# Multiplex Thinking:
# Reasoning via Token-wise Branch-and-Merge

**Yao Tang**[1*] **Li Dong**[2] **Yaru Hao**[2] **Qingxiu Dong**[2] **Furu Wei**[2] **Jiatao Gu**[1]

[1]University of Pennsylvania    [2]Microsoft Research

[1]{tangyao, jgu32}@seas.upenn.edu

Large language models often solve complex reasoning tasks more effectively with Chain-of-Thought (CoT), but at the cost of long, low-bandwidth token sequences. Humans, by contrast, often reason softly by maintaining a distribution over plausible next steps. Motivated by this, we propose **Multiplex Thinking**, a stochastic soft reasoning mechanism that, at each thinking step, samples $K$ candidate tokens and aggregates their embeddings into a single continuous *multiplex token*. This preserves the vocabulary embedding prior and the sampling dynamics of standard discrete generation, while inducing a tractable probability distribution over multiplex rollouts. Consequently, multiplex trajectories can be directly optimized with on-policy reinforcement learning (RL). Importantly, multiplex thinking is self-adaptive: when the model is confident, the multiplex token is nearly discrete and behaves like standard CoT; when it is uncertain, it compactly represents multiple plausible next steps without increasing sequence length. Across challenging math reasoning benchmarks, multiplex thinking consistently outperforms strong discrete CoT and RL baselines from Pass@1 through Pass@1024, while producing shorter sequences. The code and checkpoints are available at github.com/GMLR-Penn/Multiplex-Thinking.

## 1 INTRODUCTION

Large Language Models (LLMs) have exhibited exceptional reasoning capabilities on a wide range of complex tasks, especially in mathematics and logical problem solving (Cobbe et al., 2021; Lightman et al., 2023). A simple and effective way to elicit such behavior is *chain-of-thought* (CoT) prompting (Wei et al., 2022), which encourages the model to generate intermediate reasoning steps before producing the final answer. Beyond prompting, reinforcement learning (RL) can further improve reasoning by optimizing the model over diverse CoT rollouts using outcome- or process-level rewards (Guo et al., 2025), steering probability mass toward higher-reward reasoning trajectories.

However, both CoT prompting and RL on CoT rollouts are costly because they require generating long sequences of discrete tokens. Each rollout corresponds to a full, explicit reasoning trace, and exploring alternatives often resembles depth-first search (DFS): each sampled trace commits to a single trajectory before branching to others (Zhu et al., 2025). This cost motivates *continuous* reasoning tokens that can compactly encode a "superposition" over multiple candidate reasoning paths within a single token and decode in a more breadth-first search (BFS)-like manner (Zhu et al., 2025).

While existing continuous token approaches can reduce token cost, they are typically *deterministic*: given the next-token logits, they map the distribution to a single continuous vector, e.g., a hidden-state token (Hao et al., 2025) or a probability-weighted embedding mixture (Zhang et al., 2025). Determinism collapses the token-level policy distribution, making decoded rollouts identical and thus limiting exploration. This characteristic is fundamentally misaligned with RL, where on-policy stochastic rollouts are crucial to enable LLMs to effectively learn from *trial-and-error*s. There-
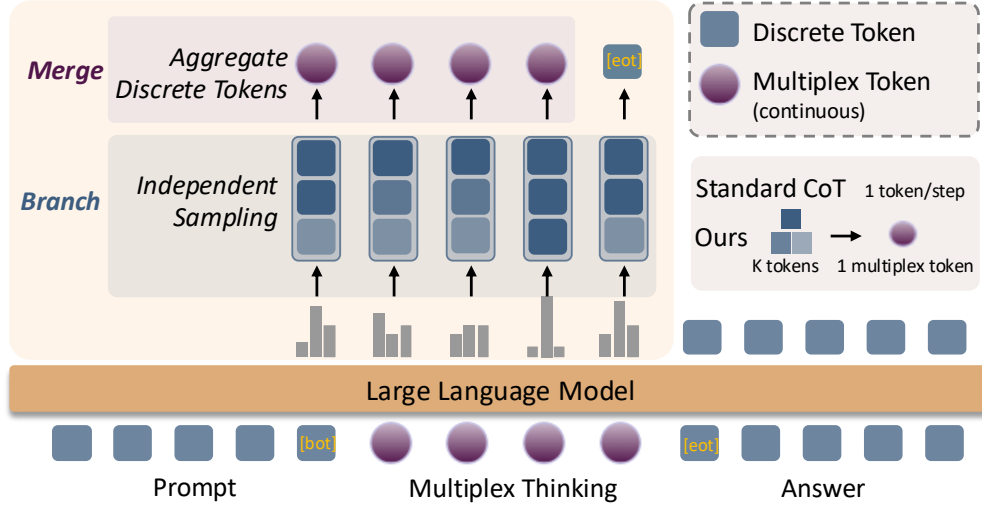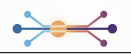
---

Figure 1: **Illustration of Multiplex Thinking**. The model first generates an initial probability distribution conditioned on the prompt and the begin of thinking token `[bot]`. Instead of committing to a single discrete token or a deterministic soft average, we conduct token-wise branching and merging by independently sampling $K$ discrete tokens and aggregate them into a continuous multiplex token. When one of the sampled discrete tokens is the end of thinking token `[eot]`, LLM continues conduct discrete decoding to give the answer. The design of sampling-based continuous thinking bridges the gap between continuous representation and stochastic discrete sampling, allowing for effective on-policy exploration and further RL training.
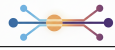
fore, we posit that continuous reasoning tokens should be *stochastic and sampling-based*, retaining discrete sampling dynamics while operating in a compact continuous space.

To fill this gap, we propose **Multiplex Thinking**, a discrete-sampling-based continuous reasoning paradigm. At each reasoning step, multiplex thinking samples $K$ independent tokens over model's token distribution, maps them to their vocabulary embeddings, and aggregates the embeddings into a single multiplex thinking token. When the logits are highly peaked (low entropy), the sampled tokens are likely to coincide, and the multiplex token effectively collapses to a standard single-token embedding. In contrast, when the logits have high entropy and exploration is desirable, the $K$ samples are more likely to differ, allowing the multiplex token to encode richer, multi-faceted information within one step. Crucially, multiplex thinking preserves both the vocabulary embedding prior and the sampling behavior of discrete token generation, while introducing a continuous-like reasoning representation. Because the sampled tokens at each step are independently from the underlying token distribution, the probability of a multiplex token is the product of the probabilities of its sampled tokens. This factorization allows us to explicitly model the probability of an entire multiplex thinking rollout and thereby optimize multiplex thinking directly with RL.

Empirically, we show that multiplex thinking consistently improves performance from Pass@1 to Pass@1024 across a range of challenging math reasoning benchmarks, surpassing strong discrete CoT and RL baselines. We further find that multiplex thinking achieves higher accuracy while maintaining better token efficiency: responses are shorter on average, since a single multiplex token can encode richer information than a standard discrete token.

Our main contributions can be summarized as:

- We introduce multiplex thinking, a token-efficient reasoning scheme that compresses multiple discrete CoT samples into continuous *multiplex* tokens while retaining stochastic exploration and probabilistic semantics.

- We formalize multiplex thinking as inducing a well-defined probability distribution over complete reasoning trajectories, enabling direct RL optimization over multiplex rollouts without paying the full token cost of long discrete CoT.

- We demonstrate consistent gains over strong discrete CoT and RL baselines across sampling budgets (Pass@1–Pass@1024), and provide analysis showing improved accuracy with shorter trajectories by compressing high-entropy reasoning steps.

