings of the AAI Conference on Artificial Intelligence*, 2024.
- Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2024.
- Bu, D., Huang, W., Han, A., Nitanda, A., Wong, H.-S., Zhang, Q., and Suzuki, T. Provable benefit of curriculum in transformer tree-reasoning post-training. *arXiv preprint arXiv:2511.07372*, 2025.
- Bu, D., Huang, W., Han, A., Wong, H.-S., Zhang, Q., Suzuki, T., and Nitanda, A. Dprm: A plug-in doob h transform-induced token-ordering module for diffusion language models. *arXiv preprint arXiv:2604.24357*, 2026.
- Cai, X.-Q., Wang, W., Liu, F., Liu, T., Niu, G., and Sugiyama, M. Reinforcement learning with verifiable yet noisy rewards under imperfect verifiers. *arXiv preprint arXiv:2510.00915*, 2025.
- Cao, Y., He, Y., Wu, D., Chen, H.-Y., Fan, J., and Liu, H. Transformers simulate mle for sequence generation in bayesian networks. *arXiv preprint arXiv:2501.02547*, 2025.
- Chen, A., Shwartz-Ziv, R., Cho, K., Leavitt, M. L., and Saphra, N. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms. *arXiv preprint arXiv:2309.07311*, 2023a.
- Chen, X., Li, T., and Zou, D. On the mechanism of reasoning pattern selection in reinforcement learning for language models. *arXiv preprint arXiv:2506.04695*, 2025.
- Chen, Y., Yuille, A., and Zhou, Z. Which layer is learning faster? a systematic exploration of layer-wise convergence rate for deep neural networks. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Chopin, N., Fulop, A., Heng, J., and Thiery, A. H. Computational doob h-transforms for online filtering of discretely observed diffusions. In *International Conference on Machine Learning*, pp. 5904–5923. PMLR, 2023.
- Clark, K., Vicol, P., Swersky, K., and Fleet, D. J. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2024.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Csiszár, I. I-divergence geometry of probability distributions and minimization problems. *The annals of probability*, pp. 146–158, 1975.
- Cui, G., Zhang, Y., Chen, J., Yuan, L., Wang, Z., Zuo, Y., Li, H., Fan, Y., Chen, H., Chen, W., Liu, Z., Peng, H., Bai, L., Ouyang, W., Cheng, Y., Zhou, B., and Ding, N. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.
- Deng, W., Ren, Y., Li, M., Sutherland, D. J., Li, X., and Thramopoulos, C. On the effect of negative gradient in group relative deep reinforcement optimization. *arXiv preprint arXiv:2505.18830*, 2025.
- Ding, Y., Lu, S., Lu, Y., Nowicki, T. J., and Gao, J. Global optimality of in-context markovian dynamics learning. <https://openreview.net/forum?id=HuBFimORiz>, 2025.

- 495 Dong, Y., Baker, B., Banino, A., Rae, J., Weber, T., and
 496 Nematzadeh, A. Raft: Leveraging ranking for fine-tuning
 497 language models. *arXiv preprint arXiv:2310.01377*,
 498 2023.
- 499
 500 Edelman, E., Tsilivis, N., Edelman, B. L., Eran Malach,
 501 and Goel, S. The evolution of statistical induction heads:
 502 in-context learning Markov chains. *Advances in Neural
 503 Information Processing Systems*, 2024.
- 504
 505 Fan, M., Han, W., Wang, D., Chen, C., Zhang, Z., and Zhou,
 506 J. When sharpening becomes collapse: Sampling bias
 507 and semantic coupling in rl with verifiable rewards. *arXiv
 508 preprint arXiv:2601.15609*, 2026.
- 509
 510 Fan, Y., Watkins, O., Du, Y., Liu, H., Ryu, M., Boutilier, C.,
 511 Abbeel, P., Ghavamzadeh, M., Lee, K., and Lee, K. Dpoc:
 512 Reinforcement learning for fine-tuning text-to-image dif-
 513 fusion models. In *Advances in Neural Information Pro-
 514 cessing Systems*, pp. 79858–79885, 2023.
- 515
 516 Foster, D. J., Mhammedi, Z., and Rohatgi, D. Is a good
 517 foundation necessary for efficient reinforcement learning?
 518 the computational role of the base model in exploration.
 519 *arXiv preprint arXiv:2503.07453*, 2025.
- 520
 521 Gai, J., Zeng, G., Zhang, H., and Raghunathan, A. Differ-
 522 ential smoothing mitigates sharpening and improves llm
 523 reasoning. *arXiv preprint arXiv:2511.19942*, 2025.
- 524
 525 Gandhi, K., Chakravarthy, A., Singh, A., Lile, N., and
 526 Goodman, N. D. Cognitive behaviors that enable self-
 527 improving reasoners, or, four habits of highly effective
 528 stars. *arXiv preprint arXiv:2503.01307*, 2025.
- 529
 530 Golowich, N., Chen, F., Rohatgi, D., Singhal, R., Domingo-
 531 Enrich, C., Foster, D. J., and Krishnamurthy, A. Reject,
 532 resample, repeat: Understanding parallel reasoning in lan-
 533 guage model inference. *arXiv preprint arXiv:2603.07887*,
 534 2026.
- 535
 536 Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R.,
 537 Zhu, Q., Ma, S., Wang, P., Bi, X., et al. DeepSeek-R1:
 538 incentivizing reasoning capability in LLMs via reinfor-
 539 cement learning. *arXiv preprint arXiv:2501.12948*, 2025a.
- 540
 541 Guo, Y., Wu, Y., Zhang, X., et al. Right is not enough: The
 542 pitfalls of outcome supervision in training llms for math
 543 reasoning. *arXiv preprint arXiv:2506.06877*, 2025b.
- 544
 545 Hazan, E. Introduction to online convex optimization. *arXiv
 546 preprint arXiv:1909.05207*, 2023.
- 547
 548 He, A., Fried, D., and Welleck, S. Rewarding the un-
 549 likely: Lifting grpo beyond distribution sharpening. *arXiv
 preprint arXiv:2506.02355*, 2025.
- Heng, J., De Bortoli, V., and Doucet, A. Diffusion
 schrödinger bridges for bayesian computation. *Statistical
 Science*, 39(1):90–99, 2024.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang,
 S., Wang, L., and Chen, W. Lora: Low-rank adaptation of
 large language models. *arXiv preprint arXiv:2106.09685*,
 2021.
- Ildiz, M. E., Huang, Y., Li, Y., Rawat, A. S., and Oymak,
 S. From self-attention to Markov models: unveiling the
 dynamics of generative transformers. In *International
 Conference on Machine Learning*, 2024.
- Ji, Z. and Telgarsky, M. Risk and parameter convergence
 of logistic regression. *arXiv preprint arXiv:1803.07300*,
 2019.
- Kalimeris, D., Kaplun, G., Nakkiran, P., Edelman, B., Yang,
 T., Barak, B., and Zhang, H. Sgd on neural networks
 learns functions of increasing complexity. In *Advances
 in Neural Information Processing Systems*, volume 32,
 2019.
- Kawata, R., Oko, K., Nitanda, A., and Suzuki, T. Direct
 distributional optimization for provable alignment of dif-
 fusion models. In *The Thirteenth International Confer-
 ence on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Nvw2szDdmI>.
- Kim, J., Wu, D., Lee, J., and Suzuki, T. Metastable dy-
 namics of chain-of-thought reasoning: Provable ben-
 efits of search, rl and distillation. *arXiv preprint
 arXiv:2502.01694*, 2025.
- Levine, S. Reinforcement learning and control as probabilis-
 tic inference: Tutorial. *arXiv preprint arXiv:1805.00909*,
 2018.
- Li, B., Wang, Y., Ding, Y., Lochab, A., Grama, A., and
 Zhang, R. Addressing performance saturation for llm
 rl via precise entropy curve control. *arXiv preprint
 arXiv:2604.26326*, 2026a.
- Li, H., He, Z., Tian, S., Wen, J., and Li, A. Martingale
 foresight sampling: A principled approach to inference-
 time llm decoding. *arXiv preprint arXiv:2601.15482*,
 2026b.
- Li, J. and Ng, H. T. The hallucination dilemma: Factuality-
 aware reinforcement learning for large reasoning models.
arXiv preprint arXiv:2505.24630, 2025.
- Li, M. Z., Agrawal, K. K., Ghosh, A., Teru, K. K., Lajoie, G.,
 and Richards, B. A. Tracing the representation geometry
 of language models from pretraining to post-training. In
High-dimensional Learning Dynamics 2025, 2025a.

- 550 Li, Y., Lin, Z., Zhang, S., Fu, Q., Chen, B., Lou, J.-G., and
 551 Chen, W. Making language models better reasoners with
 552 step-aware verifier. In *Proceedings of the 61st Annual*
 553 *Meeting of the Association for Computational Linguistics*
 554 *(Volume 1: Long Papers)*, pp. 5315–5333, 2023.
- 555 Li, Z., Chen, C., Xu, T., Qin, Z., Xiao, J., Luo, Z.-Q., and
 556 Sun, R. Preserving diversity in supervised fine-tuning of
 557 large language models. In *The Thirteenth International*
 558 *Conference on Learning Representations*, 2025b. URL
 559 [https://openreview.net/forum?id=NQEE](https://openreview.net/forum?id=NQEE7B7bSw)
 560 [7B7bSw](https://openreview.net/forum?id=NQEE7B7bSw).
- 561 Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker,
 562 B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and
 563 Cobbe, K. Let’s verify step by step. *arXiv preprint*
 564 *arXiv:2305.20050*, 2023.
- 565 Ling, W., Yogatama, D., Dyer, C., and Blunsom, P. Pro-
 566 gram induction by rationale generation: Learning to solve
 567 and explain algebraic word problems. *arXiv preprint*
 568 *arXiv:1705.04146*, 2017.
- 569 Makuva, A. V., Bondaschi, M., Girish, A., Nagle, A., Jaggi,
 570 M., Kim, H., and Gastpar, M. Attention with Markov:
 571 A framework for principled analysis of transformers via
 572 Markov chains. *arXiv preprint arXiv:2402.04161*, 2024.
- 573 Marion, P., Korba, A., Bartlett, P., Blondel, M., Bor-
 574 toli, V. D., Doucet, A., Llinares-López, F., Paquette,
 575 C., and Berthet, Q. Implicit diffusion: Efficient opti-
 576 mization through stochastic sampling. *arXiv preprint*
 577 *arXiv:arxiv.org/abs/2402.05468*, 2024.
- 578 Molina, J., Petrache, M., Costabal, F. S., and Courdurier,
 579 M. Understanding the dynamics of the frequency bias in
 580 neural networks. *arXiv preprint arXiv:2405.14957*, 2024.
- 581 Nichani, E., Damian, A., and Lee, J. D. How transform-
 582 ers learn causal structure with gradient descent. *arXiv*
 583 *preprint arXiv:2402.14735*, 2024.
- 584 OpenAI. Spinning up: proximal policy optimization (PPO),
 585 2018. URL [https://spinningup.openai.com/](https://spinningup.openai.com/en/latest/algorithms/ppo.html)
 586 [en/latest/algorithms/ppo.html](https://spinningup.openai.com/en/latest/algorithms/ppo.html). Accessed:
 587 2025-01-26.
- 588 Parashar, S., Gui, S., Li, X., Ling, H., Vemuri, S., Olson, B.,
 589 Li, E., Zhang, Y., Caverlee, J., Kalathil, D., et al. Curricu-
 590 lum reinforcement learning from easy to hard tasks im-
 591 proves llm reasoning. *arXiv preprint arXiv:2506.06632*,
 592 2025.
- 593 Qiu, X., Gan, Y., Hayes, C. F., Liang, Q., Meyerson, E.,
 594 Hodjat, B., and Miiikkulainen, R. Evolution strategies
 595 at scale: Llm fine-tuning beyond reinforcement learning.
 596 *arXiv preprint arXiv:2509.24372*, 2025.
- 597 Rafailov, R., Hashimoto, T., Zhang, C., Xiao, Y., Li, X.,
 598 Madaan, A., Leng, Y., He, Y., Zhou, Y., Singh, A.,
 599 et al. Direct preference optimization: Your language
 600 model is secretly a reward model. In *International*
 601 *Conference on Learning Representations (ICLR)*, 2024.
 602 arXiv:2305.18290.
- 603 Rajaraman, N., Jiao, J., and Ramchandran, K. An anal-
 604 ysis of tokenization: Transformers under markov data.
 In *Advances in Neural Information Processing Systems*,
 volume 37, pp. 62503–62556. Curran Associates, Inc.,
 2024.
- Ren, Y. and Sutherland, D. J. Learning dynamics of llm
 finetuning. *arXiv preprint arXiv:2407.10490*, 2025.
- Rogers, L. C. G. and Williams, D. *Diffusions, Markov*
processes, and martingales: Itô calculus, volume 2. Cam-
 bridge university press, 2000.
- Sanford, C., Fatemi, B., Hall, E., Tsitsulin, A., Kazemi, M.,
 Halcrow, J., Perozzi, B., and Mirrokni, V. Understanding
 transformer reasoning capabilities via graph algorithms.
arXiv preprint arXiv:2405.18512, 2024.
- Schmied, T., Bornschein, J., Grau-Moya, J., Wulfmeier, M.,
 and Pascanu, R. Llm’s are greedy agents: Effects of rl
 fine-tuning on decision-making abilities. *arXiv preprint*
arXiv:2504.16078, 2025.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
 Klimov, O. Proximal policy optimization algorithms.
arXiv preprint arXiv:1707.06347, 2017.
- Setlur, A., Garg, S., Geng, X., Garg, N., Smith, V., and
 Kumar, A. Rl on incorrect synthetic data scales the effi-
 ciency of llm math reasoning by eight-fold. *arXiv preprint*
arXiv:2406.14532, 2024.
- Setlur, A., Nagpal, C., Fisch, A., Geng, X., Eisenstein, J.,
 Agarwal, R., Agarwal, A., Berant, J., and Kumar, A. Re-
 warding progress: Scaling automated process verifiers for
 LLM reasoning. In *The Thirteenth International Confer-*
ence on Learning Representations, 2025a. URL <https://openreview.net/forum?id=A6Y7AqlzLW>.
- Setlur, A., Rajaraman, N., Levine, S., and Kumar, A. Scaling
 test-time compute without verification or rl is suboptimal.
arXiv preprint arXiv:2502.12118, 2025b.
- Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Netra-
 palli, P. The pitfalls of simplicity bias in neural networks.
 In *Advances in Neural Information Processing Systems*,
 volume 33, pp. 9573–9585, 2020.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
 H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Push-
 ing the limits of mathematical reasoning in open language
 models. *arXiv preprint arXiv:2402.03300*, 2024.

- 605 Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L.,
 606 Van Den Driessche, G., Schrittwieser, J., Antonoglou, I.,
 607 Panneershelvam, V., Lanctot, M., et al. Mastering the
 608 game of go with deep neural networks and tree search.
 609 *Nature*, 529(7587):484–489, 2016.
- 610 Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-
 611 time compute optimally can be more effective than scal-
 612 ing model parameters. *arXiv preprint arXiv:2408.03314*,
 613 2024.
- 614 Sun, Y., Yin, Z., Huang, X., Qiu, X., and Zhao, H. Error
 615 classification of large language models on math word
 616 problems: A dynamically adaptive framework. *arXiv*
 617 *preprint arXiv:2501.15581*, 2025.
- 618 Tian, Y. Composing global optimizers to reasoning tasks
 619 via algebraic objects in neural nets. *arXiv preprint*
 620 *arXiv:2410.01779*, 2024.
- 621 Tong, Y., Zhang, X., Wang, R., Wu, R., and He, J. Dart-
 622 math: Difficulty-aware rejection tuning for mathemat-
 623 ical problem-solving. *arXiv preprint arXiv:2407.13690*,
 624 2024.
- 625 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,
 626 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E.,
 627 Azhar, F., et al. Llama: Open and efficient foundation lan-
 628 guage models. *arXiv preprint arXiv:2302.13971*, 2023.
- 629 Tuyls, J., Foster, D. J., Krishnamurthy, A., and Ash,
 630 J. T. Representation-based exploration for language
 631 models: From test-time to post-training. *arXiv preprint*
 632 *arXiv:2510.11686*, 2025.
- 633 Uehara, M., Zhao, Y., Black, K., Hajiramezanali, E., Scalia,
 634 G., Diamant, N. L., Tseng, A. M., Biancalani, T., and
 635 Levine, S. Fine-tuning of continuous-time diffusion
 636 models as entropy-regularized control. *arXiv preprint*
 637 *arXiv:2402.15194*, 2024.
- 638 Valle-Perez, G., Camargo, C. Q., and Louis, A. A. Deep
 639 learning generalizes because the parameter-function map
 640 is biased towards simple functions. *arXiv preprint*
 641 *arXiv:1805.08522*, 2018.
- 642 Verdun, C. M., Oesterling, A., Lakkaraju, H., and Calmon,
 643 F. P. Soft best-of-n sampling for model alignment. *arXiv*
 644 *preprint arXiv:2505.03156*, 2025.
- 645 Wang, P., Li, L., Shao, Z., Xu, R. X., Dai, D., Li, Y., Chen,
 646 D., Wu, Y., and Sui, Z. Math-shepherd: Verify and
 647 reinforce llms step-by-step without human annotations.
 648 *arXiv preprint arXiv:2312.08935*, 2024.
- 649 Wang, S., Yu, L., Gao, C., Zheng, C., Liu, S., Lu, R., Dang,
 650 K., Chen, X., Yang, J., Zhang, Z., Liu, Y., Yang, A.,
 651 Zhao, A., Yue, Y., Song, S., Yu, B., Huang, G., and Lin,
 652 J. Beyond the 80/20 rule: High-entropy minority tokens
 653 drive effective reinforcement learning for llm reasoning.
 654 *arXiv preprint arXiv:2506.01939*, 2025a.
- 655 Wang, Y., Yang, Q., Zeng, Z., Ren, L., Liu, L., Peng, B.,
 656 Cheng, H., He, X., Wang, K., Gao, J., Chen, W., Wang,
 657 S., Du, S. S., and Shen, Y. Reinforcement learning for
 658 reasoning in large language models with one training
 659 example. *arXiv preprint arXiv:2504.20571*, 2025b.
- 660 Wen, X., Liu, Z., Zheng, S., Xu, Z., Ye, S., Wu, Z., Liang,
 661 X., Wang, Y., Li, J., Miao, Z., Bian, J., and Yang, M.
 662 Reinforcement learning with verifiable rewards implicitly
 663 incentivizes correct reasoning in base llms. *arXiv preprint*
 664 *arXiv:2506.14245*, 2025.
- 665 Williams, R. J. Simple statistical gradient-following algo-
 666 rithms for connectionist reinforcement learning. *Machine*
 667 *learning*, 8(3-4):229–256, 1992.
- 668 Wu, F., Xuan, W., Lu, X., Liu, M., Dong, Y., Harchaoui, Z.,
 669 and Choi, Y. The invisible leash: Why rlvr may or may
 670 not escape its origin. *arXiv preprint arXiv:2507.14843*,
 671 2025a.
- 672 Wu, Y., Mirhoseini, A., and Tambe, T. On the role of
 673 temperature sampling in test-time scaling. *arXiv preprint*
 674 *arXiv:2510.02611*, 2025b.
- 675 Xin, H., Ren, Z., Song, J., Shao, Z., Zhao, W., Wang, H.,
 676 Liu, B., Zhang, L., Lu, X., Du, Q., et al. DeepSeek-
 677 Prover-v1.5: harnessing proof assistant feedback for rein-
 678 forcement learning and Monte-Carlo tree search. *arXiv*
 679 *preprint arXiv:2408.08152*, 2024.
- 680 Xiong, W., Yao, J., Xu, Y., Pang, B., Wang, L., Sahoo,
 681 D., Li, J., Jiang, N., Zhang, T., Xiong, C., and Dong, H.
 682 A minimalist approach to llm reasoning: from rejection
 683 sampling to reinforce. *arXiv preprint arXiv:2504.11343*,
 684 2025.
- 685 Xu, K., Li, J., Zhang, M., Du, S. S., Kawarabayashi, K.-i.,
 686 and Jegelka, S. What can neural networks reason about?
 687 *arXiv preprint arXiv:1905.13211*, 2019.
- 688 Xu, Y., Dong, H., Wang, L., Xiong, C., and Li, J. Reward
 689 models identify consistency, not causality. *arXiv preprint*
 690 *arXiv:2502.14619*, 2025.
- 691 Yang, K., Swope, A., Gu, A., Chalamala, R., Song, P., Yu,
 692 S., Godil, S., Prenger, R., and Anandkumar, A. Lean-
 693 Dojo: Theorem proving with retrieval-augmented lan-
 694 guage models. In *Neural Information Processing Systems*
 695 (*NeurIPS*), 2023.
- 696 Yang, Y., Gan, E., Dziugaite, G. K., and Mirzasoleiman,
 697 B. Identifying spurious biases early in training through
 698 the lens of simplicity bias. In *International Conference*

- 660 *on Artificial Intelligence and Statistics*, pp. 2953–2961.
 661 PMLR, 2024.
- 662 Yao, J., Wang, R., and Zhang, T. PRL: Process reward learn-
 663 ing improves llms’ reasoning ability and broadens the
 664 reasoning boundary. *arXiv preprint arXiv:2601.10201*,
 665 2026.
- 667 Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y., Fan,
 668 T., Liu, G., Liu, L., Liu, X., Lin, H., Lin, Z., Ma, B.,
 669 Sheng, G., Tong, Y., Zhang, C., Zhang, M., Zhang, W.,
 670 Zhu, H., Zhu, J., Chen, J., Chen, J., Wang, C., Yu, H.,
 671 Dai, W., Song, Y., Wei, X., Zhou, H., Liu, J., Ma, W.-Y.,
 672 Zhang, Y.-Q., Yan, L., Qiao, M., Wu, Y., and Wang, M.
 673 Dapo: An open-source llm reinforcement learning system
 674 at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- 676 Yuan, Z., Yuan, H., Li, C., Dong, G., Lu, K., Tan, C., Zhou,
 677 C., and Zhou, J. Scaling relationship on learning math-
 678 ematical reasoning with large language models. *arXiv*
 679 *preprint arXiv:2308.01825*, 2023.
- 680 Yue, Y., Chen, Z., Lu, R., Zhao, A., Wang, Z., Yue, Y., Song,
 681 S., and Huang, G. Does reinforcement learning really
 682 incentivize reasoning capacity in llms beyond the base
 683 model? *arXiv preprint arXiv:arxiv.org/abs/2504.13837*,
 684 2025.
- 686 Zekri, O., Odonnat, A., Benechehab, A., Bleistein, L.,
 687 Boullé, N., and Redko, I. Large language models as
 688 Markov chains. *arXiv preprint arXiv:2410.02724*, 2024.
- 690 Zhang, R., Arora, D., Mei, S., and Zanette, A. Speed-rl:
 691 Faster training of reasoning models via online curriculum
 692 learning. *arXiv preprint arXiv:2506.09016*, 2025a.
- 693 Zhang, S., Wang, Y., Liu, Y., Liu, T., Grabowski, P., Ie,
 694 E., Wang, Z., and Li, Y. Beyond markovian: Reflective
 695 exploration via bayes-adaptive rl for llm reasoning. *arXiv*
 696 *preprint arXiv:2505.20561*, 2025b.
- 698 Zhao, W., Zhu, Q., Zeng, X., Mi, F., Shang, L., and Fung, Y.
 699 R. M. Entropy centroids as intrinsic rewards for test-time
 700 scaling. *arXiv preprint arXiv:2604.26173*, 2026a.
- 702 Zhao, Z., Land, S., Bikel, D. M., and Alshikh, W. Short-
 703 hand for thought: Compressing llm reasoning via entropy-
 704 guided supertokens. *arXiv preprint arXiv:2604.26355*,
 705 2026b.
- 706 Ziebart, B. D. *Maximum Entropy Inverse Reinforcement*
 707 *Learning*. PhD thesis, Carnegie Mellon University, Pitts-
 708 burgh, PA, 2008.
- 710 Zuo, B. and Zhu, Y. Strategic scaling of test-time com-
 711 pute: A bandit learning approach. *arXiv preprint*
 712 *arXiv:2506.12721*, 2025.
- 713
 714

A. Additional Related Work

LLMs as Markov Processes. A growing body of work has drawn connections between large language models (LLMs) and Markovian dynamics. Zekri et al. (2024) established a theoretical equivalence between next-token prediction in LLMs and finite-state Markov chains, deriving scaling laws for in-context learning when prompted with such chains. Nichani et al. (2024) demonstrated that disentangled transformers are capable of learning Markov chains in context. Ildiz et al. (2024) studied how a single self-attention layer can simulate context-conditioned Markov chains, while Ding et al. (2025) showed that multi-layer transformers can approximate preconditioned gradient descent over Markovian distributions. Edelman et al. (2024) analyzed the distinct phases of training as transformers learn Markov chains, and Makkuva et al. (2024) investigated the function landscape of single-layer transformers on Markovian data, revealing challenges in learning higher-order chains. Rajaraman et al. (2024) proved that constant-depth transformers can learn k -order Markov processes when the next-token distribution depends on the previous k tokens. Furthermore, Cao et al. (2025) showed that transformers can simulate the maximum likelihood estimation (MLE) algorithm for learning Bayesian networks, which subsume Markov chains as a special case. Despite these advances, most prior works focus on modeling sequential variable dependencies, without abstracting the structure to chain-of-thought (CoT) reasoning. The most relevant exception is the recent work of Kim et al. (2025), which investigates CoT processes under metastable Markov chain assumptions. They show the necessity of search, RL-based finetuning, and distillation to navigate sparse transition spaces, also under a softmax modeling assumption. Their proposed algorithm is tailored specifically for such metastable settings. Instead, motivated by real-world multi-task, tree-structured reasoning tasks with binary (0-1) rewards, our work aims to theoretically compare the intrinsic biases of RL-based finetuning and inference sampling, and to connect these with recent discussions on the squeezing effect, the benefits of reasoning diversity, and the inherent limitations of RL-based fine-tuning.

Spectral Bias. The study of spectral bias in deep learning is extensive, with many works showing that neural networks tend to learn low-frequency or simple patterns with high signal-to-noise ratio first (Arpit et al., 2017; Valle-Perez et al., 2018; Kalimeris et al., 2019; Chen et al., 2023a; Abbe et al., 2023; Molina et al., 2024). Edelman et al. (2024) demonstrated that this simplicity bias during training can delay convergence to the correct solution in Markov chain learning. Chen et al. (2023b) observed that shallow layers in neural networks prioritize fitting lower-order functions, while Allen-Zhu & Li (2023) showed that this tendency in shallow networks can lead to drastically increased sample complexity due to their bias toward low-order polynomials. Tian (2024) examined simplicity bias from the perspective of algebraic structure learning. Other works have highlighted potential downsides: Shah et al. (2020); Yang et al. (2024) showed that such biases can be detrimental, causing models to overlook important features or be misled by spurious correlations. Recent work Ren & Sutherland (2025) identified the *squeezing effect* of Direct Preference Optimization: probability mass becomes increasingly concentrated on the output that was most confident prior to the update. A following-up work Deng et al. (2025) identified similar phenomenon of GRPO. Separately, Li et al. (2025b) analyzes the nature of the cross-entropy loss, showing that it systematically shifts probability mass from non-target tokens to target tokens—regardless of the quality of the non-target options—ultimately leading to distribution collapse during finetuning. However, prior work has not systematically characterized how this squeezing effect influences fine-tuning dynamics. In our study, under the Tree-structured Markov Chain (TMC) framework and a linear softmax model, we show that binary outcome rewards can potentially amplify this effect, favoring simple reasoning paths during fine-tuning and contributing to the model’s inductive bias.

Distribution Sharpening, Entropy Structure, and Diversity Collapse. A rapidly growing line of work studies whether RLVR expands reasoning capabilities or mainly sharpens the base model’s existing reasoning distribution. Empirical studies suggest that RLVR often improves finite-sample precision while shrinking empirical support, thereby missing underrepresented correct answers (Wu et al., 2025a). In formal theorem proving, He et al. (2025) identify a rank bias in GRPO: high-probability correct trajectories are preferentially reinforced whereas rare correct trajectories are neglected. Concurrent works further formalize diversity collapse through selection/reinforcement bias (Gai et al., 2025) or finite-batch sampling bias and semantic coupling (Fan et al., 2026). These works motivate algorithmic corrections such as unlikeliness rewards, differential smoothing, inverse-success advantage calibration, and diversity-aware sampling. Complementarily, a token-level entropy perspective shows that only a minority of high-entropy “forking” tokens drive much of RLVR’s reasoning improvement (Wang et al., 2025a), while entropy-control methods such as Entrocraft explicitly shape the entropy curve to mitigate long-run performance saturation (Li et al., 2026a). Recent test-time and interpretability studies also reveal structured uncertainty in reasoning traces: Zhao et al. (2026a) use high-entropy segment centroids as intrinsic rewards for response selection, whereas Zhao et al. (2026b) distinguish low-entropy structural tokens from higher-entropy problem-specific tokens for reasoning compression and diagnosis. Our work differs from these algorithmic and empirical studies by giving a multi-task TMC account of why post-training reweights existing reasoning paths toward high-probability/easy CoTs and can

770 suppress rare-but-valid hard CoTs.

