

DISC: DYNAMIC DECOMPOSITION IMPROVES LLM INFERENCE SCALING

Anonymous authors

Paper under double-blind review

ABSTRACT

Inference scaling methods often rely on decomposing problems into steps (or groups of tokens), followed by sampling and selecting the best next steps. However, these steps and their sizes are often predetermined or manually designed based on domain knowledge. We propose dynamic decomposition, a method that adaptively and automatically fractions solution and reasoning traces into manageable steps during inference. By more effectively allocating compute – particularly through subdividing challenging steps and prioritizing their sampling – dynamic decomposition significantly improves inference efficiency. Experiments on benchmarks such as APPS, MATH, and LiveCodeBench demonstrate that dynamic decomposition outperforms static approaches, including token-level, sentence-level, and single-step decompositions. These findings highlight the potential of dynamic decomposition to improve a wide range of inference scaling techniques.

1 INTRODUCTION

Scaling inference efficiency remains a fundamental challenge for large language models (LLMs). Many existing approaches improve inference by decomposing problems into smaller steps and systematically exploring different solutions (Feng et al., 2023; Zeng et al., 2024; Wu et al., 2024; Nori et al., 2024; Snell et al., 2024; Brown et al., 2024; Gandhi et al., 2024; Lee et al., 2025; Light et al., 2024a; Wang et al., 2025).

Some decomposition methods rely on domain-specific heuristics and hand-crafted rules (Yao et al., 2024; Zelikman et al., 2023; Zhou et al., 2022). However, manually partitioning problems or designing task-specific heuristics is costly and lacks generalization. Moreover, identifying critical steps for an LLM can be non-trivial for humans. As shown in Sec. 3.5, LLMs may assign importance to seemingly trivial words (e.g., *therefore* or *which*), which, while counterintuitive to humans, play a crucial role in autoregressive generation (Lin et al., 2025). Other approaches employ fixed, uniform step sizes, such as token- or sentence-level decomposition (Feng et al., 2023; Guo et al., 2025). All these methods rely on **static decomposition strategies**, where step sizes are predefined or determined via heuristics. Such rigidity wastes compute on steps that are easy for the LLM (but potentially difficult for humans) while undersampling more challenging steps.

To overcome these limitations, we propose DISC (Dynamic decomposition Improves Scaling Compute), a recursive inference algorithm that dynamically partitions solution steps based on difficulty. Unlike prior methods, DISC **adapts decomposition granularity** during inference based on both the available budget and problem complexity, ensuring finer granularity for more difficult steps. By leveraging the autoregressive nature of LLMs, DISC efficiently **locates difficult steps** through binary partitioning, focusing compute on challenging regions rather than wasting resources on trivial steps. DISC is *generalizable* and requires *no human supervision, domain-specific heuristics, prompt engineering, or process annotations*, making it widely applicable across tasks.

Our main contributions are:

- We introduce DISC, a method for recursive partitioning and decomposing solutions during inference *without human supervision, domain-specific heuristics, or process reward models*.
- We demonstrate how DISC integrates decomposition with inference-time search, **allocating compute to high-impact, difficult steps**.
- We show that DISC improves inference scaling in terms of both **sample efficiency** and **token efficiency**.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

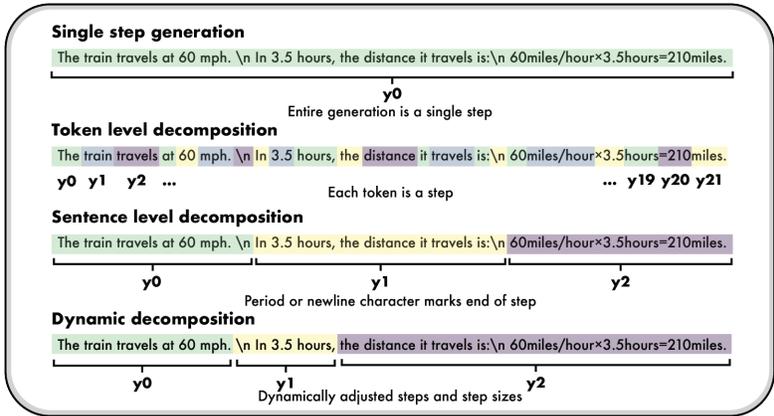


Figure 1: Comparison of different automatic decomposition methods based on step size determination.

- We provide insights into how LLMs reason and plan by identifying **critical steps** in their generation process.

2 PRELIMINARIES

2.1 PROBLEM SETTING

We consider a reasoning and code generation setting where a dataset $\mathcal{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ consists of problem prompts \mathbf{x} , and a reward model $R : \mathcal{X} \cdot \mathcal{Y} \rightarrow [0, 1]$ evaluates generated solutions $\mathbf{y} \in \mathcal{Y}$. This includes program synthesis, where correctness is verified using ground-truth tests (Chen et al., 2021; Austin et al., 2021), and mathematical reasoning, where solutions are validated numerically (Hendrycks et al., 2021a; Cobbe et al., 2021). The reward model can be a ground-truth verifier, a trained heuristic (Zhang et al., 2024), self-consistency (Wang et al., 2023a), or an LLM-as-a-judge (Zheng et al., 2023). Since our focus is on step decomposition rather than verification, we use the ground-truth reward model where available. We assume access to a pretrained language model π that generates text autoregressively. A generated response \mathbf{y} consists of both the final solution and the reasoning chain leading to it, and can be represented as a sequence of tokens $\mathbf{y} = (y_0, \dots, y_{L_y})$. Additionally, solutions can be partitioned into **solution steps** $\mathbf{y} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$, where each step \mathbf{y}_i is a contiguous string of tokens. A **partial solution** up to step k is defined as $\mathbf{y}_{1..k} := \mathbf{y}_1 \cdot \mathbf{y}_2 \cdot \dots \cdot \mathbf{y}_k$, and its **rollout** or **completion**, denoted $\mathbf{y}_{1..k+}$, is the continuation generated by π until an end-of-sequence token (EOS). The **size** of a solution step, $|\mathbf{y}_i|$, refers to its length in tokens or characters.

2.2 PRIOR AUTOMATIC DECOMPOSITION METHODS

Single-step generation. In a single-step generation, the entire solution is generated in one pass from the prompt to the EOS token, treating it as a single action. This approach underlies the widely used inference scaling method **best of n** (BoN) (Cobbe et al., 2021; Lightman et al., 2023; Snell et al., 2024; Liang et al., 2024), where n complete solutions are sampled, and the highest-scoring one is selected. Single-step generation also plays a role in alignment and fine-tuning methods such as **DPO** (Rafailov et al., 2024) and **RLOO** (Ahmadian et al., 2024).

Token-level decomposition. At the opposite end of the spectrum, token-level decomposition treats each atomic token as an individual step. While this approach dramatically increases search complexity, it enables fine-grained search that can yield higher performance gains given sufficient compute (Feng et al., 2023).

Newline and sentence-level decomposition. A commonly used decomposition method segments LLM generations into sentences or lines based on delimiters such as periods or newlines (Hao et al., 2023; Feng et al., 2023; Yao et al., 2024). Typically, each newline corresponds to a new paragraph, equation, or line of code, which often encapsulates a distinct reasoning step.

See App. C for more discussion of prior methods and related works.

Problem: Automatic and scalable decomposition

Existing decomposition methods are task-specific, manual, and static, limiting their adaptability and scalability.

2.3 STEP SAMPLING

Inference-time scaling methods must balance exploration at the current step with exploration of future steps. We implement a simple dynamic sampling process, referred to as **negative binomial sampling**, where we continue sampling completions until the *sum of their rewards exceeds a predefined threshold* σ . More formally, the number of samples M drawn from a partial solution $\mathbf{y}_{1\dots k}$ is the smallest integer satisfying $\sum_{i=1}^M R(\mathbf{x} \cdot \mathbf{y}_{1\dots k+}^{(i)}) \geq \sigma$. Here, $\mathbf{y}_{1\dots k+}^{(i)}$ represents the i -th sampled completion from the partial solution. This process ensures efficient allocation of compute by dynamically adjusting the number of samples per step. It achieves this by either: (a) *continuing to sample until a sufficiently high-reward completion is found*, or (b) *stopping early when additional sampling is unlikely to yield significant improvements*, thus redirecting the compute to future steps. The stopping criterion is governed by σ : when accumulated reward from completions surpasses σ , the method assumes further sampling is unnecessary. For a fair comparison, we apply this sampling method *uniformly across all decomposition methods* in our experiments.

2.4 INFERENCE SCALING METHODS AND DECOMPOSITION

Since our study focuses on decomposition rather than search, we primarily use **greedy step search** as the search method. In greedy step search, multiple candidate steps are sampled at each iteration, but only the highest-scoring step is retained, while the rest are discarded. The process then repeats, conditioning future steps on the best step found so far. We also perform ablation studies comparing **Monte Carlo Tree Search** (MCTS) (Feng et al., 2023; Light et al., 2024b) and **beam search** (Xie et al., 2024), two commonly used inference scaling methods. These comparisons, presented in Sec. 4.2, highlight how different search strategies interact with decomposition. Additional details on MCTS and beam search are provided in App. F.

3 METHODOLOGY

3.1 DISC ALGORITHM

The DISC algorithm employs recursive binary decomposition to iteratively break down complex solutions into smaller, more manageable steps. Given a problem prompt \mathbf{x} , the algorithm outputs a decomposition of a solution, $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$, such that the concatenation $\mathbf{y}_{1\dots K}$ forms a complete solution to \mathbf{x} .

The algorithm operates in two key stages:

1. **Solution sampling.** Starting from \mathbf{x} and $\mathbf{y}_0 = \emptyset$, the algorithm generates complete solutions $\mathbf{y}_k \sim \pi(\cdot | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)})$ using a policy π like in single step generation. The best solution, \mathbf{y}_k^* , is selected based on the reward model R .
2. **Recursive partitioning.** The selected solution \mathbf{y}_k^* is partitioned into two segments, $\mathbf{y}_k^* = \mathbf{y}_a^* \cdot \mathbf{y}_b^*$, based on a predefined partition fraction α , where $|\mathbf{y}_a^*| \approx \alpha |\mathbf{y}_k^*|$. For each part, a priority metric h is estimated: $\hat{h}(\mathbf{y}_a^* | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)})$ and $\hat{h}(\mathbf{y}_b^* | \mathbf{x} \cdot \mathbf{y}_{1\dots a}^*)$, usually through rollouts of the step using π . The part with the lower priority is further partitioned.

- If $\hat{h}(\mathbf{y}_a^* | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)}) \geq \hat{h}(\mathbf{y}_b^* | \mathbf{x} \cdot \mathbf{y}_{1\dots a}^*)$, additional samples are sampled for \mathbf{y}_b^* , with the process repeating on the new best solution, $\mathbf{y}_b'^*$.
- Conversely, if $\hat{h}(\mathbf{y}_a^* | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)}) < \hat{h}(\mathbf{y}_b^* | \mathbf{x} \cdot \mathbf{y}_{1\dots a}^*)$, the first segment \mathbf{y}_a^* is further partitioned. The first step corresponds to the α fraction of $\mathbf{y}_a'^*$, with the remaining part of the full solution forming the second step.

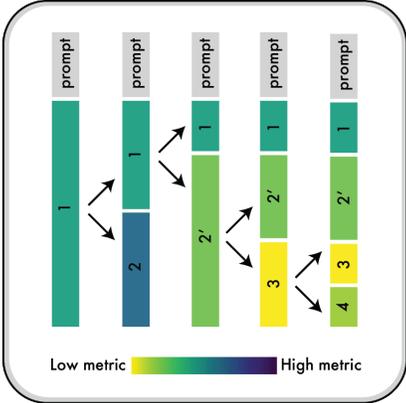


Figure 2: DISC recursively partitions the most challenging step – i.e., the one with the lowest priority metric h – to progressively refine and find the optimal solution.

This recursive process is illustrated in Fig. 2. The pseudocode for DISC is provided in Alg. 1, with an annotated Python implementation in App. A.