## 2 BACKGROUND

**Chain-of-Thought (CoT) Reasoning**  In standard CoT reasoning, LLMs are prompted to output intermediate thinking tokens before predicting the final answers to solve problems. Given a language model $\pi_\theta$ with a vocabulary $V$ and the embedding layer $E \in R^{|V| \times d}$, the $k$-th token in the vocabulary can be embedded as $E[k]$, also denoted as $e(k)$. With an input question $q = (q_1, q_2, \ldots, q_L)$, language model $\pi_\theta$ first outputs the thinking sequence $t$ and then outputs the final answer $y$, where $t_i \sim \pi_\theta(e(q), e(t_{<i}))$ and $y_i \sim \pi_\theta(e(q), e(t), e(y_{<i}))$. While effective, standard CoT samples a discrete token at each reasoning step. This abandons the rich information of the distribution over the whole vocabulary, which further motivates research on reasoning over continous tokens to better preserve the information over the vocabulary.

**Soft Thinking**  Soft Thinking (Zhang et al., 2025) enhances LLM performance without fine-tuning by replacing discrete thinking tokens $t$ with continuous *concept tokens* $c$. At the $i$-th thinking step, it constructs a concept token by using the model's next-token distribution $p_i = \pi_\theta(\cdot \mid e(q), c_{<i})$ as weights to aggregate token embeddings over the vocabulary $\mathcal{V}$:

$$c_i = \sum_{k \in V} p_i(k) e(k).$$

Although Soft Thinking effectively compresses reasoning information into continuous vectors, it is inherently *deterministic*: given a context, the mapping from the logit distribution to $c_i$ is fixed. This lack of stochasticity prevents the model from exploring diverse reasoning paths, thereby limiting its potential to be optimized via reinforcement learning objectives that rely on trial-and-error.

**Reinforcement Learning with Verifiable Rewards (RLVR)**  RLVR trains language models on tasks with verifiable answers with a reinforcement learning objective (Lambert et al., 2025). Given a verifiable dataset consisting of question-answer pairs $\mathcal{D} = \{(q, y^\star)\}$ and an answer sampled from a language model $y \sim \pi_\theta(\cdot \mid q)$, a verifiable reward function $v$ can provide a reward $r = v(y, a)$ based on the ground-truth answer $y^\star$ and the sampled answer $y$. The model is trained to maximize the reward, given by:

$$\mathcal{J}_{\text{RLVR}} = \mathbb{E}_{(q, y^\star) \sim \mathcal{D}, y \sim \pi_\theta(\cdot|q)}[v(y, y^\star)].$$
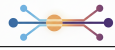
## 3 MULTIPLEX THINKING

In this section, we propose **multiplex thinking**, a reasoning paradigm effectively combining the information density of continuous representations with the probabilistic structure of discrete sampling. Based on multiplex thinking trajectories, we further introduce the reinforcement learning objective and provide entropy analysis of multiplex tokens.

### 3.1 FORMULATION

Given a question $q$ sampled from dataset $\mathcal{D}$, language model $\pi_\theta$ samples a multiplex thinking trace $c = (c_1, c_2, \ldots, c_L)$ followed by a final answer $y$. Different from standard decoding where a discrete token is sampled at each step, multiplex thinking constructs a mutliplex token with width $K$ by independently sampling $K$ times and aggregating the sampled discrete tokens. Formally, at reasoning step $i$, we independently sample $K$ discrete tokens $k_{i,1}, k_{i,2}, \ldots, k_{i,K}$ from the model's distribution $\pi_\theta(\cdot \mid e(q), c_{<i})$. We aggregate these discrete samples by averaging their one-hot vectors:

$$s_i = \frac{1}{K} \sum_{j=1}^{K} z_{i,j},$$

where $z_{i,j}$ is the one-hot vector corresponding to the discrete token $k_{i,j}$. When $K = 1$, $s_i = z_{i,1}$ collapses to a one-hot vector, and the multiplex token degenerates to a standard discrete token

sampled from $\pi_\theta$. When $K \to \infty$, the empirical distribution of $K$ i.i.d. samples converges to the model's LM head distribution $\pi_\theta(\cdot \mid e(q), c_{<i})$.

To obtain a continuous representation, we map $s_i$ through the embedding matrix $E \in \mathbb{R}^{V \times d}$, and we further define the continuous multiplex token by applying a vocabulary-space weighting $w_i \in \mathbb{R}^V$ to $s_i$:

$$c_i = E^\top (s_i \odot w_i),$$

where $\odot$ denotes element-wise multiplication. We consider two choices of $w_i$. **(i) Uniform averaging:** $w_i[v] = 1$ for all vocabulary indices $v$, which recovers $c_i = E^\top s_i$ (the averaged embedding of sampled tokens). **(ii) LM-head reweighting:** we set $w_i[v] = K \cdot \frac{\mathbf{1}[s_i[v]>0] \cdot \pi_\theta(v|e(q),c_{<i})}{\sum_{u=1}^V \mathbf{1}[s_i[u]>0] \cdot \pi_\theta(u|e(q),c_{<i})}$, i.e., we only reweight tokens that appear in the sampled set and scale them according to the model's LM-head probabilities. Empirically, we find that uniform averaging over samples and LM-head reweighting lead to comparable performance in Section 5.5. In our experiments, we adopt reweighting by default as it more directly reflects the model's confidence over the sampled candidates.

By independently sampling multiplex discrete tokens and aggregating these sampled tokens into a continuous representation, this design allows each $c_i$ to capture a stochastic ensemble of potential reasoning paths. When the distribution is sharp with low entropy, the samples collapse to the same token, reverting to standard discrete behavior. Conversely, high entropy distributions result in a diverse mixture, encoding exploration within a single continuous vector.

The probability of generating a specific multiplex token factorizes due to the independence assumption. Consequently, the log-probability of the entire reasoning trace $c$ is the sum of the log-probabilities of all constituent discrete samples:

$$\log \pi(c|e(q)) = \sum_{i=1}^{|c|} \sum_{j=1}^K \log \pi_\theta(k_{i,j}|e(q), c_{<i}).$$

## 3.2 Reinforcement Learning Objective

Leveraging the factorization above, we can directly optimize the model using Reinforcement Learning. We aim to maximize the expected reward of the generated answer $y$. The objective function is defined as:

$$\mathcal{J}_{\text{RL}}(\theta) = \mathbb{E}_{\substack{(q,y^\star) \sim \mathcal{D}, \\ c \sim \pi_\theta(\cdot|q), \\ y \sim \pi_\theta(\cdot|q,c)}} \left[ (\log \pi_\theta(c|e(q)) + \log \pi_\theta(y|e(q), c)) \cdot v(y, y^\star) \right].$$

This objective performs on-policy reinforcement learning over the joint generation process of the multiplex thinking trace $c$ and the final answer $y$,
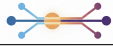
## 3.3 Entropy of Multiplex Token

To quantify the exploration capability of multiplex thinking, we compare the entropy of the multiplex token against the standard discrete token. In standard Chain-of-Thought (CoT), the discrete token $t_i$ is sampled from the policy $\pi_\theta(\cdot \mid q, t_{<i})$ at the decoding step $i$. The entropy of this single-step sampling is given by the standard Shannon entropy:

$$H_{\text{CoT}}(i) = - \sum_{v \in V} \pi_\theta(v \mid q, t_{<i}) \log \pi_\theta(v \mid q, t_{<i}).$$

In contrast, a multiplex token is constructed by independently sampling $K$ tokens $\mathcal{K}_i = \{k_{i,1}, \ldots, k_{i,K}\}$ from the distribution $\pi_\theta(\cdot \mid q, c_{<i})$. We conceptualize the generation of a multiplex token as a single *integrated action* that selects a composite outcome $(k_{i,1}, \ldots, k_{i,K})$ from the augmented state space $|\mathcal{V}|^K$. Under the assumption of sampling independence, we treat the multiplex token $c_i$ as a unified random variable and analyze the entropy of $c_i$ as a joint entropy, which is the sum of individual sampling entropies:

$$H(\mathcal{K}_i) = K \cdot H(\pi_\theta(q, c_{<i})).$$

We could observe that the entropy scales linearly with $K$, which corresponds to an exponential expansion of the effective exploration volume from $|\mathcal{V}|$ to $|\mathcal{V}|^K$. While standard CoT commits to a single discrete path, Multiplex Thinking leverages the high-capacity continuous space to encode a 'superposition' of $K$ paths simultaneously (Zhu et al., 2025). This allows the model to defer discrete decisions and retain probabilistic diversity within the reasoning trace. This property is particularly advantageous for reinforcement learning, as it provides a richer training signal. We also provide empirical analysis in Experiments 5.4.

