

EquivPruner: Boosting Efficiency and Quality in LLM-Based Search via Action Pruning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) excel at complex reasoning through search algorithms, yet current strategies often suffer from massive token consumption due to redundant exploration of semantically equivalent steps. Existing semantic similarity methods struggle to accurately identify such equivalence in domain-specific contexts like mathematical reasoning. To address this, we propose *EquivPruner*, a simple yet effective approach that identifies and prunes semantically equivalent actions during LLM reasoning search. We also introduce MathEquiv, the first dataset we created for mathematical statement equivalence, which enables the training of a lightweight equivalence detector. Extensive experiments across various models and tasks demonstrate that *EquivPruner* significantly reduces token consumption, improving searching efficiency and often bolstering reasoning accuracy. For instance, when applied to Qwen2.5-Math-7B-Instruct on GSM8K, *EquivPruner* reduced token consumption by 48.1% while also improving accuracy. Our code is available at <https://anonymous.4open.science/r/EquivPruner-2364/>.

1 Introduction

Large Language Models (LLMs) are increasingly demonstrating remarkable capabilities, yet their performance scaling during pretraining faces potential constraints due to data limitations (Lightman et al., 2023). Consequently, enhancing LLM capabilities at inference time has become a critical research frontier (Snell et al., 2024). A prominent direction involves leveraging search algorithms, particularly reward-guided tree search, to improve complex reasoning (Ke et al., 2025). These methods typically expand the search space by generating multiple reasoning steps (e.g., via chain-of-thought prompting (Wei et al., 2022)) and employ techniques like beam search (Kang et al., 2024) or Monte Carlo Tree Search (MCTS) (Chen et al.,

2024; Zhang et al., 2024a) to navigate this space and identify high-quality solutions.

However, current search strategies exhibit significant inefficiencies (Damani et al., 2024). A common practice involves sampling multiple candidate reasoning steps and exploring them, often allocating computational resources uniformly across these candidates (Yao et al., 2023; Long, 2023; Besta et al., 2024). This approach overlooks the potential semantic equivalence among textually distinct candidates. Treating semantically identical reasoning steps as unique branches leads to redundant exploration of the search space, incurring substantial computational overhead through unnecessary token generation and processing. Moreover, for search algorithms that incorporate preference learning based on intermediate steps (e.g., Xie et al., 2024; Jiang et al., 2024), learning preferences from pairs of equivalent steps may provide noisy or conflicting signals, hindering the learning of effective reasoning policies. This challenge is particularly acute in domains like mathematical reasoning, where numerous textual formulations can represent the same underlying logical operation or state. Addressing this redundancy via standard Semantic Textual Similarity (STS) techniques (Majumder et al., 2016) proves inadequate as illustrated in Figure 1. Existing embedding models, such as SBERT (Reimers and Gurevych, 2019), predominantly trained on general text, often fail to capture the nuanced structural and logical equivalence specific to mathematical statements. Even domain-specific models like MathBERT (Peng et al., 2021), which enhance mathematical text representation, along with other embedding models MATH-Similarity (Steinfeldt and Mihaljević, 2024), lack optimization for identifying functional equivalence between mathematical sentences. This limitation is further exacerbated by the lack of specialized benchmark datasets designed for mathematical statement equivalence. Although large-scale generative models can achieve

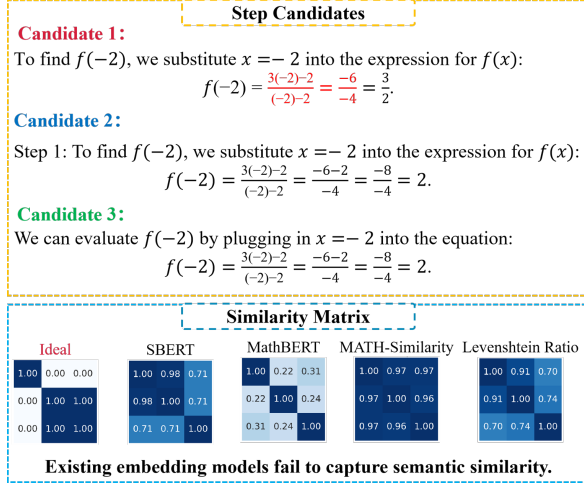


Figure 1: Illustration of the mathematical statement equivalence challenge during reasoning search. Given multiple candidate steps generated by an LLM, standard methods like embedding similarity or Levenshtein Ratio may incorrectly assess candidate 1 and candidate 2 as highly similar due to surface features, while failing to recognize the true semantic equivalence between candidate 2 and candidate 3, which represent the identical logical operation.