The **priority metric** h serves as the central heuristic for determining which solution steps to prioritize. It estimates the “difficulty” or “potential for improvement” of a step y_k given the context $x \cdot y_{1\dots(k-1)}$, computed via rollouts of the policy π . Specifically, $h(y_k | x \cdot y_{1\dots(k-1)})$ is estimated by sampling continuations and evaluating their outcomes.

In practice, estimating h and generating new samples occur simultaneously, as both rely on rollout-based computations (Sec. 3.2). Unlike standard decomposition methods, DISC does not process steps in strict temporal order, resembling goal-directed planning (Parascandolo et al., 2020) and *backtracking*.

Intuition and Benefits. DISC partitions difficult steps into smaller, simpler sub-steps, allocating additional resources to refine them. This approach is particularly effective for inference scaling, where computing must be used judiciously. A binary decomposition strategy enables fast identification of difficult or high-potential steps.

Key Insight: Recursive partitioning

Top-down, recursive partitioning means that we can both efficiently locate critical steps and also dynamically determine step sizes based on our budget.

Dynamic Compute Allocation. A key advantage of DISC is its ability to prioritize challenging or high-potential steps, improving solution quality while minimizing compute waste. By iteratively refining steps with low priority scores, DISC adaptively allocates more resources to difficult steps and less to simpler ones, optimizing inference efficiency.

Key Insight: Adaptive compute allocation

DISC dynamically allocates inference compute to harder steps, optimizing solution quality and resource efficiency.

3.2 PRIORITY METRIC

We consider two intuitive priority metrics for step selection: **Q-value priority** and **Z-score priority**, which are visualized in Fig. 4.

Q-value based priority (DISC-Q). Given a partial solution $y_{1\dots k}^* = y_{1\dots(k-1)}^* \cdot y_a^* \cdot y_b^*$, we aim to prioritize either y_a^* or y_b^* for refinement. The core intuition behind Q-value prioritization is that *steps with lower Q-values indicate areas needing refinement*, directing compute toward the most challenging parts of the solution. More formally, we define the Q-priority metric under policy π as:

$$h_Q(y_k^* | x \cdot y_{1\dots(k-1)}) = \mathbb{E}_{y_k} [Q^\pi(y_k | x \cdot y_{1\dots(k-1)})].$$

Here, y_k represents alternative steps sampled from π , and $Q^\pi(y_k | x \cdot y_{1\dots(k-1)})$ denotes the Q-value of y_k , conditioned on the partial solution $x \cdot y_{1\dots(k-1)}$. Equivalently, $h_Q(y_k^*)$ can be interpreted as the value function $V^\pi(x \cdot y_{1\dots(k-1)})$, where V^π represents the expected reward achievable from the given partial solution. The expectation is taken over sampled candidates $y_k \sim \pi(\cdot | x \cdot y_{1\dots(k-1)})$, constrained by $|y_k| = |y_k^*|$. This metric helps identify *difficult steps* that the LLM is likely to get wrong, guiding partitioning toward the most critical refinements.

To estimate h_Q for the second step y_b^* , we sample y_b from π , generate rollouts $y_{(k+1)+}$, compute rewards, and average the outcomes. For the previous step y_a^* , we reuse rollouts from earlier partitioning, as the mean of these previously generated samples provides an unbiased estimate of the Q-priority metric:

$$h_Q(y_a^* | x \cdot y_{1\dots(k-1)}) = \mathbb{E} [R(y_{1\dots(k-1)+}^*)].$$

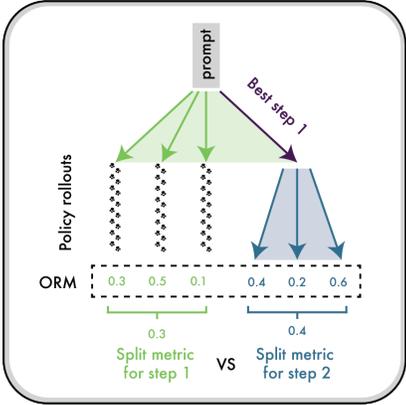


Figure 3: **Priority metric estimation.** We can estimate the priority metric \hat{h} using Monte Carlo rollouts of the LLM policy π for each step.

By leveraging existing rollouts, DISC avoids redundant sampling, improving computational efficiency. Once the rollouts for $\mathbf{y}^{*(k+1)}$ are available, the best completion $\mathbf{y}^{*(k+1)+}$ is selected as the next step to partition. This dual use of rollouts optimizes both metric estimation and inference sampling.

Takeaway: Combining step priority estimation and inference-time search
 By integrating LLM policy sampling for both metric estimation and search, we can significantly enhance computational efficiency.

Z-score based priority (DISC-Z). To allocate more compute to steps with *higher potential for improvement*, we estimate the probability of sampling a better step given existing samples. Assuming the Q-values of sampled steps follow a normal distribution, we model this probability using the cumulative distribution function (CDF). Given a mean μ_r and standard deviation σ_r of sampled Q-values, the probability of sampling better than the best-observed step \mathbf{y}_k^* with Q-value q^* is:

$$1 - \text{CDF}\left(\frac{q^* - \mu_r}{\sigma_r}\right) = 1 - \text{CDF}(z^*),$$

where z^* is the Z-score. Since the CDF is monotonic, we compare steps based on their Z-scores.

We formally define the Z-score priority metric as:

$$h_Z(\mathbf{y}_k^* | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)}) = \frac{q^* - \mathbb{E}[Q^\pi(\mathbf{y}_k | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)})]}{\text{Std}[Q^\pi(\mathbf{y}_k | \mathbf{x} \cdot \mathbf{y}_{1\dots(k-1)})]}.$$

Since $q^* - \mu_r$ represents the advantage of step \mathbf{y}_k^* , the Z-score metric can be interpreted as a *standard deviation-scaled advantage*. Lower h_Z values indicate steps with greater room for improvement, guiding decomposition toward those with higher variance in performance.

3.3 DISC AND SEARCH METHODS

DISC can also be used to enhance **Monte Carlo Tree Search (MCTS)** and other inference scaling and search methods. Recall that in our partition step, we greedily partition the best solution step \mathbf{y}^* . Instead of greedily partitioning the best step, we can partition the top k best steps instead, and select which step to partition using the upper confidence tree (UCT) formula. We can also use **beam search** to prune out steps we do not want to partition further. We explain MCTS and beam search in detail in App. F and present results of combining DISC with search in Sec. 4.9.

Takeaway: Dynamic step sizes can improve search
 DISC can enhance search based inference scaling methods by determining what step size to search across.

3.4 A MOTIVATING EXAMPLE ON DISC-Z

We use the Wiener process $W(t)$ as an example where there are intractably many actions and steps. Suppose we start at $t = 0$ with $W(0) = 0$. At each round k , the algorithm can choose one of the two options:

1. samples a trajectory and observe the final value $W(T)$ at time $t = T$, as the reward signal. Denote the whole trajectory as $w_k(\cdot)$.
2. chooses one trajectory from the previous rounds (denoted as $w_s(t)$ for round s), and time t_0 ; then sample a trajectory at $t = t_0$ with $W(t_0) = w_s(t_0)$. Denote the concatenated trajectory as $w_k(\cdot)$ with $w_k(t) = w_s(t)$ when $t \leq t_0$.

Note that we are only able to observe the final reward $W(t)$. At any intermediate time $t \in (0, T)$, the current value $W(t)$ is not observable. The goal is to design an algorithm that can reach the highest reward among the K trajectories. Formally speaking, we aim to maximize the maximum:

$$\max_{k \in K} w_k(T).$$

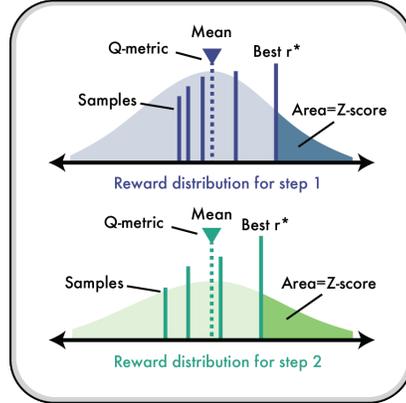


Figure 4: **Different priority metrics.** Sampled rewards of different steps can be visualized as a distribution.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

One naive solution is to call option 1 for K times and return the best-of- K reward, each following:

$$W(T) \sim \mathcal{N}(0, T).$$

Alternatively, suppose there is a promising path $w(\cdot)$ with a high final reward $w(T) = R$. It is natural to consider starting at some midpoint αT ($0 < \alpha < 1$) and perform more completions to obtain an even higher reward than R . The reward distribution sampled this way is

$$W'(T) \sim \mathcal{N}(w(\alpha T), (1 - \alpha)T).$$

The remaining question is which α we should choose. One option is to maximize the probability that the newly sampled reward is higher than R :

$$\mathbb{P}(W'(T) > R) = 1 - \Phi\left(\frac{R - w(\alpha T)}{\sqrt{(1 - \alpha)T}}\right).$$

3.5 EXAMPLE DECOMPOSITION

With sample budget 100, the decomposition of a MATH problem is as follows, where color indicates the value of the priority metric h of each step (yellow low, purple high).

DISC example decomposition

Let the length of the rectangle be l and the width of the rectangle be w . Since the perimeter of the rectangle is 24 inches, we have that $2l + 2w = 24$, so $l + w = 12$. We wish to maximize the area of the rectangle, which is $A = lw$. Let $l = 12 - w$ and plug into the area:

$$A = (12 - w)w \Rightarrow A = 12w - w^2$$

Now, we differentiate A with respect to w :

$$A'(w) = 12 - 2w$$

We wish to maximize A , so we set $A'(w) = 0$, and solve for w :

$$12 - 2w = 0 \Rightarrow w = 6$$

Since $l = 12 - w$, we have that $l = 12 - 6 = 6$. Therefore, the area of the rectangle is $A = lw = 6 \cdot 6 = \boxed{36}$.

Once the LLM generates the first three steps, the rest is easy. Interestingly, ‘which’ is an important decision point which helps decide how the LLM will complete the solution.

Takeaway: Autoregressive models require autoregressive decomposition

While words such as ‘which’, ‘therefore’, etc. may not seem like important steps to humans, they actually represent important steps for autoregressive LLMs which are trained on next token prediction.

4 EXPERIMENTAL RESULTS

4.1 BENCHMARKS

We evaluate DISC on three benchmarks: **APPS**, **MATH**, and **LiveCodeBench**, to assess its impact on inference scaling for both coding and reasoning. **APPS** (Hendrycks et al., 2021a) consists of 5000 competitive programming problems across three difficulty levels, with the competition-level subset being the hardest. We evaluate on a 200-problem subset due to computational constraints. **MATH** (Hendrycks et al., 2021b) comprises 12,500 math problems. Since the ground-truth verifier provides only binary rewards, we use a pretrained ORM (Xiong et al., 2024), trained via the method in (Wang et al., 2024b), with Llama-3.1-8B-Instruct as the base model. We test on a 500-problem subset (**MATH500**), identical to prior work (Wang et al., 2024b; Lightman et al., 2023). **LiveCodeBench** (Jain et al., 2024) is a continuously updated dataset from Leetcode, AtCoder, and CodeForces, ensuring LLMs have not been exposed to test problems. We evaluate on the 108 problems uploaded between 10/01/2024 and 12/01/2024 to prevent contamination.