## 4 EXPERIMENTS

In this section, we empirically evaluate multiplex thinking against discrete reasoning approaches and competitive continuous reasoning baselines. Our evaluation focuses on both Pass@1 accuracy and test-time scaling (Pass@1–Pass@1024), measuring effectiveness at small budgets as well as the exploration gains from increased rollouts. Overall, multiplex thinking achieves the best Pass@1 performance in major settings and exhibits stronger Pass@k scaling than discrete baselines.

### 4.1 EXPERIMENTAL SETUPS

**Implementation** We implement multiplex thinking on top of two open-source reasoning backbones: `DeepSeek-R1-Distill-Qwen-1.5B` and `DeepSeek-R1-Distill-Qwen-7B`. The models are optimized using Group Relative Policy Optimization (GRPO) (Shao et al., 2024). We train for 300 steps with a global batch size of 128 questions, a learning rate of $1 \times 10^{-6}$, and zero KL penalty and entropy penalty. A maximum response length of 4096 tokens is enforced for all training and evaluation stages.

During training, we generate 8 rollout samples per question with the temperature of 1.0 and the top-p of 1.0 to train LLMs with on-policy RL. During evaluation, we use a top-p of 0.95 to evaluate the Pass@1 performance, and these results are averaged on 64 runs. We also measure the Pass@$k$ performance (Chen et al., 2021) for $k \in \{1, 2, 4, \dots, 1024\}$ to fully investigate the exploration upper limit of different methods with the top-p of 1.0. Pass@$k$ measures the probability that at least one correct solution exists among $k$ sampled trajectories, serving as a proxy for the model's potential to discover a valid solution within a given exploration budget. The experiment evaluating the Pass@$k$ performance for $k \in \{1, 2, 4, \dots, 1024\}$ aims to provide a thorough empirical study on how the exploration space of multiplex thinking scales with exploration budget in comparison with discrete baselines. Results evaluating Pass@$k$ for $k \in \{1, 2, 4, \dots, 1024\}$ are computed by boostrapping for 1,000 times on a total of 1,024 runs. More implementation details are provided in Appendix A.1.

**Datasets** The training set is DeepScaleR-Preview-Dataset (Luo et al., 2025) consisting of approximately 40,000 unique problem-answer pairs. For evaluation, we use six challenging datasets: AIME 2024 (Veeraboina), AIME 2025 (Zhang & Math-AI, 2025), AMC 2023, MATH-500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), and OlympiadBench (He et al., 2024).

**Baselines** To validate the effectiveness of Multiplex Thinking, we compare against three distinct categories of methods: **Discrete CoT**: The backbone models using standard discrete Chain-of-Thought decoding without additional training. **Stochastic Soft Thinking**: A recent strong training-free continuous reasoning baseline (Wu et al., 2025). Building on the original deterministic Soft Thinking, this method injects stochasticity via the Gumbel–Softmax trick to mitigate the "greedy pitfall" identified in the original Soft Thinking and to enable exploration at test time. **Discrete RL**: The backbone models fine-tuned with GRPO on the same dataset using standard discrete tokens. This serves as the direct baseline to measure the gain from continuous exploration.

### 4.2 PASS@1 PERFORMANCE

Table 1 reports the Pass@1 accuracy across six mathematical reasoning benchmarks. Our proposed Multiplex Thinking consistently outperforms baselines, achieving the best results in 11 out of 12 experimental settings. Notably, it surpasses Discrete RL sharing the identical GRPO training setup across all the tasks. This observation validates the efficacy of exploring multiplex thinking trajecto-
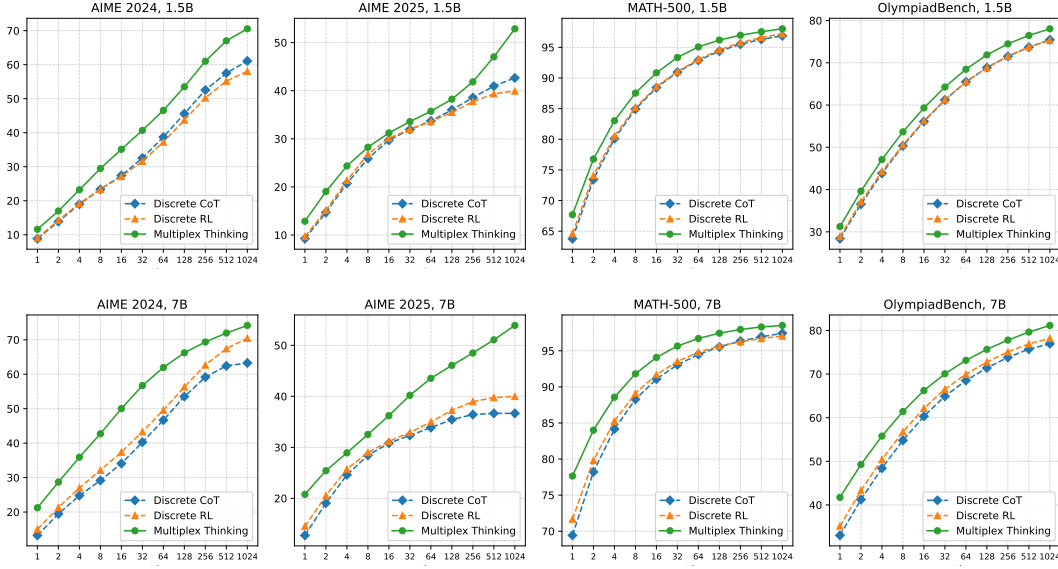
Figure 2: Pass@1–Pass@1024 performance on four representative datasets. Top row: 1.5B models; bottom row: 7B models. Full results on all six datasets are reported in Appendix A.2.

ries, proving that the performance gains stem from our unique representation brought from multiplex tokens rather than the RL training process alone.

Compared with Stochastic Soft Thinking, Multiplex Thinking demonstrates superior scaling behavior. On the 1.5B backbone, our method already demonstrates superior performance, achieving the better results on 4 out of 6 benchmarks. This advantage becomes more amplified at the 7B scale, where Multiplex Thinking dominantly get the highest scores across all six benchmarks. We postulate that the larger model capacity is essential for resolving the interference between superposed reasoning paths, thereby allowing the 7B model to fully leverage the exploration potential of multiplex trajectories.

Overall, Multiplex Thinking achieves the best Pass@1 performance in 11 out of the 12 evaluation settings spanning two model sizes and six datasets. This empirical results establish Multiplex Thinking as an effective method for advancing the reasoning capabilities of large language models.