771 **Curricula, Process Rewards, and Inference-Time Scaling.** Several works study how data selection, process supervision,
 772 and inference-time search can preserve useful exploration. Difficulty-aware filtering methods select intermediate-difficulty
 773 prompts because all-correct or all-incorrect groups provide weak learning signal (Bae et al., 2025); this is aligned with
 774 our theoretical justification for rejecting overly easy instances. Outcome-based RL can provably induce CoT-style graph
 775 traversal under suitable data distributions, but its learnability relies on sufficient mass on simple examples (Bu et al., 2025),
 776 highlighting the role of implicit curricula. For process supervision, Setlur et al. (2025a) define process rewards as progress
 777 under a prover policy, while Yao et al. (2026) derive process rewards by decomposing an entropy-regularized RL objective,
 778 turning sparse outcome rewards into step-level guidance. At inference time, PRM-guided parallel generation has been
 779 analyzed through particle filtering and SMC (Golowich et al., 2026), and martingale-based foresight decoding provides
 780 another principled route to step valuation and pruning (Li et al., 2026b). These methods are complementary to our DPRM
 781 perspective, which interprets PRM/BoN-style inference as a reweighting of base-model trajectories and studies when such
 782 reweighting overemphasizes common/easy CoTs. Bu et al. (2026) extend our PRM idea to Diffusion Language Models.
 783

784 **Process Reward Models (PRMs) & Reinforcement Learning with Verifiable Rewards (RLVR).** Process Reward Models
 785 (PRMs) and Reinforcement Learning with Verifiable Rewards (RLVR) both employ external verifiers to reward reasoning
 786 steps, with PRMs guiding inference (Lightman et al., 2023; Li et al., 2023; Snell et al., 2024) and RLVR enhancing
 787 finetuning (Wang et al., 2025b; Foster et al., 2025). Setlur et al. (2025b) show that verifier-based scaling outperforms
 788 verifier-free approaches when the reward distributions have anti-concentration and heterogeneity properties. (Foster et al.,
 789 2025) also analyzed on linear softmax model, for which they designed an algorithm that is computationally efficient, and
 790 showed the necessity of coverage within their framework. Yue et al. (2025) find RLVR’s gains limited to small k , with
 791 base models matching or surpassing it at large k , suggesting RLVR reinforces existing reasoning rather than fostering new
 792 patterns—echoing our finding that RL finetuning overfits to simpler paths due to the squeezing effect. Schmied et al. (2025)
 793 highlight RLVR’s “greediness”, favoring easy actions akin to our findings, while Yu et al. (2025)’s DAPO and Xiong et al.
 794 (2025)’s minimalist approaches counter this by rejecting overly-correct samples, promoting diverse reasoning and keeping
 795 steady entropy, whose merits are also theoretically justified in our settings. (Wang et al., 2025b) also empirically showed
 796 the critical role of promoting exploration with diverse reasoning patterns. Setlur et al. (2025a) propose a separate prover
 797 policy to enhance exploration, noting the base model’s advantage calculation limits diversity—supporting our observation of
 798 RLVR’s bias toward simpler paths. Li et al. (2025b) add that cross-entropy finetuning reduces sampling diversity, reinforcing
 799 the need for varied inference strategies. These findings collectively underscore the value of diverse reasoning, motivating
 800 our comparison of RL and PRM under binary outcome rewards.
 801

802 B. Limitations and Broader Impact

803 **Unmodeled Complexity in Large-Scale.** While our theoretical analysis introduces new perspectives on finetuning and
 804 inference-scaling under binary (0–1) outcome supervision, several limitations remain. First, the latent reasoning model
 805 and neural formulation may require further refinement to better align with practical scenarios, including: handling *varying*
 806 *reasoning depths*; incorporating *structural priors* (e.g., multi-index models); modeling with nonlinear transformers instead
 807 of a linear softmax model (per discussed in App. D.4); and analyzing parameter-efficient tuning methods like LoRA (Hu
 808 et al., 2021).
 809

810 **Reward Hacking, and the Benefit of Consistency.** Even within the TMC framework, our formulation does not fully
 811 capture challenges such as robustness to noisy rewards, hallucinations, or reward hacking. For example, in Fig. 1, the
 812 trajectories $q \rightarrow o_2^1 \rightarrow a_3$ (valid for Task 4) and $q \rightarrow o_2^2 \rightarrow a_3$ (valid for Task 5) share the same endpoints but are invalid
 813 for each other’s task, illustrating a form of reward misalignment or *hallucination*. This warrants deeper investigation. A
 814 concurrent study by Wen et al. (2025) raised a concern: rather than rewarding rare reasoning paths, they classified them as
 815 incorrect CoTs and treated common paths as *logically coherent*—which they assumed correct. They further advocated for
 816 stronger verifiers and new RLVR algorithms explicitly designed to incentivize correct reasoning paths—a perspective we
 817 share. In our Multi-task TMC (Def. 2.2), our “validity” notion is to distinguish in-correct rare paths for a task with those
 818 correct ones. We left a more detailed discussions of the pros and cons of the simplicity bias an important future direction.
 819

820 **Entropy may fail to decrease in practice—particularly when the training dynamics become unstable.** In the context of
 821 RLVR finetuning using only outcome rewards (without SFT supervision):
 822

- 823 • If the base model’s capability is too weak for the target task (i.e., pass rate is too low), the gradient variance can become
 824

excessively large, leading to chaotic updates and potential entropy increase (Li & Ng, 2025).

- Likewise, if the reward oracle is noisy (Cai et al., 2025) or unable to verify intermediate reasoning steps for difficult problems (e.g., reward hacking as discussed in the previous paragraph), the supervision signals become inconsistent, again possibly causing entropy to increase.

These situations lie outside our theoretical assumptions since our framework requires the low-probability transition edge to remain above a constant (c), and does not model oracle noise. As suggested by Li & Ng (2025) and Wen et al. (2025), incorporating teacher-forced SFT to improve the base model’s competence, and enhancing reward oracle fidelity—for example, by verifying intermediate steps using tools such as Lean4 in theorem-proving—can stabilize such finetuning processes and mitigate this phenomenon.

Faster-vs-Better Trade-Off. Moreover, although our results highlight the value of diversity—particularly when a non-negligible fraction of instances require hard-to-reason CoTs—our analysis does not quantify the additional computational cost such diversity induces. This reflects an inherent tradeoff: overfitting to simpler reasoning paths enables faster finetuning when the target is improving overall accuracy within certain iterations, while supporting diverse reasoning incurs greater complexity—a “no free lunch” scenario.

Non-Markovianity of LLM Reasoning. Markov-chain (MC) abstractions—where transition probabilities encode step difficulty—are well-established in prior theory (Xu et al., 2019; Sanford et al., 2024; Abbe et al., 2024; Besta et al., 2024; Kim et al., 2025). In particular, Kim et al. (2025) model LLM inference as a metastable MC and design algorithms showing benefits of search and distillation. Building on empirical evidence of tree-alike reasoning (Lightman et al., 2023; Snell et al., 2024; Yue et al., 2025; AI et al., 2025; Gandhi et al., 2025), and observed real-world hardness metrics (base-model *pass rates* Tong et al. (2024)), our Multi-task TMC is arguably more aligned with practice than prior work. We acknowledge that MC models cannot perfectly capture actual LLM inference, per Zhang et al. (2025b) on LLM non-Markovianity. Nonetheless, this does not diminish the value of MC-based theories: conclusions remain informative and can often be generalized to non-Markovian settings with suitable extensions.

While our findings are theoretical, they provide high-level justification for recent empirical efforts that promote reasoning diversity and reject overly easy instances, offering useful insights for future work on RL fine-tuning, PRM design, and inference strategies in LLMs. We do not anticipate any direct societal risks arising from this research.

C. Additional Experiments

C.1. Comprehensive Performance and Coverage Analysis

Building upon the empirical simulations presented in Section 5, we provide additional experimental results that further validate our theoretical findings. The following analysis examines both performance metrics (Pass@K rates) and coverage characteristics (valid CoT generation patterns) across different sampling strategies for both TASK1 and TASK2.

C.1.1. PERFORMANCE ANALYSIS

Figure 2 and Figure 3 present the Pass@30 performance for TASK1 and TASK2, respectively, across all evaluated sampling strategies. The results demonstrate several key patterns that align with our theoretical predictions:

TASK1 Performance: The performance across different strategies shows relatively consistent results, with Pass@30 rates ranging from 0.65 to 0.73. Notably, DPRM achieves the highest performance (0.73), followed closely by Reinforce-rej (e.g. RL-rej) and GRPO-KL (both at 0.72). The base model performs moderately well (0.71), while PRM-BoN shows the lowest performance (0.65). This suggests that while most strategies can achieve reasonable performance on the primary task, there are meaningful differences in their effectiveness.

TASK2 Performance: The results reveal a stark contrast, with performance ranging from 0.35 to 0.95. The base model and diversity-promoting methods (Reinforce-rej (e.g. RL-rej), GRPO-KL) achieve the highest performance (0.95), demonstrating their ability to maintain capability on secondary tasks. In contrast, standard RL fine-tuning methods (REINFORCE, RAFT, PPO) show significantly degraded performance (0.35-0.48), confirming the *forgetting* phenomenon predicted by our theoretical analysis.

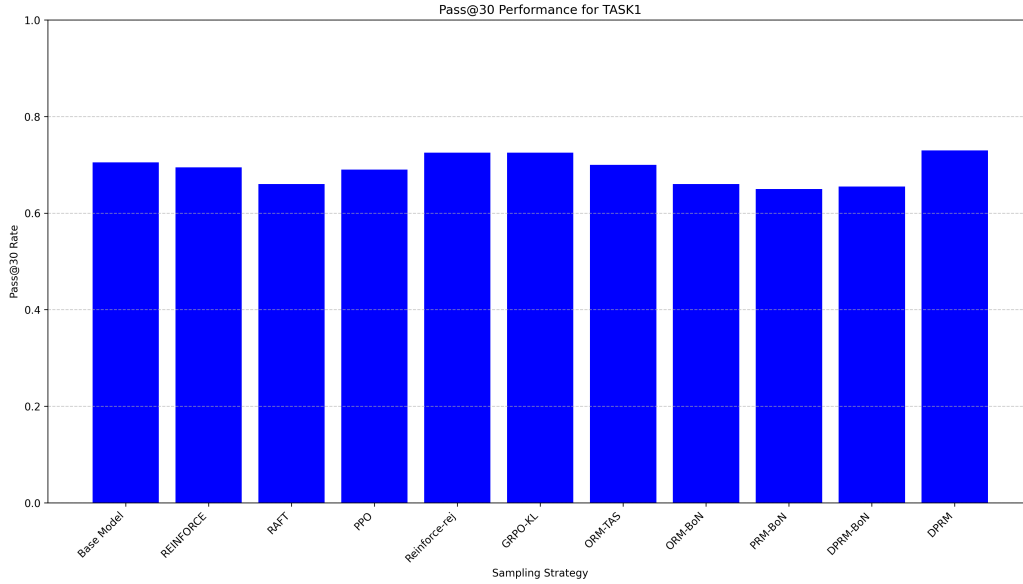


Figure 2. Pass@30 Performance for TASK1 across different sampling strategies. The results show relatively consistent performance across most methods, with DPRM achieving the highest rate of 0.73.

C.1.2. COVERAGE ANALYSIS

The coverage analysis, presented in Figure 4 and Figure 5, provides insights into the types of CoTs generated by each strategy. These stacked bar charts show the proportion of invalid, hard valid, and easy valid CoTs generated by each method.

TASK1 Coverage: The results reveal distinct patterns across different strategy categories. Standard RL fine-tuning methods (REINFORCE, RAFT, PPO) and PRM-based methods (PRM-BoN, DPRM-BoN) generate predominantly easy valid CoTs (90-98%) with minimal invalid CoTs, demonstrating strong *simplicity bias*. In contrast, diversity-promoting methods (Reinforce-rej (e.g. RL-rej), GRPO-KL) show a more balanced distribution, with substantial proportions of both easy and hard valid CoTs. The base model and ORM-based methods generate a high proportion of invalid CoTs (70-72%), indicating limited effectiveness in generating task-appropriate reasoning paths.

TASK2 Coverage: The coverage patterns for TASK2 are markedly different, reflecting the task’s increased difficulty. Most strategies generate a high proportion of invalid CoTs, with standard RL methods showing particularly poor performance (97-98% invalid). However, diversity-promoting methods (GRPO-KL, PRM-BoN, DPRM-BoN, DPRM) achieve significantly better coverage, with 45-55% valid CoTs. This demonstrates the importance of diversity-promoting mechanisms for maintaining capability across multiple tasks.

C.1.3. KEY INSIGHTS AND IMPLICATIONS

These comprehensive results provide several important insights that extend our theoretical analysis:

Simplicity Bias Confirmation: The coverage analysis clearly demonstrates the *simplicity bias* in standard RL fine-tuning methods, which overwhelmingly favor easy-to-reason CoTs while suppressing hard-to-reason alternatives. This bias is particularly pronounced in TASK1, where REINFORCE, RAFT, and PPO generate 90-95% easy valid CoTs.

Forgetting Phenomenon: The dramatic performance degradation on TASK2 for standard RL methods (from 0.70-0.72 on TASK1 to 0.35-0.48 on TASK2) provides empirical evidence for the *forgetting* phenomenon predicted by our theoretical analysis. This confirms that overfitting to the primary task can severely compromise performance on secondary tasks.

Diversity-Promoting Benefits: Methods that promote diversity (Reinforce-rej (e.g. RL-rej), GRPO-KL, DPRM variants) demonstrate superior performance on TASK2 while maintaining reasonable performance on TASK1. This validates our theoretical prediction that diversity-promoting mechanisms are crucial for multi-task scenarios.

Inference Scaling Effectiveness: The PRM-based and DPRM-based inference methods show particularly interesting

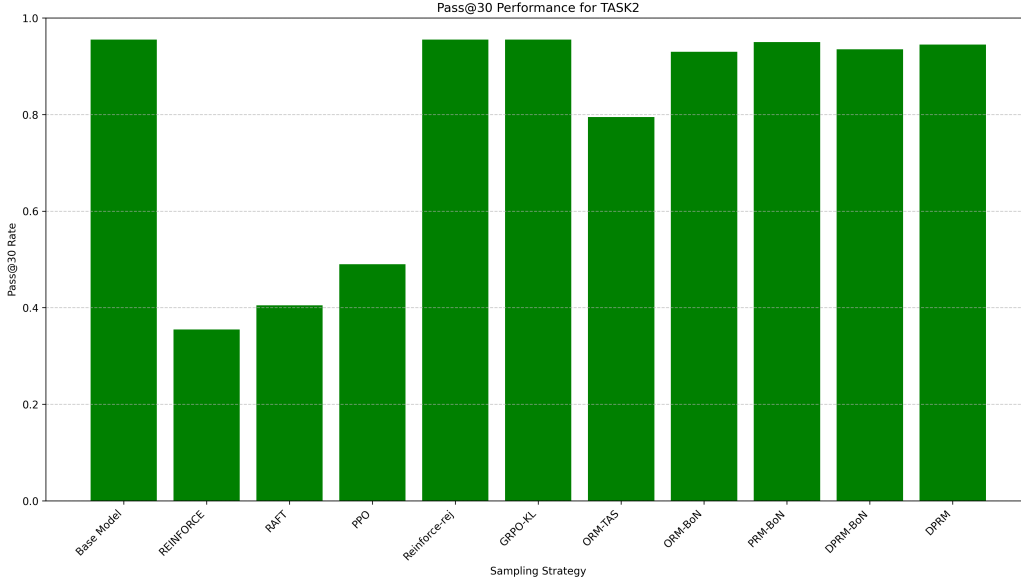


Figure 3. Pass@30 Performance for TASK2 across different sampling strategies. The results demonstrate significant performance degradation for standard RL methods (REINFORCE, RAFT, PPO) compared to diversity-promoting approaches, confirming the forgetting phenomenon.

behavior, achieving high performance on TASK1 while maintaining reasonable coverage on TASK2. This suggests that process reward models can effectively guide reasoning without the computational overhead of fine-tuning.

These results collectively support our theoretical findings and provide practical guidance for designing effective multi-task reasoning systems in large language models.

D. Details of Reward Models and Methods

D.1. Summary of Notations

We remark that in our setting, for all $l \in [L]$, $\mathbf{o}_l = e_{o_l} \in \mathbb{R}^{|S|}$ denotes the one-hot encoding of token o_l from the vocabulary. In practice, language models typically apply a softmax over the entire vocabulary to produce next-token probabilities. Hence, for simplicity, we do not distinguish between \mathbf{o}_l and o_l in notation, and treat them interchangeably throughout the paper. We summarize our notation in Table 4.

Let $\mathcal{Q}_k \subseteq S_1$ be the set of *question states* for task $k \in \mathcal{T}$. Suggest P^k is a distribution over the *question states* \mathcal{Q}_k associated with task k , denote $R_{\text{out}}^k(\mathbf{o}) = \mathbb{E}_{(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}} [R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o})]$ (Setlur et al., 2025a; 2024) as the population reward over $\mathcal{D}_{a_q}^{q,k}$ of the task tuple (q, a_q, k) .

D.2. RLVR Finetuning

REINFORCE. The classical REINFORCE algorithm (Williams, 1992) maximizes the expected reward from sampled trajectories. For mathematical reasoning, a standard approach is using 0 – 1 correctness of reasoning answer as the reward (Xiong et al., 2025; Setlur et al., 2025a). In our TMC setting, for task k and given prompt q , the REINFORCE objective is

$$\mathcal{J}_{\text{REINFORCE}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1=q \sim P^k(\mathcal{Q}_k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}, \{\mathbf{o}^i\}_{i=2}^L \sim \hat{p}_{\theta^k}(\mathbf{O}|\mathbf{o}_1)} \left[\mathbb{1}(\mathbf{o} \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}) \right], \quad (14)$$

where $\mathbf{o}_{1:L}$ denotes the trajectory sampled from the policy, and $\mathbb{1}(\mathbf{o} \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}) \in \{0, 1\}$ indicates whether the final output yields the correct answer. In our scenario, the objective would become

$$\mathcal{J}_{\text{REINFORCE}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1=q \sim P^k(\mathcal{Q}_k), \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}} [R_{\text{out}}^k(\mathbf{o})].$$

RAFT (Rejection Sampling Fine-tuning) optimizes LLMs by sampling multiple responses from a policy, using a reward

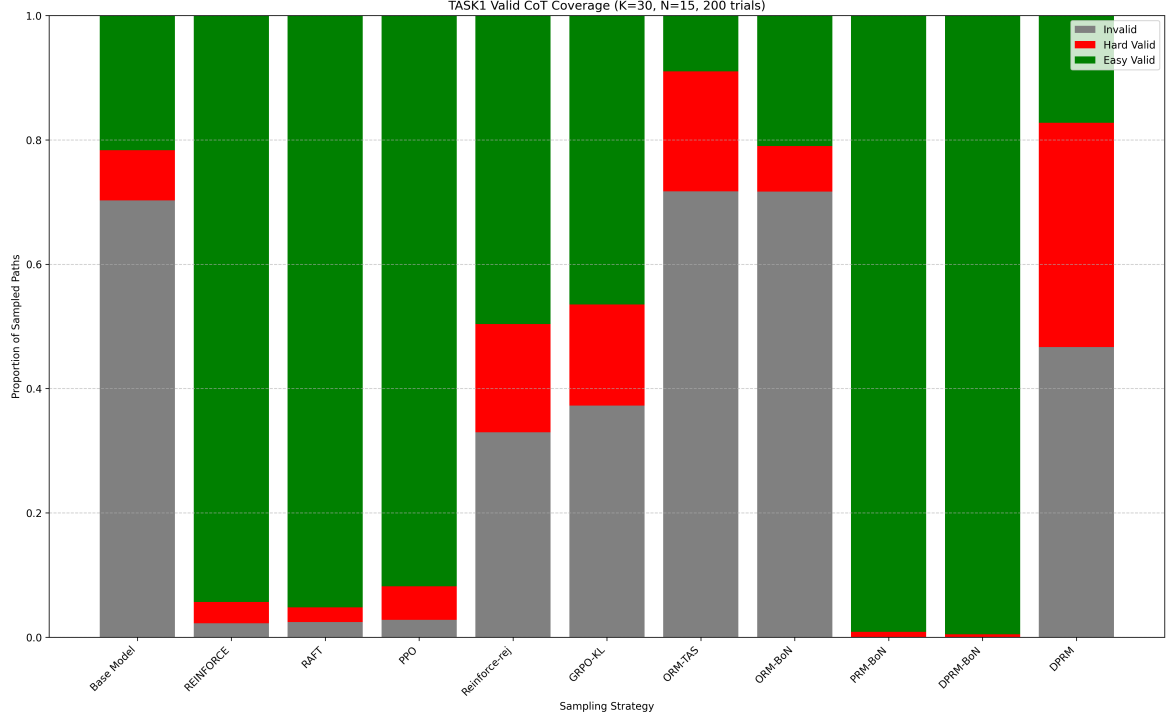


Figure 4. Valid CoT Coverage for TASK1 (K=30, N=15, 200 trials). The stacked bars show the proportion of invalid (gray), hard valid (red), and easy valid (green) CoTs generated by each strategy. Standard RL methods show strong simplicity bias with predominantly easy valid CoTs.

signal to select the best one, and then fine-tuning the policy using supervised learning on the selected best responses (Xiong et al., 2025; Dong et al., 2023). The objective is to maximize the likelihood of these high-reward outputs:

$$\mathcal{J}_{\text{RAFT}}(\theta) = \mathbb{E}_{[(q, o^*) \sim \mathcal{D}_{\text{RAFT}}]} [\log \pi_{\theta}(o^* | q)], \quad (15)$$

where $\mathcal{D}_{\text{RAFT}}$ is a dataset constructed from queries q and their corresponding best sampled responses o^* , as determined by a reward function. (Xiong et al., 2025) found that a minimal RL approach to finetune the base model is to reject both the entirely correct and incorrect responses. In our TMC case, we have

$$\mathcal{J}_{\text{RAFT}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1 \sim P^k(Q^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{\mathbf{o}_1}^{\mathbf{o}_1^k}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}(\mathbf{o}_{2:L} | \mathbf{o}_1)} \left[\sum_{l=1}^{L-1} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1} | \mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right]$$

Direct Preference Optimization (DPO) optimizes the policy directly using a dataset of human preferences, provided as pairs of preferred (o_w) and dispreferred (o_l) responses for a given prompt q (Rafailov et al., 2024). It avoids explicit reward model training or reinforcement learning, instead optimizing a loss based on the policy’s probability ratio relative to a reference policy π_{ref} :

$$\mathcal{J}_{\text{DPO}}(\theta) = \mathbb{E}_{[(q, o_w, o_l) \sim \mathcal{D}_{\text{DPO}}]} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(o_w | q)}{\pi_{\text{ref}}(o_w | q)} - \beta \log \frac{\pi_{\theta}(o_l | q)}{\pi_{\text{ref}}(o_l | q)} \right) \right], \quad (16)$$

where \mathcal{D}_{DPO} is the preference dataset, σ is the logistic sigmoid function, and β is a temperature hyperparameter that scales the difference in log-probabilities.

In our TMC setting, for task k , suppose for each prompt q , the reference model (base model \hat{p}_{θ^*} or current model \hat{p}_{old}^k) produces two candidate trajectories: a preferred one $\mathbf{o}_{1:L}^+$, and a dispreferred one $\mathbf{o}_{1:L}^-$, where $R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}_L^+) = 1 > R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}_L^-) = 0$. The DPO objective for the current policy \hat{p}_{θ^k} is:

$$\mathcal{J}_{\text{DPO}}(\theta^k) := \sum_{(q, \mathbf{o}^+, \mathbf{o}^-) \in \mathcal{D}^k} \log \sigma \left(\beta \cdot \left[\log \frac{\hat{p}_{\theta^k}(\mathbf{o}_{2:L}^+ | \mathbf{o}_1^+)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{2:L}^+ | \mathbf{o}_1^+)} - \log \frac{\hat{p}_{\theta^k}(\mathbf{o}_{2:L}^- | \mathbf{o}_1^-)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{2:L}^- | \mathbf{o}_1^-)} \right] \right), \quad (17)$$

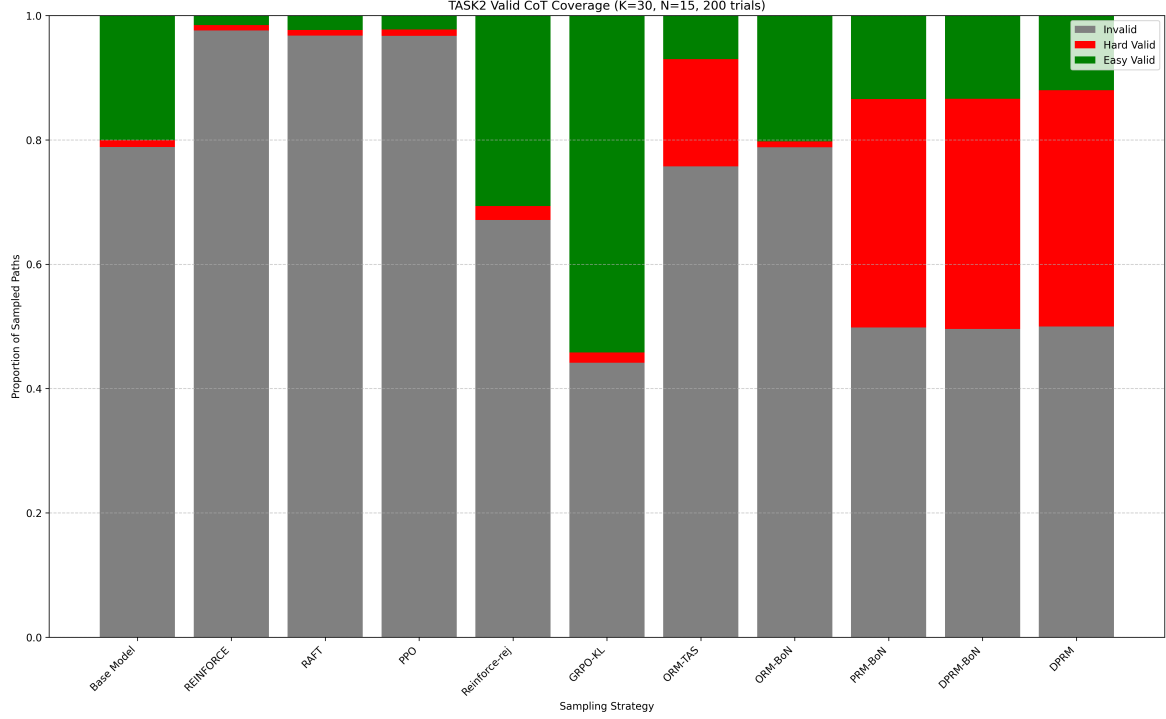


Figure 5. Valid CoT Coverage for TASK2 (K=30, N=15, 200 trials). The stacked bars show the proportion of invalid (gray), hard valid (red), and easy valid (green) CoTs generated by each strategy. Diversity-promoting methods achieve significantly better coverage compared to standard RL approaches.

where $\sigma(\cdot)$ is the sigmoid function and $\beta > 0$ is a temperature hyperparameter controlling preference sharpness. This objective promotes the likelihood ratio of preferred over dispreferred CoTs as measured under \hat{p}_{θ^k} , relative to the fixed reference \hat{p}_{θ^*} used for sampling.

Proximal Policy Optimization (PPO) (Schulman et al., 2017) optimizes LLMs by maximizing the following surrogate objective (OpenAI, 2018):

$$\mathcal{J}_{\text{PPO}}(\theta) = \mathbb{E}_{[q \sim P(Q), o \sim \pi_{\theta_{old}}(O|q)]} \frac{1}{|o|} \sum_{t=1}^{|o|} \min \left[\frac{\pi_{\theta}(o_t|q, o_{<t})}{\pi_{\theta_{old}}(o_t|q, o_{<t})} A_t, \text{clip} \left(\frac{\pi_{\theta}(o_t|q, o_{<t})}{\pi_{\theta_{old}}(o_t|q, o_{<t})}, 1 - \epsilon, 1 + \epsilon \right) A_t \right], \quad (18)$$

where A_t is the advantage computed via Generalized Advantage Estimation (GAE), requiring an *additional critic model*. ϵ is a clipping-related hyperparameter.