satisfactory performance in few-shot scenarios for such judgment tasks, their substantially higher computational complexity results in significantly slower inference speeds compared to embedding models (Brown et al., 2020). The consequent latency renders them impractical for high-throughput applications requiring real-time processing.

To overcome these limitations, we introduce *EquivPruner*, a simple yet effective approach that centers on identifying and pruning semantically equivalent actions during LLM reasoning search. We create MathEquiv, the first dataset specifically designed for mathematical statement equivalence. Leveraging this dataset, we trained a lightweight yet effective equivalence detection model. This model serves as a dynamic pruner integrated into the LLM’s search process. When the LLM generates multiple candidate reasoning steps at a given expansion point, the pruner identifies sets of semantically equivalent candidates among these siblings. For each set of equivalent steps, it retains only a single representative node for further exploration, effectively pruning the redundant branches and significantly reducing the search space.

While the proposed pruning framework is potentially generalizable, this paper focuses on its validation within mathematical reasoning due to the significant research community attention (Ke

et al., 2025) and the availability of well-developed open-source process reward models (Shao et al., 2024). We conduct extensive experiments across various models, including Mistral-7B-SFT (Shao et al., 2024) and the Qwen2.5-Math-Instruct series (Yang et al., 2024), using two widely recognized math reasoning benchmarks: GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). Our proposed *EquivPruner* demonstrates compelling improvements across these settings. For instance, when applied to Qwen2.5-Math-7B-Instruct on GSM8K—where the model already achieves a very high baseline accuracy of 96.44%—*EquivPruner* not only cuts token consumption by a substantial 48.1% but also further boosts accuracy to 96.59%. This demonstrates *EquivPruner*’s ability to significantly enhance searching efficiency.

Our main contributions are:

- To the best of our knowledge, this work is the *first* to identify and address the problem of action equivalence in LLM-based reasoning search.
- We introduce *EquivPruner*, a simple yet effective approach that centers on identifying and pruning semantically equivalent actions during LLM reasoning search.
- We release MathEquiv, the first benchmark dataset specifically designed for mathematical statement equivalence. It serves as a versatile resource applicable to a variety of mathematical tasks and scenarios.
- Extensive experiments demonstrate the effectiveness of *EquivPruner*. When applied to Qwen2.5-Math-7B-Instruct on GSM8K, *EquivPruner* not only cuts token consumption by a substantial 48.1% but also further boosts accuracy in a very high baseline.

2 Related Work

LLM Reasoning via Search Strategies Efforts to improve LLM problem-solving capabilities have moved beyond simple prompting. Chain-of-Thought prompting (Wei et al., 2022) demonstrated the value of intermediate reasoning steps. Building on this, structured search methods like Tree-of-Thoughts (Yao et al., 2023) and Graph-of-Thoughts (Besta et al., 2024) explore multiple reasoning paths, enhancing performance on complex tasks requiring exploration and backtracking.

Further advancing this direction, a particularly powerful paradigm integrates LLMs with sophisticated search algorithms. Among these, the synergy between LLMs and Monte Carlo Tree Search (MCTS) (Chen et al., 2024; Zhang et al., 2024a) has garnered significant attention for tackling complex reasoning problems. MCTS, renowned for its ability to balance exploration and exploitation in vast search spaces, becomes exceptionally potent when guided by an LLM’s generative capabilities to propose candidate steps and a reward model to estimate state values (Yao et al., 2023; Long, 2023; Besta et al., 2024). This LLM-MCTS approach, alongside other advanced search integrations like LLM-guided beam search (Chen et al., 2024), has consistently achieved state-of-the-art results in demanding areas such as science tasks (Yang et al., 2024), coding (Dainese et al., 2024; Zhang et al., 2023), and mathematical reasoning (Zhang et al., 2024b; Luo et al., 2024). However, despite the remarkable success of these advanced search strategies, a significant challenge emerges, especially prevalent in mathematical reasoning when employing methods like LLM-MCTS: the substantial token cost (Chen et al., 2024). While LLM-MCTS explores many branches effectively, it often wastes resources evaluating syntactically distinct but semantically equivalent states. This redundancy unnecessarily expands the search space, consuming tokens without yielding novel solutions, thus limiting efficiency and scalability.