Table 1: Pass@1 accuracy on six math reasoning benchmarks averaged over 64 runs. We compare discrete CoT decoding, discrete RL fine-tuning, Stochastic Soft Thinking, and our Multiplex Thinking on DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B. Multiplex Thinking achieves the best performance in most setups. The best results are **bolded** and the second best results are underlined in each column.

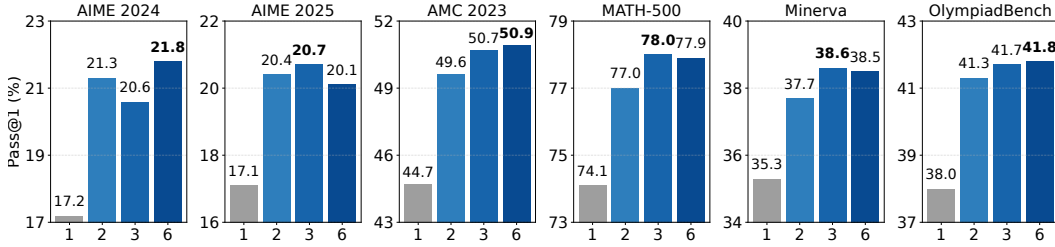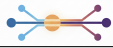| Exp Name | AIME 2024 | AIME 2025 | AMC 2023 | MATH-500 | Minerva | OlympiadBench |
|---|---|---|---|---|---|---|
| | *DeepSeek-R1-Distill-Qwen-1.5B* | | | | | |
| Discrete CoT | 10.2 | 11.3 | 36.9 | 66.3 | 23.4 | 30.5 |
| Stochastic Soft Thinking | <u>11.2</u> | **13.2** | **38.7** | <u>66.8</u> | <u>25.4</u> | 30.6 |
| Discrete RL | 10.5 | 12.0 | 38.4 | 66.7 | 24.3 | <u>31.2</u> |
| Multiplex Thinking | **11.8** | <u>12.8</u> | **38.7** | **67.5** | **26.2** | **31.3** |
| | *DeepSeek-R1-Distill-Qwen-7B* | | | | | |
| Discrete CoT | 15.7 | 16.0 | 42.4 | 71.6 | 33.3 | 35.6 |
| Stochastic Soft Thinking | <u>20.3</u> | <u>19.1</u> | <u>47.9</u> | <u>76.5</u> | <u>37.2</u> | <u>40.6</u> |
| Discrete RL | 17.2 | 17.1 | 44.7 | 74.1 | 35.3 | 38.0 |
| Multiplex Thinking | **20.6** | **19.7** | **50.7** | **78.0** | **38.6** | **41.7** |

Figure 3: Performance comparison under different multiplex widths $K$. The x-axis is the multiplex width $K$ and the y-axis is the Pass@1 performance on each dataset. The grey columns ($K = 1$) indicate the discrete RL performance and the blue columns represent Multiplex Thinking. The best results are **bolded** in each dataset.

### 4.3 TEST-TIME SCALING: FROM PASS@1 TO PASS@1024

Pass@$k$ with a large $k$ (e.g., $k = 1024$) serves as a proxy for the method's performance upper limit (Yue et al., 2025), reflecting the intrinsic exploration potential of the model. Figure 2 illustrates the Pass@$k$ performance scaling with respect to the number of sampled trajectories $k \in \{1, 2, \ldots, 1024\}$. As shown in Figure 2, Multiplex Thinking consistently achieves a higher upper bound compared to discrete baselines in most setups.

**Exploration potential on hard tasks.** The performance gap between Multiplex Thinking and discrete baselines has a trend to wide on challenging setups as $k$ increases. For instance, on AIME 2025 (7B), while the Discrete RL baseline begins to plateau around 40%, Multiplex Thinking continues to scale effectively, reaching approximately 55% at $k = 1024$. This substantial margin suggests that the continuous multiplex representation effectively expands the viable search space, enabling the model to uncover correct reasoning paths that are assigned negligible probability in the discrete token space.

**Scaling behavior differs across difficulties.** We could observe from Figure 2 that the benefits of Multiplex Thinking are difficulty-dependent. On simpler datasets like MATH-500, performance saturates quickly for all methods as accuracy approaches the ceiling. However, on tasks with sparse solution spaces (e.g., AIME 2025 and OlympiadBench), the ability of multiplex tokens to maintain superposed reasoning states proves crucial for escaping local optima, resulting in the "widening gap" trend observed in the figures.

**Sampling efficiency.** Beyond higher upper limits as exploration potential, Multiplex Thinking demonstrates superior sample efficiency. To achieve a target accuracy, our method requires significantly fewer samples than discrete baselines, directly translating to reduced test-time compute.

## 5 ANALYSIS

We analyze Multiplex Thinking through four questions that mirror the order of the following subsections. (1) **Do multiplex representations help without training?** Next, we evaluate an inference-only variant (`Multiplex Thinking-I`) to isolate the intrinsic benefit of multiplex trajectories from RL optimization (Table 2). (2) **How does multiplex width affect performance?** We first study the impact of multiplex width $K$ on Pass@1, with $K = 1$ recovering the Discrete RL baseline (Section 5.2). (3) **What is the compute trade-off between multiplex width and sequence length?** We then quantify how multiplex width can substitute for longer discrete rollouts (Table 3 and Figure 4). (4) **Which mechanisms and design choices drive the gains?** Finally, we examine training dynamics (policy entropy and response length), ablate the token aggregation strategy, and provide qualitative visualizations to illustrate how multiplex tokens encode and modulate multiple reasoning paths in practice (Section 5.4, Section 5.5, and Figure 6).
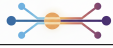
Table 2: Pass@1 (%) performance conducted on 7B backbone. The best results are **bolded** and the second best results are <u>underlined</u> in each column.

| Exp Name | AIME 2024 | AIME 2025 | AMC 2023 | MATH500 | Minerva | OlympiadBench |
|---|---|---|---|---|---|---|
| Discrete CoT | 15.7 | 16.0 | 42.4 | 71.6 | 33.3 | 35.6 |
| Stochastic Soft Thinking | 20.3 | 19.1 | 47.9 | <u>76.5</u> | <u>37.2</u> | <u>40.6</u> |
| Multiplex Thinking-I | <u>20.5</u> | <u>19.6</u> | <u>48.6</u> | 76.4 | 37.1 | <u>40.6</u> |
| Multiplex Thinking | **20.6** | **19.7** | **50.7** | **78.0** | **38.6** | **41.7** |

## 5.1 INTRINSIC CAPABILITIES OF MULTIPLEX REPRESENTATION

To disentangle the gains attributed to the multiplex representation from those yielded by Reinforcement Learning, we evaluate a training-free variant of our method, denoted as `Multiplex Thinking-I`. As shown in Table 2, we compare this inference-only baseline against the other two training-free baselines, including Discrete CoT and Stochastic Soft Thinking on the 7B model. Remarkably, applying Multiplex Thinking solely at inference time yields substantial performance gains over standard Discrete CoT. Compared to Stochastic Soft Thinking, Multiplex Thinking-I remains highly competitive, achieving better results on four out of six datasets. These empirical evidences demonstrate that the intrinsic capabilities of multiplex representation are beneficial for LLM reasoning. This inference-time superiority provides a strong starting point for further RL training; indeed, as shown in the final row of Table 2, applying RL optimization further amplifies these gains, consistently achieving the best performance across all benchmarks.

## 5.2 THE IMPACT OF TOKEN WIDTH $K$

We investigate the impact of the multiplex width $K$, the number of independently sampled discrete tokens aggregated into a multiplex token. Figure 3 compares the performance trained on the 7B backbone across varying widths $K \in \{1, 2, 3, 6\}$, where $K = 1$ corresponds to the standard Discrete RL baseline.

