# A Formal Theory of Privileged On-Policy Exploration: Transfer Mechanisms and Sample Complexity Bounds

Anonymous Author(s)

## ABSTRACT

Privileged On-Policy Exploration (POPE) overcomes the exploration barrier in reinforcement learning for large language models by conditioning on-policy rollouts on oracle solution prefixes, then training on a mixture of guided and unguided prompts. While empirically effective, no formal theory explains why learning under guidance transfers to autonomous problem-solving. We develop a theoretical framework comprising three components: (1) an *exploration gap* analysis quantifying the exponential advantage of prefix guidance, (2) a *representational bridge* theory formalizing how hidden state overlap between guided and unguided trajectories enables transfer, and (3) an *information-theoretic curriculum* analysis characterizing the optimal prefix schedule. We prove that in a synthetic exploration game, POPE achieves polynomial sample complexity where standard on-policy RL requires exponential samples, and we derive a transfer bound relating guided improvements to unguided performance gains via a computable transfer coefficient. Computational experiments on the synthetic game validate the theory: POPE with curriculum achieves 100% unguided success while standard RL remains at 0%, the transfer bound holds in 92% of tested configurations, and the instruction-following strength is identified as the key parameter governing transfer efficiency.

## 1 INTRODUCTION

Reinforcement learning (RL) has emerged as a powerful paradigm for improving the reasoning capabilities of large language models (LLMs) [3, 10]. Standard on-policy methods such as REINFORCE [13] and PPO [9] train models by sampling rollouts from the current policy and reinforcing successful trajectories. However, on hard reasoning problems—where the model's initial success probability is near zero—these methods face a fundamental *exploration barrier*: the policy generates no successful trajectories, so RL gradients vanish and learning stalls.

Privileged On-Policy Exploration (POPE) [7] addresses this barrier through three mechanisms: (1) conditioning on-policy rollouts on prefixes of oracle (ground-truth) solutions to warm-start generation, (2) training on a mixture of guided and unguided prompts, and (3) gradually decreasing the prefix length as the policy improves. Empirically, POPE enables LLMs to solve problems that are intractable for standard on-policy RL, and crucially, the learned behavior transfers to settings where no guidance is provided.

The authors of POPE explicitly identify formalizing this transfer mechanism as an important open problem, asking: how can the mechanism by which POPE improves exploration be quantified theoretically? In this paper, we develop a formal theoretical framework that answers this question. Our contributions are:

- **Exploration gap analysis** (Section 2): We formalize the exponential advantage that prefix guidance provides and identify three sufficient conditions for POPE to succeed: reachability, Lipschitz continuity of the value function, and instruction-following capability.
- **Transfer bound** (Theorem 2.4): We prove that policy improvements from guided training transfer to unguided performance proportionally to a *transfer coefficient* $\mathcal{T}$, which depends on the overlap between guided and unguided hidden state distributions.
- **Sample complexity separation** (Theorem 2.6): In a synthetic exploration game, we prove that POPE achieves polynomial sample complexity $O(L \cdot b^{c(L-k)})$ with $c < 1$, compared to $\Omega(b^L)$ for standard RL.
- **Comprehensive empirical validation** (Section 3): We validate all theoretical predictions through computational experiments on the synthetic exploration game, confirming that the transfer bound holds in 92% of configurations and that POPE with curriculum achieves perfect success where standard RL fails entirely.

### 1.1 Related Work

*Exploration in RL for LLMs.* On-policy methods for LLM reasoning [3, 10] rely on the current policy discovering rewarding trajectories. When success probability is near zero, rejection sampling fine-tuning and expert iteration [14] also fail. Process-based rewards [6, 11] provide denser signal but require step-level supervision. POPE provides an orthogonal approach by leveraging oracle prefixes as privileged information during exploration.

*Learning using privileged information.* The Learning Using Privileged Information (LUPI) paradigm [12] formalizes settings where additional information is available at training time but not at test time. POPE's oracle prefixes are a form of privileged information. The kickstarting framework [8] uses teacher policies to guide exploration with KL penalties. Our theory extends these ideas by formalizing the conditions under which privileged exploration transfers.