4.2 DECOMPOSITION COMPARISON

We compare DISC against three prior decomposition methods from Sec. 2.2: **TokenSplit** (token-level decomposition), **LineSplit** (newline-based decomposition), and **BoN** (treating the entire solution as a single step). Across all benchmarks, DISC achieves superior scaling and performance under both fixed token budgets (Fig. 5) and sample budgets (Fig. 10). We evaluate two key metrics: **Pass@k**,

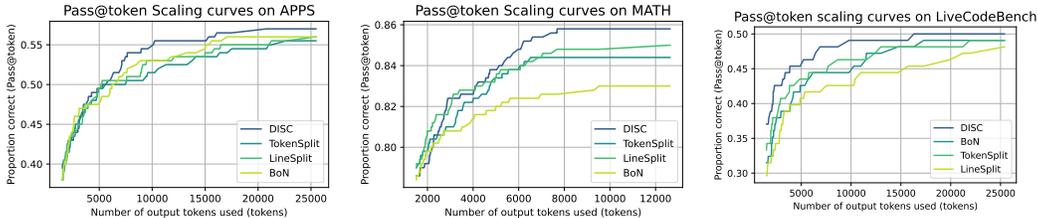


Figure 5: **Token-level comparisons across benchmarks.** (Left) APPS competition level (Middle) MATH500 (Right) LiveCodeBench. DISC achieves superior inference scaling over baselines on all three benchmarks.

the proportion of problems solved within a sample budget k , and **Pass@token**, the proportion solved within a given token budget. Notably, DISC consistently outperforms static decomposition methods on APPS, MATH, and LiveCodeBench (Fig. 5), demonstrating its ability to allocate compute adaptively for improved inference efficiency. Extended results and analyses for each benchmark are provided in App. E.1, E.4, and E.5.

4.3 DECOMPOSITION ANALYSIS AND INTERPRETATION

Our results strongly indicate that decomposition—whether line-based, token-based, or DISC—improves sample quality. Fig. 6 illustrates how the mean and variance of sampled rewards evolve with the **step number**, which represents the order in which a step is explored. Higher step numbers correspond to deeper search levels, where solutions are partitioned into finer-grained steps. As shown in Fig. 6, increasing step number correlates with higher-quality solutions, demonstrating that finer-grained decomposition improves sample quality. Additionally, Fig. 6 shows that reward variance decreases as step count increases, highlighting how decomposition enhances sampling precision.

Furthermore, DISC achieves better performance with fewer partitions under a fixed sampling budget (Fig. 7). We distinguish between **actual partitions**, the number of steps effectively explored, and **planned partitions**, the number of partitions intended by the method. Token and line split methods generate a large number of planned partitions (Fig. 7) but search over at most 15 steps due to budget constraints. In contrast, DISC dynamically adjusts the number of partitions based on available budget, efficiently identifying and focusing on critical steps.

Takeaway: Decomposition and sample quality
 Finer-grained decomposition improves sampled solution quality and reduces reward variance.

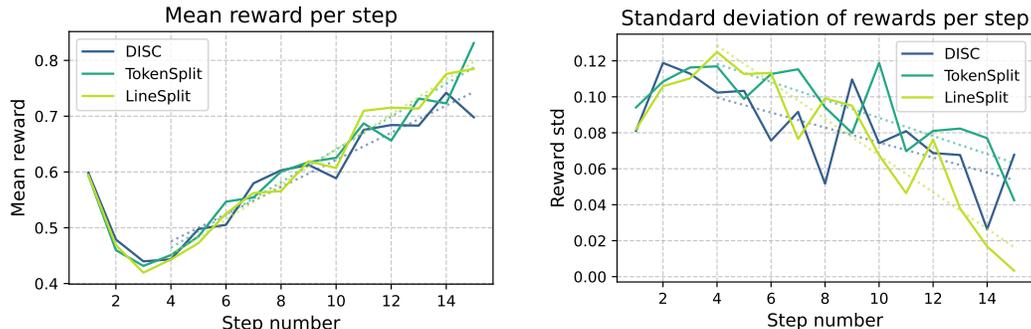


Figure 6: **Analysis of rewards per step on APPS.** (Left) Average reward per step: From step 3 onward, higher step counts strongly correlate with increased average reward, demonstrating the effectiveness of decomposition. The dip between steps 1 and 3 likely occurs because simple problems are solved early, preventing further search. (Right) Standard deviation of rewards per step: Decomposition reduces sampling variance, improving precision at deeper search depths.

4.4 INTERACTION BETWEEN TEMPERATURE AND DISC

We perform ablation studies to analyze the impact of temperature on DISC. Typically, inference scaling methods achieve optimal performance at temperatures around 0.6–0.8, as increased tempera-

378
379
380
381
382
383
384
385
386
387
388
389

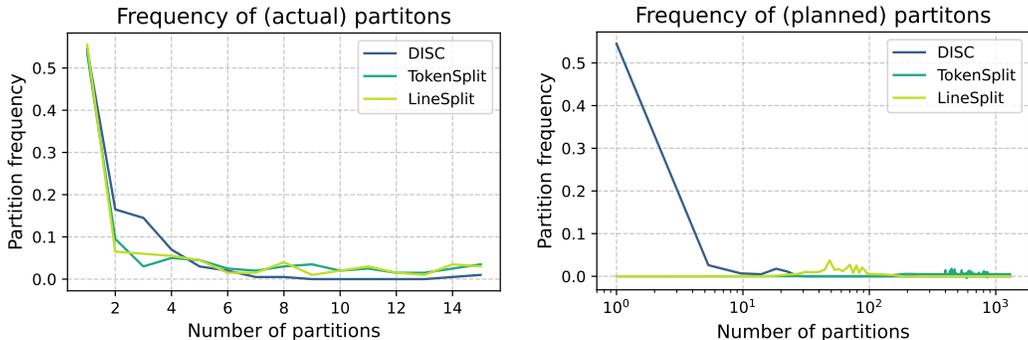


Figure 7: **Comparison of actual and planned partitions on APPS.** DISC outperforms other methods with fewer partitions by efficiently identifying critical steps. Unlike token and line split methods, which plan many partitions but search only a subset, DISC dynamically adjusts partitioning based on budget.

394
395
396
397
398
399
400
401
402
403

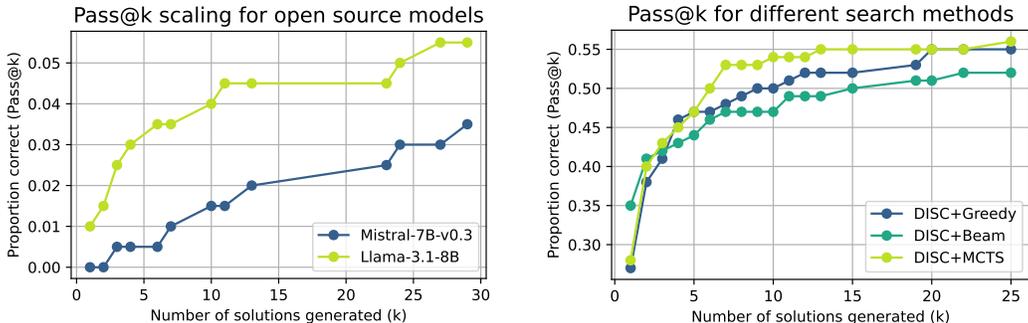


Figure 8: **Pass@k scaling on APPS.** (Left) Open-source models: DISC substantially improves performance across different LLMs, including Llama and Mistral. (Right) Search methods: MCTS scales best, followed by greedy search, then beam search (beam size 2) when combined with DISC on APPS with gpt-4o-mini.

407
408
409
410
411
412
413
414

ture promotes sample diversity (Wang et al., 2024a). Surprisingly, however, DISC performs *better at lower temperatures*, as shown in Fig. 9. This trend is in stark contrast to BoN (Fig. 16), where higher temperatures are generally beneficial. We believe this phenomenon arises because DISC depends on accurately estimating the priority metric h at each step. Lower temperatures reduce sample variance, leading to more reliable estimates of h , which in turn improves step selection. This is further supported by Fig. 14, which shows that lower temperatures yield lower standard deviations per step, indicating increased sampling consistency. Additional details and analyses can be found in App. D.1.

415

4.5 SELF-GENERATED VALIDATION TESTS

417
418
419
420
421
422
423

We also evaluate DISC in a more practical setting where a ground-truth reward model is unavailable for code generation (Chen et al., 2022; 2023b; Zhou et al., 2024). Instead of relying on predefined test cases, we prompt the LLM to generate validation test cases based on the problem prompt. In real-world applications, manually curated ground-truth test cases are often costly to obtain, making self-generated validation a more scalable approach. The results, shown in Fig. 10, indicate that DISC continues to scale better than other methods in this setting. Additional results and details are provided in App. E.3.

424

4.6 INTERACTION BETWEEN PRIORITY METRIC AND DISC

426
427
428
429
430
431

We conduct an ablation study to examine how the choice of priority metric affects DISC performance. In addition to the **Q-based** and **Z-based** priority metrics (DISC-Q and DISC-Z) introduced in Sec. 3.2, we evaluate three baselines: **DISC-R** (random step selection), **DISC-negQ**, and **DISC-negZ** (which prioritize the opposite steps of DISC-Q and DISC-Z, respectively). As shown in Fig. 9, the selection of a priority metric significantly impacts performance. Both DISC-Q and DISC-Z *significantly* outperform random selection and their inverse counterparts, demonstrating the effectiveness of their priority heuristics. Additional details and analysis are in App. D.2.

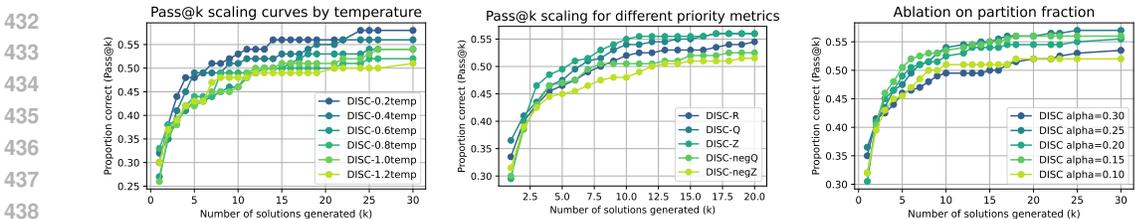


Figure 9: Analysis of factors affecting DISC performance on APPS with gpt-4o-mini. (Left) Effect of temperature: Unlike BoN and other inference scaling methods, DISC achieves higher performance at lower temperatures. (Middle) Effect of priority metrics: Both Q-based and Z-based priority metrics outperform random selection and their inverses, highlighting their effectiveness. (Right) Effect of partition fraction α : The range $0.15 \leq \alpha \leq 0.25$ appears optimal.

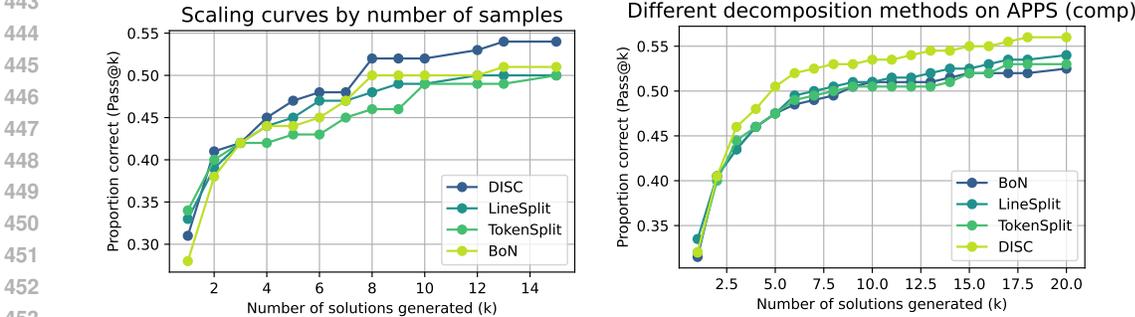


Figure 10: Comparison of Pass@k performance on APPS with ground truth tests (right) and self generated validation tests (left) using gpt-4o-mini. DISC scales more effectively in both.

4.7 ABLATION ON BASE LLM MODEL

We evaluate DISC across different LLMs, including open-source models. As shown in Fig. 8 and Fig. 21, DISC significantly enhances performance even for weaker models. Specifically, it improves Llama’s pass rate from 1% to 5.5%, a 550% relative increase, and Mistral’s from 0% to 3.5%, demonstrating substantial gains even from a nonzero baseline. Additional details and analyses are in App. D.3.

4.8 ABLATION ON PARTITION FRACTION α

We conduct an ablation study to analyze the effect of the partition fraction α on DISC performance. As shown in Fig. 9 and 25, the optimal range appears to be $0.15 \leq \alpha \leq 0.25$. Lower partition fractions ($\alpha < 0.5$) tend to perform better due to the asymmetric cost of sampling from different halves of the partition. Sampling from the first half requires generating more tokens, while the second half requires fewer, making it crucial to partition the first half more conservatively. Additional analysis are in App. D.4.