In our TMC setting, we have the advantage function as

$$A_{l+1}^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) := Q^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) - V^{\hat{p}_{\theta}, k}(\mathbf{o}_l). \quad (19)$$

Here the transition-value and state-value functions are

$$Q^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) := \mathbb{E}_{\mathbf{o}_1 = q \sim P^k(Q_k), \mathbf{o}_{l+2:L} \sim \hat{p}_{\theta}} [R_{\text{out}}^k(\mathbf{o}) | \mathbf{o}_l, \mathbf{o}_{l+1}], \quad (20)$$

$$V^{\hat{p}_{\theta}, k}(\mathbf{o}_l) := \mathbb{E}_{\mathbf{o}_1 = q \sim P^k(Q_k), \mathbf{o}_{l+1} \sim \hat{p}_{\theta}(\cdot | \mathbf{o}_l)} [Q^{\hat{p}_{\theta}, k}(\mathbf{o}_l, \mathbf{o}_{l+1})]. \quad (21)$$

Table 4. Summary of Notations

Notation	Description
S_l, S	State space at layer l ; $S = \bigcup_{l=1}^L S_l$ is the full state space.
$o, O, \mathbf{o}, o_l, \mathbf{o}_l$	O denotes a trajectory; \mathbf{o} its one-hot form; o_l, \mathbf{o}_l are step- l token and embedding.
$\mathbb{P}(\cdot \cdot)$ or $\mathbb{P}_{\text{TMC}}(\cdot \cdot), \hat{p}_{\theta}(\cdot \cdot)$	$\mathbb{P}(\cdot \cdot)$ or $\mathbb{P}_{\text{TMC}}(\cdot \cdot)$: TMC kernel in Def. 2.1; $\hat{p}_{\theta}(\cdot \cdot)$: softmax predictor based on θ .
C_{o_l}, D_{o_l}	High probability transition subset in S_{l+1} ; Non-zero probability transition subset.
$\mathcal{Q}_k, (q, a_q, k)$	$\mathcal{Q}_k \subseteq S_1$: question states in task k ; (q, a_q, k) : task tuple with $q \mapsto a_q$.
$\mathcal{D}_{a_q}^{q,k}, \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}$	Instance Distribution over task tuple (q, a_q, k) ; Correct CoTs for $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$.
$\mathcal{G}_{q,a_q}^{(k)}, \mathcal{G}_{q,a_q}^{(k), \text{easy}}, \mathcal{G}_{q,a_q}^{(k), \text{hard}}$	Valid CoTs set for (q, a_q, k) ; partitioned into easy and hard subsets.
$\mathcal{I}_{o_{l+1}, o_l}^{(k)}, \mathcal{S}_{o_l}^{(k)}, \mathcal{S}_{o_l}^{(k), \text{easy}}, \mathcal{S}_{o_l}^{(k), \text{hard}}$	Valid CoTs passing (o_l, o_{l+1}) ; subset of reachable o_{l+1} from o_l in valid CoTs; easy/hard-to-reason subsets.
$\theta^*, \theta^k, \theta^{k,(t)}$	θ^* : base model in Sec. 2.2; θ^k : task- k model; superscript (t) : iteration.
$R_{\text{out}}^k(o), R_{\text{out}}^{\hat{p}}(o)$	$R_{\text{out}}^k(o)$: Expected accuracy over $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$ for sampled CoT o , by \hat{p}_{θ^*} or \hat{p} .
$R_{\text{likelihood}}^k(o_l), R_{\text{DPRM}}^k(o_l)$	Expected accuracy of o_l ; DPRM reward in Eq.(13).
$A_{l+1}^{\hat{p}_{\theta^k}}(o_l, o_{l+1}), Q_{\hat{p}_{\theta^k}}(o_l, o_{l+1}), V^{\hat{p}_{\theta^k}}(o_l)$	RL's Advantage for task k ; Expected accuracy of state o_{l+1} and o_l .
$P_{\text{acc}}^k(o)$	Success probability of CoT o for task k .
β, λ	Temperature parameters of $\hat{p}_{\theta^k}^{\text{PO}}$ in Eq.(9) and $P_{\text{Gibbs}}^k(\mathbf{o})$ in Eq.(11).
$\text{Pass@K}_{q,k}^{\hat{p}}$	Probability that \hat{p} generates at least one correct CoT in K samples for (q, a_q, k) .
$O(\cdot), \Omega(\cdot), \Theta(\cdot)$	Standard asymptotic notation: upper, lower, and tight bounds, respectively.

The PPO objective (OpenAI, 2018) in our scenario is

$$\mathcal{J}_{\text{PPO}}(\theta^k) = \mathbb{E}_{q \sim P^k(\mathcal{Q}_k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}, \{\mathbf{o}^i\}_{i=2}^G \sim \hat{p}_{\theta^k}(O|o_1^i)} \left[\frac{1}{L} \sum_{l=1}^{L-1} \min \left[\frac{\hat{p}_{\theta^k}(o_{l+1}^i | o_l^i)}{\hat{p}_{\text{old}}^k(o_{l+1}^i | o_l^i)} A_{l+1}^{\hat{p}_{\theta^k}, k}, \right. \right. \\ \left. \left. \text{clip} \left(\frac{\hat{p}_{\theta^k}(o_{l+1}^i | o_l^i)}{\hat{p}_{\text{old}}^k(o_{l+1}^i | o_l^i)}, 1 - \epsilon, 1 + \epsilon \right) A_{l+1}^{\hat{p}_{\theta^k}, k} \right] \right]. \quad (22)$$

In our modeling setup, the advantage estimate $A_{l+1}^{\hat{p}_{\theta^k}, k}$ aims to approximate Eq.(19), the gap between the value of making a particular transition at step l , versus the expected value of acting from state o_l without knowledge of o_{l+1} .

GRPO (Shao et al., 2024), in contrast, samples a group of output trajectories $\{o^i\}_{i=1}^G$ from $\pi_{\theta_{old}}$ and optimizes:

$$\mathcal{J}_{\text{GRPO}}^k(\theta) = \mathbb{E}_{[q \sim P(Q), \{\mathbf{o}^i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]} \\ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o^i|} \sum_{t=1}^{|o^i|} \left\{ \min \left[\frac{\pi_{\theta}(o_t^i | q, o_{<t}^i)}{\pi_{\theta_{old}}(o_t^i | q, o_{<t}^i)} \hat{A}_{i,t}, \text{clip} \left(\frac{\pi_{\theta}(o_t^i | q, o_{<t}^i)}{\pi_{\theta_{old}}(o_t^i | q, o_{<t}^i)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] \right. \\ \left. - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta} || \pi_{\text{ref}}] \right\}, \quad (23)$$

where $\hat{A}_{i,t}$ is computed based on relative rewards within the sampled group, and β controls KL regularization.

In our scenario, the formulation of GRPO (Shao et al., 2024) equates

$$\mathcal{J}_{\text{GRPO}}^k(\theta^k) = \mathbb{E}_{q \sim P^k(\mathcal{Q}_k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}, \{\mathbf{o}^i\}_{i=2}^G \sim \hat{p}_{\theta^k}(O|o_1^i)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{L} \sum_{l=1}^{L-1} \left\{ \min \left[\frac{\hat{p}_{\theta^k}(o_{l+1}^i | o_l^i)}{\hat{p}_{\text{old}}^k(o_{l+1}^i | o_l^i)} \hat{A}_{i,l+1}^k, \right. \right. \right. \\ \left. \left. \text{clip} \left(\frac{\hat{p}_{\theta^k}(o_{l+1}^i | o_l^i)}{\hat{p}_{\text{old}}^k(o_{l+1}^i | o_l^i)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,l+1}^k \right] \right\} - \beta D_{\text{KL}}[\hat{p}_{\theta^k} || \hat{p}_{\theta^*}] \right], \quad (24)$$

Outcome Supervision RL with GRPO. A outcome reward model assigns scores $\mathbf{r} = \{r_1^k, \dots, r_G^k\}$ to sampled outputs, which are then normalized: $\tilde{r}_i = \frac{r_{i, \text{index}(l)}^k - \text{mean}(\mathbf{r}^k)}{\text{std}(\mathbf{r}^k)}$ within the group. The advantage is set as $\hat{A}_{i,l+1}^k = \tilde{r}_i^k, \forall l \in [L-1]$, aiming to approximate Eq. (19). Here, for the task $k \in \mathcal{T}$, if we consider an *offline* scenario, our outcome reward model is $r_{i, \text{index}(l)}^k = R_{\text{out}}^k(\cdot)$ defined in Sec. 4.

Process Supervision RL with GRPO. Instead of a single reward per output, a process reward model assigns step-wise rewards $\mathbf{R} = \{\{r_{1, \text{index}(1)}^k, \dots, r_{1, \text{index}(L)}^k\}, \dots\}$, where $\text{index}(l)$ denotes the l -th step's end token index. Rewards are

normalized: $\tilde{r}_{i,index(l)}^k = \frac{r_{i,index(l)}^k - \text{mean}(\mathbf{R})}{\text{std}(\mathbf{R})}$. The advantage is computed as:

$$\hat{A}_{i,l}^k = \sum_{index(j) \geq l} \tilde{r}_{i,index(j)}^k, \quad (25)$$

and the policy is optimized via Eq. (6). Specifically, we could adopt $r_{i,index(l)}^k = R_{\text{pro}}^k(\mathbf{o}_l^i), \forall i \in [G], l \in \{1, \dots, L\}$ in Sec. 4. However, this approach is unnatural - since $R_{\text{DPRM}}^k(\cdot)$ is designed for temperature-controlled adjusted sampling. Instead, a more common approach is to choose the $R_{\text{likelihood}}^k(\mathbf{o}_l)$ in Eq.(10) and $R_{\text{potential}}^k(\mathbf{o}_l)$.

In this work, following (Xiong et al., 2025; Yu et al., 2025), we only studied the properties of GRPO with outcome reward. However, our theorem can include the GRPO with process reward by assuming that the advantage calculated in Eq.(25) is approximating Eq.(19) accurately.

D.3. Reward-based Sampling

ORM Mode. Given an input x , the model generates an CoT trajectory $\mathbf{o}_1, \dots, \mathbf{o}_L$. Define $\mathbf{o}_l \in \mathbb{R}^{|\mathcal{S}|}$ as the one-hot vector representing \mathbf{o}_l , $\mathbf{o} = (\mathbf{o}_1, \dots, \mathbf{o}_L)^\top \in \mathbb{R}^{L \times |\mathcal{S}|}$ as the trajectory vector. An outcome reward model (ORM) $R_{\text{out}}^k(\cdot)$ assigns a scalar score based on the entire output:

$$R_{\text{out}}^k(\mathbf{o}) = f(\mathbf{o}), \quad (26)$$

where $f(\cdot)$ usually evaluates correctness, coherence, or other task-specific criteria (Shao et al., 2024; Wang et al., 2024; Li et al., 2023; Snell et al., 2024).

PRM Mode. Instead of rewarding only the final output, a *process reward model* (PRM) assigns intermediate rewards along the reasoning trajectory:

$$R_{\text{pro}}^k(\mathbf{o}_l) = g(\mathbf{o}_1, \dots, \mathbf{o}_l), \quad l \in \{1, \dots, L\}, \quad (27)$$

where $g(\cdot)$ estimates step-wise utility using heuristics, verification signals, or learned evaluation metrics (Shao et al., 2024; Snell et al., 2024; Wang et al., 2024; Li et al., 2023). Designing process rewards from outcome rewards is essential due to the high cost of human annotation. However, existing approaches are largely heuristic—either based on (i) the expected correctness of the final answer from the current state, typically via Monte Carlo rollouts (Setlur et al., 2025a; 2024; Wang et al., 2024):

$$R_{\text{likelihood}}^k(\mathbf{o}_l) = \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} [R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l], \quad (28)$$

and (ii) using binary signals to indicate whether the current state can still reach a correct solution (Snell et al., 2024; Setlur et al., 2025b):

$$R_{\text{potential}}^k(\mathbf{o}_l) = \sup_{\mathbf{o}': \mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}' \in \mathcal{T}_{\text{all}}} R_{\text{out}}^k(\mathbf{o}') = \mathbf{1} \{ \exists \mathbf{o}' \in \mathcal{T}_{\text{all}} : R_{\text{out}}^k(\mathbf{o}', a) \}, \quad (29)$$

for all $\mathbf{o}_l \in S_l, l \in \{1, \dots, L\}$. Here, \mathcal{T}_{all} is typically approximated typically by Monte Carlo rollouts.

Temperature-controlled Adjusted Sampling. Here, we consider refinubf the sampling distribution using the reward model R_{out}^k . Define the original sampling probability of a trajectory \mathbf{o} under θ^* as:

$$\hat{p}_{\theta^*}(\mathbf{o}) = \mathbb{P}_{\rho}^{\text{test}}(\mathbf{o}_1) \prod_{l=1}^{L-1} \hat{p}_{\theta^*}(\mathbf{o}_{l+1} \mid \mathbf{o}_l),$$

where $\mathbb{P}_{\rho}^{\text{test}}(\mathbf{o}_1) = \Theta(1/M_0)$ is the initial distribution over S_1 . The adjusted sampling distribution, guided by R_{out}^k , is defined as:

$$P_{\text{Gibbs}}^k(\mathbf{o}) = \frac{\hat{p}_{\theta^*}(\mathbf{o}) \exp(\lambda R_{\text{out}}^k(\mathbf{o}))}{\sum_{\mathbf{o}' \in \mathcal{T}_{\text{all}}} \hat{p}_{\theta^*}(\mathbf{o}') \exp(\lambda R_{\text{out}}^k(\mathbf{o}'))} \propto \hat{p}_{\theta^*}(\mathbf{o}) \exp(\lambda R_{\text{out}}^k(\mathbf{o})). \quad (30)$$

for a temperature parameter $\lambda > 0$, with normalization over \mathcal{T}_{all} . The estimation of the \mathcal{T}_{all} is typically through *Monte Carlo Rollout*. This discrete distribution reweights the pretrained model’s probabilities to favor trajectories with higher estimated rewards, consistent with traditional sampling literature where the exponential form amplifies the influence of the reward signal. The form $P_{\text{Gibbs}}^k(\mathbf{o}) \propto \hat{p}_{\theta^*}(\mathbf{o}) \exp(\lambda R_{\text{out}}^k(\mathbf{o}))$ mirrors soft policy sampling in RL and NLP literature (e.g., REINFORCE or importance sampling). λ controls the trade-off: large λ heavily biases toward high-reward trajectories; small λ preserves the original distribution.

D.4. Discussion on Broader Finetuning Settings

D.5. Benefit of Broader Exploration Strategies

Beyond standard RL or KL-regularized finetuning, our theoretical framework also provides insight into several *non-standard* post-training strategies that emphasize broader exploration.

One example is *Evolution Strategy (ES) finetuning* (Qiu et al., 2025), which updates parameters via isotropic perturbations,

$$\theta_{t+1} = \theta_t + \frac{\alpha}{N} \sum_{n=1}^N R(\theta_t + \sigma \epsilon_n) \epsilon_n, \quad \epsilon_n \sim \mathcal{N}(0, I),$$

where α and σ are hyperparameters and $R(\cdot)$ denotes a reward model. From the perspective of our TMC analysis, such isotropic exploration assigns comparable reward access to parameter directions corresponding to both easy-to-reason and hard-to-reason transitions at each layer l . This contrasts with advantage-based RL finetuning, which relies on the model’s own CoT sampling and therefore preferentially reinforces easy-to-reason edges with higher probability. As a result, ES does not systematically amplify already frequent easy directions and thus naturally mitigates the squeezing effect identified in our theory.

A related class of methods is *representation-based exploration finetuning* (Tuyls et al., 2025), which explicitly encourages large covariance and diversity in latent representations. Under our framework, this can be interpreted as preserving diversity across parameter directions associated with low-probability transitions, preventing them from being collapsed or suppressed by repeated advantage-driven updates. Such behavior aligns precisely with our theoretical characterization of mechanisms that counteract the squeezing of hard-to-reason CoT steps.

Taken together, these examples illustrate that our theoretical perspective not only consolidates existing intuitions about exploration, but also offers a unifying lens for understanding and motivating broader, less conventional exploration-based post-training strategies.

Case of DPO. Recall from Eq.(17) that the DPO objective is defined as:

$$\mathcal{J}_{\text{DPO}}(\theta^k) := \sum_{(q, \sigma^+, \sigma^-) \in \mathcal{D}^k} \log \sigma \left(\beta \cdot \left[\log \frac{\hat{p}_{\theta^k}(\sigma_{2:L}^+ | \sigma_1^+)}{\hat{p}_{\text{old}}^k(\sigma_{2:L}^+ | \sigma_1^+)} - \log \frac{\hat{p}_{\theta^k}(\sigma_{2:L}^- | \sigma_1^-)}{\hat{p}_{\text{old}}^k(\sigma_{2:L}^- | \sigma_1^-)} \right] \right),$$

where $R_{(\mathbf{Q}, \mathbf{A})}^k(\sigma_L^+) = 1 > R_{(\mathbf{Q}, \mathbf{A})}^k(\sigma_L^-) = 0$. As discussed in (Ren & Sutherland, 2025), DPO can exhibit a squeezing effect, and such dynamics might also apply under our TMC reasoning framework. However, DPO is not a natural fit for our setting: we are concerned with correctness rather than relative preferences over reasoning paths. As such, the data required to support DPO—pairs (q, σ^+, σ^-) indicating relative preference—is not directly meaningful in our binary (0–1) reward formulation. For this reason, while the objective form is stated for reference, we do not pursue further theoretical development of DPO in this work. Nonetheless, it may serve as a promising direction for future study of RLHF under the TMC framework with additional assumptions on preference structure.

D.6. Extension to the general (nonlinear Transformer) case

Our multi-task TMC framework recovers three empirically observed phenomena (Phenomenon 1-3 in the introduction) in nonlinear multihead Transformers. While it is common for theoretical analyses of large-scale LMs to use idealized surrogate models to distill and prove generalizable principles (e.g., Foster et al. (2025); Kim et al. (2025)), we here further clarify how our theory connects to nonlinear Transformer setting.

Linear surrogates sufficiently capture tabular latent-state transitions. Our formulation models reasoning as a discrete Markov chain—an abstraction used in several recent studies (Xu et al., 2019; Sanford et al., 2024; Abbe et al., 2024; Besta et al., 2024; Kim et al., 2025)—where the current state encodes all information for current reasoning step. Thus, global token dependencies are captured in state transitions, eliminating the need for positional entanglement. Prior work (Nichani et al., 2024; Edelman et al., 2024) has shown that transformers can successfully learn Markovian dynamics, and in our setting, the linear softmax model is already overparameterized enough to capture the TMC structure.

Extension to the general (nonlinear Transformer) case. Our multi-task TMC recovers three empirically observed phenomena (**Phenomenon 1-3** in the introduction) that also hold in nonlinear multihead Transformers:

Phenomenon 1(architecture-aware): RL-induced squeezing. Squeezing (sharpening) under RL post-training: rare-but-correct CoTs are forgotten, a behavior widely reported in math/coding systems (He et al., 2025) but previously lacking theoretical explanation. In our framework, it emerges when gradients over-reinforce easy CoTs driven by their higher advantage, also stemming from the decoupled neural representation of different states.

Nonlinear Logits. When the model’s logits deviate from the linear form in Eq.(2 and instead follow the general parameterization of Eq.(89, i.e., $\hat{p}_\theta(\cdot | \mathbf{x}) = \text{softmax}(h_\theta(\cdot, \mathbf{x}))$ for $\mathbf{x} \in \{0, 1\}^{|\mathcal{S}|}$, the fine-tuning dynamics become considerably more complex.

As noted in Remark I.4, Lemma I.2 depends on a set of extended conditions, notably the Parameter Isolation condition (Eq.(108)), which typically fails to hold in practice. In large language models (LLMs), token representations are entangled via shared parameters across layers and positions, making it impossible to isolate updates per token. This design is aligned with in-context learning (Nichani et al., 2024), where sequential dependencies are a fundamental modeling assumption.

To understand the impact of nonlinearity more concretely, we adopt a first-order approximation of the logit update at transition $o_l \rightarrow o_{l+1}$ following Proposition 1 in (Ren & Sutherland, 2025):

$$\begin{aligned} \Delta \log \hat{p}_\theta(\cdot | o_l) &= \eta [\nabla_{h(\cdot, o_l)} \log \hat{p}_\theta(\cdot | o_l)] \mathbb{E} \{ \mathcal{K}_\theta(o_l, o_l^{\text{train}}) [\nabla_{h(\cdot, o_l)} \mathcal{J}^{\text{train}}(o_{l+1}^{\text{train}}, o_l^{\text{train}})] \} + O(\eta^2 \|\nabla_\theta h_\theta(\cdot, o_l)\|_{\text{op}}) \\ &= \eta (\mathbb{I} - \mathbf{1} \hat{p}_\theta(\cdot | o_l)^\top) \mathbb{E} \{ \mathcal{K}_\theta(o_l, o_l^{\text{train}}) [\nabla_{h(\cdot, o_l)} \mathcal{J}^{\text{train}}(o_{l+1}^{\text{train}}, o_l^{\text{train}})] \} + O(\eta^2 \|\nabla_\theta h_\theta(\cdot, o_l)\|_{\text{op}}) \end{aligned}$$

where $\mathcal{J}^{\text{train}}$ is the state-wise loss function (e.g. entropy loss or expected accuracy $Q(o_{l+1}^{\text{train}}, o_l^{\text{train}})$), \mathcal{K}_θ is the empirical NTK (eNTK) defined as $\mathcal{K}_\theta = (\nabla_\theta h_\theta(\cdot, o_l) \nabla_\theta h_\theta(\cdot, o_l)^\top)$, and the expectation is taken over question states $q = o_1 \sim P^k(\mathcal{Q}^k)$, training instances $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$, and sampled CoTs $o^{\text{train}} \sim \hat{p}_{\theta^*}(\cdot | q)$. In contrast to the linear case where $\mathcal{K}_\theta = o_l o_l^\top$, the nonlinear update depends on the learned geometry of the representation space.

The squeezing effect occurs if

$$\frac{\hat{p}_{\theta^*}(o'_{l+1} | o_l)}{\hat{p}_{\theta^*}(o_{l+1} | o_l)} \leq 1 \iff \Delta \log \hat{p}_\theta(o'_{l+1}) \geq \Delta \log \hat{p}_\theta(o_{l+1}),$$

for $o'_{l+1} \in C_{o_l}$ and $o_{l+1} \in D_{o_l} \setminus C_{o_l}$. The update difference satisfies:

$$\Delta \log \hat{p}_\theta(o'_{l+1} | o_l) - \Delta \log \hat{p}_\theta(o_{l+1} | o_l) = \eta [o'_{l+1} - o_{l+1}]^\top \Theta ((\mathbb{I} - \mathbf{1} \hat{p}_\theta(\cdot | o_l)^\top) \mathbb{E} \{ \mathcal{K}_\theta(o_l, o_l^{\text{train}}) [\nabla_{h(\cdot, o_l)} \mathcal{J}^{\text{train}}(o_{l+1}^{\text{train}}, o_l^{\text{train}})] \}).$$

This shows that the relative update magnitudes—and thus the squeezing effect—depend on the eNTK structure and how different CoT representations interact. If the non-linear representations of hard and easy CoTs are highly correlated, their learning dynamics may reinforce or suppress each other, analogous to the phenomenon in (Ren & Sutherland, 2025), where learning digit 4 accelerates digit 9 but impedes unrelated classes. In our setting, this implies that whether the squeezing effect persists under nonlinearity hinges on structural coupling between CoTs in the representation space. In real-world, different reasoning patterns do have co-relations, and we left the broader investigations with certain assumptions as an important future direction.

Phenomenon 2(architecture-agnostic): Consistency bias of neural verifiers. As shown in Prop. 4.2 likelihood-based population objectives intrinsically upweight frequent patterns and downweight rare CoTs—this bias arises from the objective itself, not the network class. This explains the observed phenomenon in real-world LLMs (Xu et al., 2025).

Phenomenon 3(architecture-agnostic): Hard instances rely on rare CoTs. This is defined by base-model pass rate (Tong et al., 2024) and is independent of the underlying architecture.

E. Details and Proofs of TMDP

Remark E.1. For the reader’s high-level understanding, we here list some scenarios where the common valid reasoning patterns do not suffice for specific instance.

- **Problem type: Algebra (quadratic equations)**

Common Valid CoT: applying factorization method to solve quadratic equations.

Scenario it is not Correct: when the quadratic polynomial is irreducible over integers (e.g., $x^2 + x + 1 = 0$), factorization fails.

- 1320 • **Problem type: Geometry (triangle side relations)**
- 1321 **Common Valid CoT:** applying the Pythagorean theorem to relate side lengths of triangles.
- 1322 **Scenario it is not Correct:** when the triangle is not right-angled, Pythagoras’ theorem does not hold.
- 1323
- 1324 • **Problem type: Probability (complex event calculation)**
- 1325 **Common Valid CoT:** applying the law of total probability to compute probabilities.
- 1326 **Scenario it is not Correct:** when the partition of events is not mutually exclusive, leading to double counting.
- 1327
- 1328 • **Problem type: Number theory (modular arithmetic)**
- 1329 **Common Valid CoT:** reasoning with modular addition to check congruences.
- 1330 **Scenario it is not Correct:** when an incorrect modulus is used (e.g., reducing modulo 6 instead of 7).
- 1331
- 1332 • **Problem type: Combinatorics (counting problems)**
- 1333 **Common Valid CoT:** applying permutation and combination formulas.
- 1334 **Scenario it is not Correct:** when order vs. unordered distinction is misapplied, such as using combinations when
- 1335 permutations are required.
- 1336

1337 This view is supported by recent large-scale error analyses on real math datasets. Sun et al. (2025) construct *MWPES-300K*
 1338 (304,865 erroneous solutions across 15 LLMs and 4 datasets: SVAMP, GSM8K, AQuA, MATH) and discover that (i) error
 1339 patterns *diversify with dataset difficulty* (e.g., MATH consistently elicits more diverse error types than GSM8K/SVAMP),
 1340 indicating that simple “valid” patterns cease to be *correct* on harder instances; (ii) many failures arise from *mis-applied*
 1341 *common patterns*, such as *Assumed independence of overlapping events (AIO)*, *Misapplication of probability formulas*
 1342 *for independent events (MPI)*, *Incorrect combinatorial principles (ICP)*, *Unit/Conversion errors (UNE/FAC)*, or *algebraic*
 1343 *manipulation mistakes (MAM)*, showing that widely used CoT routes are not instance-wise reliable; and (iii) *Error-Aware*
 1344 *Prompting (EAP)* selectively diverts models from their default CoT routes and yields sizable per-category gains on hard cases
 1345 (e.g., AIO +6.1pp, MPI +6.5pp, UNE +6.5pp, FAC +13.5pp), evidencing the value of rarer, problem-specific reasoning
 1346 paths over frequent but brittle patterns.

1347 This aligns with recent findings (Xiong et al., 2025; Li et al., 2025b; Ren & Sutherland, 2025; Wang et al., 2025b)
 1348 highlighting the role of *reasoning diversity* and *entropy stability* in post-training, albeit evidence shows that post-training
 1349 and inference-scaling do not explore beyond base model’s tree-search knowledge (Yue et al., 2025; AI et al., 2025; Gandhi
 1350 et al., 2025).

1351 **Definition E.2** (Formal Version of Def. 2.2). Let $X = (X_t)_{t \geq 0}$ be a Tree-structured Markov Chain (TMC) as defined in
 1352 Def. 2.1, and let \mathcal{T} be a collection of tasks. Each task $k \in \mathcal{T}$ specifies a set of different question *state* $\mathcal{Q}_k \subset S_1$, where each
 1353 $q \in \mathcal{Q}_k$ has a corresponding unique correct answer a_q^k under task k . For $(q, a_q^k) \neq (q', a_{q'}^k) \in \mathcal{Q}_k$ we have $q \neq q'$, $a_q^k \neq a_{q'}^k$.

1354 A state tuple $(q, a_q = a_q^k, k)$ is called *common* if there exists at least one easy-to-reason chain of thought (CoT) (o_1, \dots, o_L)
 1355 from q to a_q , and *rare* otherwise. Each such state tuple is associated with a set $\mathcal{G}_{q, a_q}^{(k)} \subset S_1 \times \dots \times S_L$ of *valid* CoTs.

- 1359 1. All easy-to-reason CoTs from q to a_q belong to $\mathcal{G}_{q, a_q}^{(k)}$;
- 1360 2. These CoTs are *not valid* for any task $k' \neq k$;
- 1361 3. Hard-to-reason CoTs may or may not belong to $\mathcal{G}_{q, a_q}^{(k)}$;
- 1362 4. Every edge $(o_l \rightarrow o_{l+1})$ with non-zero transition probability appears in some valid CoT for some task;
- 1363 5. Each task state tuple (q, a_q, k) induces a QA distribution $\mathcal{D}_a^{q, k}$, and the probability that a valid CoT $o_{1:L} \in \mathcal{G}_{q, a_q}^{(k)}$ is
 1364 *correct* for a concrete instance $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_a^{q, k}$ is given by $p_{\text{acc}}^k(o)$:

$$1370 p_{\text{acc}}^k(o) = \frac{\prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l)}{\sum_{o'_{1:L} \in \mathcal{G}_{q, a_q}^{(k)}} \prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o'_{l+1} | o'_l)}. \quad (31)$$

1373 Further, we assume that the correctness of any different $o \neq o' \in \mathcal{G}_{q, a_q}^{(k)}$ for any instance $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_a^{q, k}$ is *independent*.

A task k is denoted *rare* if there is no valid easy-to-reason CoT in its $\mathcal{G}_{q,a_q}^{(k)}$ for any question-answer state pair $(a, a_q) \in \mathcal{Q}_k$, and *common* other wise.

Remark E.3. Here, the fifth condition is to provide merits for the probability distribution of the original TMC X (\mathbb{P}_{TMC}) that models after real-world LLM. Typically, the predictive distribution obtained from pretraining would match the ‘‘frequency’’ of whether a CoT be valid for certain task. That is, through the 5-th condition, we justify why the original TMC X (\mathbb{P}_{TMC}) would be equipped its distribution–driven by inherent chance to become a valid CoT for some reasoning task.