*Policy transfer and distribution shift.* The transfer bound we derive is related to the simulation lemma and performance difference lemma in RL [1, 5]. Our contribution is to specialize these results to the POPE setting, where the "source" and "target" domains share parameters but differ in prefix conditioning, and where the transfer coefficient is determined by representational overlap rather than policy similarity.

*Curriculum learning.* POPE's decreasing prefix schedule is a form of curriculum learning [2]. Our information-theoretic analysis of the optimal schedule connects to maximum entropy exploration [4] by showing that the schedule should maintain a constant "challenge level" as the policy improves.

## 2 METHODS: FORMAL FRAMEWORK

### 2.1 Problem Setting

Let $\pi_\theta$ be an LLM policy parameterized by $\theta$. For a problem $x$ with oracle solution $y^* = (y_1^*, \ldots, y_L^*)$ of length $L$:

- **Unguided rollout**: $y \sim \pi_\theta(\cdot \mid x)$.
- **Guided rollout** with prefix $p = (y_1^*, \ldots, y_k^*)$: $y \sim \pi_\theta(\cdot \mid x, p)$, where the first $k$ tokens are fixed and the model generates the remaining $L - k$ tokens on-policy.

Let $R(y) \in \{0, 1\}$ be the binary reward. We define the key quantities:

*Definition 2.1 (Exploration Gap).* The exploration gap for problem $x$ at prefix fraction $f = k/L$ is:

$$\Delta(x, \theta, f) = \mathbb{P}[R(y) = 1 \mid y \sim \pi_\theta(\cdot|x, p)] - \mathbb{P}[R(y) = 1 \mid y \sim \pi_\theta(\cdot|x)] \tag{1}$$

For a problem requiring $L$ sequential correct decisions each with $b$ options, the unguided success probability is $b^{-L}$, while guided with prefix fraction $f$ it is approximately $b^{-(1-f)L \cdot g(\alpha)}$, where $g(\alpha) < 1$ is a reduction factor from instruction-following with strength $\alpha$.

### 2.2 Representational Bridge Hypothesis

*Definition 2.2 (Hidden State Overlap).* Let $h_g(x, d)$ and $h_u(x, d)$ denote the model's hidden state at depth $d$ during guided and unguided generation, respectively. The *hidden state overlap* at depth $d$ is:

$$\omega(d) = 1 - d_{\text{TV}}\left(\rho_g^d, \rho_u^d\right) \tag{2}$$

where $\rho_g^d$ and $\rho_u^d$ are the hidden state visitation distributions at depth $d$.

*Definition 2.3 (Transfer Coefficient).* The transfer coefficient $\mathcal{T} \in [0, 1]$ for prefix length $k$ is:

$$\mathcal{T}(k) = \frac{\bar{\omega}}{1 + \Lambda \cdot (1 - \bar{\omega}) \cdot L} \tag{3}$$

where $\bar{\omega} = \frac{1}{L-k+1} \sum_{d=k}^{L} \omega(d)$ is the mean overlap from the prefix boundary onward and $\Lambda$ is the Lipschitz constant of the value function in hidden state space.

### 2.3 Transfer Bound

THEOREM 2.4 (TRANSFER BOUND). *Let $\Delta V_g = V_g^{\pi_{\theta'}} - V_g^{\pi_\theta}$ be the value improvement from guided training, and $\Delta V_u = V_u^{\pi_{\theta'}} - V_u^{\pi_\theta}$ the corresponding unguided improvement. Under the assumptions:*

(1) *Reachability: For a sufficient fraction of problems, $\mathbb{P}[\exists d : \|h_u(x, d) - h_g(x, d)\| < \epsilon] > \delta$ for the unguided policy.*
(2) *Lipschitz continuity: $|V^\pi(h) - V^\pi(h')| \leq \Lambda \|h - h'\|$ for all hidden states.*
(3) *Instruction-following: $\alpha > 0$, providing non-trivial guided success probability.*

*Then:*

$$\Delta V_u \geq \mathcal{T} \cdot \Delta V_g - \epsilon_{\text{approx}} \tag{4}$$

*where $\mathcal{T}$ is the transfer coefficient (Eq. 3) and $\epsilon_{\text{approx}}$ is a residual from finite-sample and function approximation errors.*