4.9 SEARCH AND DISC

We demonstrate that search methods such as MCTS and beam search can be combined with DISC. As shown in Fig. 41 in the Appendix, greedy search explores deeper partitions given the same search budget due to its greedy nature, while MCTS and beam search reach similar, shallower depths. However, MCTS allocates the search budget more effectively than beam search, leading to higher performance, as seen in Fig. 8. Additional details and analysis are in App. F.

5 CONCLUSION

We introduce DISC, a dynamic decomposition framework that adaptively partitions solution steps based on difficulty, improving inference scaling by directing compute toward critical steps while balancing exploration and resource allocation. DISC seamlessly integrates with search-based methods such as MCTS and beam search, further enhancing performance. It also identifies challenging steps for LLMs, aiding curriculum learning, fine-tuning, and dataset augmentation. By dynamically adjusting partitioning based on available compute, DISC enables more adaptive and efficient reasoning in large language models, with broad implications for both training and inference optimization.

REFERENCES

- 486
487
488 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,
489 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning
490 from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- 491 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
492 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
493 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 494
495 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and
496 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling.
497 *arXiv preprint arXiv:2407.21787*, 2024.
- 498 Alexandra Carpentier and Michal Valko. Extreme bandits. *Advances in Neural Information Processing*
499 *Systems*, 27, 2014.
- 500
501 Tristan Cazenave. Nested monte-carlo search. In *Twenty-First International Joint Conference on*
502 *Artificial Intelligence*, 2009.
- 503
504 Angelica Chen, David Dohan, and David So. Evoprompting: Language models for code-level neural
505 architecture search. *Advances in neural information processing systems*, 36:7787–7817, 2023a.
- 506
507 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen.
508 Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397*, 2022.
- 509
510 Jingchang Chen, Hongxuan Tang, Zheng Chu, Qianglong Chen, Zekun Wang, Ming Liu, and Bing
511 Qin. Divide-and-conquer meets consensus: Unleashing the power of functions in code generation.
512 *arXiv preprint arXiv:2405.20092*, 2024a.
- 513
514 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
515 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
516 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 517
518 Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi
519 Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu.
520 Do not think that much for $2+3=?$ on the overthinking of o1-like llms, 2024b. URL <https://arxiv.org/abs/2412.21187>.
- 521
522 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to
523 self-debug. *arXiv preprint arXiv:2304.05128*, 2023b.
- 524
525 Vincent A Cicirello and Stephen F Smith. The max k-armed bandit: A new model of exploration
526 applied to search heuristic selection. In *The Proceedings of the Twentieth National Conference on*
527 *Artificial Intelligence*, volume 3, pp. 1355–1361, 2005.
- 528
529 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
530 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
531 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 532
533 Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. Successive prompting for decom-
534 posing complex questions. *arXiv preprint arXiv:2212.04092*, 2022.
- 535
536 Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun
537 Wang. Alphazero-like tree-search can guide large language model decoding and training. *arXiv*
538 *preprint arXiv:2309.17179*, 2023.
- 539
540 Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin Liu, Winson Cheng, Archit Sharma, and
541 Noah D Goodman. Stream of search (sos): Learning to search in language. *arXiv preprint*
542 *arXiv:2404.03683*, 2024.
- 543
544 Xinyu Guan, Li Lina Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang.
545 rstar-math: Small llms can master math reasoning with self-evolved deep thinking, 2025. URL
546 <https://arxiv.org/abs/2501.04519>.

- 540 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
541 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
542 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 543
- 544 Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. Reasoning
545 with language model is planning with world model. In *Empirical Methods in Natural Language*
546 *Processing*, pp. 8154–8173, 2023.
- 547 Erik Hemberg, Stephen Moskal, and Una-May O’Reilly. Evolving code with a large language model.
548 *Genetic Programming and Evolvable Machines*, 25(2):21, 2024.
- 549
- 550 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
551 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
552 *preprint arXiv:2103.03874*, 2021a.
- 553
- 554 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
555 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
556 *preprint arXiv:2103.03874*, 2021b.
- 557 Sergio Hernández-Gutiérrez, Minttu Alakuijala, Alexander V Nikitin, and Pekka Marttinen. Recursive
558 decomposition with dependencies for generic divide-and-conquer reasoning. In *The First Workshop*
559 *on System-2 Reasoning at Scale, NeurIPS’24*, 2024.
- 560
- 561 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
562 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
563 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 564 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish
565 Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. *arXiv*
566 *preprint arXiv:2210.02406*, 2022.
- 567
- 568 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
569 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:
570 22199–22213, 2022.
- 571 Kuang-Huei Lee, Ian Fischer, Yueh-Hua Wu, Dave Marwood, Shumeet Baluja, Dale Schuurmans,
572 and Xinyun Chen. Evolving deeper llm thinking, 2025. URL [https://arxiv.org/abs/
573 2501.09891](https://arxiv.org/abs/2501.09891).
- 574
- 575 Kuang-Huei Leea, Ian Fischera, Yueh-Hua Wuc, Dave Marwood, Shumeet Baluja, Dale Schuurmans,
576 and Xinyun Chen. Evolving deeper llm thinking. *Gen*, 2:3, 2025.
- 577
- 578 Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O Stanley.
579 Evolution through large models. In *Handbook of Evolutionary Machine Learning*, pp. 331–366.
580 Springer, 2023.
- 581 Kyla H Levin, Kyle Gwilt, Emery D Berger, and Stephen N Freund. Effective llm-driven code
582 generation with pythoness. *arXiv preprint arXiv:2501.02138*, 2025.
- 583
- 584 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making
585 language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual*
586 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5315–5333,
587 2023.
- 588 Zhenwen Liang, Ye Liu, Tong Niu, Xiangliang Zhang, Yingbo Zhou, and Semih Yavuz. Improving
589 llm reasoning through scaling inference computation with collaborative verification. *arXiv preprint*
590 *arXiv:2410.05318*, 2024.
- 591
- 592 Jonathan Light, Min Cai, Weiqin Chen, Guanzhi Wang, Xiushi Chen, Wei Cheng, Yisong Yue, and
593 Ziniu Hu. Strategist: Learning strategic skills by llms via bi-level tree search. *arXiv preprint*
arXiv:2408.10635, 2024a.

- 594 Jonathan Light, Yue Wu, Yiyou Sun, Wenchao Yu, Xujiang Zhao, Ziniu Hu, Haifeng Chen, Wei
595 Cheng, et al. Scattered forest search: Smarter code space exploration with llms. *arXiv preprint*
596 *arXiv:2411.05010*, 2024b.
- 597
- 598 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
599 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
600 *arXiv:2305.20050*, 2023.
- 601
- 602 Zicheng Lin, Tian Liang, Jiahao Xu, Qiuzhi Lin, Xing Wang, Ruilin Luo, Chufan Shi, Siheng Li,
603 Yujiu Yang, and Zhaopeng Tu. Critical tokens matter: Token-level contrastive estimation enhances
604 llm’s reasoning capability, 2025. URL <https://arxiv.org/abs/2411.19943>.
- 605
- 606 Vadim Liventsev, Anastasiia Grishina, Aki Härmä, and Leon Moonen. Fully autonomous program-
607 ming with large language models. In *Proceedings of the Genetic and Evolutionary Computation*
608 *Conference*, pp. 1146–1155, 2023.
- 609
- 610 Rohin Manvi, Anikait Singh, and Stefano Ermon. Adaptive inference-time compute: Llms can predict
611 if they can do better, even mid-generation. *arXiv preprint arXiv:2410.02725*, 2024.
- 612
- 613 Harsha Nori, Naoto Usuyama, Nicholas King, Scott Mayer McKinney, Xavier Fernandes, Sheng
614 Zhang, and Eric Horvitz. From medprompt to o1: Exploration of run-time strategies for medical
615 challenge problems and beyond, 2024. URL <https://arxiv.org/abs/2411.03590>.
- 616
- 617 Giambattista Parascandolo, Lars Buesing, Josh Merel, Leonard Hasenclever, John Aslanides, Jes-
618 sica B Hamrick, Nicolas Heess, Alexander Neitz, and Theophane Weber. Divide-and-conquer
619 monte carlo tree search for goal-directed planning. *arXiv preprint arXiv:2004.11410*, 2020.
- 620
- 621 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
622 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
623 *in Neural Information Processing Systems*, 36, 2024.
- 624
- 625 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
626 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
627 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
628 *Nature*, 625(7995):468–475, 2024.
- 629
- 630 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
631 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 632
- 633 Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann
634 Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To cot or not to cot? chain-of-
635 thought helps mainly on math and symbolic reasoning, 2024. URL <https://arxiv.org/abs/2409.12183>.
- 636
- 637 Evan Wang, Federico Cassano, Catherine Wu, Yunfeng Bai, Will Song, Vaskar Nath, Ziwen Han,
638 Sean Hendryx, Summer Yue, and Hugh Zhang. Planning in natural language improves llm search
639 for code generation. *arXiv preprint arXiv:2409.03733*, 2024a.
- 640
- 641 Evan Wang, Federico Cassano, Catherine Wu, Yunfeng Bai, Will Song, Vaskar Nath, Ziwen Han,
642 Sean Hendryx, Summer Yue, and Hugh Zhang. Planning in natural language improves LLM search
643 for code generation. In *The Thirteenth International Conference on Learning Representations*,
644 2025. URL <https://openreview.net/forum?id=48WAZhwHHw>.
- 645
- 646 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang
647 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In
Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume
1: Long Papers), pp. 9426–9439, 2024b.
- Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting, 2024. URL
<https://arxiv.org/abs/2402.10200>.

- 648 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha
649 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
650 models. In *The Eleventh International Conference on Learning Representations*, 2023a. URL
651 <https://openreview.net/forum?id=1PL1NIMMrw>.
- 652 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
653 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
654 models. In *International Conference on Learning Representations*, 2023b.
- 655 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
656 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 35:
657 24824–24837, 2022.
- 658 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws: An
659 empirical analysis of compute-optimal inference for problem-solving with language models, 2024.
660 URL <https://arxiv.org/abs/2408.00724>.
- 661 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu Zhao, Min-Yen Kan, Junxian He, and Michael
662 Xie. Self-evaluation guided beam search for reasoning. *Advances in Neural Information Processing
663 Systems*, 36, 2024.
- 664 Wei Xiong, Hanning Zhang, Nan Jiang, and Tong Zhang. An implementation of generative prm.
665 <https://github.com/RLHFlow/RLHF-Reward-Modeling>, 2024.
- 666 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
667 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural
668 Information Processing Systems*, 36, 2024.
- 669 Eric Zelikman, Qian Huang, Gabriel Poesia, Noah Goodman, and Nick Haber. Parsel: Algorithmic
670 reasoning with language models by composing decompositions. *Advances in Neural Information
671 Processing Systems*, 36:31466–31523, 2023.
- 672 Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin, Bo Wang, Shimin Li, Yunhua Zhou, Qipeng Guo,
673 Xuanjing Huang, and Xipeng Qiu. Scaling of search and learning: A roadmap to reproduce o1 from
674 reinforcement learning perspective, 2024. URL <https://arxiv.org/abs/2412.14135>.
- 675 Janis Zenkner, Lukas Dierkes, Tobias Sesterhenn, and Chrisitan Bartelt. Abstractbeam: Enhancing
676 bottom-up program synthesis using library learning. *arXiv preprint arXiv:2405.17514*, 2024.
- 677 Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal.
678 Generative verifiers: Reward modeling as next-token prediction, 2024. URL <https://arxiv.org/abs/2408.15240>.
- 679 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
680 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
681 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- 682 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language
683 agent tree search unifies reasoning acting and planning in language models. *ICML*, 2024.
- 684 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,
685 Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning
686 in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- 687 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,
688 Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex
689 reasoning in large language models. In *International Conference on Learning Representations*,
690 2023.
- 691
692
693
694
695
696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A CODE IMPLEMENTATION OF DISC

