UPSKILL: MUTUAL INFORMATION SKILL LEARNING FOR STRUCTURED RESPONSE DIVERSITY IN LLMS

Anonymous authors

000

001

002003004

006

008

010 011

012

013

014

015

016

017

018

019

021

023

025 026 027

028 029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Reinforcement Learning with Verifiable Rewards (RLVR) has improved the reasoning abilities of large language models (LLMs) on mathematics and programming tasks, often by maximizing pass@1 correctness. However, optimizing singleattempt accuracy can inadvertently suppress response diversity across repeated attempts, narrowing exploration and overlooking underrepresented strategies. We introduce UpSkill, a training time method that adapts Mutual Information Skill Learning (MISL) to LLMs to induce structured response diversity: a discrete latent z selects a reproducible "strategy" that steers the token distribution toward distinct modes. We propose a novel reward that we implement within Group Relative Policy Optimization (GRPO): a token-level mutual information (MI) reward that encourages trajectory specificity to z. Experiments on GSM8K with three open-weight models, Llama 3.1-8B, Qwen 2.5-7B, and R1-Distilled-Qwen2.5-Math-1.5B show that UpSkill improves multi-attempt metrics, yielding median gains of \sim 4% in pass@k and \sim 7% in consensus@k without degrading pass@1. Additionally, we prove that improvements in pass@k are closely tied to the mutual information objective, providing a theoretical justification for UpSkill.

1 Introduction

LLMs excel at verifiable reasoning tasks such as mathematical problem solving and code generation (Guo et al., 2025). However, repeated sampling often yields highly similar outputs (Shaier et al., 2025). This is detrimental in multiattempt settings where just one correct completion solves the problem at hand, such as code generation with tests (Chen et al., 2021) or formal proofs in Lean (Trinh et al., 2024), as a lack of diversity reduces the effective number of independent attempts. Therefore, for these or other objectives evaluated by pass@k, or the probability that at least one of k completions will be correct, more deterministic output decreases the chance that some sampled attempt will succeed. Furthermore, recent work has found that post-training that optimizes single-attempt correctness suppresses response variation across attempts (Chen et al., 2025; Dang et al., 2025), creating a discrepancy between how models are trained and how they are used and evaluated.

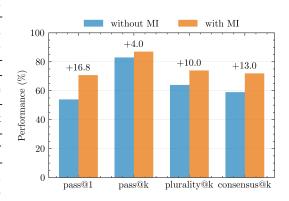


Figure 1: Token-level MISL improves multiattempt accuracy without hurting single-attempt accuracy on GSM8K for the Qwen 2.5-7B model (See Sec. 5.2)

The challenge of balancing diversity and accuracy, or exploration and exploitation (Sutton & Barto, 2015), has primarily been studied in prior works that change how decoding is done. Methods such as temperature sampling (Renze & Guven, 2024), nucleus sampling (Holtzman et al., 2020), and prompt perturbations (Shur-Ofry et al., 2024) can inject variety, but they require manual tuning (Du et al., 2025) and are brittle across domains (Shi et al., 2024; Qiang et al., 2024). Separately, prior training-time methods do not expose a controllable basis of strategies and require a complex training regime

Input

If Raymond does half as much laundry as Sarah, and Sarah does 4 times as much laundry as David, calculate the difference in the amount of laundry Raymond and David do if Sarah does 400 pounds of laundry.



Sarah does 400 pounds of laundry. David does 1/4 the laundry Sarah does, so he does 400 / 4 = 100 pounds of laundry. Raymond does 1/2 the laundry Sarah does, so he does 400 / 2 = 100 pounds of laundry. The difference is 200-100=100 pounds of

Answer: 100 pounds

laundry.

Let D be the amount of laundry David does and let R be the amount of laundry Raymond does.

Since Sarah does 4 times as much laundry as David, and Raymond does half as much laundry as Sarah, we know that 4D = 400 and R = 400 / 2

The difference is:

R - D = 400 / 2 - 400 / 4 = 200 - 100 = 100

Answer: 100 pounds



Figure 2: UpSkill is an unsupervised method for training LLMs to produce diverse responses. After training, different latent vectors z (blue boxes above) correspond to different response strategies. Because of space constraints, the figure shows summarized responses from UpSkill; we report the full responses in Appendix I.

to properly balance exploration and exploitation (Tang et al., 2025; Chen et al., 2025). We seek a training-time mechanism that (i) increases diversity in a controlled manner, (ii) produces semantically distinct and reproducible modes of reasoning, and (iii) preserves single-attempt verifiable accuracy.

We introduce UpSkill, a training-time approach that induces *structured response diversity* without prompt engineering. The key idea behind UpSkill is to introduce an input token z that structures the response, so that different values of z correspond to different responses. Formally, we will model LLM attempts on verifiable reasoning tasks as a token-level Markov decision process. We can then adopt prior work from reinforcement learning on learning *skills*, which learn a policy conditioned on a latent variable z. These methods (Eysenbach et al., 2019; Gregor et al., 2016a; Achiam et al., 2018; Sharma et al., 2020a; Florensa et al., 2017) include a loss term that maximizes the mutual information between z and the policy's behavior. Precisely, we adapt the CSF method (Zheng et al., 2024) to LLMs: the model conditions its response on a discrete latent $z \in \{1, \ldots, N\}$, and training encourages behaviors whose distribution depends strongly on z. Intuitively, each z should correspond to a reproducible strategy, and the set of strategies should span a broad range of behaviors.

The main contribution of our paper is a method for training LLMs to produce diverse responses. Our method implements mutual information skill learning by applying GRPO (Shao et al., 2024) with a novel reward term: a token-level mutual information reward, which encourages diversity in completions. Finally, we sketch a theoretical link between $\mathcal{I}(\tau; z \mid x)$ and passek: the improvement of passek after training is related to $\mathcal{I}(\tau; z \mid x)$, In summary, our contributions are as follows:

- UpSkill achieves median gains of +4% in pass@k and +7% in consensus@k on GSM8K across three open-weight models using RL fine-tuning with LoRA adapters on 2,000 problems, with preserved pass@1 accuracy.
- In an arithmetic puzzle environment, UpSkill improves pass@5 by +10% by mitigating response variation collapse and developing a collection of diverse and complementary skills.
- We prove that pass@k improvement closely corresponds to the mutual information $\mathcal{I}(\tau;z\mid x)$, showing that large improvements in multi-attempt accuracy require and are limited by sufficient mutual information.
- We provide an effective and reproducible method for token-level MI, and will release an open-source implementation focused on practical performance.

2 BACKGROUND AND RELATED WORK

2.1 Multi-attempt evaluation, redundancy, and why diversity matters

For verifiable tasks, we often consider the probability of success across *multiple* completions rather than a single attempt (Chen et al., 2025). Let x denote the input and τ a sampled completion from policy $\pi(\cdot \mid x)$. Let $Y(\tau) \in \{0,1\}$ indicate correctness under a deterministic verifier. For k attempts, the standard metric

pass@
$$k(x) = 1 - \Pr\left(\bigcap_{i=1}^{k} \{Y(\tau_i) = 0\} \mid x\right)$$
 (1)

is the complement of the joint failure probability across k i.i.d. draws $\tau_{1:k} \sim \pi(\cdot \mid x)$ (Chen et al., 2021). Letting $p = \Pr(Y(\tau) = 1 \mid x)$, we therefore have $\Pr(\text{pass}@k(x)) = 1 - (1 - p)^k$.

In practice, identical prompts with fixed decoding hyperparameters can yield strongly correlated trajectories, especially for deterministic or near-deterministic samplers (Vijayakumar et al., 2018). A useful lens is to consider an "effective number of attempts" $k_{\rm eff}$ that discounts k by a correlation term (analogous to design effects in sampling) (Kish, 1965). If completions have pairwise correlation ρ in the binary success indicators, a heuristic adjustment gives $k_{\rm eff} \approx k/(1+(k-1)\rho)$: as $\rho \to 1$, additional attempts contribute little; as $\rho \to 0$, $k_{\rm eff} \to k$. Although crude, this highlights the central point: reducing dependence among attempts is as important as raising per-attempt accuracy. Structured diversity aims to decrease redundancy so that the joint failure probability decreases faster in k. For Gaussian random variables, correlation and mutual information are closely related (as intuitively, correlated variables have information on each other) (Krafft, 2013). However, as text correctness cannot easily be framed as a Gaussian distribution, mutual information is more a natural measurement.

Beyond pass@k, plurality@k and consensus@k measure agreement among completions, examining robustness and internal consistency of the model's reasoning (Wallace et al., 2025). In many workflows, agreement acts as a proxy for confidence while still benefiting from diversity to escape shared failure modes (Hochlehnert et al., 2025).

2.2 RL on Language models: Token-Level MDPs, RLVR, and GRPO

Autoregressive LLMs can be cast as policies over a Markov decision process (MDP), where the state is the token prefix and the action is the next token (Bahdanau et al., 2017; Ouyang et al., 2022). This setup, often referred to as the token-level MDP (Zhong et al., 2025), allows reinforcement learning algorithms to directly optimize model behavior for correctness on verifiable tasks such as math or code.