Mathematical Equivalence Detection Effective detection of mathematical statement equivalence is crucial for efficient LLM-Based search tree pruning, yet current methodologies exhibit significant shortcomings. For instance, rudimentary sequence comparison metrics like Levenshtein similarity (Yujian and Bo, 2007) are fundamentally ill-suited, as they fail to capture the deep semantic and hierarchical structures inherent in mathematical language, leading to unreliable equivalence assessments. Standard Semantic Textual Similarity models, such as SBERT (Reimers and Gurevych, 2019), trained predominantly on general language corpora, are designed to capture semantic relatedness rather than strict mathematical equivalence. Even domain-specific models like MathBERT (Peng et al., 2021), which enhance mathematical text representation, along with other embedding models MATH-Similarity (Steinfeldt and Mihaljević, 2024), lack optimization for identifying

functional equivalence between mathematical sentences. Their capacity to accurately recognize semantically equivalent mathematical sentences is thereby constrained, as illustrated by the examples in Figure 1. While LLMs like GPT-4o (Hurst et al., 2024) has the ability to recognize mathematical equivalence, their complex architectures introduce significant latency. This high time overhead renders them impractical for real-time pruning scenarios. Consequently, there is an urgent need to enable efficient pruning in LLM-based search.

3 Methodology

3.1 Define Semantic Equivalence in Mathematics

Simply equating statements based on identical outcomes can be superficial and misleading, as it may overlook critical differences in conceptual articulation, structural formulation, symbolic interpretation, and methodological pathways. To address this, we propose a definition of semantic equivalence specifically attuned to these multifaceted aspects. Accordingly, in our framework, two mathematical statements are considered semantically equivalent if and only if they rigorously satisfy the following criteria:

- **Conceptual Consistency:** The statements must articulate identical mathematical concepts, definitions, or propositions without ambiguity.
- **Structural Equivalence:** Their logical formulations, encompassing assumptions, derivations, and conclusions, must be fully aligned.
- **Notational Precision:** All variables, symbols, and mathematical expressions must be employed consistently, maintaining identical meanings across the statements.
- **Methodological Congruence:** Semantic equivalence necessitates an alignment in the underlying methodology and reasoning. Statements yielding the same final result via disparate approaches are not considered fully equivalent.

Our approach to semantic equivalence thus mandates a comprehensive assessment. It scrutinizes the congruence of conceptual foundations, logical structures, notational usage, and methodological approaches. Two mathematical statements