The *independence* assumption in the 5-th condition is for technical convenience. This definition would also induces instance (\mathbf{Q}, \mathbf{A}) under task k that has no correct CoT, with probability $\prod_{o \in \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k)}} (1 - p_{\text{acc}}^k(o))$. In real-world, the situation is far more complex, and we left the consideration of theory that assumes the interaction and co-relationship of the correctness of CoTs with different difficulty level for future work.

Some examples of Valid-not-Correct:

For some task tuple (o_1, a_{o_1}, k) , denote $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}} \subseteq \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k)}$ as the subset of valid easy-to-reason CoTs inside the valid CoTs for task k , and $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}} := \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k)} \setminus \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}}$ the subset of valid hard-to-reason CoTs. For any sampled instance $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_a^{q,k}$, it has the following two scenarios:

- With probability $\prod_{o \in \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}}} (1 - p_{\text{acc}}^k(o))$, the (\mathbf{Q}, \mathbf{A}) can only be correctly solved by some valid hard-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}$.
- With probability $1 - \prod_{o \in \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}}} (1 - p_{\text{acc}}^k(o))$, the (\mathbf{Q}, \mathbf{A}) can be correctly solved by some easy-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}}$.

This division of probability space would equip bounding the pass@K performance when the model is only capable of all valid easy-to-reason CoTs. After the finetuned model is also capable of the hard-to-reason CoTs in $\widetilde{\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}}$, we turn to be interested in the following division of probability space to discuss the pass@K performance:

- With probability $\prod_{o \in \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}} \cup \widetilde{\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}}} (1 - p_{\text{acc}}^k(o))$, the (\mathbf{Q}, \mathbf{A}) can only be correctly solved by some **unlearned** valid hard-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}} \setminus \widetilde{\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}}$.
- With probability $1 - \prod_{o \in \mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}} \cup \widetilde{\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}}} (1 - p_{\text{acc}}^k(o))$, the (\mathbf{Q}, \mathbf{A}) can be correctly solved by some easy-to-reason CoTs or learned valid hard-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{easy}} \cup \widetilde{\mathcal{G}_{o_1, a_{o_1}^{(k)}}^{(k), \text{hard}}}$.

We can characterize the breadth of tasks encoded in the topology of TMC as follows.

Corollary E.4 (Cardinality of Multi-task TMC). *Let $M_0 = |S_1|$, $M_L = |S_L|$, and for each $q \in S_1$ define*

$$A(q) = \{a \in S_L : \exists \text{ easy-to-reason CoT } q \rightsquigarrow a\}, \quad R(q) = \{a \in S_L : \exists \text{ CoT } q \rightsquigarrow a\}.$$

By Def. 2.1, we have $|A(q)| = n_q = O(1)$, and $n_q \geq 2$ for all q . Define the set of all tasks as

$$\mathcal{T} = \{k : S_1 \rightarrow S_L\}, \quad \mathcal{T}_{\text{common}} = \{k : k(q) \in A(q) \text{ for all } q\}, \quad \mathcal{T}_{\text{rare}} = \mathcal{T} \setminus \mathcal{T}_{\text{common}}.$$

Then

$$|\mathcal{T}_{\text{common}}| = \prod_{q \in S_1} |A(q)| = \Theta(c^{M_0}), \quad 2 \leq c \leq \max_q n_q = O(1),$$

and

$$|\mathcal{T}| \leq \prod_{q \in S_1} |R(q)| \leq M_L^{M_0}, \quad |\mathcal{T}_{\text{rare}}| = |\mathcal{T}| - \Theta(c^{M_0}).$$

1430 In particular, although the total number of tasks grows exponentially in M_0 , the number of common tasks is exponentially
 1431 smaller whenever $M_L \gg c$.

1432
 1433 *Proof.* Proof of Lemma E.4. Each task $k \in \mathcal{T}$ is a function $k : S_1 \rightarrow S_L$, so $|\mathcal{T}| \leq \prod_{q \in S_1} |R(q)|$, where $R(q)$ contains all
 1434 reachable answers (via any CoT) from q . The set $\mathcal{T}_{\text{common}}$ consists of tasks for which $k(q) \in A(q)$ for all q , so
 1435

$$1436 |\mathcal{T}_{\text{common}}| = \prod_{q \in S_1} |A(q)| = \prod_{q \in S_1} n_q = \Theta(c^{M_0}),$$

1437
 1438
 1439 with $c \in [2, \max_q n_q] = O(1)$. The rest follows directly by subtraction. \square
 1440

1441 **Lemma E.5.** Consider a TMC $X = (X_t)_{t \geq 0}$ defined in Def. 2.1, and a specific task $k \in \mathcal{T}$ defined in Def. 2.2. Then, for
 1442 any fixed $q = o_1 \in S_1$ and corresponding correct answer $a = o_L \in S_L$, with non-trivial probability, there exists at least
 1443 one hard-to-reason CoT trajectory (i.e., a path containing at least one sparse edge) from q to a . Specifically, the probability
 1444 of having at least one such hard-to-reason trajectory, denoted $\mathbb{P}_{\text{deg}}(o_L = a | o_1 = q)$, is lower bounded as:
 1445

$$1446 \mathbb{P}_{\text{deg}}(o_L = a | o_1 = q) \geq \Theta(\epsilon \cdot M^{L-3}) \geq c \geq \Theta(M^{-2}) > 0,$$

1447
 1448 for some constant c , where $\epsilon = O(1/M^{L-2})$ is the transition probability of a sparse edge.
 1449

1450 *Proof of Lemma E.5.* Fix $q = o_1$ and $a = o_L$. Let Π be the set of all length- L trajectories $\tau = (o_1, \dots, o_L)$ with $o_L = a$.
 1451 We split $\Pi = \Pi_{\text{nd}} \cup \Pi_{\text{deg}}$ according to whether τ has zero or at least one sparse edge.
 1452

1453 By Def. 2.1, there exist $O(1)$ “easy-to-reason” trajectories from $q \in S_1$ to $a \in S_L$, each consisting entirely of high-
 1454 probability transitions C_{o_i} . Each transition $o_i \rightarrow o_{i+1}$ along these paths has probability $\Theta(1/M)$. Therefore, for a trajectory
 1455 of $L - 1$ steps, the total probability of such a path is:
 1456

$$1457 \mathbb{P}_{\text{high}} = O(1) \cdot \left(\Theta \left(\frac{1}{M} \right) \right)^{L-1} = O \left(M^{-(L-1)} \right).$$

1460 Similarly, as the number of hard-to-reason CoT is below $\Theta(M)$, given that $\mathbb{P}_{\text{sparse}} \leq \Theta(1/M^{L-2})$, we conclude the total
 1461 probability by union bound $\mathbb{P}_{\text{TMC}}(o_L = a | o_1 = q) = \Theta(M^{-(L-1)})$.
 1462
 1463 \square

1464 **Theorem E.6 (Intrinsic Properties of Multi-task TMC).** Let $X = (X_t)_{t \geq 0}$ be a Tree-structured Markov Chain (TMC) and
 1465 \mathcal{T} a set of tasks, per defined in Def. 2.1 and 2.2.
 1466

1467 1. **(Task Interference)** Let tasks $k, k' \in \mathcal{T}$ share at least one question state $q \in S_1$ or answer state $a \in S_L$, with distinct
 1468 valid QA pairs $(q, a_q) \in \mathcal{Q}_k$ and $(q', a_{q'}) \in \mathcal{Q}_{k'}$. Suppose the transition probabilities along edges in $\mathcal{G}_{q, a_q}^{(k)}$ are
 1469 amplified such that the TMC reaches a_q with probability $1 - \delta$ (where $\delta = o(M^{-L}) \ll 1$) via valid CoTs in $\mathcal{G}_{q, a_q}^{(k)}$.
 1470 Then for all shared q or a , every originally easy-to-reason CoT in $\mathcal{G}_{q', a_{q'}}^{(k')}$ must satisfy:
 1471

$$1472 \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) = o(1/M^2) \quad \exists (o_l \rightarrow o_{l+1}) \in \tau \text{ and } \tau \in \mathcal{G}_{q', a_{q'}}^{(k')},$$

1473
 1474 i.e., all such CoTs degenerate into hard-to-reason paths. Similarly, for any task $\hat{k} \neq k \in \mathcal{T}$ whose valid CoT set $\mathcal{G}_{\hat{q}, a_{\hat{q}}}^{(\hat{k})}$
 1475 has at least one easy-to-reason CoT \hat{o} sharing some transitions $\hat{o}_l \rightarrow \hat{o}_{l+1}$ with the CoTs in $\mathcal{G}_{q, a_q}^{(k)}$. Then \hat{o} becomes
 1476 hard-to-reason.
 1477

1479 2. **(Correctness Bottleneck)** Suppose the probability mass of valid hard-to-reason CoTs traveling from q to a_q for task k
 1480 in the original TMC X (\mathbb{P}_{TMC}) is Δ .
 1481

1482 Then suppose a model \hat{p} satisfies:

- 1483 • The total probability mass from q to a is $1 - C$.
- 1484 • The fraction of easy-to-reason CoTs among CoTs traveling from q to a_q is $1 - \epsilon$.

Then the expected correctness over the QA distribution $\mathcal{D}_{a_q}^{q,k}$ is upper bounded by:

$$R_{\text{out}}^k \hat{p}(\mathbf{o}) \leq \Theta((1-C)[(1-\epsilon)\frac{1}{1+\Delta M^{L-1}} + \epsilon\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}])$$

Besides, we denote the pass@K performance of model \hat{p} for task tuple (q, a_q, k) (the probability that at least succeed once among K trials) as $\text{Pass@K}_{q,k}^{\hat{p}}$:

$$\text{Pass@K}_{q,k}^{\hat{p}} := \Pr_{\substack{\{\mathbf{o}^i\}_{i \in [K]} \sim \hat{p}(O|q) \\ (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}}} \left[\bigcup_{i=1}^K \mathbb{1}(\mathbf{o}^i \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)}) \right]. \quad (32)$$

When $C = 0$, $\text{Pass@K}_{q,k}^{\hat{p}}$ is upper bounded by

$$\text{Pass@K}_{q,k}^{\hat{p}} \leq \underbrace{\Theta\left(\left[\left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_q} (1 - (1-\epsilon)^K)\right]\right)}_{\text{upper bound of pass@K of instance that cannot be solved by easy CoTs}} + \underbrace{\Theta\left(\left[\left(1 - \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_q}\right)(1 - \epsilon^K)\right]\right)}_{\text{upper bound of pass@K of instance that can be solved by some easy CoT}}.$$

If

$$\epsilon = o\left(\sqrt[K]{1 - C_{\text{Err}} / \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_q}}\right)$$

for some $C_{\text{Err}} \in (0, \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_q})$, then we have the pass@K performance upper bounded by

$$1 - \Omega(C_{\text{Err}}) = o(1),$$

with **constant error** $\Omega(C_{\text{Err}}) = \Theta(1)$.

When $C = \epsilon = 0$, we have

$$R_{\text{out}}^k \hat{p}(\mathbf{o}) \leq \Theta\left(\frac{1}{1+\Delta M^{L-1}}\right)$$

And the pass@K performance is upper bounded by

$$\text{Pass@K}_{q,k}^{\hat{p}} \leq \Theta\left(1 - \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_q}\right).$$

Proof. Proof of Thm. E.6.

1. non-negligible decay of transitions for other tasks when overfit a target task.

Fix a shared question state $q \in S_1$. By Def. 2.2(ii), the easy-to-reason CoTs in $\mathcal{G}_{q, a_q}^{(k)}$ and $\mathcal{G}_{q, a_{q'}}^{(k')}$ are disjoint. The amplification condition implies:

$$\sum_{\tau \in \mathcal{G}_{q, a_q}^{(k)}} \mathbb{P}(\tau|q) \geq 1 - \delta = 1 - o(M^{-L}).$$

Since $\sum_{\tau \in \text{all CoTs from } q} \mathbb{P}(\tau|q) = 1$, the remaining CoTs (including those in $\mathcal{G}_{q, a_{q'}}^{(k')}$) must satisfy:

$$\sum_{\tau \in \mathcal{G}_{q, a_{q'}}^{(k')}} \mathbb{P}(\tau|q) \leq \delta = o(M^{-L}) \ll 1.$$

For any easy-to-reason CoT $\tau = (q, o_2, \dots, o_L = a_{q'}) \in \mathcal{G}_{q, a_{q'}}^{(k')}$, the original transition probabilities satisfy $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) = \Theta(1/M)$ for all edges. However, since the total probability mass for τ is now $o(M^{-L})$, we have:

$$\prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) = o(M^{-L}).$$

Give $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) \leq \Theta(1/M)$, this forces at least one transition term to decay to $o(1/M^2)$. Otherwise, if any edge retained $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) = \Theta(1/M)$, the product would be $\Theta(M^{-(L-1)})$, contradicting $\mathbb{P}(\tau|q) = o(M^{-L})$.

For a shared answer state $a \in S_L$, as well as \hat{o} in other task sharing some transition with CoTs in $\mathcal{G}_{q,a_q}^{(k)}$, the same logic applies.

2. non-negligible error when only favoring easy-to-learn CoTs.

We have the total mass of the valid CoTs for task k in the original X as

$$Z \geq \underbrace{\Theta\left(\frac{1}{M^{L-1}}\right)}_{Z_{\text{easy}}} + \underbrace{\Delta}_{Z_{\text{hard}}}$$

By Def. 2.2, where the expected correctness over a QA sample is proportion to the CoT's likelihood in the original X , we can combine the components to upper bound the expected correctness:

$$R_{\text{out}}^k \hat{p}(\mathbf{o}) = \mathbb{E}_{(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}} [\mathbb{1}(\mathbf{o} \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)})] \leq \Theta\left(\left(1 - C\right)\left[1 - \epsilon\right] \frac{1}{1 + \Delta M^{L-1}} + \epsilon \frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right).$$

Especially, by the first discussion of division of probability space in Remark E.3, it is direct to deduce that the probability that one specific $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$ cannot be solved by every easy-to-reason CoT is

$$\prod_{\mathbf{o} \in \mathcal{G}_{\sigma_1, a_{\sigma_1}^k}^{(k), \text{easy}}} (1 - p_{\text{acc}}^k(\mathbf{o})) = \Theta\left(\left(1 - \frac{1}{1 + \Delta M^{L-1}}\right)^{n_q}\right) = \Theta\left(\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right)$$

When facing these instances, we have the probability of success to be **at most** ϵ when $C = 0$.

Besides, the probability that $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}$ can be solved by some easy-to-reason CoT is

$$1 - \Theta\left(\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right)$$

When facing these instances, we have the probability of success to be **at most** $1 - \epsilon$ when $C = 0$.

Therefore, collaborating with $n_q = O(1), \forall q \in S_1$, the pass@K performance (the probability that at least succeed once among K trials) is upper bounded by

$$\underbrace{\Theta\left(\left[\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q} (1 - (1 - \epsilon)^K)\right]\right)}_{\text{upper bound of pass@K of instance that cannot be solved by easy CoTs}} + \underbrace{\Theta\left(\left[\left(1 - \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right) (1 - \epsilon^K)\right]\right)}_{\text{upper bound of pass@K of instance that can be solved by some easy CoT}}).$$

If

$$\epsilon = o\left(\sqrt[K]{1 - C_{\text{Err}} / \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}}\right),$$

for some $C_{\text{Err}} \in (0, \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q})$, then we have the pass@K performance upper bounded by $1 - \Omega(C_{\text{Err}})$ with constant error $\Theta(C_{\text{Err}}) = \Theta(1)$.

When $C = \epsilon = 0$, we have $R_{\text{out}}^k(\mathbf{o}) \leq \Theta\left(\frac{1}{1 + \Delta M^{L-1}}\right)$. The pass@K performance is upper bounded by

$$\Theta\left(1 - \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right).$$

□

Lemma E.7 (Formal Version of Prop. 3.2). *Let θ^* be the base model in Eq.(2 that exact predicts the distribution of a Multi-task TMC as in Def. 2.1 and 2.2, fix a common task state tuple (q, a, k) . For any valid easy-to-reason CoT o^{easy} and hard-to-learn CoT o^{hard} that share the states o_l , and deviate at the $l + 1$ layer (i.e., $o_l^{\text{easy}} = o_l^{\text{hard}} = o_l$, $o_{l+1}^{\text{easy}} \neq o_{l+1}^{\text{hard}}$), if the total number of valid hard-to-reason CoTs is bounded by $\Theta(M)$, we have $A_{l+1}^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{easy}}) \geq c_1 > 0$, $A_{l+1}^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{hard}}) \leq -c_2 < 0$, $\forall l \in [L - 1]$ for some constants $c_1, c_2 > 0$.*

Proof of Lemma E.7. First, it is direct to see that given any valid easy-to-learn CoT $o^{\text{easy}} \in \mathcal{G}_{q, a_q}^{(k)}$ and hard-to-learn CoT $o^{\text{hard}} \in \mathcal{G}_{q, a_q}^{(k)}$ for task tuple (q, a, k) , we have

$$\frac{\prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1}^{\text{easy}} | o_l)}{\sum_{o'_{1:L} \in \mathcal{G}_{q, a_q}^{(k)}} \prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o'_{l+1} | o'_l)} / \frac{\prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1}^{\text{hard}} | o_l)}{\sum_{o'_{1:L} \in \mathcal{G}_{q, a_q}^{(k)}} \prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o'_{l+1} | o'_l)} \geq \frac{\Theta(M^{-1})}{o(M^{-1})} > \Omega(M), \quad (33)$$

by Eq.(31).

That is, the expected accuracy of o^{easy} on any instance from $\mathcal{D}_{a_q}^{q, k}$, denoted as $p_{\text{acc}}^k(o^{\text{easy}})$ is larger than the o^{hard} , denoted as $p_{\text{acc}}^k(o^{\text{hard}})$, with a ratio no less than $\Theta(M)$.

Fix l and o_l . Write

$$Q^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{easy}}) = \mathbb{E}_{q=o_l \sim P^k(\mathcal{Q}_k)} [\mathbb{1}(o_L = a_q) p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}, o \sim \hat{p}_{\theta^*}(\cdot | o_l)],$$

and

$$V^{\hat{p}_{\theta^*}, k}(o_l) = \sum_{o_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(o_{l+1} | o_l) Q^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}).$$

By definition and Lemma E.5 there are $\Theta(1)$ length- $(L - l)$ continuations each with probability $\Theta(M^{-(L-l)})$. Hence

$$Q^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{easy}}) = \Theta(M^{-(L-l)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}],$$

Also we see $o_{l+1}^{\text{easy}} \in C_{o_l}$. Then $\Pr[o_{l+1}^{\text{easy}} | o_l] = \Theta(1/M)$, thus

$$V^{\hat{p}_{\theta^*}, k}(o_l) \geq \Theta(M^{-(L-l+1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}].$$

Therefore, we have

$$A_{l+1}^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{easy}}) = Q^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{easy}}) - V^{\hat{p}_{\theta^*}, k}(o_l) \geq \Theta(M^{-(L-l+1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}] > 0$$

$o_{l+1}^{\text{hard}} \notin C_{o_l}$. Then $\Pr[o_{l+1}^{\text{hard}} | o_l] = o(M^{-2})$, and the best possible continuations contribute at most $\Theta(M^{-(L-l-1)})$ each.

Thus

$$\begin{aligned} Q^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{hard}}) &\leq o(M^{-2}) \cdot O(M^{-(L-l-1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{hard}}] \\ &= o(M^{-(L-l+1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{hard}}]. \end{aligned}$$

Therefore

$$\begin{aligned} A_{l+1}^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{hard}}) &= Q_{l+1}^{\hat{p}_{\theta^*}, k}(o_l, o_{l+1}^{\text{hard}}) - V^{\hat{p}_{\theta^*}, k}(o_l) \\ &= o(M^{-(L-l+1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{hard}}] - \Theta(M^{-(L-l+1)}) \mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}] \\ &< 0. \end{aligned}$$

Given that for a chosen TMC $\mathbb{P}_{\text{TMC}}(\cdot | \cdot) = \hat{p}_{\theta^*}(\cdot | \cdot)$, the $\mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{easy}}]$, $\mathbb{E}[p_{\text{acc}}^k(o) \mid o_{l+1} = o_{l+1}^{\text{hard}}]$ are constants.

Therefore, by choosing some positive constants $c_1 = \Theta(M^{-(L+1-l)})$, $c_2 = \Theta(M^{-(L+1-l)})$ to bound the advantages, we complete the proof. \square

Algorithm 1 Pretraining of Foundation Model (Kim et al., 2025)

```

1: set  $\theta^{(0)} = \mathbf{0}, \eta = O(M)$ ,
2:  $T_1 = \tilde{O}(M^2 c^{-2}), T_2 = \tilde{O}(M c^{-2})$ 
3: for  $t = 1, \dots, T_1$  do
4:    $\theta^{(t)} = \theta^{(t-1)} + \eta \nabla \mathbb{E}_{X_0, X_1} [\log \hat{p}_{\theta^{(t-1)}}(X_1 | X_0)]$ 
5: end for
6:  $\theta_{ij}^{(T_1)} \leftarrow -\infty$  if  $\hat{p}_{\theta}(o_{l+1} | o_l)^{(T_1)} < c_{\text{thres}}$  {thresholding}
7: for  $t - T_1 = 1, \dots, T_2$  do
8:    $\theta^{(t)} = \theta^{(t-1)} + \eta \nabla \mathbb{E}_{X_0, X_1} [\log \hat{p}_{\theta^{(t-1)}}(X_1 | X_0)]$ 
9: end for

```

F. Details and Proofs of Pretraining

Following (Kim et al., 2025), we could have the following theorem.

Theorem F.1. *Let $X_0 \sim \text{Unif}(S \setminus S_L)$ and $X_1 \sim \mathbb{P}(\cdot | X_0)$ be random samples from the TMC X in Def. 2.1. Let $X_0 \sim \text{Unif}(S \setminus S_L)$ and $X_1 \sim \mathbb{P}(\cdot | X_0)$ be random samples from the TMC X defined in Def. 2.1. For $i \in S_l$, define:*

$$D_{o_l} = \{o_{l+1} : \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) > 0\}, \quad c = \min_{o_{l+1} \in D_{o_l}} \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) > 0,$$

then the softmax predictor trained via Algorithm 1 satisfies:

1. After T iterations with $\eta = \Theta(M)$, the uniform convergence rate is:

$$\sup_{\substack{l \in [L-1], o_l \in S_l \\ o_{l+1} \in S_{l+1}}} |\hat{p}_{\theta^{(t)}}(o_{l+1} | o_l) - \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l)| \leq \tilde{O} \left(\sqrt{\frac{M}{T}} \right) \quad (34)$$

where \tilde{O} hides $\log(TM c^{-1})$ factors.

2. For threshold $c_{\text{thres}} = \Theta(1)$, after $T_1 = \tilde{\Theta}(M^2 c^{-2})$ steps:

$$\begin{cases} \hat{p}_{\theta}(o_{l+1} | o_l) = 0 & \text{if } \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) = 0 \\ \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) - \tilde{O}(c) \leq \hat{p}_{\theta}(o_{l+1} | o_l) \leq \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l) + \tilde{O}(c) & \text{otherwise} \end{cases} \quad (35)$$

3. Post-thresholding, linear convergence occurs:

$$\sup_{\substack{(o_l, o_{l+1}) \in \\ \text{supp}(\mathbb{P})}} |\hat{p}_{\theta^{(T_1+T)}}(o_{l+1} | o_l) - \mathbb{P}_{\text{TMC}}(o_{l+1} | o_l)| \leq \tilde{O}(e^{-\Omega(c^2 T)}) \quad (36)$$

Remark F.2. The logarithmic factors in \tilde{O} terms explicitly track:

- $\log(T_1) = \log(M^2 c^{-2})$ for thresholding
- $\log(c^{-1})$ for initialization dependence
- $\log M$ for high-probability transition

Since we are considering a vanilla regression setting, the proof is standard following (Kim et al., 2025; Ji & Telgarsky, 2019). For the convenience of readers, we provide the proof here.

Proof. We analyze each part of Thm. F.1 systematically.

Proof of Item 1: Uniform Convergence. Let $\mathcal{E}_{l,l+1} = \{(o_l, o_{l+1}) : o_l \in S_l, o_{l+1} \in S_{l+1}\}$ denote all potential transitions. For each $(o_l, o_{l+1}) \in \mathcal{E}_{l,l+1}$, define the parameter error $\Delta_{o_l, o_{l+1}}^{(t)} = \hat{p}_{\theta^{(t)}}(o_{l+1}|o_l) - \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l)$. For a given state o_l , the cross-entropy loss is:

$$L_{o_l}(\theta) = - \sum_{o_{l+1} \in S_{l+1}} \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) \log \hat{p}_{\theta}(o_{l+1}|o_l)$$

where the model's predicted probability is:

$$\hat{p}_{\theta}(o_{l+1}|o_l) = \frac{e^{\theta_{o_l, o_{l+1}}}}{\sum_{o'_{l+1}} e^{\theta_{o_l, o'_{l+1}}}}$$

Similar to Lemma I.1, the gradient component for parameter $\theta_{o_l, o_{l+1}}$ is:

$$\nabla_{\theta_{o_l, o_{l+1}}} L_{o_l} = - \sum_{o'_{l+1}} \mathbb{P}_{\text{TMC}}(o'_{l+1}|o_l) \nabla_{\theta_{o_l, o_{l+1}}} \log \hat{p}_{\theta}(o'_{l+1}|o_l) \quad (37)$$

$$= - \sum_{o'_{l+1}} \mathbb{P}_{\text{TMC}}(o'_{l+1}|o_l) \frac{\nabla_{\theta_{o_l, o_{l+1}}} \hat{p}_{\theta}(o'_{l+1}|o_l)}{\hat{p}_{\theta}(o'_{l+1}|o_l)} \quad (38)$$

$$= -\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) \frac{\nabla \hat{p}_{\theta}(o_{l+1}|o_l)}{\hat{p}_{\theta}(o_{l+1}|o_l)} - \sum_{o'_{l+1} \neq o_{l+1}} \mathbb{P}_{\text{TMC}}(o'_{l+1}|o_l) \frac{\nabla \hat{p}_{\theta}(o'_{l+1}|o_l)}{\hat{p}_{\theta}(o'_{l+1}|o_l)} \quad (39)$$

Using the softmax derivative property:

$$\nabla_{\theta_{o_l, o_{l+1}}} \hat{p}_{\theta}(o'_{l+1}|o_l) = \begin{cases} \hat{p}_{\theta}(o_{l+1}|o_l)(1 - \hat{p}_{\theta}(o_{l+1}|o_l)) & \text{if } o'_{l+1} = o_{l+1} \\ -\hat{p}_{\theta}(o_{l+1}|o_l)\hat{p}_{\theta}(o'_{l+1}|o_l) & \text{if } o'_{l+1} \neq o_{l+1} \end{cases} \quad (40)$$

Substituting these derivatives yields:

$$\nabla_{\theta_{o_l, o_{l+1}}} L_{o_l} = -\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l)(1 - \hat{p}_{\theta}(o_{l+1}|o_l)) + \sum_{o'_{l+1} \neq o_{l+1}} \mathbb{P}_{\text{TMC}}(o'_{l+1}|o_l) \hat{p}_{\theta}(o_{l+1}|o_l) \quad (41)$$

$$= -\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) + \hat{p}_{\theta}(o_{l+1}|o_l) \underbrace{\sum_{o'_{l+1}} \mathbb{P}_{\text{TMC}}(o'_{l+1}|o_l)}_{=1} \quad (42)$$

$$= \hat{p}_{\theta}(o_{l+1}|o_l) - \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) \quad (43)$$

Then the gradient descent update rule is:

$$\theta_{o_l, o_{l+1}}^{(t)} = \theta_{o_l, o_{l+1}}^{(t-1)} + \eta (\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) - \hat{p}_{\theta^{(t-1)}}(o_{l+1}|o_l))$$

This corresponds to the classical softmax parameter updates. The key challenge lies in the heterogeneous transition probabilities:

- For $o_{l+1} \in C_{o_l}$: $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) = \Theta(1/M)$, with $|C_{o_l}| \leq M$
- For $o_{l+1} \in D_{o_l} \setminus C_{o_l}$: $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) \geq c$ but $o(1/M)$
- For $o_{l+1} \notin D_{o_l}$: $\mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) = 0$

Phase 1 - *High-probability edges*: Let $M_0 = |S_1| = \Theta(M)$. The initial parameters $\theta^{(0)} = 0$ yield uniform distribution:

$$\hat{p}^{(0)}(\mathbf{o}_{l+1}|\mathbf{o}_l) = \frac{1}{|S_{l+1}|} \leq \frac{1}{M_0} = O(1/M)$$

For $\mathbf{o}_{l+1} \in C_{o_l}$, the initial error is $\Theta(1/M) - O(1/M) = \Theta(1/M)$. Each gradient step updates \hat{p} by $\eta \cdot \Theta(1/M)$. To reach ϵ -accuracy for these edges, we need $T \geq \Omega(M^2/\epsilon^2)$.