Python implementation of DISC

```

756
757
758
759
760
761 def dynamic_decomposition(problem, model, reward_model, split_str, complete_solution, fraction,
762     solution_budget, split_metric, stop_threshold=-float("inf"), stop_sum_score=1.0, stop_if_solved=
763     False, ):
764     """
765     Decomposes the solution using a dynamic binary search approach
766
767     Args:
768     problem (Problem): The problem to solve
769     model (Model): The model to use for generation
770     reward_model (function): The reward model to use for scoring
771     split_str (function): The function to use for splitting a string
772     complete_solution (function): The function to use for completing a solution
773     fraction (float): The fraction to split the string
774     solution_budget (int): The maximum number of solutions to generate
775     split_metric (function): The metric to use for splitting
776     stop_threshold (float): The threshold to stop splitting
777     stop_sum_score (float): The sum score to stop generating completions
778     stop_if_solved (bool): Whether to stop if the problem is solved
779     """
780
781     # Initialize results and decomposition steps
782     decomp_return = {
783         "generated_solutions": [],
784         "decomposition": []
785     }
786
787     while len(decomp_return["generated_solutions"]) < solution_budget:
788         # Combine all previous steps into an intermediate solution
789         intermediate_solution = "".join([step["step_str"] for step in decomp_return["decomposition"]])
790         new_scores = []
791         best_solution = None
792         best_completion = None
793         best_score = -float("inf")
794         sum_score = 0.0
795
796         # 1) Generate completions until we generate enough samples to estimate the split metric
797         while sum_score < stop_sum_score:
798             proposed_completion = complete_solution(problem, intermediate_solution, model)
799             proposed_solution = intermediate_solution + proposed_completion
800             decomp_return["generated_solutions"].append(proposed_solution)
801
802             # Update scores
803             proposed_score = reward_model(proposed_solution)
804             new_scores.append(proposed_score)
805             sum_score += proposed_score
806
807             # Track the best solution
808             if proposed_score > best_score:
809                 best_solution = proposed_solution
810                 best_score = proposed_score
811                 best_completion = proposed_completion
812
813             # Stop early if problem is solved
814             if stop_if_solved and proposed_score >= 1.0:
815                 decomp_return["decomposition"].append({"step_str": proposed_completion})
816                 return decomp_return
817
818         new_metric = split_metric(new_scores)
819         last_metric = decomp_return["decomposition"][-1]["metric"] if decomp_return["decomposition"]
820             else None
821
822         # Determine the split target. We always split the step with the highest metric
823         is_split_new_step = last_metric is None or new_metric >= last_metric
824         split_target = decomp_return["decomposition"][-1]["step_str"] if not is_split_new_step else
825             best_completion
826
827         # 3) Attempt to split the target
828         split_result = split_str(split_target, fraction)
829         if not split_result: # If we can't split the target, we're done
830             decomp_return["decomposition"].append({"step_str": best_completion, "metric": new_metric})
831             return decomp_return
832
833         # Update decomposition based on split
834         part1, part2 = split_result
835         if is_split_new_step:
836             decomp_return["decomposition"].append({"step_str": part1, "metric": new_metric})
837             # Stopping condition based on threshold
838             if new_metric < stop_threshold:
839                 decomp_return["decomposition"].append({"step_str": part2})
840                 return decomp_return
841         else:
842             decomp_return["decomposition"][-1] = {"step_str": part1, "metric": last_metric}
843
844     return decomp_return

```

B PSEUDOCODE FOR DISC

Algorithm 1 Dynamic Decomposition

Input: Problem instance x , reward model r , partition function f , LLM policy model π , partition fraction α , solution budget B , priority metric h , metric stopping precision θ , sampling stopping threshold σ , inference mode $\mathbf{b}_{\text{inference}}$ **Output:** Final decomposition D Initialize $D \leftarrow \{\text{generated_solutions} : \emptyset, \text{decomposition} : \emptyset\}$ # Decompose the solution recursively until we reach the desired precision θ or run out of budget B $|D.\text{generated_solutions}| < B$
 $y_{\text{intermediate}} \leftarrow \text{Concatenate}([\text{step_step_str} \forall \text{step} \in D.\text{decomposition}])$ \triangleright Concatenate previous steps to form intermediate solution $R_{\text{new}} \leftarrow \emptyset$ \triangleright
Record rewards of completions $\text{best}.y_{\text{final}} \leftarrow \text{None}$, $\text{best}.y_{\text{completion}} \leftarrow \text{None}$, $\text{best}.r \leftarrow -\infty$ \triangleright Track the best completion # Step 1:
Generate completions until we have enough samples to estimate the splitting metric. Here we use a geometric sampling distribution $\text{sum}(R_{\text{new}}) < \sigma$ $y_{\text{completion}} \leftarrow \pi(\cdot | x, y_{\text{intermediate}})$ $y_{\text{proposed}} \leftarrow y_{\text{intermediate}} \oplus y_{\text{completion}}$ Append y_{proposed} to $D.\text{generated_solutions}$ $r_{\text{proposed}} \leftarrow r(y_{\text{proposed}})$ Append r_{proposed} to R_{new} $r_{\text{proposed}} > \text{best}.r$ $\text{best}.y_{\text{final}} \leftarrow y_{\text{proposed}}$, $\text{best}.y_{\text{completion}} \leftarrow y_{\text{completion}}$ $\text{best}.r \leftarrow r_{\text{proposed}}$ $\mathbf{b}_{\text{inference}}$ and $r_{\text{proposed}} = 1.0$ Append $\{\text{step_str} : y_{\text{completion}}\}$ to $D.\text{decomposition}$ **Return** D \triangleright Exit if problem is solved # Step 2: Compute splitting metric $z_{\text{new}} \leftarrow h(R_{\text{new}})$ $z_{\text{last}} \leftarrow D.\text{decomposition}[-1].z$ **if** $D.\text{decomposition} \neq \emptyset$ **else** $-\infty$ # Step 3: Split the step with the higher metric $\mathbf{b}_{\text{split new step}} \leftarrow z_{\text{new}} \geq z_{\text{last}}$ $y_{\text{target step}} \leftarrow \text{best}.y_{\text{completion}}$ **if** $\mathbf{b}_{\text{split new step}}$ **else** $D.\text{decomposition}[-1].\text{step_str}$ $y_1, y_2 \leftarrow f(y_{\text{target step}}, \alpha)$ $y_1 = \text{None}$ or $y_2 = \text{None}$ Append $\{\text{step_str} : y_{\text{completion}}, \text{metric} : z_{\text{new}}\}$ to $D.\text{decomposition}$ **Return** D \triangleright Exit if we cannot do a finer split $\mathbf{b}_{\text{split new step}}$ Append $\{\text{step_str} : y_1, \text{metric} : z_{\text{new}}\}$ to $D.\text{decomposition}$ \triangleright Add new step $z_{\text{new}} < \theta$ Append $\{\text{step_str} : y_2\}$ to $D.\text{decomposition}$ **Return** D \triangleright Exit if all metrics are smaller than precision $D.\text{decomposition}[-1] \leftarrow \{\text{step_str} : y_1, \text{metric} : z_{\text{last}}\}$ \triangleright Split last step **Return** D

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

C RELATED WORK

Inference scaling. Inference scaling has emerged as a dominant paradigm, driven by the introduction of o1- and r1-like chain-of-thought reasoning models (Snell et al., 2024; Brown et al., 2024; Manvi et al., 2024; Leea et al., 2025). Several works examine the trade-off between inference compute and training compute (Guan et al., 2025; Chen et al., 2024b). LLM inference often relies on decomposing complex problems into intermediate reasoning steps, as seen in chain-of-thought (CoT) prompting (Wei et al., 2022; Sprague et al., 2024; Wang & Zhou, 2024) and its variants (Kojima et al., 2022; Zhou et al., 2023; Wang et al., 2023b; Li et al., 2023). We extend inference scaling by introducing a new approach for adaptive compute allocation (Manvi et al., 2024).

LLM reasoning and code generation. LLM reasoning and code generation are central tasks for inference scaling. Evolutionary inference scaling methods have been explored in program generation (Liventsev et al., 2023; Chen et al., 2023a; Romera-Paredes et al., 2024; Lehman et al., 2023; Hemberg et al., 2024). Domain-specific decomposition strategies have been applied in code generation, such as function-based decomposition (Chen et al., 2024a; Zenkner et al., 2024; Levin et al., 2025). More broadly, decomposition often involves prompting LLMs to generate subtask completions (Hernández-Gutiérrez et al., 2024; Khot et al., 2022; Dua et al., 2022), which differs from methods that refine a single LLM generation.

Reinforcement learning and Monte Carlo methods. Unlike standard RL, our setting resembles a search problem where the goal is to identify the single highest-reward path. Cazenave (2009) demonstrated that nested Monte Carlo search can accelerate optimal pathfinding. Under the bandit setting, this can be formulated as identifying the arm with the highest *maximum* reward rather than the highest mean reward (Cicirello & Smith, 2005; Carpentier & Valko, 2014).

D ABLATION STUDIES

D.1 ABLATION ON TEMPERATURE

We conduct an ablation study to analyze the effects of temperature on DISC and BoN. Temperature controls the randomness of token sampling in autoregressive models, influencing both exploration and consistency. Higher temperatures encourage more diverse outputs, whereas lower temperatures yield more deterministic generations. To examine its impact, we evaluate DISC and BoN on a 100-problem subset of APPS (the first 100 problems) using gpt-4o-mini.

Fig. 11 presents the Pass@token scaling curve for DISC across different temperatures. The results indicate that lower temperatures lead to improved performance, as DISC benefits from more deterministic step selection. Unlike BoN, which relies on broad solution sampling, DISC dynamically refines steps, making stable token probabilities advantageous.

Fig. 12 illustrates the frequency of actual partitions made by DISC at different temperatures. As temperature increases, the number of partitions fluctuates more, suggesting that high temperature introduces instability in step selection. Lower temperatures provide more structured decomposition, reducing unnecessary subdivisions.

In Fig. 13, we visualize the mean reward per step. The trend shows a linear increase in reward as step number grows, demonstrating that deeper decomposition results in progressively better solutions. This reinforces that DISC effectively allocates computation towards refining difficult steps.

The mean standard deviation per step is shown in Fig. 14. Lower temperatures yield lower standard deviations, confirming that DISC benefits from reduced variability in sample quality. This consistency allows for more reliable prioritization of difficult steps, enhancing overall inference efficiency.

For comparison, Fig. 16 and Fig. 15 display Pass@token and Pass@k scaling curves for BoN across different temperatures. Unlike DISC, BoN achieves peak performance at a temperature around 0.6-0.8, balancing diversity and consistency. Higher temperatures increase exploration but degrade precision, while lower temperatures hinder sample diversity, reducing the probability of obtaining high-quality completions.

These findings highlight the fundamental difference between DISC and BoN: DISC benefits from lower variance and stable decomposition, while BoN relies on broader exploration facilitated by moderate temperature settings. As a result, optimal temperature settings differ significantly between these methods, with DISC favoring deterministic sampling and BoN requiring a balance between diversity and coherence.