Reinforcement Learning from Verifiable Rewards (RLVR) leverages automatically checkable signals (e.g., exact numeric answers, unit tests) as rewards, with the goal being to improve the pass rate of policies while ensuring the new policy remains close to a base model (Liu et al., 2023; Dou et al., 2024). Let π_{θ} denote the trainable policy and $\pi_{\rm base}$ a frozen reference. A common form of the per-trajectory reward is

$$r_{\text{RLVR}}(\tau) = r_{\text{correctness}}(\tau) - \beta D_{\text{KL}}(\pi_{\theta}(\cdot \mid x) \| \pi_{\text{base}}(\cdot \mid x)),$$
 (2)

where $\beta > 0$ controls deviation from the base model (Xiong et al., 2024).

Group Relative Policy Optimization (GRPO) (Shao et al., 2024) adapts PPO-style updates to reasoning by sampling multiple completions per prompt x as a group. Within-group baselines reduce variance and increase the relative difference between completion rewards. Concretely, for each x one draws C trajectories $\{\tau_i\}_{i=1}^C$, computes verifiable rewards and a group baseline (e.g., a rank or meannormalized signal), and updates π_θ with clipped policy ratios as in PPO (Schulman et al., 2017). GRPO typically improves passel on math/code under RLVR (Shao et al., 2024). However, absent any explicit term for diversity, it can reduce variation across attempts as the policy sharpens around locally high-reward regions (Dang et al., 2025).

As some intuition for this distribution change, suppose that the model is attempting to predict the correct answer in a setting where it believes that the answer is Yes with probability 70% and No with probability 30%. Cross-entropy loss encourages a model to predict the correct distribution of

70% Yes and 30% No; on the other hand, GRPO training would cause the model to collapse its output distribution towards predicting 100% Yes, as it maximizes the pass@1. Empirically, this can shrink the entropy of the completion distribution and heighten redundancy among attempts, limiting pass@k improvements even as pass@1 increases (Dang et al., 2025).

2.3 MUTUAL INFORMATION AND SKILL DISCOVERY

Maximizing mutual information (MI) between latent variables and observed behavior has been a recurring tool for learning structured, controllable representations (Tishby et al., 2000; Kingma & Welling, 2013; Stratos & Wiseman, 2020).

In generative modeling, InfoGAN (Chen et al., 2016) augments GAN training with a variational lower bound on $\mathcal{I}(c;x)$ to make latent codes c predictably control semantic factors (e.g., stroke thickness for MNIST). In variational autoencoders, InfoVAE (Zhao et al., 2018b) adds an explicit MI term to counteract posterior collapse and preserve informative latents even with expressive decoders.

In sequential decision making, MI has been used to discover diverse, reusable behaviors without external rewards. Early work such as VIC (Gregor et al., 2016b) and DIAYN (Eysenbach et al., 2018) maximizes $\mathcal{I}(s;z)$ or $\mathcal{I}(\tau;z)$, encouraging skills z whose rollouts visit different parts of state or trajectory space and remain identifiable from observations. InfoGAIL (Li et al., 2017) extends this to imitation learning by maximizing MI between a latent intention and trajectories to capture multimodal expert behavior. Subsequent methods bias the MI objective toward long-horizon distinctiveness to avoid trivial short-term variation (Sharma et al., 2020b; Hansen et al., 2021). Additional related work on MI is available in Appendix B.

Unsupervised skill discovery in RL can be viewed as maximizing the MI between a latent "skill" variable and observed trajectories (Gregor et al., 2016b; Eysenbach et al., 2018). Let $z \in \mathcal{Z}$ index a skill and let τ denote a trajectory. These methods maximize

$$\max_{\pi} \mathcal{I}(\tau; z \mid x) = \mathbb{E}\left[\log p_{\pi}(\tau \mid x, z) - \log p_{\pi}(\tau \mid x)\right] = \mathcal{H}(\tau \mid x) - \mathcal{H}(\tau \mid x, z). \tag{3}$$

This decomposition clarifies the pressure on the policy: (i) to increase marginal entropy $\mathcal{H}(\tau \mid x)$ so that trajectories cover more of the solution space; and (ii) to decrease conditional entropy $\mathcal{H}(\tau \mid x, z)$ so that each z induces a reproducible, stable mode. The net effect is a set of distinct, consistent behaviors indexed by z that together span diverse solution strategies.

Our setting is closest in spirit to unsupervised skill discovery (e.g., DIAYN, VIC) and to mutual information-based skill learning (Zheng et al., 2024), but differs in applying these techniques to language models with RLVR training. We develop an approach for maximizing MI tailored to LLM reasoning, and also connect pass@k performance with the mutual information objective.

2.4 OTHER TECHNIQUES FOR RESPONSE DIVERSIFICATION

Beyond MI-based training, several other approaches aim to increase output diversity.

At inference time, decoding-time diversification alters sampling: increasing temperature, switching to nucleus/top-k sampling (Holtzman et al., 2020; Fan et al., 2018), or perturbing prompts (Qiang et al., 2024). While simple, these approaches face limitations: (i) they often fail to explore qualitatively distinct solution paths (Nguyen et al., 2025; Renze & Guven, 2024); (ii) they require domain-specific tuning (Wiher et al., 2022); and (iii) they can trade off against correctness and coherence (Nguyen et al., 2025). Prompt-cycling can inject domain knowledge (e.g., "try algebra" vs. "try geometry"), but it burdens users with prompt engineering and saturates well below human diversity (Shur-Ofry et al., 2024).

Determinantal point processes (DPPs) provide another path by rewarding sets of outputs that span high-volume regions in embedding space (Kulesza, 2012; Vijayakumar et al., 2018; Meister et al., 2023; Wang et al., 2024). In reinforcement learning, determinant-based rewards can encourage agents to explore trajectories that span complementary regions of state space (Ash et al., 2021; Zhao et al., 2024). Compared to MI, which directly couples a latent z with trajectories to ensure reproducible modes, DPP-based diversity is distribution-free: it treats a set of samples as diverse if they occupy a high-volume region in representation space, regardless of whether the same diversity is reproducible under repeated sampling.

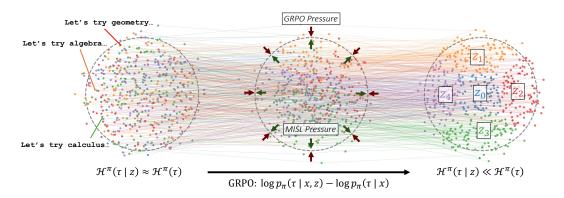


Figure 3: Example illustration of how the MISL reward improves pass@k performance. Before MISL (left), the trajectory distribution is independent of the latents z, so the conditional entropy is close to the marginal. MISL training prevents distribution collapse due to pass@1 training (middle). Adding the token-level MI reward (right) yields well-separated clusters indexed by z, reducing conditional entropy while preserving high marginal entropy. At inference, fixing different z values produces consistent and diverse solution strategies.

Finally, training-time diversification has also been studied through explicit pass@k-based objectives. Tang et al. (2025) proposed an unbiased estimator for generic k-attempt objectives, showing overall improved model efficacy. Extending this, Chen et al. (2025) argue that simply training on pass@1 falls victim to a local maximum of over-exploitation and reduced exploration. They find that pass@k training naturally focusing optimization efforts on harder problems producing significant improvements in both pass@k and pass@1. Outside of verifiable domains, DivPO (Lanchantin et al., 2025) alters preference optimization by contrasting diverse high-quality responses with common low-quality ones using a predefined diversity objective, yielding large diversity gains on creative and instruction-following tasks.

As we provide an orthogonal method to improve pass@k and diversity, our approach may complement that of Chen et al. (2025), Tang et al. (2025), and Lanchantin et al. (2025).

3 OPTIMIZING LLM DIVERSITY WITH MUTUAL INFORMATION

Given an input x and an autoregressive policy $\pi(\cdot \mid x)$ that produces a completion (trajectory) $\tau = (y_1, \ldots, y_T)$, we introduce a *discrete* latent $z \in \{1, \ldots, N\}$ via a lightweight prompt prefix (e.g., Strategy $\{z\}$ |), yielding conditional policies $\pi(\cdot \mid x, z)$. During training, z is drawn uniformly at random from the set $\{1, \ldots, N\}$. At inference, one selects $k \leq N$ distinct values of z and generates one completion per value, producing k semantically distinct attempts.

3.1 OBJECTIVE

We would like to encourage *structured response diversity* by maximizing the conditional mutual information $\mathcal{I}(\tau; z \mid x)$. Intuitively, maximizing mutual information makes the outputs of different strategies distinguishable, ensuring that each z induces a reliably different mode. By querying each strategy once, we obtain k semantically distinct attempts. Formally, this corresponds to maximizing

$$\max_{\pi} \mathcal{I}(\tau; z \mid x) = \mathbb{E}\Big[\log p_{\pi}(\tau \mid x, z) - \log p_{\pi}(\tau \mid x)\Big], \tag{4}$$

which increases the overall entropy of trajectories while reducing the conditional entropy within each z-mode, ensuring diverse yet reproducible strategies. The term $p_{\pi}(\tau \mid x) = \frac{1}{N} \sum_{z'=1}^{N} p_{\pi}(\tau \mid x, z')$ is a uniform mixture over skills. Maximizing mutual information encourages (i) high marginal entropy of trajectories, promoting broad coverage; and (ii) low conditional entropy given z, so that each response is distinct and determined by z. Figure 3 provides an overview of the relevant dynamics.