Phase 2 - *Low-probability edges*: For $\mathbf{o}_{l+1} \in D_{o_l} \setminus C_{o_l}$, the signal-to-noise ratio is weaker. The gradient signal is $\mathbb{P}_{\text{TMC}}(\mathbf{o}_{l+1}|\mathbf{o}_l) - \hat{p} \geq c - O(1/M)$. Using the regret bound for online gradient descent (Theorem 3.1 in Hazan (2023)):

$$\sum_{t=1}^T (\hat{p}^{(t)} - \mathbb{P})^2 \leq O\left(\frac{\log |S_{l+1}|}{\eta} + \eta T c^2\right)$$

Optimizing η yields $T \geq \tilde{\Omega}(M/(c^2\epsilon^2))$ for ϵ -accuracy. Combining both phases via union bound over $O(M)$ edges per layer and $O(L) = O(1)$ layers gives:

$$\sup |\Delta^{(T)}| \leq O\left(\sqrt{\frac{M \log T}{T}} \cdot \log\left(\frac{TM}{c}\right)\right)$$

This matches Equation (34) after constant absorption.

Proof of Item 2: Support Recovery via Thresholding. After $T_1 = \tilde{O}(KM^2c^{-2})$ iterations:

- *Zero-probability edges*: For $\mathbf{o}_{l+1} \notin D_{o_l}$, the true probability $\mathbb{P} = 0$. The empirical estimate satisfies:

$$\hat{p}^{(T_1)}(\mathbf{o}_{l+1}|\mathbf{o}_l) \leq \sqrt{\frac{2 \log(1/\delta)}{T_1}} + O\left(\frac{1}{M}\right)$$

via Azuma's inequality for martingales. Setting $\delta = c_{\text{thres}}c$ and $T_1 \geq \tilde{\Omega}(M^2c^{-2})$ ensures $\hat{p} \leq c_{\text{thres}}c$.

- *Non-zero edges*: From Item 1, for $\mathbf{o}_{l+1} \in D_{o_l}$:

$$|\hat{p}^{(T_1)} - \mathbb{P}| \leq O\left(\sqrt{\frac{M}{T_1}} \log T_1\right) = o(c)$$

Thus $\hat{p}^{(T_1)} \geq \mathbb{P} - o(c) \geq c - o(c) > c_{\text{thres}}c$ for proper $c_{\text{thres}} < 1$.

Thresholding at $c_{\text{thres}}c$ thus exactly recovers the support while maintaining Equation (35).

Proof of Item 3: Linear Convergence. Post-thresholding, the parameter space restricts to D_{o_l} edges. The Hessian of L_{pre} becomes:

$$\nabla^2 L_{\text{pre}}(\boldsymbol{\theta})_{(o_l, o_{l+1}), (o_l, o'_{l+1})} = \text{Cov}(\hat{p}_{\boldsymbol{\theta}}(\mathbf{o}_{l+1}|\mathbf{o}_l), \hat{p}_{\boldsymbol{\theta}}(\mathbf{o}'_{l+1}|\mathbf{o}_l))$$

Under the TMC structure, the Fisher information matrix $I(\boldsymbol{\theta})$ satisfies $\lambda_{\min}(I) \geq \Omega(c^2)$ since all active transitions have probability $\geq c$. By Theorem 4.1 in Ji & Telgarsky (2019), gradient descent on strongly convex objectives achieves:

$$\|\Delta^{(T_1+t)}\|^2 \leq \exp(-\Omega(c^2t)) \|\Delta^{(T_1)}\|^2$$

Given $\|\Delta^{(T_1)}\| = O(\sqrt{M/T_1} \log T_1) = O(\log c^{-1})$ from Item 2, we obtain Equation (36).

□

G. Details and Proofs of RLVR Finetuning

During the gradient update, for any $o_l \in S_l, l \in [L - 1]$ that appears as the transition in the valid CoT set $\cup_{q \in \mathcal{Q}_k} \mathcal{G}_{q, a_q}^{(k)}$ for task $k \in \mathcal{T}$, we define the following notations (summarized in Table 4)

- $\mathcal{I}_{o_{l+1}, o_l}^{(k)} \subseteq \cup_{q \in \mathcal{Q}_k} \mathcal{G}_{q, a_q}^{(k)}$ as the subset of valid CoTs satisfies $\forall o^i \in \mathcal{I}_{o_{l+1}, o_l}^{(k)}, o_l^i = o_l, o_{l+1}^i = o_{l+1}$.
- $\mathcal{S}_{o_l}^{(k)} := \{o_{l+1} \in S_l \mid \mathcal{I}_{o_{l+1}, o_l}^{(k)} \neq \emptyset\} \subseteq S_{l+1}$ as the subset of $l + 1$ -th layer states collecting the states such that for any valid CoT o for task k passing $o_l', o_{l+1} \in \mathcal{S}_{o_l'}^{(k)}$.
- $\mathcal{S}_{o_l}^{(k), \text{easy}} \subseteq \mathcal{S}_{o_l}^{(k)}$ contains the subset of $l + 1$ -th layer's states in $\mathcal{S}_{o_l}^{(k)}$ passed by at least one easy-to-reason CoTs (in the original TMC \mathbb{P}_{TMC}) for task k , and $\mathcal{S}_{o_l}^{(k), \text{hard}} := \mathcal{S}_{o_l}^{(k)} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}$ contains the $l + 1$ -th layer's states only passed by valid hard-to-reason CoTs for task k .
- $\mathcal{G}_{o_1, a_{o_1}^k}^{(k), \text{easy}} \subseteq \mathcal{G}_{o_1, a_{o_1}^k}^{(k)}$ as the subset of valid easy-to-reason CoTs inside the valid CoTs for task k , and $\mathcal{G}_{o_1, a_{o_1}^k}^{(k), \text{hard}} := \mathcal{G}_{o_1, a_{o_1}^k}^{(k)} \setminus \mathcal{G}_{o_1, a_{o_1}^k}^{(k), \text{easy}}$ the subset of valid hard-to-reason CoTs.
- $\theta^{k, (t)}$ be the finetuned model for the task k at post-training iteration t
- For any CoT $o \in \mathcal{G}_{o_1, a_{o_1}^k}^{(k)}$, the probability that o is correct for a sampled instance $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{o_1}}^{o_1, k}$ is given by $p_{\text{acc}}^k(o)$, which, per Condition (iv) in Def. 2.2, is proportional to its likelihood among $\mathcal{G}_{o_1, a_{o_1}^k}^{(k)}$:

$$p_{\text{acc}}^k(o) = \frac{\prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1} \mid o_l)}{\sum_{o_{l+1} \in \mathcal{G}_{o_1, a_{o_1}^k}^{(k)}} \prod_{l=1}^{L-1} \mathbb{P}_{\text{TMC}}(o_{l+1} \mid o_l)},$$

where $\mathbb{P}_{\text{TMC}}(\cdot \mid o_l) = \hat{p}_{\theta^*}(\cdot \mid o_l)$, $o_l \in S_l$, is the transition kernel of our Multi-task TMC.

- The gradient update objectives $\mathcal{J}_{\text{REINFORCE}}(\theta^k)$ and $\mathcal{J}_{\text{RAFT}}(\theta^k)$ represent

$$\begin{aligned} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{o_1}^k}^{o_1, k}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O \mid \mathbf{o}_1)} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right], \\ \mathcal{J}_{\text{RAFT}}(\theta^k) &= \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{o_1}^k}^{o_1, k}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O \mid \mathbf{o}_1)} \left[\sum_{l=1}^{L-1} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right]. \end{aligned} \quad (44)$$

- The objective of RL-rej $\mathcal{J}_{\text{Rein-rej}}(\theta^k)$, in our case, is training using the REINFORCE objective on a online distorted data distribution $\mathcal{D}_{\text{rej}, (t)}^{o_1, k}$, formally

$$\mathcal{J}_{\text{Rein-rej}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{\text{rej}, (t)}^{o_1, k}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O \mid \mathbf{o}_1)} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right]. \quad (45)$$

Here, $\mathcal{D}_{\text{rej}, (t)}^{o_1, k} := \{(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{o_1}^k}^{o_1, k} \mid \Pr_{\mathbf{o} \sim \hat{p}_{\theta^k, (t)}(\cdot \mid \mathbf{o}_1)} [\cup_{i=1}^G \mathbb{1}(\mathbf{o}^i \notin \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)})] = \Theta(1)\}$, where $G = O(1)$ is the (offset) time of parallel experiments. That is, the algorithm rejects samples with $\Pr_{\mathbf{o} \sim \hat{p}_{\theta^k, (t)}(\cdot \mid \mathbf{o}_1)} [\cap_{i=1}^G \mathbb{1}(\mathbf{o}^i \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)})] = \Theta(1)$, which represents instances that the model confidently predicts its correct CoTs of G times in parallel. Therefore, idealistically, instance sampled from $\mathcal{D}_{\text{rej}, (t)}^{o_1, k}$ would have some correct CoTs that is not well-learned by the current model $\theta^{k, (t)}$.

Theorem G.1 (Squeezing Effect and Merits of Rejecting Correct (Full Version of Thm. 3.1 and Cor. 3.5)). *Let θ^* be the base model in Eq.(2 that exact predicts the distribution of a Multi-task TMC as in Def. 2.1 and 2.2, and θ^k the current model to be finetuned from θ^* for task $k \in \mathcal{T}$. Denote the task tuples of task $k \in \mathcal{T}$ as (q, a_q^k, k) , where $a_q^k \in S_L$ is the sole answer state under task k . Assume for each (o_1, a_{o_1}, k) under task k , the number of hard-to-reason CoTs from o_1 to a_{o_1} is bounded by $O(M)$. Let the question distribution during finetuning of task k be $P^k(\mathcal{Q}^k)$ (i.e., $o_1 \sim P^k(\mathcal{Q}^k)$). Then, when finetuning the base model using REINFORCE and RAFT objectives in Eq.(44), we have*

1. **Squeezing Effect & Difference of Logit Update.** For any different state pair $\mathbf{o}'_{l+1} \neq \mathbf{o}_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}$ denoting two $l+1$ -th states in some valid hard-to-reason and easy-to-reason CoT sharing the l -th state \mathbf{o}_l for task k , we have

$$\begin{aligned} \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) > 0, \\ \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) < 0. \\ \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^k} \mathcal{J}_{\text{RAFT}}(\theta^k) > 0, \\ \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^k} \mathcal{J}_{\text{RAFT}}(\theta^k) < 0. \end{aligned} \quad (46)$$

In addition, we have the difference of the logits' update $\Delta h_{\theta^k}(\mathbf{o}'_{l+1}, \mathbf{o}_l) - \Delta h_{\theta^k}(\mathbf{o}_{l+1}, \mathbf{o}_l)$ in Reinforce as:

$$\begin{aligned} \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} - \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} &< \eta [\hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) - \hat{p}_{\theta^k}(\mathbf{o}_{l+1} | \mathbf{o}_l)] \Pr_{\substack{\mathbf{o}'_l \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \\ &\cdot [\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{\mathbf{o}}_l = \mathbf{o}_l \\ \hat{\mathbf{o}}_L = a_{\hat{\mathbf{o}}_1}^k}]] - (\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \\ &\cdot \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{\mathbf{o}}_l = \mathbf{o}_l \\ \tilde{\mathbf{o}}_L = a_{\tilde{\mathbf{o}}_1}^k}])] \\ &< 0. \end{aligned} \quad (47)$$

Also, for RAFT, the difference of logits update $\Delta h_{\theta^k}(\mathbf{o}'_{l+1}, \mathbf{o}_l) - \Delta h_{\theta^k}(\mathbf{o}_{l+1}, \mathbf{o}_l)$ is

$$\begin{aligned} \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} - \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} &< \eta [\hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l)(1 + \log \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l)) - \hat{p}_{\theta^k}(\mathbf{o}_{l+1} | \mathbf{o}_l)(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1} | \mathbf{o}_l))] \\ &\cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{\mathbf{o}}_l = \mathbf{o}_l \\ \hat{\mathbf{o}}_L = a_{\hat{\mathbf{o}}_1}^k}]] - (\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{\mathbf{o}}_l = \mathbf{o}_l \\ \tilde{\mathbf{o}}_L = a_{\tilde{\mathbf{o}}_1}^k}]] \right] \\ &\cdot \Pr_{\substack{\mathbf{o}'_l \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \\ &< 0. \end{aligned} \quad (48)$$

2. Convergence of Finetuning & Constant Error of Pass@K.

For $\forall \epsilon \in (0, 1/2)$, there exists $T \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$, for $t \geq T$, the probability that $\hat{p}_{\theta^{k,(t)}}(\cdot | \mathbf{o}_1)$ (trained by REINFORCE or RAFT) reach the $a_{\mathbf{o}_1}$ is converged:

$$\Pr_{\substack{\mathbf{o}'_l \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^{k,(t)}}(\cdot | \mathbf{o}_1)}} [o'_L = a_{\mathbf{o}'_1}^k] \geq \Pr_{\substack{\mathbf{o}'_l \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^{k,(t)}}(\cdot | \mathbf{o}_1)}} [o'_L = a_{\mathbf{o}'_1}^k, \mathbf{o}'_L \in \mathcal{G}_{\mathbf{o}_1, a_{\mathbf{o}'_1}^k}^{(k), \text{easy}}] \geq 1 - o(\epsilon). \quad (49)$$

Then, suggest the probability mass of valid hard-to-reason CoTs traveling from some $\mathbf{o}_1 \sim P(\mathcal{Q}_k)$ to $a_{\mathbf{o}_1}$ for task k in the original TMC X (\mathbb{P}_{TMC}) is Δ . Then it holds that

$$\text{Rex}_{\mathbf{o}_1, k}^{\hat{p}_{\theta^{k,(t)}}}(\mathbf{o}) \leq \Theta((1 - \epsilon) \frac{1}{1 + \Delta M^{L-1}} + \epsilon \frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}). \quad (50)$$

Further, the pass@K performance $\text{Pass@K}_{q,k}^{\hat{p}} := \Pr_{\substack{\{\mathbf{o}^i\}_{i \in [K]} \sim \hat{p}(\mathcal{O}|q) \\ (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_q}^{q,k}}} [\bigcup_{i=1}^K \mathbb{1}(\mathbf{o}^i \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)})]$ is upper bounded by

$$\text{Pass@K}_{\mathbf{o}_1, k}^{\hat{p}_{\theta^{k,(t)}}} \leq \underbrace{\Theta\left(\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q} (1 - (1 - \epsilon)^K)\right)}_{\text{Solved by hard CoTs}} + \underbrace{\Theta\left(\left(1 - \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right) (1 - \epsilon^K)\right)}_{\text{Solved by some easy CoTs}}. \quad (51)$$

When $\epsilon = o(\sqrt[\kappa]{1 - C_{\text{Err}} / (\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}})^{n_q}}) \rightarrow 0$, the pass@K performance suffer from constant error: $1 - \text{Pass@K}_{\hat{p}_{\theta^k, (t)}} = \Theta(1)$.

3. **Curriculum Learning of RL-rej.** For any $\epsilon \in (0, 1/2)$, suppose setting $G = 1$ in $\mathcal{D}_{\text{rej}, (t)}^{\sigma_1, k}$ (Eq.(45) excludes all (\mathbf{Q}, \mathbf{A}) pairs containing correct CoTs that the current model $\hat{p}_{\theta^k, (t)}$ predicts with non-trivial probability $\Theta((1 - \epsilon)/M)$. Then, optimizing Eq.(45 via RL-rej leads to the following:

(i) The model first learns the easy-to-reason CoTs within $\Theta(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ steps.

(ii) Once its predictive mass over $\mathcal{G}_{\sigma_1, a_{\sigma_1}}^{(k), \text{easy}}$ reaches $\Theta((1 - \epsilon)/M)$, learning begins on sparse edges in $\mathcal{S}_{\sigma_l}^{(k), \text{hard}} = \mathcal{S}_{\sigma_l}^{(k)} \setminus \mathcal{S}_{\sigma_l}^{(k), \text{easy}}$. Hard-to-reason CoTs in $\mathcal{G}_{\sigma_1, a_{\sigma_1}}^{(k), \text{hard}}$ are progressively learned, with those sharing more edges with $\mathcal{G}_{\sigma_1, a_{\sigma_1}}^{(k), \text{easy}}$ being learned earlier.

(iii) Let Δ denote the total probability mass of valid hard-to-reason CoTs from σ_1 to a_{σ_1} in the original TMC X (under \mathbb{P}_{TMC}). Suppose after $T_2 = \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$, there are n'_{σ_1} hard-to-reason CoTs each with likelihood ratio scale $\Theta(\rho) < 1$ in the $\mathcal{G}_{\sigma_1, a_{\sigma_1}}^{(k)}$ have been well-learned with predictive probability $\Theta((1 - \epsilon)/M)$. Then the pass@K is at the scale:

$$\text{Pass@K}_{\hat{p}_{\theta^k, (t)}} = \Theta \left[\underbrace{\left((1 - \rho)^{n'_{\sigma_1}} \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}} \right)^{n_{\sigma_1}} (1 - (1 - \epsilon)^K) \right)}_{\text{instances unsolvable by learned CoTs}} \right] + \Theta \left[\underbrace{\left(\left(1 - (1 - \rho)^{n'_{\sigma_1}} \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}} \right)^{n_{\sigma_1}} \right) (1 - \epsilon^K) \right)}_{\text{instances solvable by some learned CoT}} \right].$$

This bound tends to 1 as $(1 - \rho)^{n'_{\sigma_1}} \rightarrow 0$ and $\epsilon \rightarrow 0$, showing the superiority of RL-rej when the probability mass Δ is non-negligible.

Proof. Proof of Thm. G.1. In our proofs, we first prove the results of REINFORCE, and the results of RAFT follows directly with a more serious of squeezing effect.

Proof of Item 1: Difference of Logit Update.

Recall that by Lemma I.6, we have

$$\begin{aligned} \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\sigma_1 \sim P^k(\mathcal{Q}^k) \\ \{\sigma_{l+1} \sim \hat{p}_{\theta^k}(\cdot | \sigma_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (e_{\sigma_{l+1}, \sigma_l} - \sum_{\sigma'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\sigma'_{l+1} | \sigma_l) e_{\sigma'_{l+1}, \sigma_l}) \right], \\ \nabla_{\theta^k} \mathcal{J}_{\text{RAFT}}(\theta^k) &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\sigma_1 \sim P^k(\mathcal{Q}^k) \\ \{\sigma_{l+1} \sim \hat{p}_{\theta^k}(\cdot | \sigma_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (1 + \log \hat{p}_{\theta^k}(\sigma_{l+1} | \sigma_l)) (e_{\sigma_{l+1}, \sigma_l} - \sum_{\sigma'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\sigma'_{l+1} | \sigma_l) e_{\sigma'_{l+1}, \sigma_l}) \right]. \end{aligned} \quad (52)$$

Per conditions in our item, there are valid easy-to-reason and hard-to-reason CoTs passing σ_l .

Collaborating Eq.(52) with definitions in Item 1 and base model formula in Eq.(2), we have

$$\begin{aligned}
 \nabla_{\theta_{o_{l+1}, o_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_l)\}_{l=1}^{L-1}}} \left[\text{Rex}_{o_l, k}(\mathbf{o}) \cdot \mathbb{1}(o \in \mathcal{I}_{o_{l+1}, o_l}^k) \right. \\
 &\quad \left. - \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^k} \mathbb{1}(o \in \mathcal{I}_{o'_{l+1}, o_l}^k) \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \right] \\
 &= \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_l)\}_{l=1}^{L-1}}} \left[\mathbb{1}(o_L = a_{o_1}^k) p_{\text{acc}}^k(o) \mathbb{1}(o \in \mathcal{I}_{o_{l+1}, o_l}^k) \right. \\
 &\quad \left. - p_{\text{acc}}^k(o) \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^k} \mathbb{1}(o \in \mathcal{I}_{o'_{l+1}, o_l}^k) \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \right] \\
 &= \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_l \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_1)}} \left[\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k \right] \left[\hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{easy}}}] \right. \\
 &\quad \left. - \left(\sum_{o'_{l+1} \in \mathcal{S}_{o_l}^k} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_L = a_{\tilde{o}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}] \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \right) \right] \\
 &= \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_l \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_1)}} \left[\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k \right] \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \left[\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{easy}}}] \right. \\
 &\quad \left. - \left(\sum_{o'_{l+1} \in \mathcal{S}_{o_l}^k} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_L = a_{\tilde{o}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}] \right) \right]
 \end{aligned} \tag{53}$$

where the first equality is by Eq. (123) in Lemma I.6; the second equality is by the definition of $R_{\text{out}}^k(\cdot)$, Condition (iv) in Def. 2.2 and $p_{\text{acc}}^k(o)$ in Eq.(31); the third equality is by the condition in our item that $o_l \in S_l, l \in [L-1]$ appears as the transition in the valid CoT set as well as the definition of $\mathcal{I}_{o_{l+1}, o_l}^k$. Given that for different easy-to-reason CoT sharing o_l , the $\hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l)$ is within the same range, starting from the scale $\Theta(M^{-1})$. Therefore, it is safe to conclude that $\nabla_{\theta_{o_{l+1}, o_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) > 0$.

Similarly, we have

$$\begin{aligned}
 \nabla_{\theta_{o_{l+1}, o_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_l \sim \hat{p}_{\theta^k}(\cdot | \mathbf{o}_1)}} \left[\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k \right] \hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{hard}} | \mathbf{o}_l) \left[\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{hard}}}] \right. \\
 &\quad \left. - \left(\sum_{o'_{l+1} \in \mathcal{S}_{o_l}^k} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_L = a_{\tilde{o}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}] \right) \right]
 \end{aligned} \tag{54}$$

Notably, we have the condition in our item that $o_l \in S_l, l \in [L-1]$ appears as the transition in the valid CoT set. Then, by Eq.(31) as well as the low-probability nature of the sparse edge in Def. 2.1, similar to Eq.(33) we see that

$$\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{easy}}}] > \Omega(M \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{hard}}}]) \tag{55}$$

Therefore, given that $n_q = O(1)$ in Def. 2.1 as well as $\hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \geq \Theta(M^{-1})$, we have

$$\hat{p}_{\theta^k}(\mathbf{o}_{l+1}^{\text{easy}} | \mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{easy}}}] > \Omega(\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}_{l+1}^{\text{hard}}}]).$$

That is, the rightest term in Eq.(54) is strictly lower than 0, making $\nabla_{\theta_{o_{l+1}, o_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) < 0$, a serious squeezing effect such that $\Delta \theta_{o_{l+1}, o_l}^{k, \text{REINFORCE}} := \eta \nabla_{\theta_{o_{l+1}, o_l}^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) < 0$. The proof of RAFT is similar—the only difference in Eq.(52) is the ultra $(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1} | \mathbf{o}_l))$, where the easy-to-reason's value is larger than the hard-to-reason ones due to

the monotonicity of $\log(\cdot)$. Also, noted that $1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \in (1 + \log c, 1)$, which could be scaling as $O(1)$ such that our results of REINFORCE directly applies.

Therefore, it holds that

$$\begin{aligned}
 \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) - \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}' \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \left\{ \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \right. \\
 &\cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{o}_l = \mathbf{o}_l \\ \hat{o}_L = a_{\hat{\mathbf{o}}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right. \\
 &- \left(\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{o}_l = \mathbf{o}_l \\ \tilde{o}_L = a_{\tilde{\mathbf{o}}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right) \\
 &- \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{o}_l = \mathbf{o}_l \\ \hat{o}_L = a_{\hat{\mathbf{o}}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right. \\
 &- \left. \left. \left(\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{o}_l = \mathbf{o}_l \\ \tilde{o}_L = a_{\tilde{\mathbf{o}}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right) \right] \right\} \\
 &= (A^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}'_{l+1}^{\text{hard}}) - A^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}'_{l+1}^{\text{easy}})) \\
 &\quad + V^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l) (\hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{hard}}|\mathbf{o}_l) - \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{easy}}|\mathbf{o}_l)) \\
 &< \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}' \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \\
 &\quad \cdot [\hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{hard}}|\mathbf{o}_l) - \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{easy}}|\mathbf{o}_l)] \\
 &\quad \cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{o}_l = \mathbf{o}_l \\ \hat{o}_L = a_{\hat{\mathbf{o}}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right. \\
 &\quad \left. - \left(\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{o}_l = \mathbf{o}_l \\ \tilde{o}_L = a_{\tilde{\mathbf{o}}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right) \right]
 \end{aligned} \tag{56}$$

Then, during every update, it holds that $\Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} := \eta \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k)$ with η as the step size. Then we have

$$\begin{aligned}
 \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} - \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} &< \eta [\hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{hard}}|\mathbf{o}_l) - \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{easy}}|\mathbf{o}_l)] \cdot \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}' \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \\
 &\cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{o}_l = \mathbf{o}_l \\ \hat{o}_L = a_{\hat{\mathbf{o}}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}'_{l+1}}]] - \left(\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{o}_l = \mathbf{o}_l \\ \tilde{o}_L = a_{\tilde{\mathbf{o}}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right) \right]
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{RAFT}} - \Delta \theta_{\mathbf{o}_{l+1}, \mathbf{o}_l}^{k, \text{REINFORCE}} &< \eta [\hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{hard}}|\mathbf{o}_l) (1 + \log \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{hard}}|\mathbf{o}_l)) - \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{easy}}|\mathbf{o}_l) (1 + \log \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}^{\text{easy}}|\mathbf{o}_l))] \\
 &\cdot \Pr_{\substack{\mathbf{o}'_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}' \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_1)}} [\mathbf{o}'_l = \mathbf{o}_l, \mathbf{o}'_L = a_{\mathbf{o}'_1}^k] \\
 &\cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{\mathbf{o}}) \mid \substack{\hat{o}_l = \mathbf{o}_l \\ \hat{o}_L = a_{\hat{\mathbf{o}}_1}^k \\ \hat{o}_{l+1} = \mathbf{o}'_{l+1}}]] - \left(\sum_{\mathbf{o}'_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{\mathbf{o}}) \mid \substack{\tilde{o}_l = \mathbf{o}_l \\ \tilde{o}_L = a_{\tilde{\mathbf{o}}_1}^k \\ \tilde{o}_{l+1} = \mathbf{o}'_{l+1}}]] \right) \right]
 \end{aligned}$$

Proof of Item 2: Convergence and Failure. Per conditions in our item, there are valid easy-to-reason and hard-to-reason CoTs passing o_l . Recall that $\theta^{k,(0)} = \theta^*$ at the iteration $t = 0$ as the base model to be finetuned, $\theta^{k,(t)}$ be the finetuned model for the task k at post-training iteration t .