**Breaking the single-token bottleneck with $K \geq 2$.** As shown in Figure 3, transitioning from a single discrete token ($K = 1$) to a multiplex representation ($K \geq 2$) yields substantial gains across all benchmarks. And the performance gap between $K = 1$ and $K \geq 2$ is significant and consistent. For example, on AMC 2023, the precision jumps from 44.7% to 49.6% (+4.9%). This highlights that the primary advantage of our method stems from the paradigm shift which breaks the single-token bottleneck to enable exploration in a continuous latent space.

**Diminishing marginal gains with larger $K$.** Beyond this initial leap from $K = 1$ to $K = 2$, we observe performance continues to increase when multiplex width $K$ increases from 2 to 3 and 6 on AMC23, MATH-500, Minerva, and OlympiadBench. However, the marginal utility of widening the multiplex window gradually diminishes. As shown in Figure 3, the performance increase from $K = 2$ to $K \in \{3, 6\}$ is notable, yet the difference between $K = 3$ and $K = 6$ becomes considerably smaller. This suggests that the marginal gains from exploring additional tokens are most pronounced in the initial expansion. Therefore, a moderate width (e.g., $K = 3$ in our main experiments) is typically sufficient to capture the high-probability modes of the reasoning distribution, effectively covering the critical diverse paths. Further increasing the sample size yields diminishing returns, as the most valuable exploration directions are likely already included within the first few samples. We also provide Pass@$k$ for $k \in \{1, 2, \ldots, 1024\}$ analysis in Appendix A.2.2.

Beyond accuracy, we analyze how $K$ influences exploration behavior (policy entropy) and trajectory compactness (response length).

## 5.3 TEST-TIME COMPUTE: MULTIPLEX WIDTH V.S. SEQUENCE LENGTH

We study the trade-off between multiplex width and sequential length under different test-time compute budgets. Figure 4 shows that increasing the response length improves the performance of both discrete CoT and Multiplex Thinking-I, as expected. However, even with substantially shorter trajectories, Multiplex Thinking consistently achieves higher accuracy.
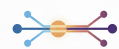
Table 3: Pass@1 accuracy averaged on six math reasoning datasets over 64 runs. We compare Multiplex Thinking-I-4k could match the performance of Discrete CoT-5k which has 25% more sequence token length budget.

| Exp Name | Averaged Accuracy (%) |
|---|---|
| Discrete CoT-4k | 35.8 |
| Discrete CoT-5k | 39.6 (+2.8) |
| Multiplex Thinking-I-4k | 40.5 (+4.7) |

Table 4: Entropy reduction ratio (%) under different multiplex width $K$, measured as the relative decrease in average policy entropy from the beginning to the end of training.

| $K$ Size | Entropy Reduction (%) |
|---|---|
| 1 (Discrete RL) | 9.44 |
| 2 | 5.82 |
| 3 | 6.03 |
| 6 | 7.09 |



Figure 4: Length scaling comparison. The y-axis is the accuracy averaged on six challenging reasoning datasets.

Figure 5: Response length dynamics under different multiplex width $K$. The x-axis is the training step and the y-axis is the averaged response length at each step.

To quantify this exchange rate between multiplex width and sequential length, we conduct a controlled experiment by scaling the inference token budget of the discrete CoT baseline from 4,096 to 5,120, while strictly constraining Multiplex Thinking-I to a 4,096-token budge, denoted as Multiplex Thinking-I-4k. As shown in Table 3, the performance of discrete CoT increases with raised sequence length budget. However, Multiplex Thinking-I-4k with a 4k limit consistently outperforms or matches the Discrete CoT-5k baseline with 20% shorter sequences. This indicates that performance gains do not solely depend on longer discrete rollouts, but can instead be achieved through richer token representations.

Further evidence is provided by the training dynamics in Figure 5, where we observe that Multiplex Thinking consistently generates trajectories with fewer tokens than the Discrete RL baseline, yet achieves superior accuracy. This observation aligns the intuition of multiplex thinking that multiplex tokens possess higher information density. Since each token encodes multiple potential paths, the model can express complex reasoning steps more compactly, effectively compressing the reasoning process in a shorter response. Importantly, increasing $K$ changes the number of sampled candidates per step, but does not require additional forward passes beyond sampling from the same logits distribution.

## 5.4 ENTROPY ANALYSIS

We analyze how multiplex tokens influence exploration during RL by measuring the *entropy reduction ratio* of the policy. Concretely, for each setting we compute $H_{\text{start}}$ as the average policy entropy over the first 10 training steps and $H_{\text{end}}$ over the last 10 steps, and report $(H_{\text{start}} - H_{\text{end}})/H_{\text{start}} \times 100$ (Table 4). A smaller reduction ratio indicates less entropy collapse and thus more sustained exploration throughout training.

As shown in Table 4, multiplex training exhibits consistently lower entropy reduction than the discrete RL baseline ($K = 1$), suggesting that multiplex tokens mitigate premature commitment to a single reasoning path. This trend is consistent with the higher Pass@k upper bounds observed for

larger $K$ (Figure 2), where maintaining exploration helps discover correct trajectories that would otherwise receive negligible probability under discrete decoding.

## 5.5 ABLATION STUDY ON TOKEN AGGREGATION STRATEGY

We conduct ablation study on a crucial component of multiplex thinking: the token aggregation strategy. We investigate the impact of the aggregation strategy used to construct the multiplex token representation. Specifically, we compare the probability-weighted summation (denoted as Multiplex Thinking-Weighted) against a simple unweighted average (denoted as Multiplex Thinking-Averaged) of the $K$ token embeddings. As detailed in Table 5, both strategies yield highly comparable empirical performance across model scales and benchmarks, and both consistently outperform the Discrete RL baseline by a significant margin. The similarity suggests that the effectiveness of Multiplex Thinking stems from the inclusion of diverse reasoning paths in the latent space rather than the specific weighting scheme used to combine them. This finding highlights the robustness of our approach, as the model can effectively learn to extract relevant features from the multiplex representation regardless of the precise linear combination coefficients.

Table 5: Pass@1 accuracy on six math reasoning datasets, averaged over 64 runs. We compare discrete RL and our Multiplex Thinking with two different token strategies, denoted as Multiplex Thinking-Averaged and Multiplex Thinking-Weighted respectively. The best results are **bolded**.

| Exp Name | AIME 2024 | AIME 2025 | AMC 2023 | MATH-500 | Minerva | OlympiadBench |
|---|---|---|---|---|---|---|
| | *DeepSeek-R1-Distill-Qwen-1.5B* | | | | | |
| Discrete RL | 10.5 | 12.0 | 38.4 | 66.7 | 24.3 | 31.2 |
| Multiplex Thinking-Averaged | 11.7 | **14.2** | 38.1 | **67.7** | **26.3** | 31.0 |
| Multiplex Thinking-Weighted | **11.8** | 12.8 | **38.7** | 67.5 | 26.2 | **31.3** |
| | *DeepSeek-R1-Distill-Qwen-7B* | | | | | |
| Discrete RL | 17.2 | 17.1 | 44.7 | 74.1 | 35.3 | 38.0 |
| Multiplex Thinking-Averaged | 19.9 | **20.0** | 49.9 | 77.6 | 38.4 | **41.8** |
| Multiplex Thinking-Weighted | **20.6** | 19.7 | **50.7** | **78.0** | **38.6** | 41.7 |

## 5.6 QUALITATIVE ANALYSIS

To better understand how Multiplex Thinking operates in practice, we visualize a representative reasoning trajectory in Figure 6. The figure displays a representative multiplex trajectory while solving a math problem. A key observation is that Multiplex Thinking effectively modulates its exploration strategy based on the uncertainty of the current reasoning state. Concretely, the trajectory alternates between *consensus* and *exploration* phases. During *consensus* steps (unboxed), the sampled candidates collapse to the same token, indicating a peaked next-token distribution and a locally stable reasoning state. In contrast, during *exploration steps* highlighted by yellow, purple, and red where divergent candidates are sampled, Multiplex Thinking compacts these alternatives into a single continuous multiplex token and continues the rollout while preserving uncertainty.