Algorithm 1 UpSkill: A method for training LLMs to produce diverse responses with mutual information.

- 1: **Inputs:** base policy π_{base} , trainable policy π , latent count N, completions per group C, weights $(\alpha_1, \alpha_2, \beta)$
- 2: repeat

- 3: Sample a minibatch of prompts $\{x\}$
- 4: **for** each x in the minibatch **do**
 - 5: Sample $z \sim \text{Unif}(\{1, \dots, N\})$; generate C completions $\{\tau_i\}_{i=1}^C$ with $\pi(\cdot \mid x, z)$
 - 6: Compute $r_{\text{corr}}(\tau_i)$, $r_{\text{TMI}}(\tau_i; x, z)$, and $\Delta_{\text{KL}}(\tau_i)$ as above
 - 7: **end for**
 - 8: Form per-sample rewards via equation 6; compute advantages; update π with GRPO
- 9: **until** convergence

3.2 Token-level mutual information reward

We now focus on implementing the mutual information as a token-level reward. For each pair (x, z), let $\{\tau_i\}_{i=1}^C$ be C completions sampled from $\pi(\cdot \mid x, z)$. We define a *per-sample* token-level score

$$r_{\text{TMI}}(\tau_i; x, z) = \sum_{t=1}^{C} \left[\log p_{\pi}(y_t \mid x, z, y_{< t}) - \log p_{\pi}(y_t \mid x, y_{< t}) \right], \tag{5}$$

where the second term is the uniform mixture

$$p_{\pi}(y_t \mid x, y_{< t}) = \frac{1}{N} \sum_{z'=1}^{N} p_{\pi}(y_t \mid x, z', y_{< t}).$$

Log-probabilities are computed by π on the realized τ_i . In our implementation the mixture is computed *exactly* across all N skills; this is feasible for the N used in our experiments (Section 5). Since $\frac{1}{C}r_{\text{TMI}}(\tau_i;x,z)$ is a Monte Carlo estimator of $\mathcal{I}(\tau;z\mid x)$, we make this our main reward term with UpSkill, with the other reward term being considered in ablation experiments. Appendix D discusses an alternative based on semantic mutual information.

3.3 COMBINED RL OBJECTIVE

Let $r_{corr}(\tau_i) \in \mathbb{R}$ denote the verifiable correctness reward (often binary) and define the per-trajectory KL penalty as

$$\Delta_{\text{KL}}(\tau_i) = \sum_{t=1}^{|\tau_i|} \log \frac{\pi(y_t \mid x, z, y_{< t})}{\pi_{\text{base}}(y_t \mid x, z, y_{< t})}.$$

The per-sample scalar reward is

$$r(\tau_i; x, z) = r_{\text{corr}}(\tau_i) - \beta \Delta_{\text{KL}}(\tau_i) + \alpha_1 r_{\text{TMI}}(\tau_i; x, z), \tag{6}$$

with $\alpha_1, \beta \geq 0$. We apply GRPO to optimize the sum of the combined rewards.

3.4 Training procedure

We fine-tune a trainable policy π_{θ} with GRPO while injecting a discrete strategy variable $z \in \{1,\ldots,N\}$. At each step, we draw a minibatch of prompts x and, for each x, sample a strategy z uniformly and generate C completions $\tau_{1:C} \sim \pi_{\theta}(\cdot \mid x,z)$ under fixed decoding. For every completion τ , we compute: (i) a verifiable correctness reward $r_{\text{corr}}(\tau)$ from the task's deterministic checker; (ii) the token-level MISL term r_{TMI} that measures how specific the trajectory is to the chosen strategy; and (iii) a KL control term toward a frozen base policy. We then update the policy with GRPO on this reward.

3.5 Inference

Given a budget of k attempts, we choose k distinct latents from $\{1, \ldots, N\}$ and generate one completion per latent under fixed decoding hyperparameters. Aggregation (e.g., majority vote) can optionally be applied. Because each completion is produced by a trained, distinct mode, conditional success probabilities remain larger than with redundant samplings, improving multi-attempt metrics.

4 THEORETICAL CONNECTION BETWEEN PASS@K IMPROVEMENT AND MUTUAL INFORMATION

Our main theoretical result shows that the mutual information objective is closely tied to pass@k. In particularly, we will show that the mutual information objective is a lower bound on *improvement* in the pass@k objective, so maximizing mutual information provably results in an increased (lower bound on) pass@k. Our theoretical results will require the following assumptions:

- 1. *k*-uniform mixture model: Assume that the marginal distribution over the skills is identical to the base model.
- 2. Distributional impact: Let a_z be the probability of success of strategy z and a be the probability of success of the base model. Assume that for all $x \in \mathcal{X}$ there exists $\eta > 0$ such that for all $z \in [k]$, $|a_z a| \ge \eta \delta(\pi_{M,z}(\cdot \mid x), \pi_B(\cdot \mid x))$, where δ is the total variation distance.

The second assumption says that the distribution shifts induced by UpSkill correspond to different problem approaches, and, as a result, will have different probabilities of success. A more precise definition of k-uniform mixture models, additional justification for the assumptions, and the statement and proof of the lemma is in Appendix C.

Lemma 1. Let pass@k be the pass@k score of the base model on prompt x and pass@k $_M$ be the pass@k score of the mixture model on prompt x. Under the above assumptions, we show that:

$$1 - \exp\left(-C_1\eta^2\mathcal{I}(\tau;z\mid x)^2\right) \leq \frac{\mathrm{pass@k}_M - \mathrm{pass@k}_B}{1 - \mathrm{pass@k}_B} \leq 1 - \exp\left(-C_2\mathcal{I}(\tau;z\mid x)\right)$$

where C_1 depends on k and C_2 depends on k and $\max_z a_z$.

The quantity in the middle can be interpreted as the fraction of possible improvement from the base model that is realized by the mixture model. Since monotonically increasing functions in $\mathcal{I}(\tau;z\mid x)$ provide both lower and upper bounds on how much the mixture model improves over the base model, it makes sense to optimize directly for the mutual information. UpSkill explicitly increases $\mathcal{I}(\tau;z\mid x)$ during training, ensuring diversity across skills and giving a guaranteed improvement in pass@k over the base model.

5 EXPERIMENTS

We evaluate whether conditioning on a discrete latent z and training with our token-level mutual-information reward (equation 5) improves multi-attempt metrics. We present results in two settings: (1) a controlled arithmetic environment that allows for fully verifiable evaluation and direct inspection of distributional effects, and (2) the GSM8K benchmark across three open-weight models. We then report ablations on the number of strategies N and on adding a semantic-MI surrogate.

5.1 ARITHMETIC ENVIRONMENT

The arithmetic environment consists of prompts with three single-digit integers and a latent skill index $z \in \{0, \dots, N-1\}$. A small transformer model chooses one of the integers to be a target and is required to produce a simple arithmetic expression with the other two digits that evaluates to this target. We use an automatically-verified correctness reward and training uses GRPO without a KL penalty. Full details of the environment, model, and training procedure are provided in Appendix E.

Figure 4 illustrates that in the control condition (α_1 =0), training quickly collapses to a single deterministic strategy: by the end of training, pass@1 and pass@5 are indistinguishable (0.793), offering no benefit from multiple attempts. In contrast, with UpSkill (α_1 =0.5, MI reward cap at 1.0), the model maintains diverse trajectories, yielding substantially higher multi-attempt accuracy (pass@5=0.897) despite a lower single-attempt accuracy (pass@1=0.390). This difference aligns with the entropy dynamics: UpSkill preserves broad output distributions (token entropy 0.723 \rightarrow 0.797), sustaining diverse strategies and higher pass@5. By contrast, the control run—while substantially improving pass@1—collapses to near-deterministic outputs (entropy 0.723 \rightarrow 0.030), leaving pass@5 identical to pass@1 (see Appendix E.5 for more details).

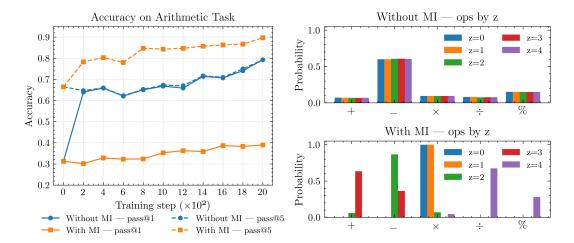


Figure 4: Arithmetic environment results. Training curves show that under GRPO alone (blue), pass@1 and pass@5 converge together, indicating that multiple attempts provide little benefit. With MISL (orange; N=5), pass@5 improves substantially while pass@1 remains modest, demonstrating that different latents yield complementary solutions. Operator distributions further highlight this effect: without MISL, they are nearly identical across z, reflecting a lack of specialization, whereas with MISL, distinct latents focus on different operators, producing diverse strategies that drive multi-attempt gains.

Figure 4 illustrates that under UpSkill, different z values yield distinct distributions over operators and digits, whereas the control produces nearly identical distributions across z. On this small scale environment, we can directly observe the learned strategies, and notably z=1 and z=2 converge to risky yet common modes, whereas other values of z cover the remaining operations, improving multi-attempt success with a strategy infeasible for optimizing a pass@1 objective. The distributions over the first digit are available in Appendix E.7.