For any $w \in (c_w/L^2, o(1/L))$ for some small positive constant $c_w > 0$, we consider the finetuning dynamics during:

$$\sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{hard}}} \hat{p}_{\theta^{k,(t)}}(o'_{l+1}|o_l) \leq w, \quad \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^{k,(t)}}(o'_{l+1}|o_l) \geq 1 - w. \quad (57)$$

Then for REINFORCE we have the lower bound over the $\Delta\theta_{o'_{l+1}, o_l}^{k,(t)}$ for any $o_{l+1}^{\text{easy}} \in \mathcal{S}_{o_l}^{(k)}$ as

$$\begin{aligned} \Delta\theta_{o'_{l+1}, o_l}^{k,(t)} &\geq \eta \Pr_{\substack{o'_1 \sim P^k(\mathcal{Q}^k) \\ o' \sim \hat{p}_{\theta^k}(\cdot|o_1)}} [o'_l = o_l, o'_L = a_{o'_1}^k] \hat{p}_{\theta^k}(o_{l+1}^{\text{easy}}|o_l) \\ &\quad \cdot \left[\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{easy}}}] - \left(\sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k)}} \hat{p}_{\theta^k}(o'_{l+1}|o_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_l = o_l \\ \tilde{o}_L = a_{\tilde{o}_1}^k \\ \tilde{o}_{l+1} = o'_{l+1}}] \right) \right] \\ &\geq \Theta\left(\frac{\eta}{ML-1} \cdot \frac{1}{M} \cdot [\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{easy}}}]w - \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{hard}}}]w\right) \\ &\geq \Theta\left(\frac{\eta}{ML} \cdot \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{easy}}}](w - \frac{1}{M}w)\right) \\ &= \Theta\left(\frac{\eta}{ML} \cdot \frac{(M-1)w}{M} \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{easy}}}]\right) \end{aligned}$$

For a given TMC $\mathbb{P}_{\text{TMC}}(\cdot|o_l) = \hat{p}_{\theta^*}, \forall l \in [L-1]$, the values of

$$\mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{easy}}}], \quad \mathbb{E}[p_{\text{acc}}^k(\hat{o}) \mid \substack{\hat{o}_l = o_l \\ \hat{o}_L = a_{\hat{o}_1}^k \\ \hat{o}_{l+1} = o_{l+1}^{\text{hard}}}]$$

are indeed deterministic positive constants within $(0, 1)$, for any valid CoTs $o_{l+1}^{\text{easy}} \in \mathcal{S}_{o_l}^{(k)}$ and $o_{l+1}^{\text{hard}} \in \mathcal{S}_{o_l}^{(k)}$ passing o_l . That is, we could omit it in $O(1)$:

$$\Delta\theta_{o'_{l+1}, o_l}^{k,(t)} \geq \Theta\left(\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}\right). \quad (58)$$

Similarly, recall that

$$D_{o_l} = \{o_{l+1} : \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) > 0\}, \quad c = \min_{o_{l+1} \in D_{o_l}} \mathbb{P}_{\text{TMC}}(o_{l+1}|o_l) > 0,$$

For neuron $o_{l+1}^{\text{other}} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k)}$, then by similar derivations we have

$$\Delta\theta_{o_{l+1}^{\text{other}}, o_l}^{k,(t)} \leq -\Theta\left(\frac{\eta}{ML} \cdot \frac{(M+1)w}{M}\right) = -\Theta\left(\frac{\eta(M+1)w}{ML+1}\right). \quad (59)$$

Therefore, by choosing $T \geq \Omega(\eta^{-1}L^2M^L \log(ML/\epsilon))$ where $0.5 > \epsilon > wL > 0$ is a small constant, we have

$$\begin{aligned}
 \sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^{k, (t)}}(o_{l+1} | o_l) &= \frac{\sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^{k, (T)}}}{\sum_{o_{l+1} \in D_{o_l}} e^{\theta_{o_{l+1}, o_l}^{k, (T)}}} = \frac{\sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^* + \sum_{t=0}^{T-1} \Delta \theta_{o_{l+1}, o_l}^{k, (t)}}}{\sum_{o_{l+1} \in D_{o_l}} e^{\theta_{o_{l+1}, o_l}^* + \sum_{t=0}^{T-1} \Delta \theta_{o_{l+1}, o_l}^{k, (t)}}} \\
 &\geq \frac{\sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^* + T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} }{\sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^* + T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} + \sum_{o_{l+1} \in C_{o_l} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^*} + \sum_{o_{l+1} \in D_{o_l} \setminus C_{o_l}} e^{\theta_{o_{l+1}, o_l}^*}} \\
 &\geq \Theta\left(\frac{e^{T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} }{e^{T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} + (M-1) + \sum_{o_{l+1} \in D_{o_l} \setminus C_{o_l}} M^{-1}}\right) \\
 &\geq \Theta\left(\frac{e^{T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} }{e^{T[\frac{\eta}{ML} \cdot \frac{(M-1)w}{M}]} + M}\right) = \Theta\left(\frac{e^{T[\frac{\eta w}{ML}]} }{e^{T[\frac{\eta w}{ML}]} + M}\right) \\
 &\geq 1 - o\left(\frac{\epsilon}{L}\right)
 \end{aligned} \tag{60}$$

Here, the first inequality is by the negative updates in Eq.(46) and Eq.(59) as well as the update lower bound in Eq.(58); the second inequality is by dividing $\sum_{o_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} e^{\theta_{o_{l+1}, o_l}^*}$ term, $|\mathcal{S}_{o_l}^{(k), \text{easy}}| \leq n_{o_l} = O(1)$ by Def. 2.1, $M = \max_{l, o_l \in \mathcal{S}_l} |C_{o_l}|$, as well as

$$\frac{\hat{p}_{\theta^{k, (t)}}(o_{l+1} | o_l)}{\hat{p}_{\theta^{k, (t)}}(o'_{l+1} | o_l)} \geq \frac{\mathbb{P}_{\text{TMC}}(o_{l+1} | o_l)}{\mathbb{P}_{\text{TMC}}(o'_{l+1} | o_l)} \geq \Omega(M),$$

for $\forall o_{l+1} \in C_{o_l}, o'_{l+1} \in D_{o_l} \setminus C_{o_l}$; the third inequality is by the condition in our item $|D_{o_l} \setminus C_{o_l}| = O(M)$; the last inequality is by choosing $T \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$.

Then

$$\begin{aligned}
 \Pr_{\substack{o'_1 \sim P^k(\mathcal{Q}^k) \\ o' \sim \hat{p}_{\theta^{k, (t)}}(\cdot | o_1)}} [o'_L = a_{o'_1}^k] &\geq \Pr_{\substack{o'_1 \sim P^k(\mathcal{Q}^k) \\ o' \sim \hat{p}_{\theta^{k, (t)}}(\cdot | o_1)}} [o'_L = a_{o'_1}^k, o' \in \mathcal{G}_{o_1, a_{o'_1}^k}^{(k), \text{easy}}] \\
 &\geq (1 - o(\frac{\epsilon}{L}))^{L-1} = 1 - o(\epsilon).
 \end{aligned} \tag{61}$$

Therefore, suggest the probability mass of valid hard-to-reason CoTs traveling from some $o_1 \sim P(\mathcal{Q}_k)$ to a_{o_1} for task k in the original TMC X (\mathbb{P}_{TMC}) is Δ . Then by Thm. E.6, we have

$$\text{Rex}_{o_1, k}^{\hat{p}_{\theta^{k, (t)}}}(\mathbf{o}) \leq \Theta\left((1 - \epsilon) \frac{1}{1 + \Delta M^{L-1}} + \epsilon \frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right). \tag{62}$$

Also, by Thm. E.6 the pass@K performance (the probability that at least succeed once among K trials) is upper bounded by

$$\text{Pass@K}_{o_1, k}^{\hat{p}_{\theta^{k, (t)}}} \leq \underbrace{\Theta\left(\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q} (1 - (1 - \epsilon)^K)\right)}_{\text{upper bound of pass@K of instance that cannot be solved by easy CoTs}} + \underbrace{\Theta\left(\left(1 - \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right)(1 - \epsilon^K)\right)}_{\text{upper bound of pass@K of instance that can be solved by some easy CoT}}.$$

Also, by Thm. E.6 we see that, when $\epsilon = o\left(\sqrt[κ]{1 - C_{\text{Err}} / \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}}\right) \rightarrow 0$, the pass@K performance would suffer from constant error.

The proof of RAFT is similar—the only difference in Eq.(52) is the ultra $(1 + \log \hat{p}_{\theta^k}(o_{l+1} | o_l))$, where the easy-to-reason's value is larger than the hard-to-reason ones due to the monotonicity of $\log(\cdot)$. Also, noted that $1 + \log \hat{p}_{\theta^k}(o_{l+1} | o_l) \in (1 + \log c, 1)$, which could be scaling as $O(1)$ such that our results of REINFORCE directly applies.

Proof of Item 3: Curriculum Learning of RL-rej.

Per Remark E.3, we see that in our case, after $T_1 \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ for a small ϵ , RL-rej would learn valid easy-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}^k}^{(k), \text{easy}}$ well with non-trivial predictive probability $\Theta((1 - \epsilon)/M)$, and start to reject the $(\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_a^{q, k}$ that can be solved by those CoTs.

That is, there exists $T_1 = \Omega(\eta^{-1}L^2M^L \log(ML/\epsilon))$, for $t \in (0, T_1]$, the RL-rej behaves exactly the same with REINFORCE. After $t \geq T_1$ we have

$$\nabla_{\theta_{o_{l+1}^{\text{easy}}, o_l}^k} \mathcal{J}_{\text{Rein-rej}}(\theta^k) = \Pr_{\substack{o'_l \sim P^k(\mathcal{Q}^k) \\ o'_L = a_{o'_l}^k \\ o' \sim \hat{p}_{\theta^k}(\cdot | o_l)}} [o'_l = o_l, o'_L = a_{o'_l}^k] \hat{p}_{\theta^k}(o_{l+1}^{\text{easy}} | o_l) \left[0 - \left(\sum_{o'_{l+1} \in \sum_{o'_l \in \mathcal{S}_{o'_l}^{(k), \text{hard}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_l = o_l \\ \tilde{o}_L = a_{o'_l}^k \\ \tilde{o}_{l+1} = o'_{l+1}}]}] \right) \right] < 0, \quad (63)$$

for all $o_{l+1}^{\text{easy}} \in \mathcal{S}_{o_l}^{(k), \text{easy}}, \forall l \in [L-1]$, where the correctly predicted easy-to-reason CoTs are rejected according to the condition in our item. Similarly, for $o_{l+1}^{\text{hard}} \in \mathcal{S}_{o_l}^{(k), \text{hard}} = \mathcal{S}_{o_l}^{(k)} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}, \forall l \in [L-1]$ we have

$$\nabla_{\theta_{o_{l+1}^{\text{hard}}, o_l}^k} \mathcal{J}_{\text{Rein-rej}}(\theta^k) = \Pr_{\substack{o'_l \sim P^k(\mathcal{Q}^k) \\ o'_L = a_{o'_l}^k \\ o' \sim \hat{p}_{\theta^k}(\cdot | o_l)}} [o'_l = o_l, o'_L = a_{o'_l}^k] \hat{p}_{\theta^k}(o_{l+1}^{\text{hard}} | o_l) \left[\mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_l = o_l \\ \tilde{o}_L = a_{o'_l}^k \\ \tilde{o}_{l+1} = o_{l+1}^{\text{hard}}}] - \left(\sum_{o'_{l+1} \in \sum_{o'_l \in \mathcal{S}_{o'_l}^{(k), \text{hard}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l) \mathbb{E}[p_{\text{acc}}^k(\tilde{o}) \mid \substack{\tilde{o}_l = o_l \\ \tilde{o}_L = a_{o'_l}^k \\ \tilde{o}_{l+1} = o'_{l+1}}]}] \right) \right] > 0, \quad (64)$$

where the inequality is by the feeble $\hat{p}_{\theta^k}(o'_{l+1} | o_l) = o(1/M)$, $o'_{l+1} \in \sum_{o'_l \in \mathcal{S}_{o'_l}^{(k), \text{hard}}}$.

Therefore, we have

$$\Delta \theta_{o_{l+1}^{\text{easy}}, o_l}^{k, \text{Rein-rej}} := \eta \nabla_{\theta_{o_{l+1}^{\text{easy}}, o_l}^k} \mathcal{J}_{\text{Rein-rej}}(\theta^k) < 0, \quad \Delta \theta_{o_{l+1}^{\text{hard}}, o_l}^{k, \text{Rein-rej}} := \eta \nabla_{\theta_{o_{l+1}^{\text{hard}}, o_l}^k} \mathcal{J}_{\text{Rein-rej}}(\theta^k) > 0,$$

and thus the $\hat{p}_{\theta^{k, (t)}}(o_{l+1}^{\text{hard}} | o_l) / \hat{p}_{\theta^{k, (t)}}(o_{l+1}^{\text{easy}} | o_l)$ would strictly increase.

By Eq.(60), similarly there exists $T_2 = T_1 + \Theta(\eta^{-1}L^2M^L \log(ML/\epsilon))$, the predictive probability $\hat{p}_{\theta^{k, (T_2)}}$ of some valid hard-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}}^{(k), \text{hard}}$ would reach $\Theta((1-\epsilon)/M)$. Simultaneously, we see that by adjusting the learning rate to be appropriately small, the predictive probability of easy-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}}^{(k), \text{easy}}$ would not decay below the scale $\Theta((1-\epsilon)/M)$, other wise it would be re-collected in the $\mathcal{D}_{\text{rej}, (t)}^{o_1, k}$ at that iteration t according to our condition setting in item 3. That is, the model $\hat{p}_{\theta^{k, (t)}}$ would gradually increase the predictive probability of sparse edges in $\mathcal{S}_{o_l}^{(k), \text{hard}} = \mathcal{S}_{o_l}^{(k)} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}$ for $\forall l \in [L-1]$, and thus the CoTs in $\mathcal{G}_{o_1, a_{o_1}}^{(k), \text{hard}}$ sharing the most common edges with some CoTs in $\mathcal{G}_{o_1, a_{o_1}}^{(k), \text{easy}}$, would first be learned. Afterwards, more and more hard-to-reason CoTs in $\mathcal{G}_{o_1, a_{o_1}}^{(k), \text{hard}}$ is getting learned, until the point where further learning a new sparse edge in some $\mathcal{S}_{o_l}^{(k), \text{hard}}, l \in [L-1]$ will make another already learned CoT's predictive probability to be lower than the scale $o(1-\epsilon)$.

Suppose the probability mass of valid hard-to-reason CoTs traveling from o_1 to a_{o_1} for task k in the original TMC X (\mathbb{P}_{TMC}) is Δ . Suggest after $T_2 = \Omega(\eta^{-1}L^2M^L \log(ML/\epsilon))$, there are n'_{o_1} hard-to-reason CoTs each with likelihood ratio scale $\Theta(\rho) < 1$ in the $\mathcal{G}_{o_1, a_{o_1}}^{(k)}$ have been well-learned with predictive probability $\Theta((1-\epsilon)/M)$, then, similar to Thm. E.6, we have

$$\text{Pass@K}^{\hat{p}_{\theta^{k, (t)}}} = 4 \underbrace{\Theta\left(\left[(1-\rho)^{n'_{o_1}} \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_{o_1}} (1-(1-\epsilon)^K)\right]\right)}_{\text{upper bound of pass@K of instance that cannot be solved by easy CoTs}} + \underbrace{\Theta\left(\left[(1-(1-\rho)^{n'_{o_1}} \left(\frac{\Delta M^{L-1}}{1+\Delta M^{L-1}}\right)^{n_{o_1}})(1-\epsilon^K)\right]\right)}_{\text{upper bound of pass@K of instance that can be solved by some easy CoT}}$$

which would tends to 1 when $(1-\rho)^{n'_{o_1}} \rightarrow 0, \epsilon \rightarrow 0$.

RL-rej with algorithms other than REINFORCE directly follows. □

Theorem G.2 (Advantage-based Finetuning Favors Easy-to-Reason CoTs (Formal Version of Thm. 3.1)). *Let θ^* be the base model in Eq.(2 that exact predicts the distribution of a Multi-task TMC as in Def. 2.1 and 2.2, and θ^k the current model to be finetuned from θ^* for task $k \in \mathcal{T}$. Denote the task tuples of task $k \in \mathcal{T}$ as (q, a_q^k, k) , where $a_q^k \in \mathcal{S}_L$ is the sole answer state under task k . Assume for each (o_1, a_{o_1}, k) under task k , the number of hard-to-reason CoTs from o_1 to a_{o_1} is bounded by $O(M)$. Let the question distribution during finetuning of task k be $P^k(\mathcal{Q}^k)$ (i.e., $o_1 \sim P^k(\mathcal{Q}^k)$). Suppose the estimates of the RL advantage of PPO / GRPO (without the KL term) by some outer oracle or group-level normalization are accurate during the finetuning: $A_{l+1}^{\hat{p}_{\theta^k, k}}, \hat{A}_{i, l+1}^k = A_{l+1}^{\hat{p}_{\theta^k, k}}(o_l, o_{l+1})$ for any CoT o , and the \hat{p}_{old}^k is appropriately chosen that*

clip operation is *always functioning* starting from the finetuning with $\epsilon_{\text{clip}} = o(1)$ such that

$$\begin{aligned} (1 + \epsilon_{\text{clip}})A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) &\leq \frac{\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}|\mathbf{o}_l)}A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}), \text{ if } A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0, \\ (1 - \epsilon_{\text{clip}})A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) &\leq \frac{\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}|\mathbf{o}_l)}A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}), \text{ if } A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \leq 0, \end{aligned} \quad (65)$$

Then the shared form of objective as

$$\mathcal{J}_{\text{PO}} = \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{\mathbf{a}_{\mathbf{o}_1}^{1, k}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}}(\mathcal{O}|\mathbf{o}_1)} \left[\frac{1}{L} \sum_{l=1}^{L-1} (1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}})A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \right], \quad (66)$$

where $\epsilon_{\text{clip}} > 0$ is a offset clipping parameter, and by Eq.(19),

$$A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) = \mathbb{E}_{\mathbf{o}_1 = q \sim P^k(\mathcal{Q}_k), \mathbf{o}_{l+2:L} \sim \hat{p}_{\theta}} [R_{\text{out}}^k(\mathbf{o})|\mathbf{o}_{l+1}] - \mathbb{E}_{\mathbf{o}_1 = q \sim P^k(\mathcal{Q}_k), \mathbf{o}_{l+1:L} \sim \hat{p}_{\theta}} [R_{\text{out}}^k(\mathbf{o})|\mathbf{o}_l]. \quad (67)$$

Favor Easy CoTs. For any different state pair $\mathbf{o}'_{l+1} \neq \mathbf{o}_{l+1} \in \mathcal{S}_{\mathbf{o}_l}^{(k)}$ denoting two $l+1$ -th states in some valid hard-to-reason and easy-to-reason CoT sharing the l -th state \mathbf{o}_l for task k , it holds that

$$\begin{aligned} \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{PO}, \text{easy}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{\text{easy}}} \mathcal{J}_{\text{PO}}(\theta^k) > 0, \quad \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{PO}, \text{hard}} := \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{\text{hard}}} \mathcal{J}_{\text{PO}}(\theta^k) < 0, \\ \Delta \theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^{k, \text{PO}} &:= \eta \nabla_{\theta_{\mathbf{o}'_{l+1}, \mathbf{o}_l}^k} \mathcal{J}_{\text{PO}}(\theta^k) < 0, \end{aligned} \quad (68)$$

for $\forall \mathbf{o}'_{l+1} \in D_{\mathbf{o}_l} \setminus \mathcal{S}_{\mathbf{o}_l}^{(k)}$.

There exists $T \geq \Omega(\eta^{-1}L^2M^L \log(ML/\epsilon))$, for $t \geq T$, the probability that $\hat{p}_{\theta^{k, (t)}}(\cdot|\mathbf{o}_1)$ reach the $\mathbf{a}_{\mathbf{o}_1}$ is converged:

$$\Pr_{\substack{\mathbf{o}'_L \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^{k, (t)}}(\cdot|\mathbf{o}_1)}} [o'_L = a_{\mathbf{o}_1}^k] \geq \Pr_{\substack{\mathbf{o}'_L \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}'_L \sim \hat{p}_{\theta^{k, (t)}}(\cdot|\mathbf{o}_1)}} [o'_L = a_{\mathbf{o}_1}^k, \mathbf{o}' \in \mathcal{G}_{\mathbf{o}_1, a_{\mathbf{o}_1}^k}^{(k), \text{easy}}] \geq 1 - o(\epsilon). \quad (69)$$

Further, the **pass@K performance** $\text{Pass@K}_{q, k}^{\hat{p}} := \Pr_{\substack{\{\mathbf{o}^i\}_{i \in [K]} \sim \hat{p}(\mathcal{O}|\mathbf{o}_l) \\ (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{\mathbf{a}_q}^{q, k}}} [\bigcup_{i=1}^K \mathbb{1}(\mathbf{o}^i \in \mathcal{G}_{\mathbf{Q}, \mathbf{A}}^{(k)})]$ is upper bounded by

$$\text{Pass@K}_{\mathbf{o}_1, k}^{\hat{p}_{\theta^{k, (t)}}} \leq \underbrace{\Theta\left(\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q} (1 - (1 - \epsilon)^K)\right)}_{\text{Solved by hard CoTs}} + \underbrace{\Theta\left(\left(1 - \left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}\right)(1 - \epsilon^K)\right)}_{\text{Solved by some easy CoTs}}. \quad (70)$$

When $\epsilon = o\left(\sqrt[1]{1 - C_{\text{Err}}/\left(\frac{\Delta M^{L-1}}{1 + \Delta M^{L-1}}\right)^{n_q}}\right) \rightarrow 0$, the pass@K performance suffer from constant error: $1 - \text{Pass@K}_{\mathbf{o}_1, k}^{\hat{p}_{\theta^{k, (t)}}} = \Theta(1)$.

Proof. By Prop. 3.2, for each l we have

$$A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{easy}}) = a_{\text{easy}}^l \geq \Theta(M^{-(L+1-l)}) > 0, \quad A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{hard}}) = -a_{\text{hard}}^l \leq -\Theta(M^{-(L+1-l)}) < 0, \quad (71)$$

for constants $a_{\text{easy}}, a_{\text{hard}} > 0$.

Also, for $\forall \mathbf{o}'_{l+1} \in D_{\mathbf{o}_l} \setminus \mathcal{S}_{\mathbf{o}_l}^{(k)}$, by definition it directly holds that

$$A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}'_{l+1}) < A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}^{\text{hard}}) = -a_{\text{hard}}^l < -\Theta(M^{-(L+1-l)}) < 0. \quad (72)$$

Therefore, by Lemma I.6 as well as the property of TMC, it holds that

$$\nabla_{\theta^k} \mathcal{J}_{\text{PO}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[(1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}})A_{l+1}^{\hat{p}_{\theta^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot (e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in D_{\mathbf{o}_l}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l)e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}) \right].$$

Therefore, collaborating with Eq.(71) and Eq.(72), for $t = 0$ where $\hat{p}_{\theta^k, (0)} = \hat{p}_{\theta^*}$, we directly have

$$\begin{aligned} \Delta \theta_{o'_{l+1}, o_l}^{k, \text{PO, easy}} &:= \eta \nabla_{\theta_{o'_{l+1}, o_l}^k} \mathcal{J}_{\text{PO}}(\theta^k) > 0, & \Delta \theta_{o'_{l+1}, o_l}^{k, \text{PO, hard}} &:= \eta \nabla_{\theta_{o'_{l+1}, o_l}^k} \mathcal{J}_{\text{PO}}(\theta^k) < 0, \\ \Delta \theta_{o'_{l+1}, o_l}^{k, \text{PO}} &:= \eta \nabla_{\theta_{o'_{l+1}, o_l}^k} \mathcal{J}_{\text{PO}}(\theta^k) < 0, \end{aligned} \quad (73)$$

for $\forall o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k)}$. Indeed, following the proof strategies in Lemma I.6, we directly see that when the transitions of $o_l \rightarrow o_{l+1}^{\text{easy}}$ is further strengthened and the transitions of $o_l \rightarrow o_{l+1}^{\text{hard}}$ is further weaken, the $A_{l+1}^{\hat{p}_{\theta^k, (t)}, k}(o_l, o_{l+1}^{\text{easy}})$ is strictly increasing along the iterations, and $A_{l+1}^{\hat{p}_{\theta^k, (t)}, k}(o_l, o_{l+1}^{\text{hard}})$, $A_{l+1}^{\hat{p}_{\theta^k, (t)}, k}(o_l, o'_{l+1})$, $\forall o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k)}$ is strictly decreasing. This makes Eq.(71), Eq.(72) and Eq.(68) hold during the finetuning iterations.

Specifically, for any different state pair $o_{l+1}^{\text{hard}} \neq o_{l+1}^{\text{easy}} \in \mathcal{S}_{o_l}^{(k)}$ and $\forall o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k)}$, it holds that

$$\begin{aligned} \Delta \theta_{o'_{l+1}, o_l}^{k, (t), \text{easy}} &\geq \Theta(\eta M^{-(L+1-l)}(1 - \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l) \\ &\quad + \sum_{o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l))) > 0, \\ \Delta \theta_{o'_{l+1}, o_l}^{k, (t), \text{hard}} &\leq -\Theta(\eta M^{-(L+1-l)}(1 + \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l) \\ &\quad - \sum_{o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l))) < 0, \\ \Delta \theta_{o'_{l+1}, o_l}^{k, (t)} &\leq -\Theta(\eta M^{-(L+1-l)}(1 + \sum_{o'_{l+1} \in \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l) \\ &\quad - \sum_{o'_{l+1} \in D_{o_l} \setminus \mathcal{S}_{o_l}^{(k), \text{easy}}} \hat{p}_{\theta^k}(o'_{l+1} | o_l))) < 0, \end{aligned}$$

where the inequalities is by Eq.(71), Eq.(72), as well as $\epsilon_{\text{clip}} = o(1)$.

Similar to the techniques in Thm. G.1, given that $M^{-(L+1-l)} > M^{-L+1}$ and $p_{\text{acc}}^k \leq 1$, after $T \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ iterations, the remaining proofs and results follows as in Thm. G.1.

□

Remark G.3. To simplify the discussion of the policy gradient case and avoid the non-convexity of $\min\{\cdot\}$, we assume the clip operation with $\epsilon_{\text{clip}} = o(1)$ and Eq. (65). However, our results still hold without this assumption.

Specifically, when the min does not select the clipped term, we instead encounter:

$$\begin{aligned} \nabla_{\theta^k} \left[\frac{\hat{p}_{\theta^k}(o_{l+1} | o_l)}{\hat{p}_{\text{old}}(o_{l+1} | o_l)} \hat{p}_{\theta^k}(o_{l+1} | o_l) \right] &= 2 \frac{\hat{p}_{\theta^k}(o_{l+1} | o_l)}{\hat{p}_{\text{old}}(o_{l+1} | o_l)} \nabla_{\theta^k} \hat{p}_{\theta^k}(o_{l+1} | o_l) \\ &= 2 \frac{\hat{p}_{\theta^k}(o_{l+1} | o_l)^2}{\hat{p}_{\text{old}}(o_{l+1} | o_l)} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(o_{l+1} | o_l) \\ &= \mathbb{E} \left[2 \frac{\hat{p}_{\theta^k}(o_{l+1} | o_l)}{\hat{p}_{\text{old}}(o_{l+1} | o_l)} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(o_{l+1} | o_l) \right], \end{aligned} \quad (74)$$

instead of

$$\begin{aligned} \nabla_{\theta^k} [(1 \pm \epsilon_{\text{clip}}) \hat{p}_{\theta^k}(o_{l+1} | o_l)] &= (1 \pm \epsilon_{\text{clip}}) \hat{p}_{\theta^k}(o_{l+1} | o_l) \nabla_{\theta^k} \log \hat{p}_{\theta^k}(o_{l+1} | o_l) \\ &= (1 \pm \epsilon_{\text{clip}}) \mathbb{E} [\log \hat{p}_{\theta^k}(o_{l+1} | o_l)]. \end{aligned}$$

Since clearly $\hat{p}_{\theta^k}(o_{l+1}^{\text{easy}} | o_l) > \hat{p}_{\theta^k}(o_{l+1}^{\text{hard}} | o_l)$, Eq. (74) shows that the gradient magnitude for easy edges dominates that of sparse ones. Thus, the squeezing effect persists even without the assumption. We adopt the assumption in our theorem purely to reduce discussion complexity.

Lemma G.4. [Detailed Version of Lemma 3.4] Let θ^* be the base model in Eq.(2 that exact predicts the distribution of a Multi-task TMC as in Def. 2.1 and 2.2, and θ^k the current model to be finetuned from θ^* for task $k \in \mathcal{T}$. Suppose the estimates of RL advantage by GRPO through group-level normalization is accurate as $A_{l+1}^{\hat{p}_{\theta^*,k}}(\mathbf{o}_l, \mathbf{o}_{l+1})$ for any CoT o . The optimal step-wise sampling distribution of the KL-regularized GRPO objective in Eq.(6) is:

$$\hat{p}_{\theta^k}^{\text{PO}}(\mathbf{o}_{l+1}|\mathbf{o}_l) \propto \hat{p}_{\theta^*}(\mathbf{o}_{l+1}|\mathbf{o}_l) \exp\left(\hat{r} \frac{A_{l+1}^{\hat{p}_{\theta^*,k}}(\mathbf{o}_l, \mathbf{o}_{l+1})}{\beta}\right), \quad (75)$$

where $\hat{r} \leq \max\{1 + \epsilon_{\text{clip}}, c^{-1}, \Theta(M)\}$.

Proof. This result is standard in RL and distribution optimization literature (Ziebart, 2008; Levine, 2018; Foster et al., 2025; Kawata et al., 2025; Fan et al., 2023; Black et al., 2024; Clark et al., 2024; Uehara et al., 2024). The proofs mirror the proof of Corollary 4.3 in Sec. H, and we therefore omit their full proofs for brevity. \square

Corollary G.5 (Full Version of Corollary 3.5). Let θ^* be the base model in Eq.(2 that exactly predicts the distribution of a Multi-task TMC as in Defs. 2.1 and 2.2. For any target task $k \in \mathcal{T}$, consider the following two categories of instances:

1. Instances $(\mathcal{Q}, \mathbf{A})$ whose correct CoTs only lie in $\mathcal{G}_{q, \alpha_q^k}^{(k), \text{hard}}$.
2. Instances $(\mathcal{Q}, \mathbf{A})$ sampled from another task $k' \neq k$.

For PPO/GRPO without KL regularization that satisfy the conditions in Thm. G.2, the pass@K upper bound for these instances after $T \geq \Omega(\eta^{-1} L^2 M^L \log(ML/\epsilon))$ is $(1 - (1 - \epsilon)^K)$.

In contrast, for the optimal sampler $\hat{p}_{\theta^k}^{\text{PO}}$ in Eq. (75), for any ϵ' satisfying $1/N_{\mathbf{o}_1} > \epsilon' \geq \epsilon > 0$, denote $\hat{p}_{\theta^{k,(\epsilon)}}^k$ as the PPO/GRPO in Thm. 3.1 with ϵ , if

$$\beta > \frac{2\hat{r}(L-1)}{\ln\left(\frac{1}{\epsilon' \prod_{l=1}^{L-1} |D_{\mathbf{o}_l}|}\right)},$$

then the pass@K performance of $\hat{p}_{\theta^k}^{\text{PO}}$ is strictly better than that of PPO/GRPO without KL regularization under the same conditions:

1. **Capable of Hard CoTs:** For instance $(\mathcal{Q}, \mathbf{A})$ with only some hard-to-reason CoTs correct:

$$\mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^{\text{PO}}(\cdot|\mathbf{o}_1)} \left[R_{(\mathcal{Q}, \mathbf{A})}^k(\mathbf{o}) \right] \geq \epsilon' \geq \epsilon \geq \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^{k,(\epsilon)}}^k(\cdot|\mathbf{o}_1)} \left[R_{(\mathcal{Q}, \mathbf{A})}^k(\mathbf{o}) \right].$$

2. **Preserve Multi-task:** For instance $(\mathcal{Q}, \mathbf{A})$ belonging to untargeted task $k' \neq k$:

$$\mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^{\text{PO}}(\cdot|\mathbf{o}_1)} \left[R_{(\mathcal{Q}, \mathbf{A})}^{k'}(\mathbf{o}) \right] \geq \epsilon' \geq \epsilon \geq \mathbb{E}_{\mathbf{o}_{2:L} \sim \hat{p}_{\theta^{k,(\epsilon)}}^k(\cdot|\mathbf{o}_1)} \left[R_{(\mathcal{Q}, \mathbf{A})}^{k'}(\mathbf{o}) \right].$$

Proof. It suffices to prove that with a large β , any non-zero transition within the TMC is larger than $\epsilon' < N_{\mathbf{o}_1}^{-1} := |D_{\mathbf{o}_1}|$.