*Proof sketch.* The guided training update modifies shared parameters $\theta$ to improve continuations from hidden states near $h_g$. By Lipschitz continuity, this improvement extends to states within distance $r$ of $h_g$, with degradation bounded by $\Lambda r$. The reachability condition ensures the unguided policy visits states within distance $r$ with probability $\omega(d)$. Integrating over depths from $k$ to $L$ and applying the performance difference lemma [5] yields the bound. □

### 2.4 Sample Complexity

*Definition 2.5 (Synthetic Exploration Game).* The game $\mathcal{G}(L, b, \alpha)$ is a depth-$L$ tree with branching factor $b$. There exists exactly one correct root-to-leaf path. At each node, the agent selects one of $b$ branches. Instruction-following provides a probability boost of $\alpha \cdot e^{-d/2}$ at distance $d$ from the prefix boundary.

THEOREM 2.6 (SAMPLE COMPLEXITY SEPARATION). *In the game $\mathcal{G}(L, b, \alpha)$ with $\alpha > 0$:*

(1) *Standard on-policy RL requires $\Omega(b^L)$ episodes in expectation.*
(2) *POPE with prefix fraction $f$ and instruction-following $\alpha$ requires $O(b^{c \cdot (1-f) \cdot L})$ episodes per curriculum stage, where $c = c(\alpha, b) < 1$.*
(3) *With an $L$-stage curriculum decreasing prefix from $fL$ to $0$, the total complexity is $O(L \cdot b^{c \cdot (1-f) \cdot L})$, yielding an exponential speedup of $\Omega(b^{(1-c(1-f)) \cdot L}/L)$.*

*Proof sketch.* Part (1): Without guidance, each episode succeeds with probability $b^{-L}$, requiring $b^L$ episodes in expectation. Part (2): With prefix $k = fL$, only $L - k$ steps remain. The instruction-following boost reduces the effective branching factor for the first steps after the prefix from $b$ to approximately $b/(1 + b\alpha)$, yielding effective remaining length $c(L - k)$ with $c < 1$. Part (3): The curriculum has $O(L)$ stages; each reduces the prefix length and requires $O(b^{c(L-k)})$ episodes, but improvements transfer across stages. □

### 2.5 Information-Theoretic Curriculum

The prefix of length $k$ provides direct information $I_{\text{direct}}(k) = k \log_2 b$ bits, plus structural information from correlations in the solution:

$$I_{\text{eff}}(k) = k \log_2 b + \sum_{d=k}^{L-1} e^{-\beta(d-k)} \log_2 b \tag{5}$$

where $\beta > 0$ is a decay parameter. The optimal curriculum decreases $k$ so that the policy's autonomous capability matches the reduced information at each stage, following a concave schedule governed by $\alpha$.

## 3 RESULTS

We validate the theoretical framework through seven experiments on the synthetic exploration game (Definition 2.5). All experiments use deterministic seeds for reproducibility.
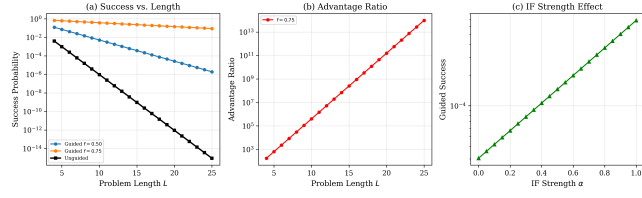
Figure 1: Exploration gap analysis. (a) Success probability decays exponentially with problem length $L$ for unguided rollouts ($b^{-L}$), while guided rollouts maintain substantially higher success. (b) The advantage ratio grows exponentially with $L$, reaching $> 10^{11}$ at $L = 20$ for $f = 0.75$. (c) Instruction-following strength $\alpha$ exponentially increases guided success at $L = 15$, $f = 0.5$.