Notably, the highlighted exploration steps correspond to higher-entropy positions where multiple plausible continuations compete. Multiplex Thinking explicitly retains these alternatives via multiplex aggregation, enabling branching behavior at the very tokens that act as decision points. This is consistent with recent findings that high-entropy minority tokens function as critical forks in CoT reasoning and account for most gains in RLVR (Wang et al., 2025). We also provide a full trajectory in Appendix A.4.

## 6 RELATED WORKS

**Discrete Reasoning** Chain-of-Thought (CoT) prompting (Wei et al., 2022) has become the standard training-free method to elicit reasoning in LLMs. To further enhance these capabilities, methods like STaR (Zelikman et al., 2022) and ReST (Gulcehre et al., 2023) utilize iterative fine-tuning on self-generated rationales. More recently, DeepSeek-R1 (Guo et al., 2025) demonstrated that large-scale RL with verifiable rewards can incentivize the reasoning capability in LLMs. However, discrete RL
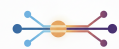
```
I/The need to figure out how many different ways the/these six cars can
stack up at/so the intersection such/so that all three lanes are occupied/
used. The \n/Hmm cars are distinguishable, meaning/and each one/car is
unique, and the/they since order in which they approach/arrive matters
is fixed ./.\n\n So \n/Hmm, the car/different first/cars car arrives/
gets first ,/, then the second, and so on up/until to the sixth ./car \n Hmm
First. Let/So me 's try/break start to break/parse this down. Since/
First the each car has three choices, without if/I might/could guess
think at/initially first that/it there/the total number of possible/
ways ways/arrangements they can choose/stack move their/the lanes
is ./just 3^6 ./, But/Let Because wait, that/the would 's/includes be/
include the case without/if the all/they order problem didn of/in 't
matter ,/or but right/here since ,/, the since the order cars does in/
matter because the/each cars are/pull distinguish coming/able and they/
arrive pull arrive pull/come one sequentially after at/another. So/Hmm
However , maybe maybe that each considering arrangement car is a just/
permutation sequence with where repetition specific ?/?\n\n So \n ,
yes/yeah thinking , so ./I each ./that total each/three number without
of/is should ways possible/possibilities without is/should any the
restrictions would/is be indeed 3 multiplied ^6 ,/. which is 729 .\n\n .
```

Figure 6: **Visualization of a representative multiplex thinking trajectory** ($K = 3$). At each step, the model samples $K = 3$ independent discrete tokens. The boxes summarize the sample outcomes:(i) when all three samples are distinct, the tokens are shown with yellow , purple , and red boxes to indicate high diversity; (ii) when two samples agree, a yellow box marks the majority token (sampled twice) and a purple box marks the minority token (sampled once); (iii) plain text (no box) indicates complete consensus, where all three samples are identical.

methods suffer from high computation cost of generating long sequences by conducting a depth-first style decoding to reach solutions (Zhu et al., 2025).

**Continuous Reasoning** A line of studies (Yang et al., 2024; Biran et al., 2024; Hao et al., 2025) define the hidden states of the transformer as latent reasoning steps. Yang et al. (2024) showed intermediate reasoning steps can be decoded from transformer hidden states. COCONUT (Hao et al., 2025) trained Large Language Models (LLMs) by using the last hidden states of the transformer as input embeddings and showed continuous chain of thought (CoT) could outperform discrete CoT on logical reasoning tasks. However, using hidden states as embeddings can suffer from representational misalignment (Zhang et al., 2025) when input embeddings and the prediction head become decoupled in larger LLMs, and full-model retraining can induce catastrophic forgetting (Xu et al., 2025). Another line of study(Zhang et al., 2025) proposed to utilize probability-weighted mixture of token embeddings as continuous tokens. This preserves the embedding prior and avoid extensive retraining.Both families of continuous tokens are inherently deterministic. Some recent works attempt to introduce stochasticity into continuous reasoning by injecting external noise Gaussian noise (Butt et al., 2025) or Gumbel noise (Wu et al., 2025) into the logits. However, none of these methods have explored using multiple independent sampling to form a single aggregated token representation at each decoding step that could achieve the best of discrete sampling and continuous representation.

**Parallel Reasoning** Parallel reasoning improves problem solving by exploring multiple reasoning paths and aggregating their outcomes. Representative methods include self-consistency (Wang et al., 2023) and Best-of-N (BoN) selection using outcome rewards (Cobbe et al., 2021) or process rewards (Lightman et al., 2023), as well as search-based approaches that branch over intermediate steps such as such as Tree-of-Thought (Yao et al., 2023) and adaptive parallel reasoning methods (Pan et al., 2025; Lian et al., 2025). A common drawback of parallel reasoning methods is that computation scales roughly linearly with the number of sampled paths, since each path requires generating a full sequence. Our proposed multiplex thinking serves as a complementary dimension to existing parallel reasoning strategies because it changes the per-step token distribution rather than the

outer-loop sampling budget. It can be seamlessly integrated into frameworks like Self-Consistency or BoN to further push the boundaries of reasoning.
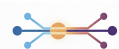
## 7 CONCLUSION

In this work, we introduced Multiplex Thinking, a novel framework that bridges the gap between discrete Chain-of-Thought and continuous reasoning representations. By aggregating multiple independent discrete tokens in continuous representations, our method achieves higher information density while preserving the probabilistic sampling required for effective reinforcement learning. Empirical evaluations across challenging mathematical benchmarks demonstrate that Multiplex Thinking consistently outperforms strong discrete CoT and RL baselines from Pass@1 to Pass@1024. Crucially, our analysis reveals that this performance gain comes with improved token efficiency, as the model learns to compress complex reasoning steps into shorter trajectories. These findings suggest that multiplex thinking is a promising direction for scaling test-time compute, offering a scalable path toward more capable and efficient reasoning models.

## ACKNOWLEDGEMENTS

## REFERENCES

Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. Hopping too late: Exploring the limitations of large language models on multi-hop queries, 2024. URL `https://arxiv.org/abs/2406.12775`.

Natasha Butt, Ariel Kwiatkowski, Ismail Labiad, Julia Kempe, and Yann Ollivier. Soft tokens, hard truths, 2025. URL `https://arxiv.org/abs/2509.19170`.

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 Mc-Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL `https://arxiv.org/abs/2107.03374`.

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.

Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.

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.

Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space, 2025. URL `https://arxiv.org/abs/2412.06769`.
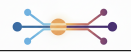
Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3828–3850, 2024.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL `https://arxiv.org/abs/2411.15124`.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857, 2022.

Long Lian, Sida Wang, Felix Juefei-Xu, Tsu-Jui Fu, Xiuyu Li, Adam Yala, Trevor Darrell, Alane Suhr, Yuandong Tian, and Xi Victoria Lin. Threadweaver: Adaptive threading for efficient parallel reasoning in language models, 2025. URL `https://arxiv.org/abs/2512.07843`.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.

Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. `https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL\-19681902c1468005bed8ca303013a4e2`, 2025. Notion Blog.