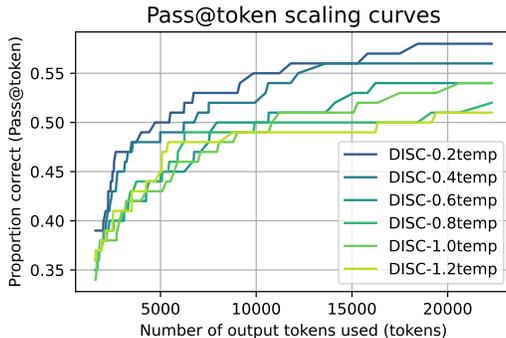


Figure 11: **Pass@token scaling curve for different temperatures on APPS using gpt-4o-mini.** The lower the temperature, the stronger the DISC performance.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

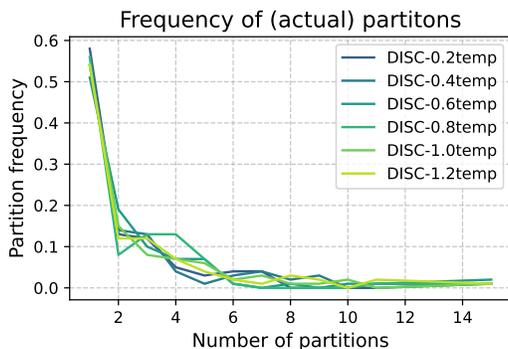


Figure 12: Partition frequency of DISC with different temperatures on APPS using gpt-4o-mini

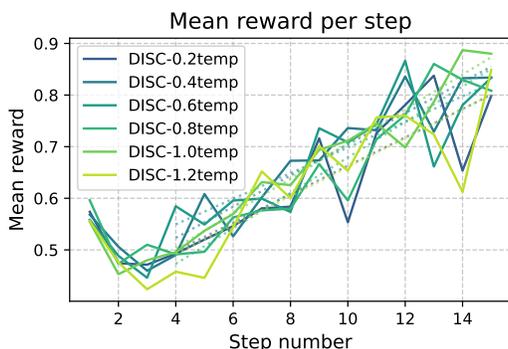


Figure 13: Mean reward per step of DISC with different temperatures on APPS using gpt-4o-mini. The mean reward scales linearly with step number.

D.2 ABLATION ON PRIORITY METRIC h

We analyze the effect of different priority metrics on DISC performance. We evaluate DISC using the first 200 competition-level APPS problems with gpt-4o-mini, setting the temperature to 0.8 for all experiments. The priority metric determines which steps are refined during recursive decomposition, impacting both efficiency and final solution quality.

Fig. 17 presents a token-level comparison of different priority metrics. Both DISC-Q and DISC-Z significantly outperform random selection and their inverse counterparts, demonstrating the importance of prioritizing high-value steps.

Fig. 18 illustrates the partition frequency under different priority metrics. We observe that effective metrics such as DISC-Q and DISC-Z lead to fewer, more meaningful partitions, whereas suboptimal strategies result in excessive, redundant partitioning.

The relationship between mean reward and step number is shown in Fig. 19. All tested metrics exhibit a strong correlation between increasing step depth and mean reward, indicating that decomposition progressively refines solutions. However, DISC-Q and DISC-Z achieve higher reward gains at earlier stages, suggesting that they prioritize the most impactful refinements.

Finally, Fig. 20 reports the standard deviation of rewards per step. Lower standard deviation suggests more stable solution quality, a property that DISC-Q and DISC-Z maintain better than random selection methods. This highlights their effectiveness in identifying and refining challenging steps efficiently.

Overall, these results confirm that choosing an appropriate priority metric is crucial for DISC. While DISC-Q and DISC-Z consistently enhance inference efficiency and quality, random or inverse strategies lead to poorer performance due to misallocation of compute resources.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

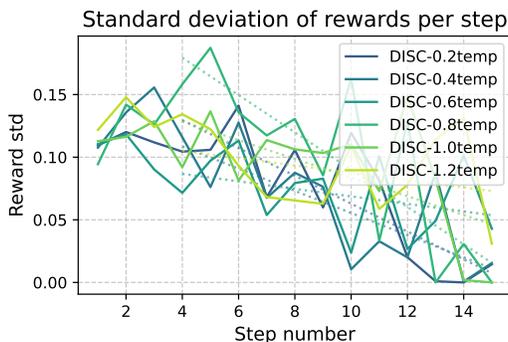


Figure 14: Mean standard deviation per step of DISC with different temperatures on APPS using gpt-4o-mini. Lower temperature means lower average standard deviation.

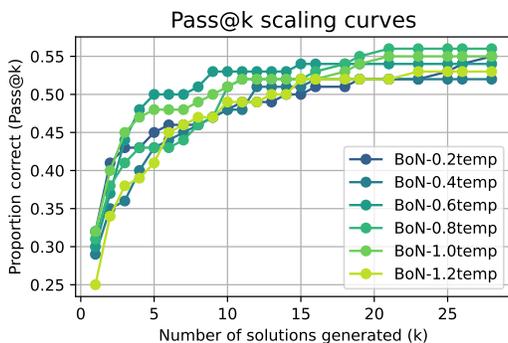


Figure 15: Pass@k scaling curve for different temperatures on APPS using gpt-4o-mini for BoN. A temperature around 0.6-0.8 leads to the best performance and balance between diversity and consistency.

D.3 MODEL ABLATION

We investigate how different LLMs perform when used with DISC on 200 competition-level APPS problems, given a sample budget of 30. The groundtruth reward model was used to evaluate correctness, and all models were set to a temperature of 0.8. Due to the challenging nature of the benchmark, open-source models struggled to achieve strong performance independently. However, when paired with DISC, their performance significantly improved.

Figure 21 presents the Pass@token scaling curve for open-source models using DISC. The results demonstrate that DISC substantially enhances the capabilities of these models, closing the gap between them and proprietary alternatives.

Figure 22 visualizes the partition frequency of DISC with different open-source models. Compared to their standalone performance, the use of DISC led to more structured and effective decomposition, highlighting its adaptability to different architectures.

The mean reward per step is shown in Figure 23. Similar to prior findings, we observe that deeper decomposition leads to increasingly higher rewards. Notably, even lower-capacity models benefit from DISC’s ability to iteratively refine their solutions.

Finally, Figure 24 presents the mean standard deviation per step. With DISC, the variance in performance is significantly reduced, resulting in more stable and reliable inference.

Overall, these findings emphasize that DISC is a robust framework capable of enhancing inference performance across diverse LLMs, particularly those with limited standalone capabilities.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

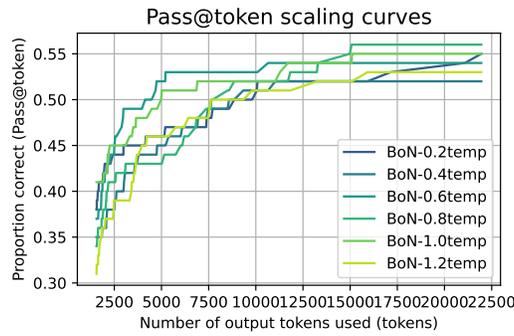


Figure 16: **Pass@token scaling curve for different temperatures on APPS using gpt-4o-mini for BoN.** A temperature around 0.6-0.8 leads to the best performance and balance between diversity and consistency.

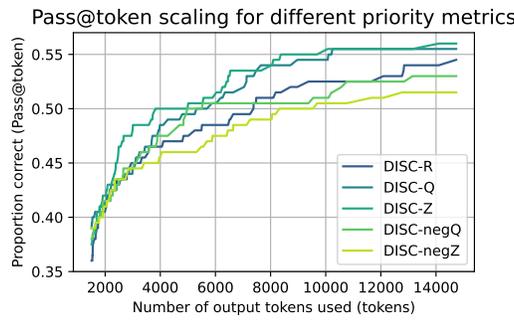


Figure 17: **Token level comparison of different priority metrics on DISC in the APPS setting with gpt-4o-mini.** Both Q and Z based priority metrics perform well.

D.4 ABLATION ON PARTITION FRACTION α

We include some more analysis on the partition fraction here.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

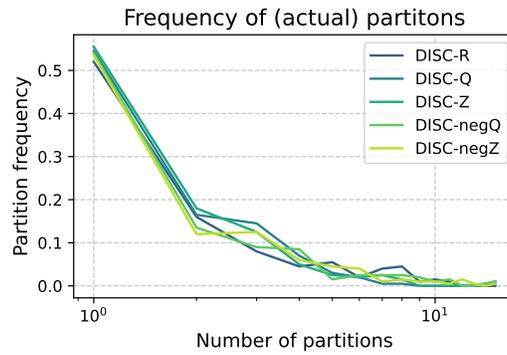


Figure 18: **Partition frequency of DISC with different priority metrics on APPS using gpt-4o-mini**

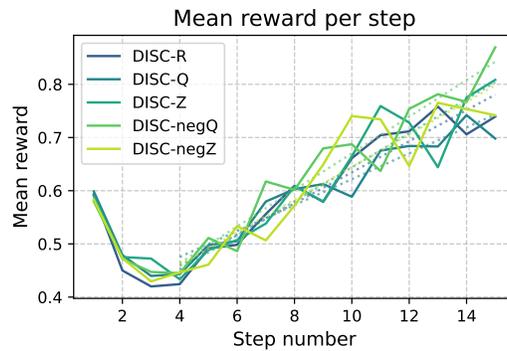


Figure 19: **Mean reward per step of DISC with different priority metrics on APPS using gpt-4o-mini.** All metrics display strong correlation between step depth and the mean reward.

E EXPERIMENTAL RESULTS EXTENDED

E.1 APPS

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

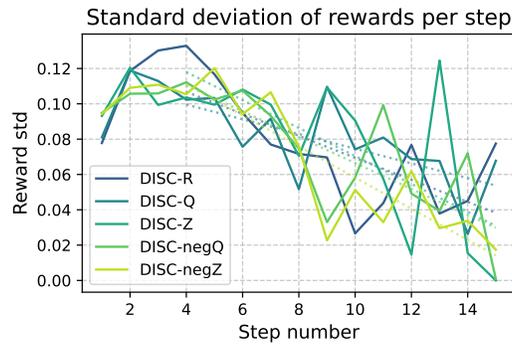


Figure 20: **Mean standard deviation per step of DISC with different priority metrics on APPS using gpt-4o-mini.** All metrics display correlation between step depth and the standard deviation.

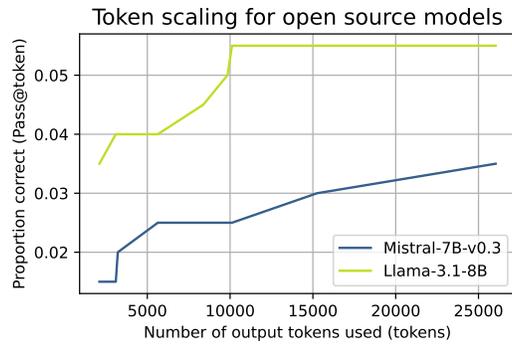


Figure 21: **Pass@token scaling curve for open source models with DISC on APPS.** DISC also demonstrates strong performance gains with open source models.

E.2 ADDITIONAL EXAMPLES

Below is another computed decomposition for the same problem as the one shown in the main text:

Let x be the length and y be the width of the rectangle. Since the perimeter is 24 inches, we have $2x + 2y = 24 \Rightarrow x + y = 12$. Therefore, we must maximize the area xy subject to this constraint.

We use the method of Lagrange multipliers. The Lagrangian is

$$\mathcal{L}(x, y, \lambda) = xy - \lambda(x + y - 12).$$

We differentiate with respect to x , y , and λ to obtain

$$\frac{\partial \mathcal{L}}{\partial x} = y - \lambda = 0, \quad \frac{\partial \mathcal{L}}{\partial y} = x - \lambda = 0, \quad \frac{\partial \mathcal{L}}{\partial \lambda} = x + y - 12 = 0.$$

We find that $x = y = \lambda$, so $x = y = \frac{12}{2} = 6$. Therefore, the maximum area of the rectangle is $6 \cdot 6 = \boxed{36}$.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

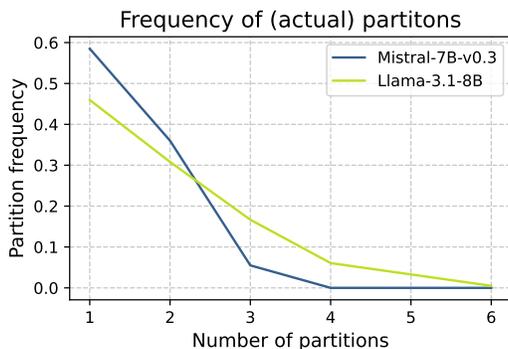


Figure 22: **Partition frequency of DISC with open source models on APPS**

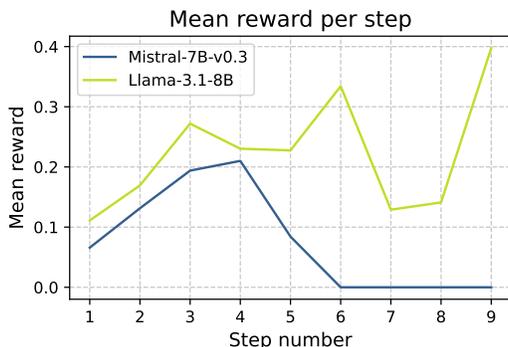


Figure 23: **Mean reward per step of DISC with open source models on APPS**

E.3 APPS WITH SELF-GENERATED VALIDATION TESTS

We examine DISC performance on APPS when using self-generated validation tests. All methods utilized the same set of self-generated validation tests to ensure fair comparisons. Each problem received 5-10 validation tests, with the exact number determined dynamically by the LLM. We evaluated a subset of 100 APPS problems, generating samples until the sample budget was exhausted or a correct solution was found.