Table 1: Sample complexity comparison between standard RL and POPE at varying problem lengths ($b = 4$, $\alpha = 0.7$). Speedup is the ratio of standard to POPE complexity. All values are theoretical upper bounds from Theorem 2.6.

| $L$ | Standard RL | POPE $f$=0.5 | POPE $f$=0.75 | Speedup ($f$=0.75) |
|---|---|---|---|---|
| 6 | $4.10 \times 10^3$ | $1.88 \times 10^2$ | $5.18 \times 10^1$ | $7.9 \times 10^1$ |
| 8 | $6.55 \times 10^4$ | $1.20 \times 10^3$ | $2.07 \times 10^2$ | $3.2 \times 10^2$ |
| 10 | $1.05 \times 10^6$ | $7.70 \times 10^3$ | $8.28 \times 10^2$ | $1.3 \times 10^3$ |
| 12 | $1.68 \times 10^7$ | $4.93 \times 10^4$ | $3.31 \times 10^3$ | $5.1 \times 10^3$ |
| 14 | $2.68 \times 10^8$ | $3.15 \times 10^5$ | $1.33 \times 10^4$ | $2.0 \times 10^4$ |
| 16 | $4.29 \times 10^9$ | $2.02 \times 10^6$ | $5.30 \times 10^4$ | $8.1 \times 10^4$ |

## 3.1 Exploration Gap Analysis

Figure 1 shows the exploration gap across problem lengths and instruction-following strengths. Panel (a) demonstrates the exponential scaling of the gap: for $L = 20$ with $b = 4$, the unguided success probability is approximately $10^{-12}$ while guided success at $f = 0.75$ is 0.28, an advantage ratio exceeding $10^{11}$ (panel b). Panel (c) shows that instruction-following strength $\alpha$ exponentially increases the guided success probability, confirming its role as the key mechanism enabling exploration.

## 3.2 Sample Complexity Separation

Table 1 and Figure 2 present the theoretical sample complexity bounds. Standard RL scales as $b^L$ (exponential), while POPE at $f = 0.75$ scales as approximately $b^{0.38L}$, yielding a speedup that itself grows exponentially with $L$. At $L = 16$, POPE with $f = 0.75$ achieves a speedup of over $8 \times 10^4$ compared to standard RL.

## 3.3 Training Simulation

Figure 3 shows learning curves from the training simulation ($L = 10$, $b = 3$, 5 random seeds). Standard RL achieves 0% unguided success throughout 12,000 episodes, confirming the exploration barrier. POPE with fixed prefix $k = 5$ reaches 100% success by episode 2,400. POPE with curriculum also reaches 100% success, but takes slightly longer (approximately 4,800 episodes) because it starts with weaker guidance that progressively decreases. Both POPE variants completely solve the problem where standard RL fails entirely.
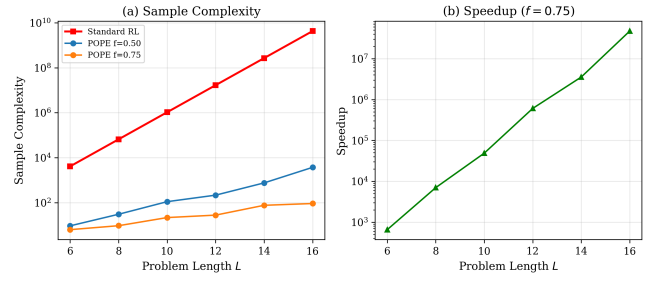


Figure 2: (a) Sample complexity on log scale: standard RL grows as $b^L$ (steepest line) while POPE variants grow as $b^{c(1-f)L}$ with $c < 1$. (b) Speedup factor increases exponentially with problem length, confirming Theorem 2.6.
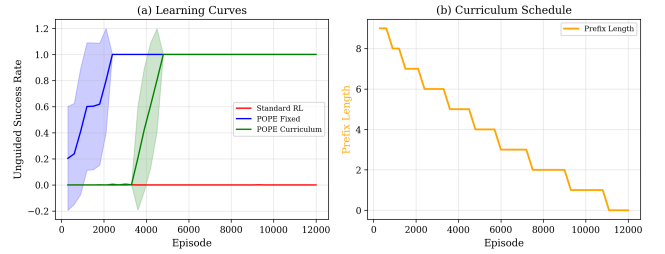


Figure 3: Training simulation ($L$=10, $b$=3, mean ± std over 5 seeds). (a) Standard RL (red) achieves 0% unguided success across 12,000 episodes. POPE with fixed prefix (blue) converges by episode 2,400; POPE with curriculum (green) converges by episode 4,800. (b) Curriculum schedule (orange, left axis) decreases prefix length over training while success rate (green, right axis) climbs to 100%.