Jiayi Pan, Xiuyu Li, Long Lian, Charlie Snell, Yifei Zhou, Adam Yala, Trevor Darrell, Kurt Keutzer, and Alane Suhr. Learning adaptive parallel reasoning with language models, 2025. URL `https://arxiv.org/abs/2504.15466`.

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.

Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems*, pp. 1279–1297, 2025.

Hemish Veeraboina. Aime problem set 1983–2024. Kaggle dataset. URL `https://www.kaggle.com/datasets/hemishveeraboina/aime-problem-set-1983-2024`. Accessed: 2025-12-17.

Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen Yu, Gao Huang, and Junyang Lin. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning, 2025. URL `https://arxiv.org/abs/2506.01939`.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_VjQlMeSB_J.

Junhong Wu, Jinliang Lu, Zixuan Ren, Gangqiang Hu, Zhi Wu, Dai Dai, and Hua Wu. Llms are single-threaded reasoners: Demystifying the working mechanism of soft thinking, 2025. URL https://arxiv.org/abs/2508.03440.

Yige Xu, Xu Guo, Zhiwei Zeng, and Chunyan Miao. Softcot: Soft chain-of-thought for efficient reasoning with llms. *arXiv preprint arXiv:2502.12134*, 2025.

Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do large language models latently perform multi-hop reasoning? *arXiv preprint arXiv:2402.16837*, 2024.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023. URL https://arxiv.org/abs/2305.10601.

Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?, 2025. URL https://arxiv.org/abs/2504.13837.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. STar: Bootstrapping reasoning with reasoning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_3ELRdg2sgI.

Yifan Zhang and Team Math-AI. American invitational mathematics examination (aime) 2025, 2025.

Zhen Zhang, Xuehai He, Weixiang Yan, Ao Shen, Chenyang Zhao, Shuohang Wang, Yelong Shen, and Xin Eric Wang. Soft thinking: Unlocking the reasoning potential of llms in continuous concept space, 2025. URL https://arxiv.org/abs/2505.15778.

Hanlin Zhu, Shibo Hao, Zhiting Hu, Jiantao Jiao, Stuart Russell, and Yuandong Tian. Reasoning by superposition: A theoretical perspective on chain of continuous thought, 2025. URL https://arxiv.org/abs/2505.12514.

# A  APPENDIX

## A.1  IMPLEMENTATION DETAILS

We implement our framework based on verl (Sheng et al., 2025) and SGLang. Our codebase references the implementation design of Soft Thinking (Zhang et al., 2025). We use the the version 0.4.9.post6 of SGLang. All experiments are conducted on $8\times$ NVIDIA DGX B200 GPUs[2].

**Stopping Criteria.** In multiplex thinking, the transition from the thinking phase to the final answer generation is triggered when the discrete token with the highest probability corresponds to the special token `</think>`. We explicitly adopt this strategy over training-free heuristics, such as monitoring consecutive low-entropy tokens for early stopping in other works. We observed that such heuristics introduce artificial patterns that the model tends to exploit during RL optimization, resulting in training instability and the generation of incoherent content. By removing these handcrafted constraints, we allow the RL objective to naturally regulate the thinking process. Since indefinite thinking without producing a final answer yields low rewards, the model autonomously learns to generate the stop token at appropriate steps to maximize the return. This approach helps reduce the reward hacking associated with rule-based stopping criteria.

### A.1.1  HYPER-PARAMETERS

Table 6 provides a detailed summary of the hyper-parameters used for GRPO training of the discrete RL and multiplex thinking.

| Params | Values |
|---|---|
| Batch size | 128 |
| PPO mini batch size | 128 |
| Rollout number | 8 |
| Learning rate | $10^{-6}$ |
| Sampling temperature | 1.0 |
| Sampling top $p$ | 1.0 |
| Multiplex width $K$ | 3 |
| Max prompt length | 1024 |
| Max response length | 4096 |
| Entropy loss coefficient | 0 |
| KL loss coefficient | 0 |
| Model data type | bfloat16 |

Table 6: Hyper-parameters used in GRPO training.

## A.2  ADDITIONAL RESULTS

In this section, we present the complete experimental results that were omitted from the main body due to space constraints.

### A.2.1  FULL PASS@1–PASS@1024 RESULTS

Figure 7 illustrates the detailed Pass@$k$ performance (from $k = 1$ to 1024) across all six benchmarks for both the 1.5B and 7B model scales, supplementing the representative results discussed in the main text.

---

[2]Computations are conducted using bfloat16 precision on the NVIDIA Blackwell architecture to balance numerical stability and computational throughput.

Figure 7: Pass@1 to Pass@1024 performance spanning all tasks and all models.

### A.2.2 PASS@1–PASS@1024 RESULTS COMPARING DIFFERENT MULTIPLEX WIDTH $K$

In Section 5.2, we report the ablation study on multiplex width $k$ with the Pass@1 performance using top-p of 0.95. We extend this analysis by presenting the full Pass@$k$ trajectories for $k \in \{1, 2, 4, \dots, 1024\}$ in Figure 8 on four representative datasets.

Figure 8: K size ablation.



(a) Training score dynamics

(b) Validation score

Figure 9: Training dynamics comparison of multiplex thinking and discrete reasoning.

From Figure 8, we could observe that there is a significant performance gap between the baseline Discrete RL ($K = 1$) and the multiplex thinking variants ($K \geq 2$). This empirical observation strongly corroborates our analysis in Section 5.2. While the transition from discrete RL to multiplex thinking overcomes the limitation of depth-first style search, scaling the multiplex width further results in diminishing returns, with the performance curves for $K = 2, 3$, and 6 remaining closely clustered.

### A.3 TRAINING DYNAMICS

We provide the training score and validation score in Figure 9. The trained model is validated every 25 training steps and we use the Pass@4 metric on MATH-500 for validation.

### A.4 FULL TRAJECTORY EXAMPLE

To qualitatively illustrate the on-policy exploration of Multiplex Thinking, we visualize a full reasoning trajectory below.

so that all three lanes are occupied. The cars are distinguishable, meaning each one is unique, and the order in which they approach matters is fixed.\n\n So, the first car arrives first, then the second, and so on up to the sixth.\n Let me try to break this down. Since each car has three choices, I might guess initially that the total number of possible arrangements they can choose their lanes is just 3^6. But wait, that would includes be the case without the order problem didn't matter, but right here since the order does matter because each cars are distinguishable and they arrive one sequentially after another. So However, maybe that each arrangement is a permutation with repetition?\n\n So, yes yeah thinking, so I each total number without of is should ways possible possibilities without any restrictions would be indeed 3 multiplied ^6, which is 729.\n\n But the problem specifies that all three lanes must be occupied. So, I need to subtract the cases scenarios where one or more lanes are empty.\n\n This sounds seems feels like an inclusion-exclusion problem principle. In general, inclusion-exclusion is a used principle counting method where you we can subtract calculate find the number total of elements in a union set of overlapping sets by including adding the excluding the sizes of various intersections. So, applying I that this here idea, I perhaps need can to subtract the all cases number where one at only lane is empty, but add then add back in the cases where two lanes are empty because since those we I were have cases subtracted too many times.\n\n Let me recall write structure this that out formal.\n\n:\n\n First Total, total the number of ways without : any restrictions restrictions: is 3^6 = .\n\n 729.\n\nNow, subtract the compute the number cases of ways where at one least one lane is empty.\n\n There are three lanes, so the number of ways where the a left-turn lane is empty is, 2^6, because since each car can has only either choose go between straight -going or and right-turn lanes. Similarly, the number same of ways where the straight-going lane is empty is 2^6, and the number of ways where the right-turn lane is empty is also 2^6.\n\nSo, subtract each ing these, : we it get have :\n\n 3 7*(2^6). ) = 3*64 = 192.\n\n But wait now, now inclusion by we've subtracted too much because the when cases where two lanes are empty have been subtracted multiple three times. So For Each, we need have to add them those back in.\n\n The How many ways are there where two lanes are empty? For Well each example, if both the left-turn and straight-going lanes are empty, all then cars six must would go right straight. There Similarly's only one 1 way for each that pair such lane of lanes. Since There there are C(3,2) = 3 pairs such ways of lanes, ( each contributing with 1 way ^, ( where so we adding we add back. 3 * ( *1 = ^ 3.\n\n So Therefore, applying putting inclusion-exclusion, the total number of ways valid where all three lanes are occupied is:\n\n Total 3 = 3^6 − 3*(2^6) + 3*(1 ^ ) ^6 )\n\n ) )\n Calcul Pl ating that :\n\n:\n\n3^6 = 729\n 3 2 * ( *2^6) = 3 1*64 = 192\n3*(1^6) = 3*1 = 3\nnSo ,\n\n Total = 729 − 192 + 3\n Let Calcul Comput me compute that step :\n\n by step :\n\n .\n\n729 − 192 = is 537\n537 + 3 = 540\nSo, the is there total number is of ways is the 540.\n\n ? Wait Hmm
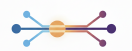
, but hold me think again. Is this the correct approach? Alternatively, another way to think about it is using the principle of inclusion-exclusion for counting surjective functions. Since we want all three lanes occupied have at least one car, this is equivalent to counting the number of onto functions from the set of six cars to the three lanes. The formula for this is indeed:

3! * S(6,3), where S(6,3) is the Stirling numbers of the second kind, which count the number of ways to partition a set of $n$ objects into $k$ non-empty subsets.

But wait, in our case, the cars are distinguishable, and the lanes are distinguishable, so the number of onto functions is indeed 3! * S(6,3). So let me recall that the formula for Stirling numbers of the second kind S(n,k) = S(n-1,k-1) + k*S(n-1,k). So, perhaps I can compute S(6,3).

But perhaps instead we can compute it directly. Or instead I could use inclusion-exclusion formula as done. Wait, but actually hold on, let me make sure the formula results. The number of onto functions from a set of size $n$ to a set of size $k$ is given:

$k$! * S(n,k). So, for $n$=6 and k=3, it would be 3! * S(6,3).

But what perhaps is S(6,3)? I think S(6,3) can be calculated using inclusion-exclusion formula as well:

Alternatively, I can remember that S(n,k) is

= (1/k!) * sum_{i=0 to k} (-1)^i * C(k,i)* (k – i)^n

So for applying S(6,3): S(6,3) = (1/3!) * [C(3,0) * (3^6 – C(3,1)*2^6 + C(3,2)*1^6 – C(3,3)*0^6]\n\n Which is (1/6)[(1*729 – 3*64 + 3*1 – 1*0] = (1/6) (729 – 192 + 3) = (1/6)(540)= 90.\n\n So S(6,3) = 90.\n\n Therefore, the number of onto functions is

3! * 90 = 6 * 90 = 540. Which matches our result earlier from inclusion-exclusion earlier. So, that this's consistent.\n\n So Thus, both methods give me the 540.\n\n Wait, but let me just think through this another way to make sure.\n\n Suppose I model each arrangement of each car the six cars to one of the three lanes, with such the condition that each lane has at least one car. So, the number of such assignments is equal to the number of ways to distribute six distinguishable objects into three distinguishable boxes, with each box containing at least one object.\n\n Which is exactly again the definition inclusion of onto functions, so yeah

onto 3! * S(6,3) = 540.\n\n Alternatively, another I can think of it as follows:\n\n assigning

choosing each of the six cars to a lane, without subtracting the assignments where one or more lanes are empty. Which is what I did initially, leading me to the 540. So, seems like consistent multiple answer is 540. Wait, but maybe now just hang to double check, let me consider plug of in another smaller case. Suppose instead of six cars, we have two cars and three lanes. How many ways can they stack up so both lanes are occupied? Using the same formula: total ways without restriction: 2^2 = 4. Subtract the cases where one lane is empty: 2*(1^2) = 2.\n\nSo, total = 4 − 2 = 2. Which is correct, because the two cars possible have either first car left, second right; or first car straight, second left. Alternatively, with the number of onto functions from 2 cars to 2 lanes is 2! * S(2,2) = 2 * 1 = 2, which matches.\n\nAnother test with n cars and two lanes. Total ways without: 2^3 = 8. Subtract cases where one lane is empty: 2*(1^3) = 2. So Total onto = 8 − 2 = 6. But actually the number of ways where both lanes are occupied is 6, which is also equivalent 2! * S(3,2) = 2 * 3 = 6. Which is correct since applying inclusion formula seems to hold.\n\nSimilarly, for the case of three cars and three lanes, the number total of onto functions would be 3! * S(3,3) = 6 * 1 = 6.\n\nWhich is also correct since each car must be assigned uniquely to a lane, so permutations 3! = 6.\n\nTherefore the same formula holds. Therefore going back, for six cars and three lanes, the number is 540.\n\nTherefore, I think I can confidently be confident that the answer is 540. Wait, but just let me make sure entirely, I perhaps should think about it another arranging as assigning each car position to a lane, and then subtract computing cases where one or more lanes are empty.\n\nSo Total, similar steps as inclusion exclusion.\n\nTotal number of assignments: 3^6 = .\n\n729.\n\nNumber of assignments where left lane is empty: 2^6 = 64.\n\nSimilarly, assignments where straight-going is empty: 64.\n\nAnd assignments where right-turn is empty: 64.\n\nSo, subtract these 3*64= 192.\n\nBut now wait, we have subtracted too much because the cases where two lanes are empty have been subtracted twice times.\n\nSo how many assignments have two lanes empty?\n\nIf left two

and straight are empty: all six cars must go right: only 1 way assignment.\n\nSimilarly, left and right empty: all cars go straight: .\n\n 1 way.\n\nStraight and right empty: all cars go left: 1 way.\n\n So three such assignments, each contributing 1 .\n\n . So, we have to add them back .\n\n So, 3*1 = 3.\n\n Therefore, total number of assignments where no lanes are empty is :\n\n729 – 192 + 3 = 540.\n\n Yes, so that matches again.\n\n Alternatively, I can be confident that the 540 is the correct answer.\n\n So, the number of ways the six cars can stack up at the intersection with all three lanes occupied is 540.\n\n**Final Answer**\nThe number of ways is \boxed{540}.\n</think>\nTo determine the number of ways six distinguishable cars can stack up at a red light such that all three lanes (left-turn, straight-going, and right-turn) are occupied, we use the principle of inclusion-exclusion.\n\n1. Calculate the total number of unrestricted ways the cars can choose their lanes:\n \[\n 3^6 = 729\n \]\n\n2. Subtract the number of ways where at least one lane is empty. There are three lanes, so the number of ways for each lane to be empty is:\n \[\n 2^6 = 64\n \]\n Therefore, total ways for one lane being empty:\n \[\n 3 \times 64 = 192\n \]\n\n3. Add back the cases where two lanes are empty, since these were subtracted too many times. There are three pairs of lanes, and each pair contributes 1 way:\n \[\n 3 \times 1 = 3\n \]\n\n4. Apply the inclusion-exclusion principle:\n \[\n 729 – 192 + 3 = 540\n \]\n\nThus, the number of ways the six cars can stack up so that all three lanes are occupied is \(\boxed{540}\).<|end_of_sentence|>' (SPECIAL)