We additionally ablate the impact of starting model capabilities, the coefficient of KL-penalty, and GRPO parameters (see Appendix F for full details and results). KL penalty β discourages entropy collapse by ensuring the new policy remains close to the initial policy, thereby improving performance. On models with $\beta \in [0.05, 0.10]$, we test $\alpha_1 \in [0.1, 0.3, 0.5]$ and find that well-chosen MI-reward parameters increase pass@k by an average of 3% for the weaker base model. However, for the stronger base model, we find the opposite trend. It is always best to choose $\alpha_1 = 0$, with the best choice of $\alpha_1 \in [0.1, 0.3, 0.5]$ still leading to a 1.2% performance decrease. Our theoretical results in Lemma 1 suggest that UpSkill improvement is negatively related to pass@1 (equivalently pass@k_B) capability, and thus, although surprising, this result is in line with our theoretical analysis. We separately conjecture that β , which corresponds with a decrease in exploration from the base policy, conflicts with the mutual information incentive to explore.

5.2 GSM8K

We next evaluate on GSM8K (Cobbe et al., 2021), a dataset of grade-school arithmetic word problems. We use 2,000 training problems and a held-out set of 100 questions. All experiments are conducted in a zero-shot setting with a maximum sequence length of 1024 tokens. We train LoRA adapters (approximately 80M trainable parameters) on top of three open-weight backbones: Llama 3.1-8B (Meta AI, 2025), Qwen 2.5-7B (Qwen Team, 2025), and R1-Distilled–Qwen2.5-Math-1.5B. As before, we apply GRPO with a correctness reward and with default KL penalty, and UpSkill is applied at the token level as in Eq. (6). At inference, we fix k distinct skill indices and generate one completion per skill. More details are available in Appendix G.

Figure 5 shows our performance on the withheld evaluation set; asterisks mark p < .05 improvements. Token-level MISL improves pass@k at a significant level in R1-Distilled-Qwen and improved consensus@k at a significant level in Qwen, while pass@1 is maintained or improved. Interestingly, these results do not entirely parallel the arithmetic environment, as UpSkills improvement has

not come at the cost of pass@1. We hypothesize this is due to the existence of a larger set of correct reasoning approaches, and thus MISL does not necessarily come at a cost of pass@1. Chen et al. (2025) have separately found that pass@k training method can improve pass@1 performance. We include summarized outputs from Qwen in Figure 1.

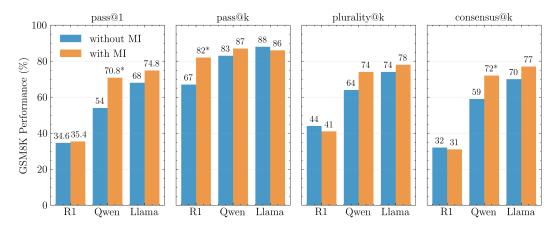


Figure 5: Performance with N=5 strategies and token MI only. We observe consistent gains on pass@1 and consensus@5. Asterisks denote p<.05 from a two-proportions z-test on the results.

5.3 GSM8k Ablations

We conduct two ablations to probe the robustness of the approach. First, increasing the number of skills beyond $N\!=\!5$ produces mixed results. Many GSM8K problems admit only a limited number of distinct solution paths, so larger N values fragment the capacity into modes that do not translate into additional gains. Second, we study the effect of replacing the token-level MI with a semantic MI, which we formally introduce in Appendix D. This semantic MI occasionally yields further improvements but introduces instability due to estimator variance in high dimensions. We provide more details in Appendix H. Our main results therefore use token-level MISL only, with semantic variants left as a direction for future work.

6 CONCLUSION

Our experiments show that UpSkill provides a simple and effective way to induce strategy-level diversity in LLMs, leading to consistent gains on multi-attempt metrics such as pass@5 and consensus@5. By conditioning on discrete latent variables, the model learns reproducible modes of reasoning that reduce redundancy across attempts and increase the likelihood of success. Beyond these empirical findings, UpSkill provides a principled training-time approach for improving response diversity, and our analysis links $\mathcal{I}(\tau;z\mid x)$ to upper bounds in pass@k improvement in training. We hope this stimulates research on robust semantic diversity signals and theoretical ties between information-theoretic objectives and multi-attempt success.

Limitations. There are a few notable limitations with our work. Assumption 1 is rather limiting and difficult to enforce empirically, so we hope to find a natural way to either incorporate it into our training method or find another assumption that can also prove Lemma 1 As future work, we hope for further experimentation across domains and with larger models, to better study how our method is affected by KL-penalty and base model performance, and to approach the theoretical guarantees from the perspective of Tsallis entropy (Furuichi, 2006).

Reproducibility. We make several efforts towards encouraging reproducibility of our work and results. We have included all experiment code in the supplementary files, along with instructions to reproduce our results. We use only models under permissive licenses. For our theoretical results, we provide a description of assumptions and the complete proof of claims in the appendix. Additionally,

for various experiments whose full results could not fit in the main paper, we include the full results in the appendix.

REFERENCES

- Joshua Achiam, Harrison Edwards, Dario Amodei, and Pieter Abbeel. Variational option discovery algorithms. *arXiv preprint arXiv:1807.10299*, 2018.
- Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. Gone Fishing: Neural Active Learning with Fisher Embeddings, December 2021. URL http://arxiv.org/abs/2106.09675. arXiv:2106.09675 [cs].
- Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. An Actor-Critic Algorithm for Sequence Prediction, March 2017. URL http://arxiv.org/abs/1607.07086.arXiv:1607.07086 [cs].
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. Mine: Mutual information neural estimation, 2021. URL https://arxiv.org/abs/1801.04062.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating Large Language Models Trained on Code, July 2021. URL http://arxiv.org/abs/2107.03374. arXiv:2107.03374 [cs].
- Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets, 2016. URL https://arxiv.org/abs/1606.03657.
- Zhipeng Chen, Xiaobo Qin, Youbin Wu, Yue Ling, Qinghao Ye, Wayne Xin Zhao, and Guang Shi. Pass@k training for adaptively balancing exploration and exploitation of large reasoning models, 2025. URL https://arxiv.org/abs/2508.10751.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Xingyu Dang, Christina Baek, J Zico Kolter, and Aditi Raghunathan. Assessing diversity collapse in reasoning. In *Scaling Self-Improving Foundation Models without Human Supervision*, 2025. URL https://openreview.net/forum?id=AMiKsHLjQh.
- Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Junjie Shan, Caishuang Huang, Wei Shen, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Xuanjing Huang, and Tao Gui. StepCoder: Improve Code Generation with Reinforcement Learning from Compiler Feedback, February 2024. URL http://arxiv.org/abs/2402.01391.arXiv:2402.01391 [cs] version: 1.
- Weihua Du, Yiming Yang, and Sean Welleck. Optimizing Temperature for Language Models with Multi-Sample Inference, February 2025. URL http://arxiv.org/abs/2502.05234.arXiv:2502.05234 [cs] version: 1.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is All You Need: Learning Skills without a Reward Function, October 2018. URL http://arxiv.org/abs/1802.06070. arXiv:1802.06070 [cs].

- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=SJx63jRqFm.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation, 2018. URL https://arxiv.org/abs/1805.04833.
 - Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical reinforcement learning. *arXiv* preprint arXiv:1704.03012, 2017.
 - Shigeru Furuichi. Information theoretical properties of Tsallis entropies. *Journal of Mathematical Physics*, 47(2):023302, February 2006. ISSN 0022-2488, 1089-7658. doi: 10.1063/1.2165744. URL http://arxiv.org/abs/cond-mat/0405600. arXiv:cond-mat/0405600.
 - Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv* preprint arXiv:1611.07507, 2016a.
 - Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational Intrinsic Control, November 2016b. URL http://arxiv.org/abs/1611.07507. arXiv:1611.07507 [cs].
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
 - Steven Hansen, Guillaume Desjardins, Kate Baumli, David Warde-Farley, Nicolas Heess, Simon Osindero, and Volodymyr Mnih. Entropic Desired Dynamics for Intrinsic Control. In *Advances in Neural Information Processing Systems*, volume 34, pp. 11436–11448. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/hash/5f7f02b7e4ade23430f345f954c938c1-Abstract.html.
 - Andreas Hochlehnert, Hardik Bhatnagar, Vishaal Udandarao, Samuel Albanie, Ameya Prabhu, and Matthias Bethge. A Sober Look at Progress in Language Model Reasoning: Pitfalls and Paths to Reproducibility, April 2025. URL http://arxiv.org/abs/2504.07086. arXiv:2504.07086 [cs].
 - Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020. URL https://arxiv.org/abs/1904.09751.
 - Vineet John, Lili Mou, Hareesh Bahuleyan, and Olga Vechtomova. Disentangled representation learning for non-parallel text style transfer, 2018. URL https://arxiv.org/abs/1808.04339.
 - Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes, December 2013. URL http://arxiv.org/abs/1312.6114. arXiv:1312.6114 [stat].
 - Leslie Kish. Survey Sampling. Wiley, 1965.
 - Peter Krafft. Correlation and mutual information. https://lips.cs.princeton.edu/correlation-and-mutual-information/, February 2013. Laboratory for Intelligent Probabilistic Systems, Princeton University Department of Computer Science.
 - Alexander Kraskov, Harald Stoegbauer, and Peter Grassberger. Estimating Mutual Information. *Physical Review E*, 69(6):066138, June 2004. ISSN 1539-3755, 1550-2376. doi: 10.1103/PhysRevE.69.066138. URL http://arxiv.org/abs/cond-mat/0305641. arXiv:cond-mat/0305641.
 - Alex Kulesza. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286, 2012. ISSN 1935-8245. doi: 10.1561/2200000044. URL http://dx.doi.org/10.1561/2200000044.
 - Jack Lanchantin, Angelica Chen, Shehzaad Dhuliawala, Ping Yu, Jason Weston, Sainbayar Sukhbaatar, and Ilia Kulikov. Diverse preference optimization, 2025. URL https://arxiv.org/abs/2501.18101.

- Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual demonstrations, 2017. URL https://arxiv.org/abs/1703.08840.
 - Jiate Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: Reinforcement Learning from Unit Test Feedback, November 2023. URL http://arxiv.org/abs/2307.04349. arXiv:2307.04349 [cs].
 - Clara Meister, Martina Forster, and Ryan Cotterell. Determinantal Beam Search, June 2023. URL http://arxiv.org/abs/2106.07400. arXiv:2106.07400 [cs].
 - Meta AI. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. Meta AI Blog, April 2025. URL https://ai.meta.com/blog/llama-4-multimodal-intelligence/. Announcement of the Llama 4 multimodal AI model family. Accessed: September 25, 2025.
 - Minh Nhat Nguyen, Andrew Baker, Clement Neo, Allen Roush, Andreas Kirsch, and Ravid Shwartz-Ziv. Turning Up the Heat: Min-p Sampling for Creative and Coherent LLM Outputs, June 2025. URL http://arxiv.org/abs/2407.01082. arXiv:2407.01082 [cs] version: 7.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, March 2022. URL http://arxiv.org/abs/2203.02155. arXiv:2203.02155 [cs].
 - Yao Qiang, Subhrangshu Nandi, Ninareh Mehrabi, Greg Ver Steeg, Anoop Kumar, Anna Rumshisky, and Aram Galstyan. Prompt Perturbation Consistency Learning for Robust Language Models. In Yvette Graham and Matthew Purver (eds.), *Findings of the Association for Computational Linguistics: EACL 2024*, pp. 1357–1370, St. Julian's, Malta, March 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.findings-eacl.91/.
 - Qwen Team. Qwen3: Think deeper, act faster. Qwen Blog, April 2025. URL https://qwenlm.github.io/blog/qwen3/. Announcement of Qwen3 large language model family. Accessed: September 25, 2025.
 - Matthew Renze and Erhan Guven. The Effect of Sampling Temperature on Problem Solving in Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 7346–7356, 2024. doi: 10.18653/v1/2024.findings-emnlp.432. URL http://arxiv.org/abs/2402.05201. arXiv:2402.05201 [cs].
 - Igal Sason and Sergio Verdú. \$f\$-divergence Inequalities. *IEEE Transactions on Information Theory*, 62(11):5973–6006, November 2016. ISSN 0018-9448, 1557-9654. doi: 10.1109/TIT.2016. 2603151. URL http://arxiv.org/abs/1508.00335. arXiv:1508.00335 [cs].
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
 - Sagi Shaier, Mario Sanz-Guerrero, and Katharina von der Wense. Asking Again and Again: Exploring LLM Robustness to Repeated Questions, March 2025. URL http://arxiv.org/abs/2412.07923. arXiv:2412.07923 [cs].
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
 - Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020a. URL https://openreview.net/forum?id=HJgLZR4KvH.
 - Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-Aware Unsupervised Discovery of Skills, February 2020b. URL http://arxiv.org/abs/1907.01657. arXiv:1907.01657 [cs].

- Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, and Wai Lam. A Thorough Examination of Decoding Methods in the Era of LLMs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 8601–8629, Miami, Florida, USA, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.489. URL https://aclanthology.org/2024.emnlp-main.489.
- Michal Shur-Ofry, Bar Horowitz-Amsalem, Adir Rahamim, and Yonatan Belinkov. Growing a Tail: Increasing Output Diversity in Large Language Models, November 2024. URL http://arxiv.org/abs/2411.02989.arXiv:2411.02989 [cs].
- Greg Ver Steeg. gregversteeg/npeet, May 2025. URL https://github.com/gregversteeg/NPEET.original-date: 2014-10-10T19:57:02Z.
- Karl Stratos and Sam Wiseman. Learning Discrete Structured Representations by Adversarially Maximizing Mutual Information, July 2020. URL http://arxiv.org/abs/2004.03991. arXiv:2004.03991 [cs].
- Richard S Sutton and Andrew G Barto. Reinforcement Learning: An Introduction. 2015.
- Yunhao Tang, Kunhao Zheng, Gabriel Synnaeve, and Rémi Munos. Optimizing language models for inference time objectives using reinforcement learning, 2025. URL https://arxiv.org/abs/2503.19595.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method, April 2000. URL http://arxiv.org/abs/physics/0004057. arXiv:physics/0004057.
- Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019. URL https://arxiv.org/abs/1807.03748.
- Ashwin K. Vijayakumar, Michael Cogswell, Ramprasath R. Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. Diverse Beam Search: Decoding Diverse Solutions from Neural Sequence Models, October 2018. URL http://arxiv.org/abs/1610.02424.arXiv:1610.02424 [cs].
- Eric Wallace, Olivia Watkins, Miles Wang, Kai Chen, and Chris Koch. Estimating Worst-Case Frontier Risks of Open-Weight LLMs, August 2025. URL http://arxiv.org/abs/2508.03153. arXiv:2508.03153 [cs] version: 1.
- Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. Diversity Measurement and Subset Selection for Instruction Tuning Datasets, February 2024. URL http://arxiv.org/abs/2402.02318. arXiv:2402.02318 [cs].
- Gian Wiher, Clara Meister, and Ryan Cotterell. On Decoding Strategies for Neural Text Generators, March 2022. URL http://arxiv.org/abs/2203.15721. arXiv:2203.15721 [cs].
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative Preference Learning from Human Feedback: Bridging Theory and Practice for RLHF under KL-Constraint, May 2024. URL http://arxiv.org/abs/2312.11456.arXiv:2312.11456 [cs].
- Jake Zhao, Yoon Kim, Kelly Zhang, Alexander M. Rush, and Yann LeCun. Adversarially regularized autoencoders, 2018a. URL https://arxiv.org/abs/1706.04223.
- Kaiyan Zhao, Yiming Wang, Yuyang Chen, Xiaoguang Niu, Yan Li, and Leong Hou U. Efficient Diversity-based Experience Replay for Deep Reinforcement Learning, October 2024. URL http://arxiv.org/abs/2410.20487. arXiv:2410.20487 [cs] version: 1.
- Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Information maximizing variational autoencoders, 2018b. URL https://arxiv.org/abs/1706.02262.

Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders, 2017. URL https://arxiv.org/abs/1703.10960.

Chongyi Zheng, Jens Tuyls, Joanne Peng, and Benjamin Eysenbach. Can a misl fly? analysis and ingredients for mutual information skill learning. *arXiv preprint arXiv:2412.08021*, 2024.

Han Zhong, Zikang Shan, Guhao Feng, Wei Xiong, Xinle Cheng, Li Zhao, Di He, Jiang Bian, and Liwei Wang. Dpo meets ppo: Reinforced token optimization for rlhf, 2025. URL https://arxiv.org/abs/2404.18922.

A USE OF LLMS

 Large language models were used in the preparation of this work for writing assistance (including polishing, improving presentation of concepts, and restructuring of text), for retrieval and discovery of related work, and for support in producing experimental code and figures. Language models were additionally used for feedback on the paper and to formalize mathematical arguments. All analysis, experimental design, and final interpretations are our own.

B EXTENDED MUTUAL INFORMATION RELATED WORK

Estimating MI reliably is challenging in high dimensions. Variational bounds (Barber–Agakov) optimize a classifier or regressor $q_{\phi}(z|\cdot)$ as a proxy for the intractable posterior (van den Oord et al., 2019). Contrastive bounds such as InfoNCE (van den Oord et al., 2019) reduce MI estimation to noise-contrastive classification and have become standard due to their stability. Neural MI estimators like MINE (Belghazi et al., 2021) directly optimize a Donsker–Varadhan bound but can suffer from bias/variance trade-offs and training instability. Nonparametric kNN estimators (KSG) (Kraskov et al., 2004) avoid parametric critics but require many samples and are sensitive to dimension, motivating careful batching and normalization when used inside policy gradients. In text generation, MI-style objectives have been used to prevent latent collapse and enable controllable generation, e.g., by encouraging informative latents in variational text models (Zhao et al., 2017; 2018a) or aligning codes with style attributes (John et al., 2018). These approaches typically maximize MI between prompts or attributes and latent variables, rather than between a discrete strategy and the full trajectory distribution, and are optimized with supervised losses rather than RL.