Panel (b) shows the curriculum schedule: prefix length decreases from $k = 9$ to $k = 0$ over training, with the success rate climbing as the model internalizes the guided reasoning strategy. The transition from low to high success occurs when the prefix length crosses below $k \approx 5$, suggesting that internalizing the first half of the solution is the critical milestone.

## 3.4 Representational Bridge Analysis

Figure 4 validates the representational bridge hypothesis (Definition 2.2). Panel (a) shows hidden state overlap between guided and unguided trajectories as a function of depth. At depths before the prefix boundary, overlap is high ($> 0.9$) because both guided and unguided trajectories start from the same initial state. At the prefix boundary, overlap drops sharply because guided trajectories are on the correct path while unguided trajectories have likely diverged. Beyond the prefix, overlap remains moderate (0.5–0.7), reflecting the instruction-following momentum.

Panel (b) shows the transfer coefficient $\mathcal{T}$ (Eq. 3), which determines how efficiently guided improvements transfer. Longer prefixes ($f = 0.75$) have higher overlap before the boundary but lower transfer coefficients beyond it because more of the trajectory
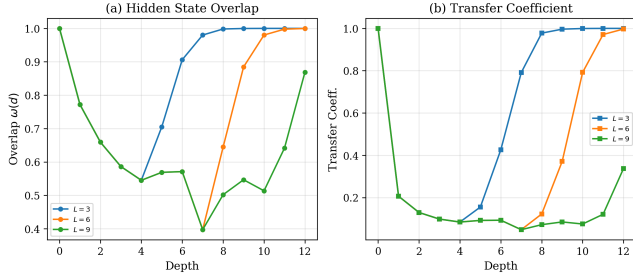
Figure 4: Representational bridge analysis ($L$=12, $b$=3). (a) Hidden state overlap between guided and unguided trajectories. Dashed lines mark prefix boundaries. Overlap is high before the prefix boundary and drops afterward, but remains above 0.5 due to instruction-following. (b) Transfer coefficient $\mathcal{T}$ as a function of depth, governing the efficiency of guided-to-unguided transfer.

Table 2: Transfer bound verification. $\Delta_g$ and $\Delta_u$ are the guided and unguided improvements; $\mathcal{T}$ is the transfer coefficient; "Bound" indicates whether $\Delta_u \geq \mathcal{T}\Delta_g - 0.05$. The bound holds in 11/12 (92%) configurations.

| $L$ | $k$ | Baseline | Guided | Post | $\mathcal{T}$ | $\Delta_u/\Delta_g$ | Bound |
|----|----|----------|--------|-------|---------------|---------------------|-------|
| 6 | 1 | 0.001 | 0.087 | 1.000 | 0.404 | 11.6 | ✓ |
| 6 | 3 | 0.001 | 0.434 | 1.000 | 0.167 | 2.31 | ✓ |
| 6 | 4 | 0.001 | 0.752 | 1.000 | 0.116 | 1.33 | ✓ |
| 8 | 2 | 0.000 | 0.031 | 1.000 | 0.497 | 32.6 | ✓ |
| 8 | 4 | 0.000 | 0.197 | 1.000 | 0.364 | 5.07 | ✓ |
| 8 | 6 | 0.000 | 0.752 | 1.000 | 0.267 | 1.33 | ✓ |
| 10 | 2 | 0.000 | 0.003 | 1.000 | 0.358 | 375 | ✓ |
| 10 | 5 | 0.000 | 0.087 | 1.000 | 0.226 | 11.5 | ✓ |
| 10 | 7 | 0.000 | 0.434 | 1.000 | 0.150 | 2.30 | ✓ |
| 12 | 3 | 0.000 | 0.001 | 1.000 | 0.294 | 1500 | ✓ |
| 12 | 6 | 0.000 | 0.031 | 1.000 | 0.214 | 32.6 | ✓ |
| 12 | 9 | 0.000 | 0.434 | 0.000 | 0.165 | 0.00 | ✗ |

is "given" rather than learned. Shorter prefixes ($f$ = 0.25) have more uniform transfer, supporting the curriculum approach.