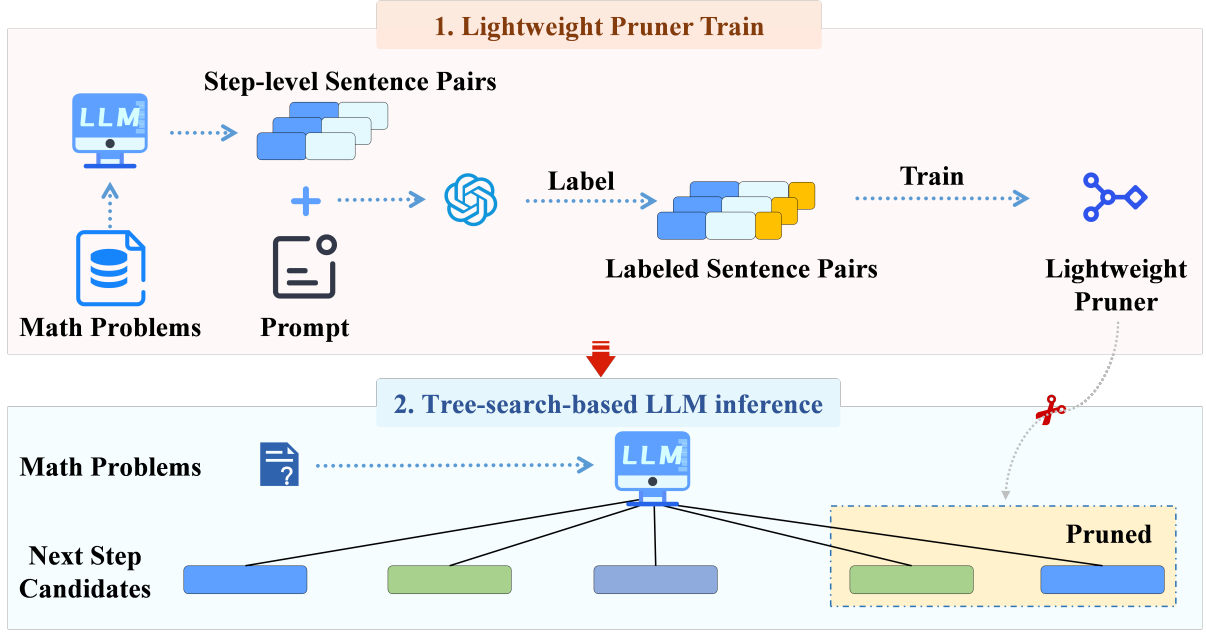


Figure 2: The EquivPruner framework. Top: Training the lightweight equivalence pruner from labeled step-level sentence pairs. Bottom: Applying the trained lightweight pruner during tree-search-based LLM inference to remove semantically equivalent candidates generated by the LLM.

are judged completely equivalent only when they demonstrate indivisible identity across all these critical facets.

3.2 The MathEquiv Dataset

Recognizing the absence of dedicated datasets for mathematical statement equivalence, we constructed and released MathEquiv to bridge this gap. The MathEquiv dataset was curated by initially employing a Step-level Beam Search algorithm (Chen et al., 2024) to gather action candidates. These candidates were subsequently formulated into step-level sentence pairs.

For the task of equivalence scoring, we implemented a five-tiered classification system. This granular approach was adopted to enhance the stability of the GPT model’s outputs, as preliminary experiments with binary classification (equivalent/non-equivalent) revealed inconsistencies in judgments. The five-tiered system yielded significantly more consistent and reliable assessments:

- **Level 4 (Exactly Equivalent):** The statements are mathematically interchangeable in all respects, exhibiting identical meaning and form.
- **Level 3 (Likely Equivalent):** Minor syntactic differences may be present, but the core mathematical content and logic align.

- **Level 2 (Indeterminable):** Insufficient information is available to make a definitive judgment regarding equivalence.
- **Level 1 (Unlikely Equivalent):** While some partial agreement may exist, critical discrepancies in logic, definition, or mathematical structure are observed.
- **Level 0 (Not Equivalent):** The statements are fundamentally distinct in their mathematical meaning, derivation, or resultant outcomes.

The MathEquiv dataset was labeled via an iterative refinement process. Initially, GPT-4o labeled a data subset, followed by human expert review. For discrepancies, the human-adjudicated label and its rationale were incorporated into GPT-4o’s prompt as few-shot examples. This cycle was repeated until model outputs for a randomly sampled subset consistently aligned with human consensus. Subsequently, the collection of few-shot examples was pruned to a minimal, representative set sufficient to maintain this model-human alignment. This iterative calibration process yielded the MathEquiv dataset, characterized by high-quality labels and an accurate assessment of semantic equivalence. The final prompt is detailed in Figure 4. The MathEquiv dataset is available at https://anonymous.4open.science/r/MathEquiv_dataset-B2B2/.

3.3 Lightweight Pruner for Tree Search

To facilitate dynamic, real-time pruning within our tree search algorithm, we developed and trained a dedicated Lightweight Pruner. The data collection process for training this pruner and its integration into the broader Tree-search-based LLM inference pipeline are illustrated in Figure 2.