Conceptually, our objective reconciles two desiderata emphasized in prior MI work: *coverage* (high marginal entropy over trajectories) and *control* (low conditional entropy given z). Whereas decoding-time diversity manipulates token entropy without guarantees about identifiable modes, MI-based diversification learns *reusable*, *reproducible* modes indexed by a small discrete latent. This makes diversity a first-class, training-time property that can be cleanly exercised at inference by selecting distinct z values.

C STATEMENT AND DERIVATION OF THEORETICAL BOUNDS

C.1 PROBLEM SETUP AND ASSUMPTIONS

Let \mathcal{X} be the set of all possible prompts. The statement of Lemma 1 applies to the general class of k-uniform mixture models.

Definition. A k-uniform mixture model (M,B) is defined to be an ordered pair of a *mixture model*, which is a set of k different policies for generating trajectories, which we will call $\pi_{M,z}(\cdot \mid x)$ for a prompt $x \in \mathcal{X}$ and strategy $z \in [k]$, along with a *base model* $\pi_B(\cdot \mid x)$ for generating trajectories subject to the condition that

$$\frac{1}{k} \sum_{z=1}^{k} \pi_{M,z}(\cdot \mid x) = \pi_B(\cdot \mid x) \,\forall \, x \in \mathcal{X}.$$

 This definition can be interpreted as follows: π_B is the trajectory distribution of the original, non-strategy conditioned language model. If we weigh each strategy as being equally important, we sample once from the mixture model by randomly choosing one strategy. In this case the joint distribution of trajectories from the mixture is

$$\pi_M(\cdot \mid x) := \frac{1}{k} \sum_{z=1}^k \pi_{M,z}(\cdot \mid x).$$

The condition essentially means that the joint distribution of trajectories over picking a strategy uniformly at random must be the same as the original distribution. Therefore, in essence the mixture model partitions π_B into k different policies that together average back to π_B .

In practice, this condition imposes undue constraints on the strategy distribution, and thus for the practical implementation, this is not enforced. Also, while in practice one may actually train N>k different strategies and then randomly sample k different strategies so that they still have equal probabilities of being selected, here we make the simplifying assumption that N=k, which is true for all of our experiments.

For ease of notation, let $a = \Pr(Y_x(\tau_{B,1}) = 1 \mid x)$ and $a_z = \Pr(Y_x(\tau_{M,z}) = 1)$ for $z \in [k]$. Because the trajectories $\tau_{M,z}$ are sampled independently, we have that

$$\pi_M(\cdot \mid x) = \pi_B(\cdot \mid x) \implies \frac{1}{k} \sum_{z=1}^k a_z = a. \tag{7}$$

We provide additional justification for the second assumption made in section 4. For the first, if the strategies differ in more than just style and contain meaningful semantic differences, we expect that the difference in success should be proportional to how different these two distributions are. The total variation distance measures this distance in trajectory space, while the constant η controls how sensitive the success probabilities are to changes in the distribution shift.

C.2 EXTENDING PASS@K TO k-UNIFORM MIXTURE MODELS

We now extend the traditional definition of pass@k to fit the setting of k-uniform mixture models to leverage the fact that we now have a natural structure for querying k different strategies by varying k. This is notably different from the setting in the consistency assumption, which can be interpreted as querying just one strategy uniformly at random.

For a given prompt x and deterministic verifier $Y_x(\tau)$ outputting 1 if τ is a valid output on prompt x and 0 otherwise, and k-uniform mixture model (M,B), define $\mathtt{pass@k}_M$ be the probability that querying each of these strategies independently exactly once (see Sec. 3.5) results in at least one correct answer, and define $\mathtt{pass@k}_B$ to be the probability that querying π_B independently k times results in at least one correct answer. Writing this out mathematically, for each $z \in [k]$ we sample $\tau_{M,z} \sim \pi_{M,z}(\cdot \mid x)$; then

$$\operatorname{pass@k}_{M} = 1 - \operatorname{Pr}\left(\bigcap_{z=1}^{k} \{Y_{x}(\tau_{M,z}) = 0\} \middle| x\right) = 1 - \prod_{z=1}^{k} \operatorname{Pr}(Y_{x}(\tau_{M,z}) = 0 \mid x). \tag{8}$$

While we use the same definition of pass@k for B as in standard literature, we include it for the sake of completeness; similarly sampling $\tau_{B,z} \sim \pi_B(\cdot \mid x)$ independently for each $z \in [k]$ we find that

$$\operatorname{pass@k}_{B} = 1 - \operatorname{Pr}\left(\bigcap_{z=1}^{k} \{Y_{x}(\tau_{B,z}) = 0\} \middle| x\right) = 1 - \operatorname{Pr}(Y_{x}(\tau_{B,1}) = 0 \mid x)^{k}$$
(9)

where we simplify the product into the RHS of equation 9 with the independence of samples. Then

$$\mathrm{pass@k}_B = 1 - (1-a)^k, \; \mathrm{pass@k}_M = 1 - \prod_{z=1}^k (1-a_z).$$

We now give the precise statement of the main result.

Lemma 2 (pass@k Improvement for k-uniform Mixture Models, Full Statement). Let $u = \max_{z \in [k]} a_z$ and let $\varphi : \mathbb{R}_{\geq 0} \mapsto \mathbb{R}$ be defined as $\varphi(x) = \frac{x \log x}{x-1}$ for $x \neq 0, 1$ and $\varphi(0) = 0, \varphi(1) = 1$. Then

$$1 - \exp\left(-\frac{k\eta^2\mathcal{I}(\tau;z\mid x)^2}{2\varphi(k)^2}\right) \leq \frac{\mathrm{pass@k}_M - \mathrm{pass@k}_B}{1 - \mathrm{pass@k}_B} \leq 1 - \exp\left(-\frac{k\mathcal{I}(\tau;z\mid x)}{4(1-u)^2}\right).$$

C.3 DERIVATION OF LOWER BOUND

 Using Taylor's Theorem on $f(y) = \log(1-y)$ gives the equations $f(a_z) = f(a) + (a_z - a)f'(a) + \frac{1}{2}(a_z - a)^2 f''(\xi_z)$ where ξ_z lies in between a_z and a, for $z \in [k]$. Summing all of these equations, the linear terms cancel due to equation 7. Then

$$\sum_{z \in [k]} \log(1 - a_z) = k \log(1 - a) + \sum_{z \in [k]} \frac{1}{2} (a_z - a)^2 f''(\xi_z).$$
 (10)

Let A be the random variable that takes value a_z for each $z \in [k]$ with probability $\frac{1}{k}$. Since $f''(y) = -\frac{1}{(1-y)^2}$ is a decreasing function in the interval [0,1), we find that $f''(\xi_z) \leq f''(0) = -1$ for all $z \in [k]$, i.e.

$$\begin{split} \sum_{z \in [k]} \log(1 - a_z) &\leq k \log(1 - a) - \sum_{z \in [k]} \frac{1}{2} (a_z - a)^2 \\ \Longrightarrow k \log(1 - a) - \sum_{z \in [k]} \log(1 - a_z) &\geq \frac{k \mathrm{Var}(A)}{2} \implies \frac{1 - \mathrm{pass@k}_B}{1 - \mathrm{pass@k}_M} \geq \exp\left(\frac{k}{2} \mathrm{Var}(A)\right) \\ \Longrightarrow \mathrm{pass@k}_M - \mathrm{pass@k}_B &= (1 - \mathrm{pass@k}_B) \left(1 - \frac{1 - \mathrm{pass@k}_M}{1 - \mathrm{pass@k}_B}\right) \\ &\geq (1 - \mathrm{pass@k}_B) \left(1 - \exp\left(-\frac{k \mathrm{Var}(A)}{2}\right)\right). \end{split}$$

We now find a lower bound for Var(A) in terms of the mutual information to finish off the proof. Let $\delta(\cdot, \cdot)$ represent the total variation distance between two distributions. We have that $Var(A) \geq \mathbb{E}[|A|]^2$ by Jensen's Inequality, so using the distributional impact assumption,

$$\operatorname{Var}(A) \ge \left(\mathbb{E}_{z \sim \operatorname{Unif}\{1, 2, \dots, k\}}[|a_z - a|]\right)^2 = \left(\frac{1}{k} \sum_{z=1}^k |a_z - a|\right)^2 \ge \left(\frac{1}{k} \sum_{z=1}^k \eta \delta(\pi_{M, z}, \pi_B)\right)^2$$

$$= \eta^2 \left(\mathbb{E}_{z \sim \text{Unif}\{1,2,\ldots,k\}} [\delta(\pi_{M,z}, \pi_B)] \right)^2.$$

Note that by the assumption that $\pi_M(\cdot \mid x) = \pi_B(\cdot \mid x)$ we have that $0 \le \pi_{M,z}(\cdot \mid x) \le k\pi_B(\cdot \mid x)$ for all $z \in [k]$. Therefore,

$$\frac{d\pi_{M,z}}{d\pi_B}(\tau) \in [0,k].$$

Applying Theorem 26 from Sason & Verdú (2016) with $\beta_2 = 0$, $\beta_1 = \frac{1}{k}$ where the constants are from the assumption on bounded likelihood ratio, we have that $D_{\text{KL}}(\pi_{M,z} \parallel \pi_B) \leq \varphi(k)\delta(\pi_{M,z},\pi_B)$. Summing over $z \in [k]$ and dividing by k yields

$$\mathcal{I}(\tau;z\mid x) = \mathbb{E}_{z\sim \mathrm{Unif}\{1,2,\ldots,k\}}[D_{\mathrm{KL}}(\pi_{M,z}\parallel\pi_B)] \leq \varphi(k)\mathbb{E}_{z\sim \mathrm{Unif}\{1,2,\ldots,k\}}[\delta(\pi_{M,z},\pi_B)].$$

Putting both bounds together, we finally find that

$$\begin{aligned} \operatorname{Var}(A) & \geq \left(\frac{\eta \mathcal{I}(\tau;z\mid x)}{\varphi(k)}\right)^2 \\ \implies \operatorname{pass@k}_M - \operatorname{pass@k}_B \geq \left(1 - \operatorname{pass@k}_B\right) \left(1 - \exp\left(-\frac{k\eta^2 \mathcal{I}(\tau;z\mid x)^2}{2\varphi(k)^2}\right)\right) \end{aligned}$$

as desired.