Figure 27 illustrates the Pass@token scaling curve, showing that DISC maintains strong scaling performance in this setting, though at a slightly lower rate compared to ground-truth verification.

Figure 28 and Figure 29 compare actual and planned partition frequencies, respectively. The results indicate that DISC continues to make structured decompositions even with self-generated validation, preserving its efficiency.

The mean reward per step, shown in Figure 30, follows a similar trend as in previous experiments, reinforcing that DISC effectively allocates compute resources for iterative refinement.

Lastly, Figure 31 demonstrates that DISC maintains lower standard deviations in performance, indicating stable quality improvements across steps.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

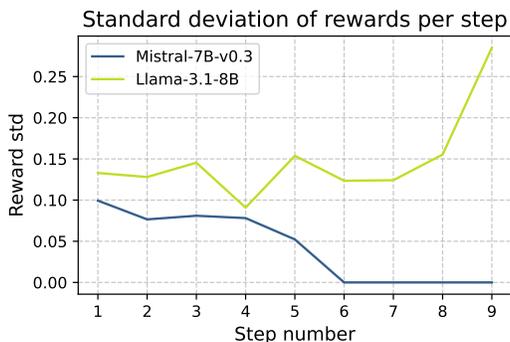


Figure 24: Mean standard deviation per step of DISC with open source models on APPS

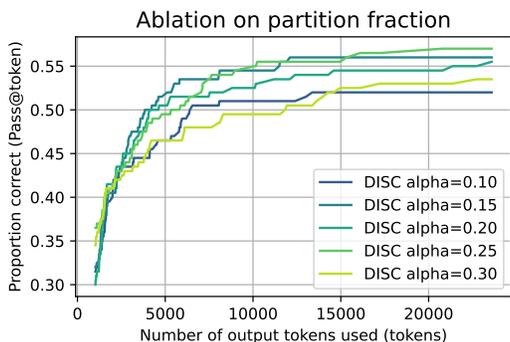


Figure 25: Token level comparison of different DISC splitting fraction α on APPS competition level. $0.15 \leq \alpha \leq 0.25$ seems to be optimal.

E.4 MATH500

Completions for MATH500 include both the reasoning steps and the final answer. Since MATH500 contains more problems than APPS200 and MATH problems tend to be relatively easier, solution quality saturates quickly. Therefore, we use a lower sample budget of 10 for these experiments.

Figure 32 presents the Pass@k performance for different decomposition methods on MATH500. We observe that all decomposition-based approaches achieve similar Pass@k performance, consistently outperforming BoN. This indicates that the structured nature of MATH problems allows multiple decomposition strategies to be effective.

Despite similar Pass@k results, the true advantage of DISC lies in its token efficiency, as shown in Figure 5. DISC significantly reduces the number of tokens required to reach correct solutions compared to alternative methods, demonstrating its ability to allocate computational effort efficiently in mathematical reasoning tasks.

Additionally, we analyze the partitioning behavior of DISC on MATH500. Figure 33 illustrates the actual partition frequency for different decomposition methods. The planned partitioning behavior, shown in Figure 34, further highlights how DISC effectively balances exploration and refinement.

Finally, we present the mean standard deviation per step in Figure 35. Lower variance suggests that DISC produces more stable and reliable decompositions over multiple runs, reinforcing its robustness in both mathematical and program synthesis domains.

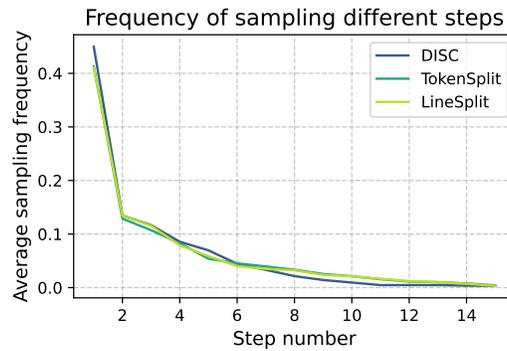


Figure 26: **Sampling frequency of each step averaged over the problems on APPS with gpt-4o-mini.** DISC seems to have preferences for spending more compute on earlier found steps.

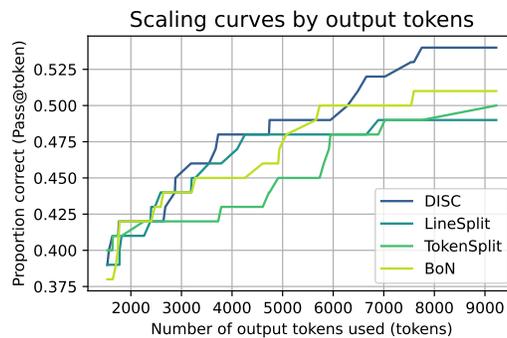


Figure 27: **Token level comparison of different decomposition methods on APPS with gpt-4o-mini and self-generated validation tests.** DISC still scales better than other methods in this setting, albeit at a lower rate.

E.5 LIVECODEBENCH

We evaluate DISC on LiveCodeBench, a benchmark designed for code generation tasks with a focus on real-world software development challenges. LiveCodeBench presents a unique set of problems requiring both reasoning and structured decomposition, making it a suitable testbed for evaluating DISC’s ability to refine and improve intermediate steps.

Figure 36 shows the Pass@k comparison of different decomposition methods on LiveCodeBench. DISC consistently scales better than other decomposition methods, highlighting its ability to refine intermediate steps more effectively in complex coding scenarios.

Figure 37 illustrates the observed partition frequency of different decomposition methods. The structured approach of DISC results in well-balanced decomposition across steps, reducing unnecessary partitioning while maintaining sufficient granularity for improved solution refinement.

Figure 38 displays the planned partition frequency across methods. DISC dynamically determines the most effective partitions based on the evolving problem state, leading to more targeted and efficient decompositions.

Finally, Figure 39 presents the mean standard deviation per step across decomposition methods. Lower variance in DISC suggests that it produces more stable and reliable decompositions, reinforcing its robustness for solving LiveCodeBench problems.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415

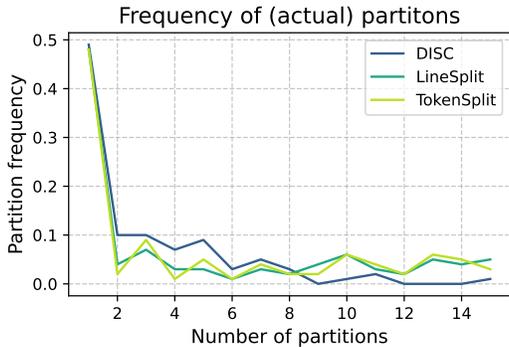


Figure 28: Actual partition frequency of different decomposition methods on APPS with gpt-4o-mini and self-generated validation tests.

1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431

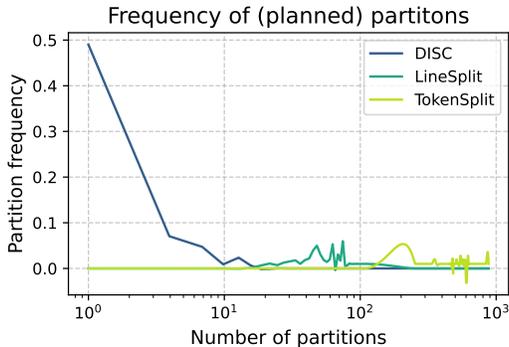


Figure 29: Planned partition frequency of different decomposition methods on APPS with gpt-4o-mini and self-generated validation tests.

F SEARCH AND SCALING

1432
1433
1434

F.1 MONTE CARLO TREE SEARCH (MCTS)

1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445

Monte Carlo Tree Search (MCTS) is a widely used algorithm for sequential decision-making in large search spaces, particularly in applications such as *game playing, planning, and inference scaling*. The algorithm builds a search tree incrementally by simulating different sequences of actions and updating estimates of state quality. A key advantage of MCTS is its ability to balance *exploration* (discovering new states) and *exploitation* (refining promising ones) using a data-driven search process. The MCTS pipeline consists of four fundamental steps: *selection, expansion, simulation, and backpropagation*.

1446

F.1.1 SELECTION

1447
1448
1449
1450
1451

Starting from the root node representing the current state s , MCTS iteratively traverses the search tree by selecting child nodes based on a *selection policy*. The most commonly used selection criterion is the *Upper Confidence Bound for Trees (UCT)*, which balances exploration and exploitation:

1452
1453
1454

$$UCT(s, d) = \widehat{Q}(s, d) + c \sqrt{\frac{\ln(\sum_b n(s, b))}{n(s, d)}}, \tag{1}$$

1455
1456
1457

where $\widehat{Q}(s, d)$ represents the estimated value of selecting action d from state s , $n(s, d)$ is the visit count for this action, and c is a hyperparameter controlling the trade-off between exploring new actions and favoring those with high past rewards.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469

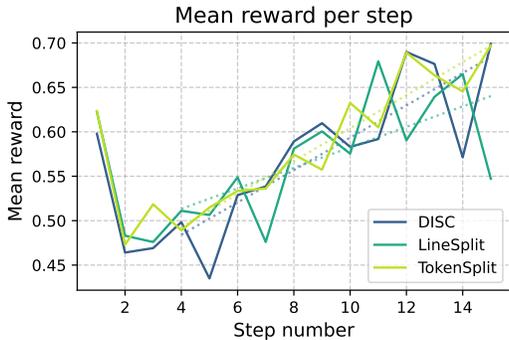


Figure 30: Mean reward per step of different decomposition methods on APPS with gpt-4o-mini and self-generated validation tests.

1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484

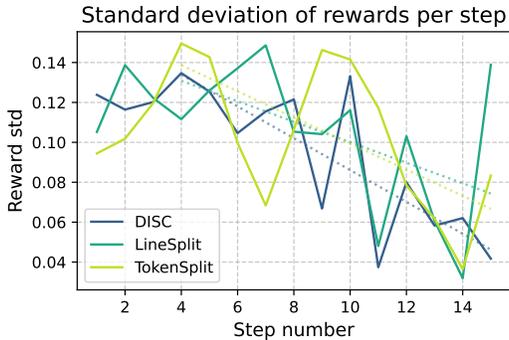


Figure 31: Mean standard deviation different decomposition methods on APPS with gpt-4o-mini and self-generated validation tests.

F.1.2 EXPANSION

Once a leaf node (a previously unexplored state) is reached, the algorithm expands the tree by *adding one or more new nodes*. These new nodes represent potential future states s' generated by sampling an action d from a predefined policy. This step broadens the search space and allows MCTS to evaluate new possibilities.

F.1.3 SIMULATION

Following expansion, the algorithm conducts a *simulation* (or rollout) from the newly added state. This step involves generating a sequence of actions according to a predefined policy until reaching a terminal state or an evaluation horizon. The outcome of the simulation, denoted as $v(s')$, provides an estimate of the quality of the new state. Depending on the application, this could represent a *game result, an optimization score, or an inference accuracy metric*.

F.1.4 BACKPROPAGATION

The final step involves *propagating the results of the simulation back up the search tree* to refine the estimated values of prior states and actions. Each node along the trajectory $\tau = [s_0, d_1, s_2, \dots, s_{-1}]$ is updated iteratively:

$$\hat{Q}(s_i, d_{i+1})^{(t+1)} \leftarrow (1 - \alpha_n)\hat{Q}(s_i, d_{i+1})^{(t)} + \alpha_n \max\{\hat{Q}(s_i, d_{i+1})^{(t)}, \hat{Q}(s_{i+1}, d_{i+2})^{(t+1)}\}, \quad (2)$$

where α_n is a learning rate that depends on the visit count, and the maximum function ensures that the best-performing trajectories are emphasized.