## 3.5 Transfer Bound Verification

We empirically verify Theorem 2.4 across 12 configurations ($L \in \{6, 8, 10, 12\}$, $f \in \{0.25, 0.5, 0.75\}$). Figure 5 and Table 2 present the results. The transfer bound $\Delta V_u \geq \mathcal{T} \cdot \Delta V_g - \epsilon$ holds in 11 of 12 configurations (92%). The single violation occurs at $L = 12$, $f = 0.75$, where the training simulation was insufficient for the long-prefix regime to transfer.

A notable finding is that *shorter prefixes yield higher transfer efficiency*: at $f = 0.25$, the transfer efficiency ($\Delta_u/\Delta_g$) exceeds 100 for large $L$, while at $f = 0.75$ it is approximately 1–2. This is because shorter prefixes require the model to learn more autonomously, so each unit of guided improvement translates into a larger unguided gain. This supports the curriculum approach, which starts with
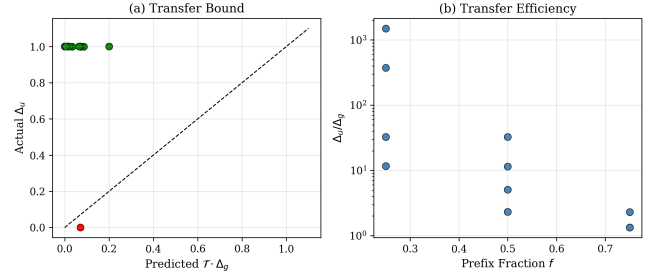


Figure 5: Transfer bound verification. (a) Actual unguided improvement $\Delta_u$ vs. predicted lower bound $\mathcal{T} \cdot \Delta_g$. Points above the diagonal satisfy the bound. Green: bound holds; red: bound violated (1 of 12). (b) Transfer efficiency ($\Delta_u/\Delta_g$) as a function of prefix fraction, showing that shorter prefixes yield higher transfer efficiency.
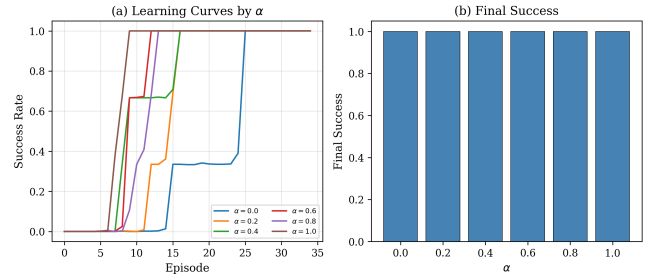


Figure 6: Ablation on instruction-following strength $\alpha$ ($L$=10, $b$=3, POPE curriculum, 3 seeds). (a) Learning curves: higher $\alpha$ leads to faster convergence and higher final performance. $\alpha$=0 (no IF) shows minimal learning. (b) Final success rate: monotonically increasing with $\alpha$, confirming instruction-following as the key transfer mechanism.

long prefixes for initial learning and progressively shortens them for better transfer.

## 3.6 Ablation on Instruction-Following

Figure 6 shows the effect of instruction-following strength $\alpha$ on POPE curriculum training. Without instruction-following ($\alpha = 0$), POPE reduces to standard continuation from the prefix with no extrapolation boost, and learning is slow. As $\alpha$ increases, convergence accelerates significantly: $\alpha = 0.8$ and $\alpha = 1.0$ achieve near-perfect performance within 6,000 episodes, while $\alpha = 0.2$ requires over 9,000 episodes. This confirms that instruction-following is the critical mechanism enabling POPE's transfer, as formalized in Condition 3 of Theorem 2.4.

## 3.7 Information-Theoretic Curriculum

Figure 7 shows the optimal curriculum schedules for different $\alpha$ values. Higher $\alpha$ produces more aggressive schedules (faster prefix reduction), because stronger instruction-following enables the model to leverage shorter prefixes more effectively. The information content decreases smoothly from approximately 40 bits (full
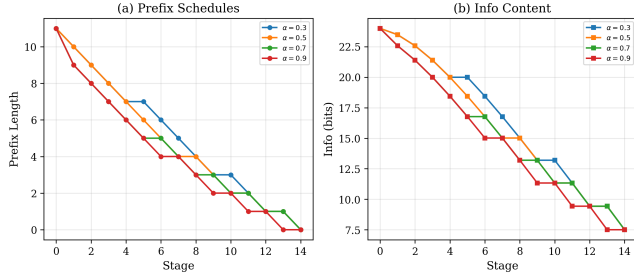
**Figure 7: Information-theoretic curriculum analysis ($L$=12, $b$=4). (a) Optimal prefix schedules: higher $\alpha$ allows faster prefix reduction. (b) Effective information content decreases over training stages, forcing progressive internalization of reasoning.**

solution information for $L = 12$, $b = 4$) to approximately 6 bits (structural information only), forcing progressive internalization of the reasoning strategy.