3.3.1 Data Complexity in Pruner Training

The MathEquiv dataset, suitable for assessing overall statement equivalence, presents specific challenges for training the Lightweight Pruner. The dataset’s step-level sentence pairs often consist of multiple sentences. A key difficulty is that step pairs labeled as non-equivalent at a macro-level may nevertheless contain sub-pairs of sentences that are semantically equivalent. This characteristic, common in data derived from intermediate mathematical problem-solving steps, can introduce ambiguity and hinder the pruner’s ability to learn fine-grained distinctions if not appropriately addressed. The true equivalence status of these sub-sentence pairs can be viewed as a latent aspect of the data.

3.3.2 Pruner Training via Expectation-Maximization (EM)

To effectively train the Lightweight Pruner amidst this data complexity, we employ the Expectation-Maximization (EM) algorithm, which is effective for handling the unobserved equivalence status of sub-sentence pairs within larger, complex training instances. The algorithm alternates between an Expectation (E) step and a Maximization (M) step:

1. E-step (Expectation Step): Given the model parameters $\theta^{(t)}$ at iteration t , the pruner predicts the equivalence probability of each sub-sentence pair in multi-sentence samples. Sub-sentence pairs with probabilities exceeding a threshold are treated as high-confidence equivalents and removed from samples to refine the dataset for the next step.

2. M-step (Maximization Step): The model parameters are updated to $\theta^{(t+1)}$ by maximizing the likelihood of the observed data, conditioned on the expectations derived in the E-step.

By training on samples that have been simplified or where latent equivalences have been accounted for, the model can better focus on learning more subtle or challenging distinctions necessary for effective pruning.

4 Experiments

In this section, we present a series of comprehensive experiments designed to validate the efficacy of EquivPruner.

4.1 MathEquiv Dataset Generation

We constructed the MathEquiv dataset for mathematical statement equivalence. The foundation of this dataset consists of 7,500 mathematical problems sourced from the MATH training set (Hendrycks et al., 2021). To prevent data leakage between training, validation, and test phases of EquivPruner, these 7,500 problems were first split into training, validation, and test sets using an 8:1:1 ratio, respectively. For each problem in these distinct sets, we generated candidate reasoning step pairs using the Qwen2.5-Math-7B-Instruct model (Yang et al., 2024) via Step-level Beam Search. These pairs were subsequently filtered based on Levenshtein distance, and a balanced sample from each set was then annotated for equivalence by GPT-4o. This process resulted in distinct training, validation, and test sets of annotated step pairs for EquivPruner. The specific parameters for step pair generation, filtering criteria, and the final dataset sizes are detailed in Appendix A.1.

4.2 Experimental Setup

4.2.1 Models and Datasets

For inference, we utilized several LLMs: Qwen2.5-Math-7B-Instruct (Yang et al., 2024), Mistral-7B-SFT (Shao et al., 2024), and Qwen2.5-Math-1.5B-Instruct (Yang et al., 2024). Given that existing open-source PRMs are predominantly tailored for mathematical reasoning, our current investigation is confined to mathematical tasks. Nevertheless, the EquivPruner framework is designed for generalizability and can be readily extended to other domains like code generation and commonsense reasoning. The Process Reward Model (PRM) employed for guiding the Monte Carlo Tree Search (MCTS) was Math-Shepherd-Mistral-7B-PRM (Shao et al., 2024). As EquivPruner was trained on data generated by Qwen2.5-Math-7B-Instruct, the Mistral-7B-SFT and Qwen2.5-Math-1.5B-Instruct models serve as out-of-distribution (OOD) models in our experiments.

Our evaluation was conducted on two standard benchmark datasets:

- **MATH** (Hendrycks et al., 2021): Featuring challenging competition-level mathematics

Table 1: Performance comparison of Vanilla MCTS and MCTS + EquivPruner across three language models on the MATH and GSM8K datasets. EquivPruner significantly reduces token consumption (Tokens, Ratio) while generally maintaining or improving accuracy (Acc, %). Best results within each model-dataset block are in **bold**.