C.4 DERIVATION OF UPPER BOUND

Let $u = \max_z a_z < 1$. In the interval $[\min(a, a_1, a_2, \dots, a_k), \max(a, a_1, a_2, \dots, a_k)]$ f'' achieves its minimum at $f''(u) = -\frac{1}{(1-u)^2}$. Then from equation 10

$$\sum_{z \in [k]} \log(1 - a_z) \ge k \log(1 - a) - \frac{1}{2(1 - u)^2} \sum_{z \in [k]} (a_z - a)^2$$

$$\implies \prod_{z=1}^k (1 - a_z) \ge (1 - a)^k \exp\left(-\frac{1}{2(1 - u)^2} \sum_{z \in [k]} (a_z - a)^2\right)$$

$$\implies 1 - \operatorname{pass@k}_M \ge (1 - \operatorname{pass@k}_B) \exp\left(-\frac{1}{2(1 - u)^2} \sum_{z \in [k]} (a_z - a)^2\right). \tag{11}$$

This places an upper bound on how much we can possibly improve pass@k compared to our original trajectory distributions. Using Pinsker's Inequality, we find that for all $i \in [k]$,

$$|a_{i} - a| = |\Pr[Y(\tau_{M,i}) = 1] - \Pr[Y(\tau_{B,i}) = 1]| = \left| \sum_{Y(\tau')=1} \pi_{M,i}(\tau') - \sum_{Y(\tau')=1} \pi_{B,i}(\tau') \right|$$

$$\leq \sum_{Y(\tau')=1} |\pi_{M,i}(\tau') - \pi_{B,i}(\tau')| \leq \sum_{\tau'} |\pi_{M,i}(\tau') - \pi_{B,i}(\tau')|$$

$$\leq \delta(\pi_{B,i}, \pi_{M,i}) \leq \sqrt{\frac{1}{2} D_{KL}(\pi_{M,i} \parallel \pi_{B,i})}$$

$$\implies (a_{i} - a)^{2} \leq \frac{1}{2} D_{KL}(\pi_{M,i} \parallel \pi_{B,i}) \implies \sum_{z \in [k]} (a_{z} - a)^{2} \leq \frac{k}{2} \cdot \frac{1}{k} \sum_{z \in [k]} D_{KL}(\pi_{M,i} \parallel \pi_{B,i})$$

$$= \frac{k}{2} \mathbb{E}_{z \sim \text{Unif}\{1,2,...,k\}} [D_{KL}(\pi_{M,i} \parallel \pi_{B})] = \frac{k}{2} \mathcal{I}(\tau; z \mid x).$$

Combining this with equation 11 yields

$$1 - \operatorname{pass@k}_{M} \ge (1 - \operatorname{pass@k}_{B}) \exp\left(-\frac{k}{4(1-u)^{2}} \mathcal{I}(\tau; z \mid x)\right). \tag{12}$$

As a result, if $\mathcal{I}(\tau; z \mid x)$ is too small, then our theoretical upper bound on improvement in pass@k between steps 0 and T will also be very small.

Rearranging equation 12 to bound the improvement $\Delta := pass@k_M - pass@k_B$, we obtain

$$\begin{split} \Delta &= (1 - \mathrm{pass@k}_B) - (1 - \mathrm{pass@k}_M) \\ &\leq (1 - \mathrm{pass@k}_B) \left(1 - \exp\left(-\frac{k}{4(1-u)^2} \mathcal{I}(\tau;z\mid x) \right) \right) \\ &\leq (1 - \mathrm{pass@k}_B) \cdot \frac{k}{4(1-u)^2} \cdot \mathcal{I}(\tau;z\mid x), \end{split}$$

where the final inequality uses $1 - e^{-x} \le x$.

D SEMANTIC MUTUAL INFORMATION REWARD

One observation that we made empirically is that token-level differences tend to reflect formatting or paraphrasing, rather than semantically distinct strategies. To bias toward more meaningful differences, we test an alternative method for measuring mutual information by embedding completions with a fixed encoder $\psi(\tau) \in \mathbb{R}^d$ and estimating the mutual information between embeddings and skills for a single prompt x:

$$\widehat{\mathcal{I}}(\psi(\tau); z \mid x) \tag{13}$$

using the KSG k-nearest-neighbor estimator (Kraskov et al., 2004), implemented with the library NPEET (Steeg, 2025). Concretely, for each x we collect the set of embeddings across strategies and samples, $\mathcal{B}(x) = \{(\psi(\tau_i^{(z)}), z) : z \in \{1, \dots, N\}, i = 1, \dots, C\}$, and apply KSG to $\mathcal{B}(x)$ to obtain a single scalar $r_{\text{SMI}}(x)$.

E ARITHMETIC ENVIRONMENT

E.1 TASK

Each problem instance consists of three integers $a,b,c\in\{0,\ldots,9\}$. A small transformer model chooses one of the integers to be a target and is required to produce a simple arithmetic expression with the other two digits that evaluates to this target. Valid operators are $\{+,-,\times,\div,\mathrm{mod}\}$, with division and modulo defined only when results are integer and denominators are nonzero. A latent skill index $z\in\{0,\ldots,N-1\}$ is provided as part of the prompt, conditioning the model on which strategy to adopt.

E.2 PROMPT FORMAT AND CONDITIONING

The input is formatted as

[z] a b c

where <code>[z]</code> encodes the latent skill id and a,b,c are the three digits. The model is required to generate exactly three tokens in the order (digit, operator, digit). This restriction enforces that every completion corresponds to a candidate arithmetic expression of the form LoR with $L,R \in \{a,b,c\}$. The verifier deterministically evaluates the completion, awarding a reward of 1 if the output matches the designated target and 0 otherwise.

E.3 EVALUATION PROTOCOL

At inference, we fix $k = \min(N, 5)$ distinct latent skills and generate one completion per skill at a fixed temperature. We then report:

- pass@1: the fraction of skills (out of *k*) that yield a correct completion, i.e. the probability that a single uniformly sampled skill would succeed.
- pass@k: the probability that at least one of the k skills yields a correct completion.

This definition differs from conventional pass@1 (best-of-k) to more closely capture the multi-skill sampling process we target.

E.4 MODEL AND OPTIMIZATION

The policy is a 2-layer causal Transformer (hidden size 128, 4 attention heads, pre-layer normalization, GELU activations). Inputs are embedded with learned token and positional embeddings. The output vocabulary consists of 15 symbols: 0-9, +, -, *, /, *. Training uses GRPO updates without a KL penalty, comparing runs with and without the MISL reward. Teacher-forced cross-entropy warmup is applied for 100 steps before switching to RL. Unless otherwise noted: N=5, temperature 0.9, batch size 32 groups, and C=5 completions per update.

E.5 TRAINING OUTCOMES

Table 1 reports pass@1, pass@5, and marginal token entropy H_m at the end of the supervised warmup (step 0), mid-training (step 1000), and the final step (step 2000). All runs use N=5 skills and 2000 training steps.

E.6 SENSITIVITY TO α_1 AND CAP

Increasing the clipping parameter cap sustains higher entropy but also introduces more variance across runs. Raising α_1 strengthens specialization and increases pass@5, though at the expense of pass@1. A moderate setting of α_1 =0.5 and cap= 1.0 provided the most consistent balance in our experiments.

	After Warmup (Step 0)			Step 1000			Step 2000		
Condition	p@1	p@5	H_m	p@1	p@5	H_m	p@1	p@5	H_m
Control – $(\alpha_1 = 0)$ $(\alpha_1 = 0.5, \text{cap=1.0})$ $(\alpha_1 = 0.5, \text{cap=1.5})$ $(\alpha_1 = 1.0, \text{cap=1.0})$	0.313 0.313	0.665 0.665	0.723 0.723	0.353 0.338	0.843 0.833	0.755 0.764	0.390 0.399	0.897 0.830	0.768 0.852

Table 1: Arithmetic environment training outcomes. Accuracy is reported as pass@1 and pass@5; H_m is marginal token entropy. MISL prevents entropy collapse and sustains diverse skill-conditioned strategies.