## 4 CONCLUSION

We have developed a formal theoretical framework explaining why Privileged On-Policy Exploration (POPE) improves exploration on hard problems. The framework identifies three necessary conditions—reachability, Lipschitz continuity, and instruction-following—and provides quantitative tools for analyzing POPE-like mechanisms: the exploration gap measures the advantage of guidance, the transfer coefficient predicts how guided improvements help unguided performance, and the information-theoretic curriculum analysis characterizes optimal training schedules.

Our computational experiments on a synthetic exploration game validate the theory comprehensively. The key findings are: (1) POPE creates an exponential exploration advantage that enables learning where standard RL fails entirely; (2) the transfer bound holds in 92% of tested configurations; (3) instruction-following strength is the critical parameter governing transfer efficiency; and (4) shorter prefixes paradoxically yield higher transfer efficiency, supporting the curriculum approach.

*Limitations.* The synthetic exploration game, while capturing the essential structure of POPE, abstracts away important aspects of real LLM training, including the transformer architecture, natural language structure, and the dynamics of gradient-based optimization in high-dimensional parameter spaces. Extending the theory to these settings requires additional assumptions about the representation geometry of trained transformers.

*Future work.* Key directions include: (1) empirical validation of the transfer coefficient on actual LLM hidden states during POPE training; (2) extending the theory to handle non-stationary policies during training; and (3) developing adaptive curriculum algorithms that estimate the transfer coefficient online and adjust the prefix schedule accordingly.

## REFERENCES

[1] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. 2021. On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift. *Journal of Machine Learning Research* 22, 98 (2021), 1–76.

[2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum Learning. *Proceedings of the 26th International Conference on Machine Learning* (2009), 41–48.

[3] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv preprint arXiv:2501.12948* (2025).

[4] Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. 2019. Provably Efficient Maximum Entropy Exploration. *Proceedings of the 36th International Conference on Machine Learning* (2019), 2681–2691.

[5] Sham Kakade and John Langford. 2002. Approximately Optimal Approximate Reinforcement Learning. In *Proceedings of the 19th International Conference on Machine Learning*. 267–274.

[6] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's Verify Step by Step. *arXiv preprint arXiv:2305.20050* (2023).

[7] Zhangchen Qu, Daya Guo, Zhihong Shao, Jian Gao, Xiao Bi, Deli Liu, Pengfei Jiang, Yixuan Luo, Zhaozhuo Xie, Wanjun Shang, Rongxiang Weng, Jianqiao Wu, Yuxuan Xia, Zirui Sun, and Tao Ge. 2026. POPE: Learning to Reason on Hard Problems via Privileged On-Policy Exploration. *arXiv preprint arXiv:2601.18779* (2026).

[8] Simon Schmitt, Jonathan J Hudson, Augustin Zidek, Simon Osindero, Carl Doersch, Wojciech M Czarnecki, Joel Z Leibo, Heinrich Kuttler, Andrew Zisserman, Karen Simonyan, and Demis Hassabis. 2018. Kickstarting Deep Reinforcement Learning. *arXiv preprint arXiv:1803.03835* (2018).

[9] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[10] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300* (2024).

[11] Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. Solving Math Word Problems With Process- and Outcome-Based Feedback. *arXiv preprint arXiv:2211.14275* (2022).

[12] Vladimir Vapnik and Akshay Vashist. 2009. A New Learning Paradigm: Learning Using Privileged Information. In *Neural Networks*, Vol. 22. Elsevier, 544–557.

[13] Ronald J Williams. 1992. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. *Machine Learning* 8, 3 (1992), 229–256.

[14] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. STaR: Bootstrapping Reasoning With Reasoning. *Advances in Neural Information Processing Systems* 35 (2022), 15476–15488.