Methods	MATH			GSM8K		
	Acc	Tokens	Ratio	Acc	Tokens	Ratio
<i>Qwen2.5-Math-7B-Instruct:</i>						
Vanilla MCTS	83.40	106773	100.00%	96.44	34826	100.00%
+ EquivPruner	84.00	74194	69.49%	96.59	18071	51.89%
<i>Mistral-7b-sft:</i>						
Vanilla MCTS	36.60	49251	100.00%	83.78	20217	100.00%
+ EquivPruner	37.40	38265	77.69%	85.06	12537	62.01%
<i>Qwen2.5-Math-1.5B-Instruct:</i>						
Vanilla MCTS	75.60	91811	100.00%	91.05	39337	100.00%
+ EquivPruner	75.60	71878	78.29%	90.75	23752	60.38%

problems. Due to computational demands, our evaluation on the MATH dataset was performed on the MATH-500 subset, identical to the test partition used in [Lightman et al. \(2023\)](#).

- **GSM8K** ([Cobbe et al., 2021](#)): Consisting of grade school mathematics word problems. Its test set has 1319 problems. Since EquivPruner was trained on data derived from MATH dataset problems, GSM8K is considered an OOD dataset.

4.2.2 Implementation Details

The EquivPruner model itself is a fine-tuned Longformer-base ([Beltagy et al., 2020](#)), chosen for its efficiency suitable for real-time pruning. During the MCTS inference phase, the determination of equivalence between two reasoning step nodes involves a two-stage process. First, the Levenshtein ratio between the steps is calculated. If the ratio is less than or equal to 0.75, the nodes are immediately considered non-equivalent, acting as a fast filter. Only if the Levenshtein ratio is greater than 0.75 is the EquivPruner model invoked to make the final equivalence prediction. This hierarchical check balances speed and accuracy in the pruning process. The maximum number of newly generated tokens by the LLMs (max_new_tokens) was set to 1024, and the generation temperature was 0.7. All experiments were conducted on NVIDIA GeForce RTX 3090 GPUs. Further details are available in [Appendix A.2](#).

4.3 Evaluation Metrics

We adopted a vanilla MCTS ([Chen et al., 2024](#)) as the baseline for comparison. The evaluation of

EquivPruner focuses on two primary aspects:

- **Effectiveness:** Measured using solution accuracy (Acc), the percentage of problems solved correctly.
- **Efficiency:** Assessed through the total number of tokens generated (Tokens) and a token ratio (Ratio), defined as the ratio of tokens generated by the EquivPruner-enhanced search to those generated by the baseline MCTS.

4.4 Main Results

Table 1 presents our main experimental findings, comparing vanilla MCTS against MCTS augmented with EquivPruner. The results consistently demonstrate that EquivPruner substantially enhances computational efficiency across different language models and datasets, primarily by reducing token generation while largely preserving or even improving solution accuracy.

Efficiency Gains EquivPruner achieves significant reductions in token counts across all configurations. For instance, with Qwen2.5-Math-7B-Instruct on GSM8K, tokens were reduced by approximately 48.11% (Ratio: 51.89%), and on MATH, by 30.51% (Ratio: 69.49%). Similar substantial token savings were observed for Mistral-7B-SFT (e.g., 37.99% reduction on GSM8K) and Qwen2.5-Math-1.5B-Instruct (e.g., 39.62% reduction on GSM8K). These figures highlight EquivPruner’s effectiveness in pruning the search space.

Accuracy Impact and Resource Optimization

Crucially, these efficiency improvements are generally accompanied by maintained or enhanced

Table 2: Performance of EquivPruner with Step-level Beam Search (SBS) using the Qwen2.5-Math-7B-Instruct model on MATH and GSM8K. EquivPruner enhances accuracy (Acc, %) by promoting diversity among selected nodes, with token counts (Tokens, Ratio) remaining largely stable.