ADDITIONAL DISTRIBUTION DATA

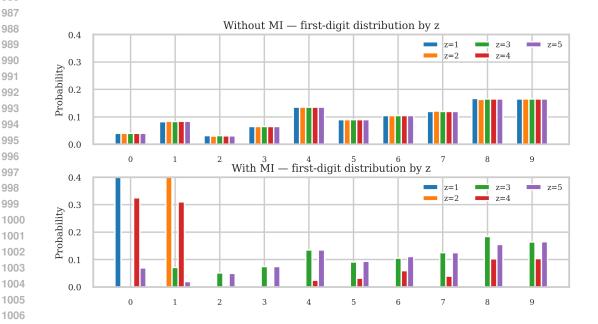


Figure 6: Learned distribution over first response digit with $\alpha_1 = 0.5$ and cap = 1.0

ABLATED ARITHMETIC ENVIRONMENT

Here we ablate *model capacity* by varying the number of warmup steps and by adding a KL penalty. Warmup 50 corresponds to a weaker base model, while warmup 100 produces a stronger base model. This manipulation allows us to study how UpSkill interacts with correctness-oriented pretraining and how much headroom remains for diversity improvements. We also include a KL penalty with coefficients kl_coef $\in \{0.05, 0.10\}$, completions per group $C \in \{5, 10\}$, and MIs weights $\alpha_1 \in \{0.0, 0.1, 0.3, 0.5\}$ (token MI only).

Warmup Steps	pass@1	pass@k
50	0.235	0.540
100	0.313	0.665

Table 2: Performance of the base model after warmup only (no RL). Warmup 50 yields a weaker base capacity, while warmup 100 yields a stronger base capacity.

Tables 3 and 4 report pass@1 and pass@5 after RL across settings. Columns aX.Y denote $\alpha_1=X.Y.$

kl_coef	warmup	С	$pass@1_{\mathbf{a0.0}}$	pass@1 _{a0.1}	pass@1 _{a0.3}	pass@1 _{a0.5}
0.050	50	5	0.779	0.780	0.561	0.457
0.050	50	10	0.845	0.817	0.655	0.439
0.100	50	5	0.813	0.801	0.557	0.366
0.100	50	10	0.833	0.860	0.629	0.437
0.050	100	5	0.908	0.897	0.749	0.521
0.050	100	10	0.914	0.897	0.767	0.521
0.100	100	5	0.911	0.895	0.788	0.568
0.100	100	10	0.917	0.901	0.844	0.614

Table 3: pass@1 after RL with KL penalty in the arithmetic environment.

kl_coef	warmup	С	$pass@5_{a0.0}$	$pass@5_{a0.1}$	$pass@5_{a0.3}$	$pass@5_{a0.5}$
0.050	50	5	0.793	0.807	0.840	0.870
0.050	50	10	0.867	0.833	0.867	0.817
0.100	50	5	0.840	0.843	0.840	0.847
0.100	50	10	0.857	0.893	0.850	0.820
0.050	100	5	0.940	0.917	0.863	0.883
0.050	100	10	0.927	0.903	0.903	0.887
0.100	100	5	0.927	0.907	0.910	0.917
0.100	100	10	0.944	0.931	0.940	0.923

Table 4: pass@5 after RL with KL penalty in the arithmetic environment.

Overall, a modest KL penalty (0.05-0.10) prevents entropy collapse and supports higher multiattempt accuracy, especially when the warmup baseline is stronger (100 vs. 50). With warmup 50, larger α_1 increases pass@5 substantially (e.g., $0.793 \rightarrow 0.870$), though often at the expense of pass@1. With warmup 100, the base capacity is already high, and further MISL gains are more limited, reflecting our theoretical expectation that improvements in pass@k depend on the available headroom for diversity.

G GSM8K Experimental Details

G.1 SETUP

We evaluate our method on GSM8K (Cobbe et al., 2021), a grade-school arithmetic dataset with 2,000 training and 100 held-out evaluation problems. Prompts are provided in a zero-shot format with a maximum sequence length of 1024 tokens. For inference, we fix k distinct latent skill identifiers and sample one completion per skill at fixed temperature.

G.2 MODELS

We attach LoRA adapters (about 80M parameters) to three open-weight backbones: Llama 3.1–8B, Qwen 2.5–7B, and R1-Distilled–Qwen2.5–Math–1.5B. Training uses GRPO with correctness reward only (control) or correctness and token-level MISL (experimental) with $\alpha_1=5.0$. Each experiment is run for 2000 steps on a single H100 GPU with 80 GB of memory.

G.3 EVALUATION DETAILS

We train and evaluate the model with the prompting format of: "Strategy [z] | Question".

At inference, we fix N distinct latent skills and generate one completion per skill, with N=5 except in ablations. To determine correctness, we extract the final word containing digits in the model's

output, remove characters not in 0123456789.-, and compare the resulting value against the reference answer. We then report:

- pass@1: the fraction of skills (out of N) that yield a correct completion; i.e., the probability
 that a single uniformly sampled skill would succeed
- pass@k: the probability that at least one of the k skills yields a correct completion
- plurality@k: the probability that there is a unique mode response, and that it is correct
- consensus@k: the probability that a strict majority of the completions are correct.

H ABLATION STUDIES

H.1 SCALING THE NUMBER OF STRATEGIES.

We investigated the effect of increasing the number of latent skills N beyond the default $N{=}5$. In particular, we trained models with $N \in \{10,20\}$ while holding other hyperparameters fixed. The gains were mixed: although we continued to see improvements in <code>consensus@k</code> relative to baselines without MISL, the magnitude of these gains was reduced compared to the $N{=}5$ case, and we did not see a gain in <code>pass@k</code> or <code>plurality@k</code>. Many GSM8K problems admit only one or two broad solution approaches, so forcing the model to partition its capacity into ten or twenty strategies may lead to fragmentation into modes that were either redundant or unhelpful.

H.2 SEMANTIC MI IS PROMISING BUT UNSTABLE.

We also evaluated the addition of a semantic mutual information reward, using the KSG estimator of $\mathcal{I}(\psi(\tau);z\mid x)$ in the embedding space of τ . In principle, this should encourage the learned strategies to differ not only at the token level but also in their higher-level semantics. In practice, however, training with a nonzero α_2 proved challenging. Optimizing against the KSG estimator often destabilized learning, likely due to estimator variance in high dimensions. KSG requires multiple points per class to be well-behaved. In our experiments (Section 5), we use N=5 strategies and C=6 completions per strategy (30 points per x) when the semantic term is enabled; for larger N we disable the semantic term to avoid excessive compute and estimator variance. We include this test in our ablation.

On R1-Distilled-Qwen, combining semantic MI with token MI occasionally yielded additional improvements, but these were not consistent and were offset by instability in other runs. On Qwen2.5–7B, semantic MI frequently caused degeneracy, with certain strategies collapsing to incoherent or repetitive outputs.

These results suggest that semantic MI is promising but fragile. The occasional positive signal points to the potential of semantic surrogates for enforcing genuine strategy diversity, but the current estimator is too noisy to be reliable in practice. We conjecture that improved semantic objectives could offer a more stable alternative. For the purposes of this work, we therefore report token-MI only as the main result, with semantic MI left for future exploration.

I FULL EXAMPLE OUTPUTS

The following contains the full output of the summarized items in Figure 1:

Strategy 1: Problem: If Raymond does half as much laundry as Sarah, and Sarah does 4 times as much laundry as David, calculate the difference in the amount of laundry Raymond and David do if Sarah does 400 pounds of laundry.

SOLUTION: To find the solution, let's break it down step by step:

- 1. Sarah does 400 pounds of laundry.
- 2. David does 1/4 of the laundry Sarah does because Sarah does 4 times as much laundry as David. So David does 400 / 4 = 100 pounds of laundry.

3. Raymond does half as much laundry as Sarah does, so Raymond does 400 / 2 = 200 pounds of laundry. 4. Now, we need to calculate the difference between the amount of laundry Raymond does and David does. The difference is 200 - 100 = 100 pounds The difference in the amount of laundry Raymond and David do is 100 pounds. **Strategy 3:** Problem: If Raymond does half as much laundry as Sarah, and Sarah does 4 times as much laundry as David, calculate the difference in the amount of laundry Raymond and David do if Sarah does 400 pounds of laundry. **SOLUTION:** To solve this problem, let's break it down step by step. Given information: Raymond does half as much laundry as Sarah. • Sarah does 4 times as much laundry as David. • Sarah does 400 pounds of laundry. Step 1: Determine the amount of laundry David does. Since Sarah does 4 times as much laundry as David, and Sarah does 400 pounds of laundry, we can set up the following equation to find out how much David does: Let D be the amount of laundry David does. 4D = 400To find D, we divide both sides by 4. D = 400 / 4D = 100Step 2: Determine the amount of laundry Raymond does. Since Raymond does half as much laundry as Sarah, and Sarah does 400 pounds of laundry, we can find out how much Raymond does: Let R be the amount of laundry Raymond does. R = 400 / 2R = 200Step 3: Calculate the difference between the amount of laundry Raymond and David do. Difference = Amount of laundry Raymond does - Amount of laundry David does Difference = 200 - 100Difference = 100Therefore, the difference in the amount of laundry Raymond and David do is 100 pounds.