By the definition of the advantage function in Eq.(19, we have

$$-1 \leq A_{l+1}^{\hat{p}_{\theta^*,k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \leq 1.$$

Therefore, from Eq.(75, the minimum sampling probability over any edge in $D_{\mathbf{o}_l}$ is

$$\hat{p}_{\theta^k}^{\text{PO}}(\mathbf{o}_{l+1}|\mathbf{o}_l) \geq \frac{e^{-\frac{\hat{r}}{\beta}}}{|D_{\mathbf{o}_l}| e^{\frac{\hat{r}}{\beta}}} = \frac{e^{-\frac{2\hat{r}}{\beta}}}{|D_{\mathbf{o}_l}|}.$$

Hence, for any trajectory of length L , the probability of sampling a specific terminal state \mathbf{o}_L from any starting state \mathbf{o}_1 whose \mathbf{o}_{l+1} transitions are in $D_{\mathbf{o}_l}$ is lower bounded by

$$\prod_{l=1}^{L-1} \frac{e^{-\frac{2\hat{r}}{\beta}}}{|D_{\mathbf{o}_l}|} = e^{-\frac{2\hat{r}(L-1)}{\beta}} \cdot \prod_{l=1}^{L-1} \frac{1}{|D_{\mathbf{o}_l}|}.$$

Define $C := \prod_{l=1}^{L-1} \frac{1}{|D_{o_l}|}$. We seek the condition on β such that this probability is at least ϵ' , i.e.,

$$C \cdot e^{-\frac{2\hat{r}(L-1)}{\beta}} \geq \epsilon.$$

Dividing both sides by C and taking logarithms yields

$$-\frac{2\hat{r}(L-1)}{\beta} \geq \ln\left(\frac{\epsilon'}{C}\right), \quad \text{so} \quad \beta \geq \frac{2\hat{r}(L-1)}{\ln\left(\frac{1}{\epsilon' C}\right)}.$$

Substituting $C = \prod_{l=1}^{L-1} \frac{1}{|D_{o_l}|}$, we obtain the desired bound:

$$\beta > \frac{2\hat{r}(L-1)}{\ln\left(\frac{1}{\epsilon' \prod_{l=1}^{L-1} |D_{o_l}|}\right)}.$$

That is, the probability of the path $(o_1, o_2, \dots, o_L), o_{l+1} \in D_{o_l}, \forall l \in [L-1]$ is larger than ϵ' . This ensure that the model is more capable of sampling valid hard-to-reason CoTs for current task as well as valid CoTs for other tasks, as long as the path with transition probability larger than zero ($c > 0$) in Def. 2.1. \square

H. Details and Proofs of Reward-based Sampling

Lemma H.1 (BoN/BS with Ground-true Signal Oracle). *Let θ^* be the base model in Eq.(2 that exactly predicts the distribution of a Multi-task TMC as defined in Definitions 2.1 and 2.2. Under task tuple (q, a, k) , consider the ORMs $R_{\mathbf{Q}, \mathbf{A}}^k(\cdot)$, as well as the PRM given in Eqs. 10. For any target task $k \in \mathcal{T}$ and instance distribution $\mathcal{D}_{a_q}^{a, k}$, if the total number of valid hard-to-reason CoTs is $\Theta(M)$, then during pass@K sampling:*

- ORM/PRM-based BoN or BS achieves success probability $\Theta(1)$ on task k ;
- ORM/PRM-based BoN or BS fails on any other task $k' \neq k$.

Proof. Consider task tuple (q, a, k) and an instance $(\mathbf{Q}, \mathbf{A}) \in \mathcal{D}_{a_q}^{a, k}$ that is solvable, i.e., it admits at least one valid CoT in $\mathcal{G}_{q, a_q}^{(k)}$. Since the base model θ^* assigns $\Theta(c^{L-1})$ sampling probability to a correct CoT, the success probability of ORM-based BoN using the ground-truth reward $R_{\mathbf{Q}, \mathbf{A}}^k(\cdot)$ satisfies:

$$\text{pass@K} = \Theta\left(1 - (1 - c^{L-1})^{NK}\right) = \Theta(1),$$

where the final equality holds for sufficiently large K .

For ORM-based BoN under outcome-population reward $R_{\text{out}}^k(\cdot)$, the CoT credit depends on relative likelihood. Consider the worst case where there is exactly one correct CoT with success probability $\Theta(c^{L-1})$, while each incorrect but valid CoT has sampling probability $\Theta(1/M^{L-1})$ (by Lemma E.5), and dominates $R_{\text{out}}^k(\cdot)$. Then the probability of sampling the correct CoT at least once in N attempts, while avoiding any misleading CoTs, is:

$$\Theta\left(\left[1 - \frac{1}{M^{L-1}}\right]^N \cdot [1 - (1 - c^{L-1})^N]\right).$$

Hence, the pass@K success probability is lower bounded by:

$$\Theta\left(1 - \left(1 - \left[1 - \frac{1}{M^{L-1}}\right]^N \cdot [1 - (1 - c^{L-1})^N]\right)^K\right) = \Theta(1),$$

again holding when K is large.

Now consider PRM-based BoN under Eq.(10. At each step, the minimal success probability is:

$$\Theta\left(\left[1 - \frac{1}{M}\right]^N \cdot [1 - (1 - c)^N]\right),$$

so across $L - 1$ steps, the overall probability is:

$$\Theta \left(\left[1 - \frac{1}{M} \right]^{N(L-1)} \cdot \left[1 - (1-c)^{N(L-1)} \right] \right),$$

and the corresponding pass@K is lower bounded by:

$$\Theta \left(1 - \left(1 - \left[1 - \frac{1}{M} \right]^{N(L-1)} \cdot \left[1 - (1-c)^{N(L-1)} \right] \right)^K \right) = \Theta(1).$$

Now consider any different task $k' \neq k$. By Definitions 2.1 and 2.2, the oracle rewards $R_{\mathbf{Q}, \mathbf{A}}^k(\cdot)$, as well as the PRMs in Eqs. 10, all assign zero credit to instances sampled from k' . Therefore, all ORM/PRM-based BoN or BS strategies fail on task k' .

For Beam Search (BS), the result follows by analogous arguments since BS depends on the same reward signals layer-wise. \square

Proof of Thm. 4.1. Fix any instance (\mathbf{Q}, \mathbf{A}) of task (o_1, a, k) and assume the premise of the theorem: all correct CoTs are hard-to-reason and there exists at least one depth $l^* \in [L]$ at which the hard CoTs diverge from a valid easy-to-reason CoT (“sparse edge”). Let \mathbf{o}^{easy} denote one such easy CoT and \mathbf{o}^{hard} any hard CoT. By Prop. 4.2, population-level ORM and PRM scores strictly prefer the easy branch whenever they differ:

$$R_{\text{out}}^k(\mathbf{o}^{\text{easy}}) > R_{\text{out}}^k(\mathbf{o}^{\text{hard}}), \quad R_{\text{likelihood}}^k(\mathbf{o}_l^{\text{easy}}) > R_{\text{likelihood}}^k(\mathbf{o}_l^{\text{hard}}) \quad \text{for all } l \text{ with } \mathbf{o}_l^{\text{easy}} \neq \mathbf{o}_l^{\text{hard}}. \quad (\text{A})$$

We analyze (i) and (ii)&(iii) separately. Throughout, M is the per-node branching factor and L is the CoT length. We take the conservative lower bounds that (a) at each node a particular child has sampling probability at least $1/M$, and (b) samples across the N trials are i.i.d.

(i) ORM + BoN. Best-of- N (BoN) first draws N full trajectories (CoTs) i.i.d. from the generator and then selects the one with the largest ORM score $R_{\text{out}}^k(\cdot)$. By (A), if among the N samples there exists at least one \mathbf{o}^{easy} , BoN will select an easy CoT, hence it will fail under the theorem’s premise (easy branch is valid but leads away from any correct hard solution due to the sparse-edge divergence).

We bound the probability that at least one \mathbf{o}^{easy} appears among N samples. Consider any fixed easy CoT \mathbf{o}^{easy} that agrees with \mathbf{o}^{hard} on the prefix up to (but excluding) l^* and then takes a different child at l^* . A conservative lower bound on the probability of sampling this *specific* easy CoT in one draw is

$$p_{\text{traj}} \geq \left(\frac{1}{M} \right)^{L-1} = \frac{1}{M^{L-1}},$$

since at $L - 1$ branching decisions (excluding the terminal) we multiply the minimal per-step mass $1/M$. Hence the probability that none of the N i.i.d. draws equals this easy trajectory is

$$(1 - p_{\text{traj}})^N \leq \left(1 - \frac{1}{M^{L-1}} \right)^N = \left(\frac{M^L - M}{M^L} \right)^N.$$

Therefore, with probability at least $1 - (1 - 1/M^{L-1})^N$ an easy CoT appears among the N draws, and by (A) BoN selects it and thus fails. Imposing

$$\left(1 - \frac{1}{M^{L-1}} \right)^N \leq \epsilon \iff N \geq \frac{\log(\epsilon)}{\log\left(\frac{M^L - M}{M^L}\right)},$$

ensures that the failure probability is at least $1 - \epsilon$, which proves the first bullet.

(ii) PRM + BoN (step-wise) and (iii) PRM + Beam Search (width N , beam $B \geq 1$). PRM-based inference expands *partial* CoTs and uses the local PRM score $R_{\text{likelihood}}^k(\mathbf{o}_l)$ to select among candidates. Consider the first divergence depth l^* . In each expansion round at depth l^* , the procedure proposes N children i.i.d. (BoN: propose and take the best child by PRM; Beam: propose N and keep the top- B by PRM). Let

$$p_{\text{child}} \geq \frac{1}{M}$$

be the conservative lower bound that a given proposal at depth l^* takes the (PRM-favored) easy child rather than the hard sparse edge. Thus the probability that *none* of the N proposals includes the easy child at that step is

$$(1 - p_{\text{child}})^N \leq \left(1 - \frac{1}{M}\right)^N = \left(\frac{M-1}{M}\right)^N.$$

Consequently, with probability at least $1 - (1 - 1/M)^N$ the easy child appears among the N proposals at depth l^* . By (A), PRM strictly prefers that easy child over the hard child at depth l^* , so:

- *PRM + BoN (step-wise)*: the chosen next token is the easy child, irrevocably steering the trajectory onto the easy branch. Repeating this argument at later depths where branches differ keeps the easy path strictly preferred, so the final selection is easy and the method fails under the theorem’s premise.
- *PRM + Beam Search*: since $B \geq 1$, any PRM-strictly-better easy child is ranked above the hard child and therefore included in the beam at depth l^* ; by standard beam monotonicity with strictly better local scores at each subsequent divergence, the easy branch remains in the top- B and dominates the final selection, hence failure.

Imposing

$$\left(1 - \frac{1}{M}\right)^N \leq \epsilon \iff N \geq \frac{\log(\epsilon)}{\log\left(\frac{M-1}{M}\right)},$$

ensures that an easy child appears at the first divergence step with probability at least $1 - \epsilon$, and by the PRM preference this forces selection of the easy branch, completing the second bullet.

Conclusion. In all cases, Prop. 4.2 ensures a strict scoring advantage for the easy branch whenever it is present among candidates; the displayed lower bounds control the probability that such an easy candidate *does* appear given N proposals. Choosing N to satisfy

$$\left(1 - \frac{1}{M^{L-1}}\right)^N \leq \epsilon \quad (\text{ORM + BoN}), \quad \left(1 - \frac{1}{M}\right)^N \leq \epsilon \quad (\text{PRM + BoN/BS}),$$

yields failure probability at least $1 - \epsilon$ for (i) and for (ii)&(iii), respectively. \square

Proof. Heuristic Proof of Corollary 4.3. Let $(\Omega, \mathcal{F}, \mu)$ be a base measure space where $\hat{p}_{\theta^*} \ll \mu$ with Radon-Nikodym derivative $d\hat{p}_{\theta^*}/d\mu > 0$ μ -a.e. We consider the optimization over absolutely continuous measures $P_{\text{new}}^k \ll \hat{p}_{\theta^*}$.

The objective functional can be written as:

$$J(P_{\text{new}}^k) = \mathbb{E}_{P_{\text{new}}^k} [R(\mathbf{o})] - \frac{1}{\lambda} D_{KL}(P_{\text{new}}^k \parallel \hat{p}_{\theta^*}) \quad (76)$$

where $R(\mathbf{o}) := R_{\text{out}}^k(\mathbf{o})$. We require:

- (C1) $R \in L^1(\hat{p}_{\theta^*})$ (finite expected reward)
- (C2) $\exists \epsilon > 0$ s.t. $\hat{p}_{\theta^*} \geq \epsilon \mu$ -a.e. (strict positivity)

High-levelly, the remaining proof is convex optimization in probability space. Define the Lagrangian with measure-theoretic notation:

$$\mathcal{L}(P, \eta) = \int R dP - \frac{1}{\lambda} \int \log\left(\frac{dP}{d\hat{p}_{\theta^*}}\right) dP + \eta \left(1 - \int dP\right) \quad (77)$$

Require:

- (C3) $P \in \mathcal{P}(\Omega)$, the space of probability measures absolutely continuous to μ
- (C4) $\log(dP/d\hat{p}_{\theta^*}) \in L^1(P)$ (finite KL divergence)

For $P \in \mathcal{P}(\Omega)$, consider variation $P_\epsilon = P + \epsilon Q$ where Q is a signed measure with $\int dQ = 0$. The Gâteaux derivative is:

$$\left. \frac{d}{d\epsilon} \mathcal{L}(P_\epsilon, \eta) \right|_{\epsilon=0} = \int R dQ - \frac{1}{\lambda} \int \left(\log \frac{dP}{d\hat{p}_{\theta^*}} + 1 \right) dQ - \eta \int dQ \quad (78)$$

For optimality, this must vanish for all admissible Q , requiring:

$$R(\mathbf{o}) - \frac{1}{\lambda} \left(\log \frac{dP}{d\hat{p}_{\theta^*}}(\mathbf{o}) + 1 \right) - \eta = 0 \quad P\text{-a.s.} \quad (79)$$

Rearranging gives:

$$\log \frac{dP}{d\hat{p}_{\theta^*}} = \lambda R(\mathbf{o}) - (1 + \lambda\eta) \quad (80)$$

Exponentiating both sides:

$$dP = \hat{p}_{\theta^*}(\mathbf{o}) \exp(\lambda R(\mathbf{o})) \exp(-1 - \lambda\eta) d\mu(\mathbf{o}) \quad (81)$$

Normalization requires:

$$\exp(1 + \lambda\eta) = \int \hat{p}_{\theta^*} \exp(\lambda R) d\mu =: Z \quad (82)$$

Thus the optimal measure is:

$$dP_{\text{adjusted}}^k = \frac{1}{Z} \hat{p}_{\theta^*} \exp(\lambda R) d\mu \quad (83)$$

First verify $P_{\text{adjusted}}^k \in \mathcal{P}(\Omega)$:

- Absolute continuity: Immediate from $\hat{p}_{\theta^*} \ll \mu$ and $Z^{-1} \exp(\lambda R) > 0$
- Integrability: By (C1) and $\exp(\lambda R) \leq \exp(\lambda \|R\|_\infty) < \infty$ from $R \leq 1$

Second, confirm stationarity. For any $Q \in T_{P_{\text{adjusted}}^k} \mathcal{P}(\Omega)$ (tangent space):

$$d\mathcal{L}(P_{\text{adjusted}}^k, \eta)(Q) = \int \underbrace{\left[R - \frac{1}{\lambda} \left(\log \frac{dP_{\text{adjusted}}^k}{d\hat{p}_{\theta^*}} + 1 \right) - \eta \right]}_{=0} dQ = 0 \quad (84)$$

Substitute P_{adjusted}^k into J :

$$\begin{aligned} J(P_{\text{adjusted}}^k) &= \mathbb{E}_{P_{\text{adjusted}}^k} [R] - \frac{1}{\lambda} \mathbb{E}_{P_{\text{adjusted}}^k} \left[\log \frac{P_{\text{adjusted}}^k}{\hat{p}_{\theta^*}} \right] \\ &= \mathbb{E}_{P_{\text{adjusted}}^k} [R] - \frac{1}{\lambda} (\lambda \mathbb{E}[R] - \log Z) \\ &= \frac{1}{\lambda} \log Z \end{aligned}$$

By Gibbs' inequality, this maximizes the trade-off between expected reward and KL regularization.

To validate our conditions required, we summarized:

- (C1): Holds as $\|R\|_\infty \leq 1$ by assumption
- (C2): Guaranteed by model construction $\hat{p}_{\theta^*} = \text{softmax}(\cdot) > 0$
- (C3): Inherited from base measure μ
- (C4): Satisfied because $D_{KL}(P_{\text{adjusted}}^k \parallel \hat{p}_{\theta^*}) = \log Z - \lambda \mathbb{E}[R] < \infty$

2640 Thus under these conditions, P_{adjusted}^k is the unique maximizer of $J(P_{\text{new}}^k)$ in $\mathcal{P}(\Omega)$. □

2641

2642 *Proof.* Proof of the legitimacy of Def.4.4. To show $h_k(\cdot)$ is a harmonic function, let us verify

2643

$$2644 \quad 1 = \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}'_{l+1} | \mathbf{o}_l) = \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}'_{l+1})}{h_k(\mathbf{o}_l)}.$$

2645

2646 By the fact that

$$2647 \quad \begin{aligned} 2648 \quad h_k(\mathbf{o}_l) &= \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o})) | \mathbf{o}_l] \\ 2649 \quad &= \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}_{\mathbf{o}'_{l+2:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}'_{l+1}] \end{aligned} \quad (85)$$

2650

2651 we see that

2652

$$2653 \quad \begin{aligned} 2654 \quad &\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}'_{l+1})}{h_k(\mathbf{o}_l)} \\ 2655 \quad &= \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \frac{\mathbb{E}_{\mathbf{o}'_{l+2:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}'_{l+1}]}{\mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o})) | \mathbf{o}_l]} \\ 2656 \quad &= \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \frac{\mathbb{E}_{\mathbf{o}'_{l+2:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}'_{l+1}]}{\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \mathbb{E}_{\mathbf{o}'_{l+2:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}'_{l+1}]} = 1 \end{aligned}$$

2657

2658

2659 Recall the definition of our **DPRM**:

$$2660 \quad R_{\text{DPRM}}^k(\mathbf{o}_l) = \frac{1}{\lambda} \log \left(\mathbb{E}_{\mathbf{o}'_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}_l] \right), \quad (86)$$

2661

2662

$$2663 \quad \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) = \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)} = \frac{\hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1}))}{Z_l(\mathbf{o}_l)},$$

2664 where $Z_l(\mathbf{o}_l) = \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}'_{l+1} | \mathbf{o}_l) \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}'_{l+1}))$. Collaborating with Eq.(87) as well as the definition of

2665

2666 $R_{\text{DPRM}}^k(\mathbf{o}_l)$, we can equate:

2667

$$2668 \quad \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_l)) = h_k(\mathbf{o}_l) = \mathbb{E}_{\mathbf{o}'_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o}')) | \mathbf{o}_l].$$

2669

2670 Therefore, it holds that

2671

$$2672 \quad \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) = \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \cdot \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)} \propto \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1})).$$

2673

2674 Recall:

2675

$$2676 \quad h_k(\mathbf{o}_l) = \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o})) | \mathbf{o}_l],$$

2677

2678 so $h_k(\mathbf{o}_L) = \exp(\lambda R_{\text{out}}^k(\mathbf{o}))$ and $h_k(\mathbf{o}_0) = Z := \sum_{\mathbf{o}' \in \mathcal{T}_{\text{all}}} \hat{p}_{\theta^*}(\mathbf{o}') \exp(\lambda R_{\text{out}}^k(\mathbf{o}'))$. The h-transformed transition is:

2679

$$2680 \quad \hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) = \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)}, \quad (87)$$

2681

2682 yielding:

2683

$$2684 \quad P_{\text{DPRM}}^k(\mathbf{o}) = \prod_{l=1}^{L-1} \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)} = \hat{p}_{\theta^*}(\mathbf{o}) \frac{h_k(\mathbf{o}_L)}{h_k(\mathbf{o}_0)} = P_{\text{Gibbs}}^k(\mathbf{o}). \quad (88)$$

2685

2686 The proof is completed.

2687

2688

2689 □

Corollary H.2. For the task $k \in \mathcal{T}$, let θ^* be the pretrained Foundation Model from Thm. F.1, and $R_{\text{out}}^k(\cdot)$ be the ORM. Consider the ORM-equipped and DPRM-equipped adjusted sampling distributions defined in Corollary 4.3.

- As the temperature parameter $\lambda \rightarrow \infty$, we have the following situation

1. **ORM-Equipped Adjusted Sampling:** The distribution in Eq.(11) converges to:

$$P_{\text{Gibbs}}^k(\mathbf{o}) \xrightarrow{\lambda \rightarrow \infty} \begin{cases} 1, & \text{if } \mathbf{o} = \arg \max_{\mathbf{o}' \in \mathcal{T}_{\text{all}}} R_{\text{out}}^k(\mathbf{o}') \\ 0, & \text{otherwise} \end{cases}$$

akin to a ORM-based BoN with $R_{\text{out}}^k(\cdot)$.

2. **DPRM-Equipped Adjusted Sampling:** The step-wise distribution (13) with $R_{\text{DPRM}}^k(\mathbf{o}_{l+1}) = \frac{1}{\lambda} \log h_k(\mathbf{o}_{l+1})$ converges to:

$$\hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) \xrightarrow{\lambda \rightarrow \infty} \begin{cases} 1, & \text{if } \mathbf{o}_{l+1} = \arg \max_{\mathbf{o}' \in S_{l+1}} R_{\text{likelihood}}^k(\mathbf{o}') \\ 0, & \text{otherwise} \end{cases}$$

akin to a PRM-based BoN with $R_{\text{likelihood}}^k(\cdot)$.

- When the temperature parameter $\lambda > 0$, each step $l \in \{0, \dots, L-1\}$ satisfies:

$$\arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} R_{\text{DPRM}}^k(\mathbf{o}_l) = \arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} R_{\text{likelihood}}^k(\mathbf{o}_l),$$

where $S_l^{\text{BoN}} = \{\mathbf{o}_l^1, \dots, \mathbf{o}_l^N\}$ denotes the N candidates sampled by the base model \hat{p}_{θ^*} . Therefore, using $\lambda > 0$ with **BoN**, **Beam Search** or **Lookahead Search** equates to prior PRM methods employing the same search strategies.

Proof. Proof of Corollary 4.5. We analyze the asymptotic behavior of the sampling distributions as $\lambda \rightarrow \infty$.

For part (1), consider the ORM-equipped adjusted sampling distribution:

$$P_{\text{Gibbs}}^k(\mathbf{o}) = \frac{P_{\theta^*}(\mathbf{o}) \exp(\lambda R_{\text{out}}^k(\mathbf{o}))}{\sum_{\mathbf{o}' \in \mathcal{T}_{\text{all}}} P_{\theta^*}(\mathbf{o}') \exp(\lambda R_{\text{out}}^k(\mathbf{o}'))},$$

where \mathcal{T}_{all} is the set of all possible trajectories. Let $\mathbf{o}^* = \arg \max_{\mathbf{o}' \in \mathcal{T}_{\text{all}}} R_{\text{out}}^k(\mathbf{o}')$, with maximum reward $R_{\text{out}}^k(\mathbf{o}^*)$. As $\lambda \rightarrow \infty$, the term $\exp(\lambda R_{\text{out}}^k(\mathbf{o}))$ dominates for \mathbf{o} with the largest $R_{\text{out}}^k(\mathbf{o})$. For $\mathbf{o} \neq \mathbf{o}^*$, if $R_{\text{out}}^k(\mathbf{o}) < R_{\text{out}}^k(\mathbf{o}^*)$, then

$$\frac{\exp(\lambda R_{\text{out}}^k(\mathbf{o}))}{\exp(\lambda R_{\text{out}}^k(\mathbf{o}^*))} = \exp(\lambda(R_{\text{out}}^k(\mathbf{o}) - R_{\text{out}}^k(\mathbf{o}^*))) \rightarrow 0,$$

since $R_{\text{out}}^k(\mathbf{o}) - R_{\text{out}}^k(\mathbf{o}^*) < 0$. Assuming $R_{\text{out}}^k(\mathbf{o})$ has a unique maximum (or summing over all maximizers if not unique), the denominator is dominated by $P_{\theta^*}(\mathbf{o}^*) \exp(\lambda R_{\text{out}}^k(\mathbf{o}^*))$. Thus,

$$P_{\text{Gibbs}}^k(\mathbf{o}) \rightarrow \begin{cases} 1, & \text{if } \mathbf{o} = \mathbf{o}^*, \\ 0, & \text{otherwise,} \end{cases}$$

which matches the behavior of BoN Sampling, where the trajectory with the highest $R_{\text{out}}^k(\mathbf{o})$ is selected.

For part (2), consider the DPRM-equipped step-wise distribution:

$$\hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) = \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \frac{h_k(\mathbf{o}_{l+1})}{h_k(\mathbf{o}_l)},$$

with $h_k(\mathbf{o}_{l+1}) = \mathbb{E}_{\mathbf{o}_{l+2:L} \sim \hat{p}_{\theta^*}} [\exp(\lambda R_{\text{out}}^k(\mathbf{o})) | \mathbf{o}_{l+1}]$, and $R_{\text{DPRM}}^k(\mathbf{o}_{l+1}) = \frac{1}{\lambda} \log h_k(\mathbf{o}_{l+1})$. Substituting h_k , we get

$$\hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} | \mathbf{o}_l) = \hat{p}_{\theta^*}(\mathbf{o}_{l+1} | \mathbf{o}_l) \exp(\lambda(R_{\text{DPRM}}^k(\mathbf{o}_{l+1}) - R_{\text{DPRM}}^k(\mathbf{o}_l))) \frac{1}{Z},$$

where $Z = \sum_{\mathbf{o}_{l+1} \in S_{l+1}} \hat{p}_{\theta^*}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) \exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1}))$ is the normalizing constant. Let $\mathbf{o}_{l+1}^* = \arg \max_{\mathbf{o}' \in S_{l+1}} R_{\text{DPRM}}^k(\mathbf{o}')$. As $\lambda \rightarrow \infty$, the term $\exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1}))$ dominates for $\mathbf{o}_{l+1} = \mathbf{o}_{l+1}^*$. For $\mathbf{o}_{l+1} \neq \mathbf{o}_{l+1}^*$, if $R_{\text{DPRM}}^k(\mathbf{o}_{l+1}) < R_{\text{DPRM}}^k(\mathbf{o}_{l+1}^*)$, then

$$\frac{\exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1}))}{\exp(\lambda R_{\text{DPRM}}^k(\mathbf{o}_{l+1}^*))} = \exp(\lambda(R_{\text{DPRM}}^k(\mathbf{o}_{l+1}) - R_{\text{DPRM}}^k(\mathbf{o}_{l+1}^*))) \rightarrow 0.$$

Thus, the distribution concentrates on \mathbf{o}_{l+1}^* :

$$\hat{p}_{\theta}^{\text{new},k}(\mathbf{o}_{l+1} \mid \mathbf{o}_l) \rightarrow \begin{cases} 1, & \text{if } \mathbf{o}_{l+1} = \mathbf{o}_{l+1}^*, \\ 0, & \text{otherwise,} \end{cases}$$

which mimics BoN Sampling by selecting the state with the highest $R_{\text{likelihood}}^k$ by our last item.