MCTS has been widely adopted in inference scaling techniques due to its ability to *efficiently allocate computational resources*, focusing more on *high-reward states* while avoiding unnecessary

1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

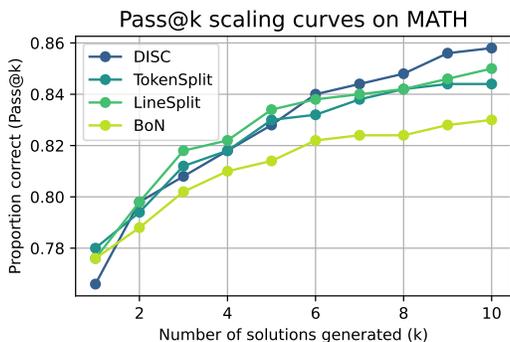


Figure 32: Pass@k performance comparison for different decomposition methods on MATH500. DISC consistently outperforms BoN across different sampling budgets.

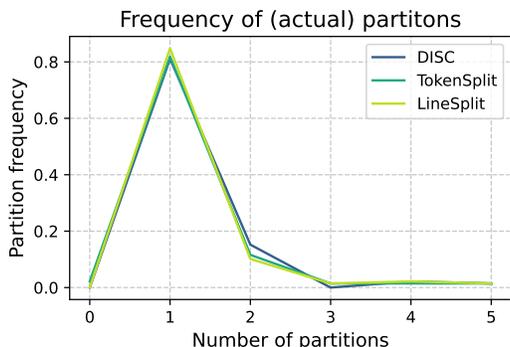


Figure 33: Observed partition frequency of different decomposition methods on MATH500. DISC effectively segments problems into meaningful subcomponents.

exploration of unpromising regions. In later sections, we explore how MCTS can be combined with *dynamic decomposition* to further optimize inference scaling.

F.1.5 COMBINING DYNAMIC DECOMPOSITION WITH MCTS

MCTS can be enhanced by integrating *dynamic decomposition*, where each node in the search tree represents a decomposition of the problem into steps. Instead of treating states as atomic decisions, we recursively decompose reasoning steps, dynamically adjusting granularity based on difficulty.

In this framework:

- Each node in the MCTS tree represents a partial decomposition of the problem, with child nodes corresponding to alternative step partitions.
- Branching occurs by generating candidate next steps using dynamic decomposition, allowing finer steps for complex regions while maintaining efficiency for simpler ones.
- The selection step prioritizes nodes that represent more promising decompositions, dynamically refining challenging areas through recursive subdivision.
- The backpropagation step ensures that decompositions leading to high-quality solutions are reinforced, helping the search tree converge toward optimal inference paths.

By integrating dynamic decomposition with MCTS, we efficiently allocate compute to the most critical reasoning steps, improving inference quality while maintaining computational efficiency.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577

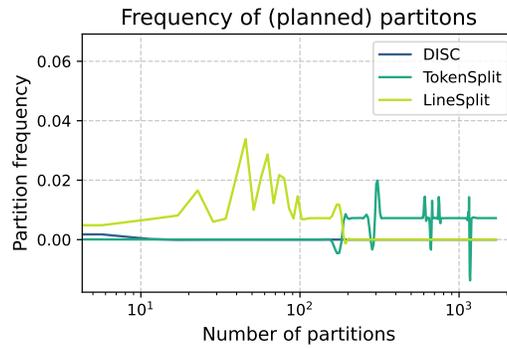


Figure 34: **Planned partitioning strategy of different decomposition methods on MATH500. DISC ’s structured approach leads to more efficient problem breakdowns.**

1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592

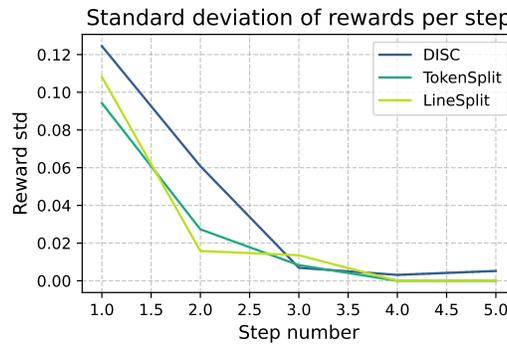


Figure 35: **Mean standard deviation per step for different decomposition methods on MATH500. Lower variance in DISC suggests more stable and reliable problem-solving steps.**

1593
1594
1595
1596
1597

F.2 BEAM SEARCH

1599
1600
1601
1602
1603

Beam search is a heuristic search algorithm commonly used in inference tasks where computational efficiency is a priority. Unlike exhaustive search methods, beam search maintains only the top k best candidates at each step, making it an effective strategy for structured prediction problems and sequential decision-making.

1604
1605
1606

At each iteration:

1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

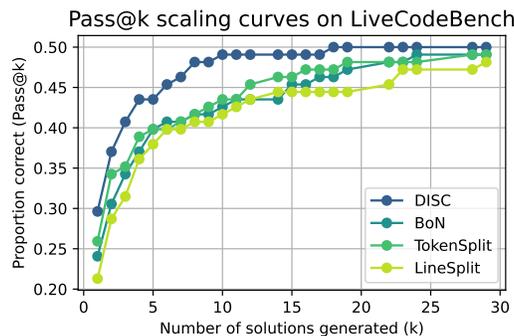
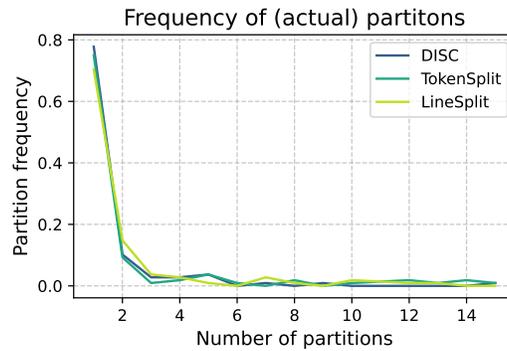


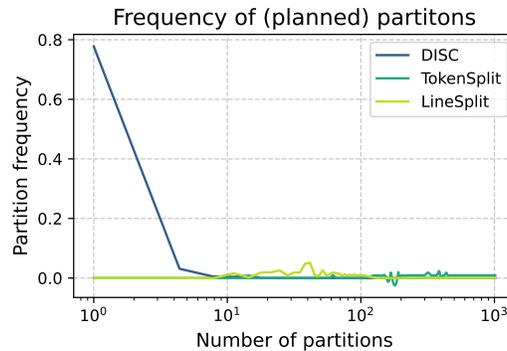
Figure 36: **Pass@k performance comparison for different decomposition methods on Live-CodeBench. DISC consistently outperforms other methods in structured problem refinement.**

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631



1632 **Figure 37: Observed partition frequency of different decomposition methods on LiveCodeBench.**
1633 DISC effectively balances problem segmentation while avoiding excessive partitioning.

1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647



1648 **Figure 38: Planned partitioning strategy of different decomposition methods on LiveCodeBench.**
1649 DISC dynamically adapts its partitioning to optimize search efficiency.

1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660

- The algorithm selects the k most promising partitions from the previous step based on an evaluation metric.
- Each selected partition is expanded by generating possible next-step samples.
- The newly generated partitions are ranked, and only the top k candidates are retained for the next iteration.
- This process continues until a stopping criterion is met, such as reaching a predefined depth or finding a sufficiently high-quality solution.

1661
1662
1663
1664
1665

Beam search provides a computationally efficient way to explore structured solution spaces while maintaining high-quality search trajectories. By integrating beam search with dynamic decomposition, we ensure that inference computation is allocated efficiently, focusing on the most promising reasoning paths at each step.

1666
1667

F.3 ADDITIONAL RESULTS AND ANALYSIS

1668
1669
1670
1671

Experiments comparing different search methods were conducted on a 100-problem subset of the APPS dataset (first 100 problems) using GPT-4o-mini. All methods used a temperature of 0.2, with $\alpha = 0.15$, Q priority metric, and $\sigma = 1.0$.

1672
1673

Token-level comparison: As shown in Figure 40, MCTS scales best among the tested methods, demonstrating superior efficiency in identifying promising partitions. Greedy search follows closely, while beam search exhibits the slowest scaling.

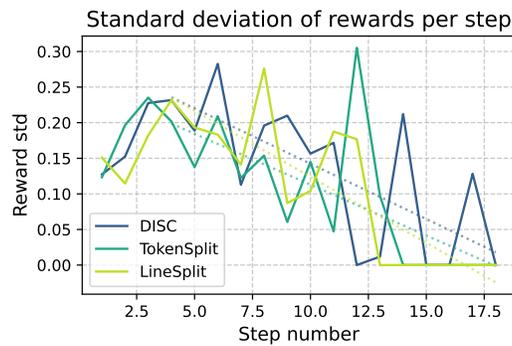


Figure 39: **Mean standard deviation per step for different decomposition methods on LiveCodeBench.** Lower variance in DISC suggests more stable and reliable problem-solving steps.

Partition frequency analysis: Figure 41 reveals that greedy search explores to greater depths within the same sampling budget. This suggests that greedy search prioritizes deep refinements, whereas MCTS and beam search balance depth with breadth.

Step variance analysis: Figure 42 illustrates that all search methods display decreasing standard deviation with increasing search depth. This trend indicates that deeper searches converge towards stable, high-quality partitions, reinforcing the benefits of dynamic decomposition.

These results highlight the trade-offs between search methods: MCTS offers robust exploration-exploitation balance, greedy search favors depth-first refinement, and beam search provides a structured yet computationally constrained approach. The integration of dynamic decomposition further enhances these search strategies by adaptively allocating computational resources to critical reasoning steps.

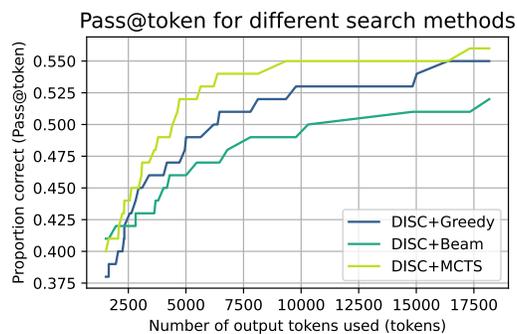


Figure 40: **Token level comparison of different decomposition search methods combined with DISC on APPS with gpt-4o-mini.** MCTS scales best, followed by greedy search, followed by beam search.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

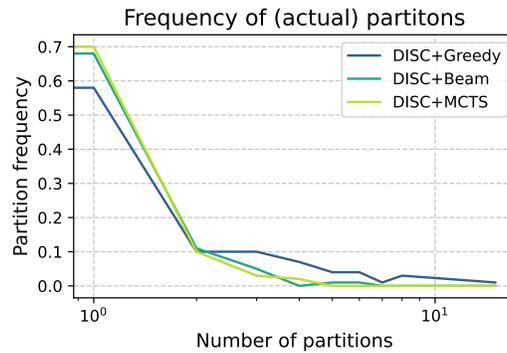


Figure 41: **Actual partition frequency of different decomposition search methods combined with DISC on APPS with gpt-4o-mini.** Greedy is able to search to higher depths given the same sampling budget.

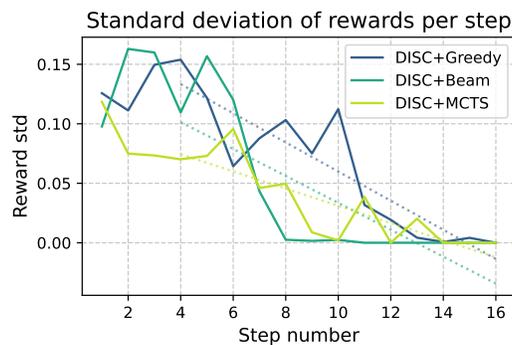


Figure 42: **Mean standard deviation of different decomposition search methods combined with DISC on APPS with gpt-4o-mini.** All search methods display decreasing standard deviation with search depth.