For the last item, for step $l \in \{0, \dots, L-1\}$, the DPRM with $\lambda > 0$ is given as $R_{\text{DPRM}}^k(\mathbf{o}_l) = \log \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} \left[\exp \left(R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l \right) \right]$. Since log is strictly increasing, we have

$$\arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} R_{\text{DPRM}}^k(\mathbf{o}_l) = \arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} \left[\exp \left(R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l \right) \right].$$

Similarly, since exp is strictly increasing, the arg max over $\mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} \left[\exp \left(R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l \right) \right]$ is equivalent to the arg max over $\mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} \left[R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l \right]$. Thus, given that $R_{\text{likelihood}}^k(\mathbf{o}_l) = \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} \left[R_{\text{out}}^k(\mathbf{o}) \mid \mathbf{o}_l \right]$ it holds that

$$\arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} R_{\text{DPRM}}^k(\mathbf{o}_l) = \arg \max_{\mathbf{o}_l \in S_l^{\text{BoN}}} \mathbb{E}_{\mathbf{o}_{l+1:L} \sim \hat{p}_{\theta^*}} R_{\text{likelihood}}^k(\mathbf{o}_l).$$

This shows that BoN Sampling with R_{pro}^k maximizes the expected outcome reward, aligning with prior methods. The equivalence for Beam Search follows similarly by replacing the sampling strategy with the respective search method, as they also maximize $R_{\text{pro}}^k(\mathbf{o}_l)$. This completes the proof. \square

Proof. Proof of Cor. 4.6. The proof directly follows the proof of Cor. G.5. \square

Corollary H.3 (Extension: Comparison with Ground-true Oracle). *Let θ^* be the base model in Eq.(2 that exactly predicts the distribution of a Multi-task TMC as in Definitions 2.1 and 2.2. Under task tuple $(q, a, k) \in S_1 \times S_L \times \mathcal{T}$, consider the ORMs $R_{\mathbf{Q}, \mathbf{A}}^k(\cdot)$ and $R_{\text{out}}^k(\cdot)$, and the PRMs of Eqs. 10. For any target task k with instance distribution $\mathcal{D}_{a,q}^{a,k}$, suppose the number of hard-to-reason CoTs is $\Theta(M)$ and the number of nonzero-probability CoTs from q to S_L is N_q . Then under pass@K sampling:*

1. DPRM is More Capable of Hard CoTs. *If a specific hard CoT has sampling probability $p = o(M^{-(L-1)})$ under the base model, then for any BoN budget*

$$N = O\left(\frac{\log(1-N_q^{-1})}{\log(1-p)}\right),$$

there exists $\lambda = o\left(\ln \frac{(1-p)^N}{(N_q-1)(1-(1-p)^N)}\right)$ such that DPRM with temperature λ achieves strictly higher pass@K than ORM-based or PRM-based BoN (or BS).

2. Preserve Multi-task. *For any $\varepsilon > 0$, if*

$$K = \Omega\left(\frac{\ln \varepsilon}{\ln\left(\frac{(N_q-1)e^\lambda}{1+(N_q-1)e^\lambda}\right)}\right),$$

then DPRM with $\lambda > 0$ attains pass@K $\geq 1 - \varepsilon$ on any other task $k' \neq k$.

In both cases, adjusting the temperature $\lambda > 0$ controls the pass@K performance.

Proof. The arguments parallel those in Cor. G.5, so we focus on the comparison of pass@K success probabilities.

(i) Hard-CoT capability. Under ORM-based BoN with ground-truth reward, the success probability for the unique hard CoT is

$$1 - (1 - p)^N.$$

Under DPRM (Eq.(11)), every valid CoT—including the correct one—has sampling probability at least

$$\frac{1}{N_q - M + Me^\lambda} \geq \frac{1}{1 + (N_q - 1)e^\lambda}.$$

Choosing $\lambda = o\left(\ln \frac{(1-p)^N}{(N_q-1)(1-(1-p)^N)}\right)$ ensures $\frac{1}{1+(N_q-1)e^\lambda} \geq 1-(1-p)^N$, so DPRM outperforms ORM. Similarly, when the budget of PRM-based BoN (or BS) in pass@K is limited and $\lambda \rightarrow 0$ would achieve more satisfactory success probability.

(ii) Multi-task preservation. For any other task k' , DPRM still assigns probability at least $\frac{1}{1+(N_q-1)e^\lambda}$ to each valid CoT. Thus, with

$$K = \Omega\left(\frac{\ln \varepsilon}{\ln\left(\frac{(N_q-1)e^\lambda}{1+(N_q-1)e^\lambda}\right)}\right),$$

the pass@K guarantee $1 - \left(1 - \frac{1}{1+(N_q-1)e^\lambda}\right)^K \geq 1 - \varepsilon$ holds, completing the proof. \square

I. Auxiliary Lemmas

Lemma I.1. Let θ^* be the base model

$$\hat{p}_\theta(\cdot|\mathbf{x}) = \text{softmax}(h_\theta(\cdot, \mathbf{x})), \quad \mathbf{x} \in \{0, 1\}^{|S|}. \quad (89)$$

Then for $\forall \mathbf{o}_l \in S_l, \mathbf{o}_{l+1} \in S_{l+1}$

$$\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) = \nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l). \quad (90)$$

Further, if the base model is Eq.(2), we have

$$\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) = e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}, \quad (91)$$

where $e_{\mathbf{o}_{l+1}, \mathbf{o}_l} := \mathbf{o}_{l+1} \mathbf{o}_l^\top \in \{0, 1\}^{|S| \times |S|}$ is the one-hot matrix with only the position corresponding to $(\mathbf{o}_{l+1}, \mathbf{o}_l)$ is 1 and 0 elsewhere.

Proof. By Eq. (89), we have

$$\begin{aligned} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) &= \nabla_{\theta^k} \log \frac{e^{h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l)}}{\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)}} \\ &= \nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) - \nabla_{\theta^k} \log \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)} \\ &= \nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) - \frac{\nabla_{\theta^k} \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)}}{\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)}} \\ &= \nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) - \frac{\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)} \nabla_{\theta^k} h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)}{\sum_{\mathbf{o}'_{l+1} \in S_{l+1}} e^{h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l)}} \\ &= \nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_\theta(\mathbf{o}'_{l+1}, \mathbf{o}_l). \end{aligned} \quad (92)$$

Besides, if the base model is $\hat{p}_\theta(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2), we have

$$\nabla_{\theta^k} h_\theta(\mathbf{o}_{l+1}, \mathbf{o}_l) = \nabla_{\theta^k} \langle \theta_{\mathbf{o}_{l+1}, \cdot}, \mathbf{o}_l \rangle = e_{\mathbf{o}_{l+1}, \mathbf{o}_l}, \quad (93)$$

Dragging Eq.(93) into Eq.(90), we could obtain Eq.(91).

The proof is completed. \square

Lemma I.2 (Policy Gradient for REINFORCE & RAFT under TMC). *Let θ^* be the base model in Eq.(2) that exact predicts the distribution of Multi-task TMC as in Def. 2.1 and 2.2, and θ^k the current model to be finetuned from θ^* for task $k \in \mathcal{T}$. The gradient of the REINFORCE objective for task k is given by:*

$$\nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right], \quad (94)$$

$$\nabla_{\theta^k} \mathcal{J}_{\text{RAFT}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right], \quad (95)$$

$$\left[(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right], \quad (96)$$

where

$$\mathcal{J}_{\text{REINFORCE}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{\mathbf{o}_1}^{01, k}}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O|\mathbf{o}_1)} \left[R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right], \quad (97)$$

$$\mathcal{J}_{\text{RAFT}}(\theta^k) = \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim \mathcal{D}_{a_{\mathbf{o}_1}^{01, k}}, \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O|\mathbf{o}_1)} \left[\sum_{l=1}^{L-1} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{(\mathbf{Q}, \mathbf{A})}^k(\mathbf{o}) \right], \quad (98)$$

Remark I.3. In the main text, Eq.(8 contains a typo: the summation term “ $\sum_{l=1}^{L-1}$ ” inside the expectation is omitted. The correct formulation is provided in Eq.(98). Additionally, the formal versions of Eq.(8 are given as Eq.(94 and Eq.(95), respectively.

Proof. For any complete trajectory $\mathbf{o} = (\mathbf{o}_1, \dots, \mathbf{o}_L)$:

$$\hat{p}_{\theta^k}(\mathbf{o}) = P^k(\mathcal{Q}^k) \prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \quad (99)$$

where $P^k(\mathcal{Q}^k)$ is the initial state distribution (parameter-independent by Def. 2.1). By the property of TMC, we have

$$\begin{aligned} \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \nabla_{\theta^k} \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O|\mathbf{o}_1)} \left[R_{\text{out}}^k(\mathbf{o}) \right] \\ &\stackrel{(1)}{=} \nabla_{\theta^k} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}) \left[P^k(q) \prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] d\mathbf{o}_{1:L} \\ &\stackrel{(2)}{=} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}) P^k(q) \nabla_{\theta^k} \left[\prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] d\mathbf{o}_{1:L} \\ &\stackrel{(3)}{=} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}) P^k(q) \left[\prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] \sum_{l=1}^{L-1} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) d\mathbf{o}_{1:L} \quad (100) \\ &\stackrel{(4)}{=} \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O|\mathbf{o}_1)} \left[R_{\text{out}}^k(\mathbf{o}) \sum_{l=1}^{L-1} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] \\ &\stackrel{(5)}{=} \sum_{l=1}^{L-1} \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), \mathbf{o}_{2:L} \sim \hat{p}_{\theta^k}^k(O|\mathbf{o}_1)} \left[\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right]. \end{aligned}$$

Step (1) expands the expectation as an integral over trajectories using the MDP's joint distribution $P^k(q) \prod_{t=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t)$;

Step (2) applies the Leibniz interchange under Markovian policy structure:

$$\nabla_{\theta^k} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}) \hat{p}_{\theta^k}(\mathbf{o}_{1:L}) d\mathbf{o}_{1:L} = \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}) \nabla_{\theta^k} \hat{p}_{\theta^k}(\mathbf{o}_{1:L}) d\mathbf{o}_{1:L} \quad (\text{a.s.}) \quad (101)$$

valid when: (i) *Policy Gradient Dominance*: $\exists h \in L^1(\mu)$ such that $\|R_{\text{out}}^k(\mathbf{o}) \nabla_{\theta^k} \hat{p}_{\theta^k}(\mathbf{o}_{1:L})\| \leq h(\mathbf{o}_{1:L}) \forall \theta^k \in \Theta^k$ where $\Theta^k = \mathbb{R}^{|S| \times |S|}$ denotes the parameter space; (ii) *Parameterized Measure Continuity*: The map $\theta^k \mapsto \sqrt{\hat{p}_{\theta^k}(\mathbf{o}_{1:L})}$ is $W^{1,1}$ -

continuous with: $\lim_{\|v\| \rightarrow 0} \mathbb{E}_{\mu} \left[\left\| \frac{\sqrt{\hat{p}_{\theta^k+v}(\mathbf{o}_{1:L})} - \sqrt{\hat{p}_{\theta^k}(\mathbf{o}_{1:L})}}{\|v\|} \right\|^2 \right] < \infty$, which are all satisfied under our case since $\|R_{\text{out}}^k(\cdot)\|_{\infty} = O(1)$

and $\hat{p}_{\theta^k}(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2);

Step (3) decomposes using Markovian parameter isolation:

$$\nabla_{\theta^k} \prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) = \sum_{l=1}^{L-1} \left(\prod_{m=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{m+1}|\mathbf{o}_m) \right) \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \quad (102)$$

valid under: (i) *Disjoint Parameter Control*: $\theta^k = \bigsqcup_{l=1}^{L-1} \theta_l^k$ where $\theta_l^k \cap \theta_{l'}^k = \emptyset$ for $l' \neq l$, with each $\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) = f_l(\mathbf{o}_{l+1}|\mathbf{o}_l; \theta_l^k)$ and $\frac{\partial f_l}{\partial \theta_l^k} \equiv 0$; (ii) *Log-Smoothness*: $\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) > 0$ μ -a.e. and $\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \in L^2(\hat{p}_{\theta^k} \otimes \mu)$;

(iii) *Sequential Fubini Condition*: $\int_{\mathcal{O}^L} \prod_{l=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) d\mathbf{o}_{1:L} = \prod_{l=1}^{L-1} \int_{\mathcal{O}} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) d\mathbf{o}_{l+1}$ in terms of total variation norm, which are all easily verified under our $\hat{p}_{\theta^k}(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2);

Step (4) rewrites the integral as $\mathbb{E}_{\mathbf{o}_{1:L} \sim \hat{p}_{\theta^k}} [R_{\text{out}}^k(\mathbf{o}) \sum_{l=1}^{L-1} \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)]$;

Step (5) exchanges summation and expectation via Fubini's theorem, valid when $\mathbb{E}[\|R_{\text{out}}^k \nabla_{\theta^k} \log \hat{p}_{\theta^k}\|] < \infty$, which obviously hold in our setting.

Similarly, we have

$$\begin{aligned} \nabla_{\theta^k} \mathcal{J}_{\text{RAFT}}(\theta^k) &= \nabla_{\theta^k} \mathbb{E}_{\substack{\mathbf{o}_1=q \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}_{t+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_t)}} \left[\sum_{l=1}^{L-1} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \cdot R_{\text{out}}^k(\mathbf{o}_{1:L}) \right] \\ &\stackrel{(1)}{=} \sum_{l=1}^{L-1} \nabla_{\theta^k} \left(\int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot P^k(\mathbf{o}_1) \cdot \prod_{\substack{t=1 \\ t \neq l}}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \cdot \left[\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] d\mathbf{o}_{1:L} \right) \\ &\stackrel{(2)}{=} \sum_{l=1}^{L-1} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot P^k(\mathbf{o}_1) \cdot \left(\prod_{\substack{t=1 \\ t \neq l}}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \right) \cdot \nabla_{\theta^k} \left[\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] d\mathbf{o}_{1:L} \\ &\stackrel{(3)}{=} \sum_{l=1}^{L-1} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot P^k(\mathbf{o}_1) \left(\prod_{t=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \right) \cdot \frac{\nabla_{\theta^k} \left[\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right]}{\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)} d\mathbf{o}_{1:L} \\ &\stackrel{(4)}{=} \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}_{t+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_t)}} \left[R_{\text{out}}^k(\mathbf{o}_{1:L}) \left(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right) \cdot \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] \\ &\stackrel{(5)}{=} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \mathbf{o}_{t+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_t)}} \left[R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot \sum_{l=1}^{L-1} \left(1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right) \cdot \nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] \end{aligned} \quad (103)$$

Step (1) expands the expectation using the Markovianity's factorized structure $P^k(\mathbf{o}_1) \prod_{t=1}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t)$, isolating the l -th transition's $\hat{p} \log \hat{p}$ term while keeping others as standard transitions, which is legitimate under our $\hat{p}_{\theta^k}(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2);

Step (2) enforces parameter-localized differentiation through:

$$\sum_{l=1}^{L-1} \nabla_{\theta^k} \int F_l d\mathbf{o} = \sum_{l=1}^{L-1} \int_{\mathcal{O}^L} R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot P^k(\mathbf{o}_1) \cdot \left(\prod_{\substack{t=1 \\ t \neq l}}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \right) \cdot \nabla_{\theta^k} \left[\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] d\mathbf{o}_{1:L} \quad (\text{a.s.}) \quad (104)$$

where $F_l = R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot P^k(\mathbf{o}_1) \cdot \prod_{\substack{t=1 \\ t \neq l}}^{L-1} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \cdot [\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)]$, valid when:

(i) *Architectural Parameter Isolation*: Policy parameters partition as $\theta^k = \bigsqcup_{l=1}^{L-1} \theta_l^k$ with:

$$\frac{\partial}{\partial \theta_m^k} \hat{p}_{\theta^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) = \begin{cases} \nabla_{\theta_l^k} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) & t = l \text{ and } m = l \\ 0 & \text{otherwise} \end{cases} \quad (105)$$

which is satisfied as $\hat{p}_{\theta^k}(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2); (ii) *Localized Dominance*: $\exists h_l \in L^1(\mu_l)$ where μ_l is the base measure on $(\mathbf{o}_l, \mathbf{o}_{l+1})$, such that:

$$|R_{\text{out}}^k(\mathbf{o}_{1:L}) \cdot \nabla_{\theta^k} [\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)]| \leq h_l(\mathbf{o}_l, \mathbf{o}_{l+1}), \quad (106)$$

which clearly held under our conditions; (iii) *Decoupled Integration*: For each l ,

$$\int_{\mathcal{O}^L} F_l d\mathbf{o} = \int_{\mathbf{o}_1} P^k \int_{\mathbf{o}_{l+1}} [\hat{p}_{\theta^k} \log \hat{p}_{\theta^k}] \left(\prod_{\substack{t=1 \\ t \neq l}}^{L-1} \int_{\mathbf{o}_{t+1}} \hat{p}_{\theta^k} d\mathbf{o}_{t+1} \right) d\mathbf{o}_l d\mathbf{o}_{l+1} \quad (107)$$

with $\prod_{\substack{t=1 \\ t \neq l}}^{L-1} \int \hat{p}_{\theta^k} d\mathbf{o}_{t+1} = 1$ μ -a.e. This condition holds apparently under our model $\hat{p}_{\theta^k}(\cdot|x) = \text{softmax}(\langle \theta, x \rangle)$ by Eq.(2),

which linearly isolates each states; (iv) *Transition Differentiability*: Each $\theta_l^k \mapsto \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)$ is Fréchet differentiable with:

$$\mathbb{E} \left[\left\| \frac{\nabla_{\theta_l^k} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)}{\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)} \right\|^2 \right] < \infty \text{ which holds in our softmax model;}$$

Step (3) multiplies one $\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)$ in the front and divide it subsequently;

Step (4) uses the chain rule:

$$\nabla_{\theta^k} (\hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) = (1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) \nabla_{\theta^k} \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l),$$

and reconstructs the expectation by recognizing $\prod_{t=1}^{L-1} \hat{p}_t = \hat{p}_{\theta^k}(\mathbf{o}_{1:L})/P^k(\mathbf{o}_1)$, with cross-terms vanishing due to $\mathbb{E}_{\mathbf{o}_{m+1} \sim \hat{p}_m} [f(\mathbf{o}_l)] = \mathbb{E}[f(\mathbf{o}_l)]$ for $m \neq l$; Step (5) applies Fubini's theorem to exchange summation and expectation, valid by the fact $\mathbb{E} \left[\sum_{l=1}^{L-1} |(1 + \log \hat{p}_l) \nabla \log \hat{p}_l \text{Rex}| \right] < \infty$ in our case. \square

Remark I.4. When the base model is no longer in the linear form in Eq.(2), but a general form in Eq.(89) with $\hat{p}_{\theta}(\cdot|x) = \text{softmax}(h_{\theta}(\cdot, \mathbf{x}))$, $\mathbf{x} \in \{0, 1\}^{|S|}$, the conclusions still holds when

- **Architectural Conditions**

- *Parameter Isolation*: $\theta = \bigsqcup_{l=1}^{L-1} \theta_l$ where $\theta_l \cap \theta_{l'} = \emptyset$ for $l \neq l'$, with:

$$h_{\theta}(\mathbf{o}_{l+1}|\mathbf{o}_l) = h_l(\mathbf{o}_{l+1}|\mathbf{o}_l; \theta_l), \quad \frac{\partial h_l}{\partial \theta_{l'}} \equiv 0 \quad \forall l' \neq l \quad (108)$$

- *Module Independence*: Each $h_l(\cdot; \theta_l)$ uses distinct computational subgraphs without parameter sharing across l

- **Smoothness & Differentiability**

- *Lipschitz Continuity*: $\exists C_l > 0$ s.t.

$$\|h_l(\cdot; \theta_l + \Delta\theta) - h_l(\cdot; \theta_l)\|_{\infty} \leq C_l \|\Delta\theta\|_2 \quad \forall \theta_l \quad (109)$$

– *Twice Differentiability*: $h_l \in C^2(\Theta_l)$ with bounded Hessians:

$$\mathbb{E}_{\mathbf{o}_l} [\|\nabla_{\boldsymbol{\theta}_l}^2 h_l\|_{\text{op}}^2] < \infty \quad (110)$$

• **Gradient Control**

– *Bounded Logits*: $\exists C < \infty$ s.t.

$$\max_a |h_l(a, \mathbf{o}_l)| \leq C \quad \forall \mathbf{o}_l, l \quad (111)$$

– *Gradient Norm Bound*:

$$\mathbb{E}_{\mathbf{o}_l} [\|\nabla_{\boldsymbol{\theta}_l} h_l\|_2^2] \leq B_l < \infty \quad \forall l \quad (112)$$

• **Probability Regularity**

– *Strict Positivity*: $\exists \epsilon > 0$ s.t.

$$\hat{p}_{\boldsymbol{\theta}}(\mathbf{o}_{l+1}|\mathbf{o}_l) \geq \epsilon \quad \mu\text{-a.e. } \forall l \quad (113)$$

– *Measure Consistency*:

$$\int_{\mathcal{O}} \hat{p}_{\boldsymbol{\theta}}(\mathbf{o}_{l+1}|\mathbf{o}_l) d\mathbf{o}_{l+1} = 1 \quad \forall \mathbf{o}_l, l \quad (114)$$

These conditions guarantee: 1. Leibniz rule applicability; through Lipschitz continuity 2. Fubini’s theorem validity via measure consistency; 3. Gradient dominance via bounded logits; 4. Policy smoothness via C^2 differentiability; 5. Numerical stability through strict positivity.

Lemma I.5 (Policy Gradient for PO (Eq.(66)) under TMC). *Let $\boldsymbol{\theta}^*$ be the base model in Eq.(2) that exact predicts the distribution of Multi-task TMC as in Def. 2.1 and 2.2, and $\boldsymbol{\theta}^k$ the current model to be finetuned from $\boldsymbol{\theta}^*$ for task $k \in \mathcal{T}$. Suggest the accurate $A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1})$ is available from some outer oracle, the clip operation is always active, and Eq.(65) holds. The gradient of the PO objective for task k is given by:*

$$\nabla_{\boldsymbol{\theta}^k} \mathcal{J}_{\text{PO}}(\boldsymbol{\theta}^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{t+1} \sim \hat{p}_{\boldsymbol{\theta}^k}(\cdot|\mathbf{o}_t)\}_{t=1}^{L-1}}} \left[(1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot \nabla_{\boldsymbol{\theta}^k} \log \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right], \quad (115)$$

$$\quad (116)$$

where $r_{l+1} = \frac{\hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)}{\hat{p}_{\text{old}}^k(\mathbf{o}_{l+1}|\mathbf{o}_l)}$. By the condition that the clip operation is always active, we have

$$\mathcal{J}_{\text{PO}} = \mathbb{E}_{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k), (\mathbf{Q}, \mathbf{A}) \sim D_{\mathbf{a}_{\mathbf{o}_1}}^{\mathbf{o}_1, k}, \mathbf{o}_{2:L} \sim \hat{p}_{\boldsymbol{\theta}^k}(\cdot|\mathbf{o}_1)} \left[\frac{1}{L} \sum_{l=1}^{L-1} (1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \right], \quad (117)$$

Proof. It holds that

$$\begin{aligned} \nabla_{\boldsymbol{\theta}^k} \mathcal{J}_{\text{PO}} &= \nabla_{\boldsymbol{\theta}^k} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k \\ \mathbf{o}_{t+1} \sim \hat{p}_{\boldsymbol{\theta}^k}}} \left[\frac{1}{L} \sum_{l=1}^{L-1} (1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \right] \\ &\stackrel{(1)}{=} \frac{1}{L} \sum_{l=1}^{L-1} \nabla_{\boldsymbol{\theta}^k} \int_{\mathcal{O}^L} (1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) P^k(\mathbf{o}_1) \prod_{t=1}^{L-1} \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) d\mathbf{o}_{1:L} \\ &\stackrel{(2)}{=} \frac{1}{L} \sum_{l=1}^{L-1} \int_{\mathcal{O}^L} (1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot P^k(\mathbf{o}_1) \prod_{t=1}^{L-1} \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{t+1}|\mathbf{o}_t) \cdot \nabla_{\boldsymbol{\theta}^k} \log \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) d\mathbf{o}_{1:L} \\ &\stackrel{(3)}{=} \frac{1}{L} \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k \\ \mathbf{o}_{t+1} \sim \hat{p}_{\boldsymbol{\theta}^k}}} \left[(1 + (2\mathbb{1}(A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1)\epsilon_{\text{clip}}) A_{l+1}^{\hat{p}_{\boldsymbol{\theta}^k, k}}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot \nabla_{\boldsymbol{\theta}^k} \log \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) \right] \end{aligned} \quad (118)$$

The methodologies follow Lemma I.5.

Step (1) expands the expectation using the MDP factorization $P^k(\mathbf{o}_1) \prod_{t=1}^{L-1} \hat{p}_{\boldsymbol{\theta}^k}(\mathbf{o}_{t+1}|\mathbf{o}_t)$, noting that \hat{p}_{old} is treated as fixed behavioral policy;

Step (2) applies parameter-localized differentiation through:

$$\nabla_{\theta^k} \prod_{t=1}^{L-1} \hat{p}_t = \prod_{\substack{t=1 \\ t \neq l}}^{L-1} \hat{p}_{\text{old},t} \cdot \nabla_{\theta^k} \hat{p}_l \quad (119)$$

with conditions similar in Lemma I.2.

Step (3) reconstructs the expectation by recognizing $\prod_{t=1}^{L-1} \hat{p}_t = \hat{p}_{\theta^k}(\mathbf{o}_{1:L})/P^k(\mathbf{o}_1)$, leveraging the Markov property.

Key Conditions Inherited from REINFORCE/RAFT in Lemma I.2: 1. *Parameter Isolation:* $\theta^k = \bigsqcup_{l=1}^{L-1} \theta_l^k$ with disjoint subparameters 2. *Policy Smoothness:* $\hat{p}_{\theta^k} \in C^2(\Theta)$ with bounded Hessians 3. *Measure Consistency:* $\prod_{t \neq l} \int \hat{p}_t d\mathbf{o}_{t+1} = 1$ μ -a.e. 4. *Advantage Regularity:* $A_{l+1}^{\hat{p}_{\theta^k}, k}(\mathbf{o}_l, \mathbf{o}_{l+1})$ is $\sigma(\mathbf{o}_{1:l+1})$ -measurable and bounded.

□

Based on the policy gradient results, the logit update lemma is provided as below.

Lemma I.6. *Let θ^* be the base model in Eq.(2) that exact behave like a Multi-task TMC as in Def. 2.1 and 2.2, and θ^k the current model to be finetuned from θ^* for task $k \in \mathcal{T}$. Then*

$$\nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (\nabla_{\theta^k} h_{\theta}(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_{\theta}(\mathbf{o}'_{l+1}, \mathbf{o}_l)) \right], \quad (120)$$

$$\nabla_{\theta^k} \mathcal{J}_{\text{RAFT}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) (\nabla_{\theta^k} h_{\theta}(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_{\theta}(\mathbf{o}'_{l+1}, \mathbf{o}_l)) \right], \quad (121)$$

$$\nabla_{\theta^k} \mathcal{J}_{\text{PO}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[(1 + 2\mathbb{1}(A_{l+1}^{\hat{p}_{\theta^k}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1) \epsilon_{\text{clip}} A_{l+1}^{\hat{p}_{\theta^k}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot (\nabla_{\theta^k} h_{\theta}(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_{\theta}(\mathbf{o}'_{l+1}, \mathbf{o}_l)) \right]. \quad (122)$$

Further, we have

$$\nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}) \right], \quad (123)$$

$$\nabla_{\theta^k} \mathcal{J}_{\text{RAFT}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (1 + \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l)) (e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}) \right], \quad (124)$$

$$\nabla_{\theta^k} \mathcal{J}_{\text{PO}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[(1 + 2\mathbb{1}(A_{l+1}^{\hat{p}_{\theta^k}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \geq 0) - 1) \epsilon_{\text{clip}} A_{l+1}^{\hat{p}_{\theta^k}, k}(\mathbf{o}_l, \mathbf{o}_{l+1}) \cdot (e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}) \right]. \quad (125)$$

Proof. By Eq. (94) and Eq.(89), we have

$$\begin{aligned} \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right] \\ &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (\nabla_{\theta^k} h_{\theta}(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_{\theta}(\mathbf{o}'_{l+1}, \mathbf{o}_l)) \right]. \end{aligned} \quad (126)$$

Similarly By Eq. (96), Eq.(89), we obtain the results of RAFT. Given Eq.(2), we have

$$\begin{aligned} \nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[\nabla_{\theta^k} \log \hat{p}_{\theta^k}(\mathbf{o}_{l+1}|\mathbf{o}_l) R_{\text{out}}^k(\mathbf{o}) \right] \\ &= \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (\nabla_{\theta^k} h_{\theta}(\mathbf{o}_{l+1}, \mathbf{o}_l) - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) \nabla_{\theta^k} h_{\theta}(\mathbf{o}'_{l+1}, \mathbf{o}_l)) \right]. \end{aligned} \quad (127)$$

Given Eq.(2), by Eq.(91) we have

$$\nabla_{\theta^k} \mathcal{J}_{\text{REINFORCE}}(\theta^k) = \sum_{l=1}^{L-1} \mathbb{E}_{\substack{\mathbf{o}_1 \sim P^k(\mathcal{Q}^k) \\ \{\mathbf{o}_{l+1} \sim \hat{p}_{\theta^k}(\cdot|\mathbf{o}_l)\}_{l=1}^{L-1}}} \left[R_{\text{out}}^k(\mathbf{o}) \cdot (e_{\mathbf{o}_{l+1}, \mathbf{o}_l} - \sum_{\mathbf{o}'_{l+1} \in S_{l+1}} \hat{p}_{\theta^k}(\mathbf{o}'_{l+1}|\mathbf{o}_l) e_{\mathbf{o}'_{l+1}, \mathbf{o}_l}) \right].$$

The results of RAFT and PO follows.

□