Methods	Acc	Tokens	Ratio
MATH:			
SBS	82.00	21341	100.00%
+ EquivPruner	82.20	20952	98.18%
GSM8K:			
SBS	96.06	8004	100.00%
+ EquivPruner	96.13	7927	99.04%

accuracy. Qwen2.5-Math-7B-Instruct saw accuracy gains of +0.60% on MATH and +0.15% on GSM8K. With Mistral-7B-SFT, an OOD model relative to EquivPruner’s training data source, accuracy improved by +0.80% on MATH and +1.28% on GSM8K (also an OOD dataset for EquivPruner). This suggests that by eliminating redundant explorations, EquivPruner enables MCTS to allocate its search resources more effectively. For Qwen2.5-Math-1.5B-Instruct (another OOD model), accuracy was maintained on MATH and saw a minor dip of -0.30% on GSM8K, which is a reasonable trade-off given the nearly 40% token reduction.

Generalization The positive outcomes on OOD models (Mistral-7B-SFT, Qwen2.5-Math-1.5B-Instruct) and the OOD dataset (GSM8K) underscore EquivPruner’s generalization capabilities. It effectively identifies and removes equivalent reasoning steps, allowing MCTS to conduct a more focused and efficient search across varied settings.

4.5 Effectiveness in Step-level Beam Search

To demonstrate its versatility beyond MCTS, we evaluated EquivPruner with Step-level Beam Search (SBS) (Chen et al., 2024) using the Qwen2.5-Math-7B-Instruct model. Unlike MCTS, SBS does not construct an extensive search tree; instead, it dynamically selects top- k child nodes during expansion. Given this mechanism, applying EquivPruner to SBS is not primarily aimed at reducing the total number of generated tokens, as SBS inherently limits the breadth of the search. Instead, our hypothesis is that EquivPruner can enhance the **quality** of the search by eliminating redundant nodes **before** the top- k selection occurs. This process ensures that the k chosen candidates are more diverse, potentially leading to the discov-

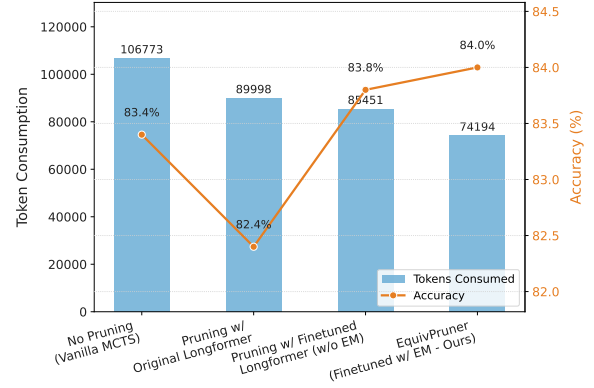


Figure 3: Ablation study of EquivPruner components. The plot illustrates the impact of different pruning strategies within a MCTS framework on Token Consumption (bars, left y-axis) and Accuracy (line, right y-axis).

ery of more effective reasoning paths and thereby improving overall task performance.

The results in Table 2 validate this. On MATH, EquivPruner increased accuracy from 82.00% to 82.20% (+0.20%), and on GSM8K from 96.06% to 96.13% (+0.07%). Concurrently, token counts remained largely unchanged, with ratios of 98.18% on MATH and 99.04% on GSM8K. These findings suggest that even in search algorithms like SBS where token generation is already constrained, EquivPruner can still offer benefits. By ensuring that the limited slots in the beam are occupied by semantically distinct reasoning steps, EquivPruner promotes a more diverse and potentially more fruitful exploration of the solution space. This demonstrates that EquivPruner is a versatile component that can enhance different types of search strategies in LLM-based reasoning by improving the quality and diversity of explored paths.

4.6 Ablation Study

To investigate the individual contributions of the key components of our EquivPruner—specifically, the fine-tuning process and the use of the EM algorithm—we conducted an ablation study. The experiments were performed using the Qwen2.5-Math-7B-Instruct model on the MATH dataset. We compare our full method, EquivPruner (Finetuned w/ EM), against three variants: (1) No Pruning (vanilla MCTS baseline); (2) Pruning w/ Original Longformer (using a pre-trained Longformer-base without task-specific fine-tuning for equivalence); and (3) Pruning w/ Finetuned Longformer (w/o EM) (standard supervised fine-tuning without the EM algorithm).

The results in Figure 3 demonstrate the impact of each component. Using the Original Longformer-base for pruning (Setting 2) reduces tokens (106773 to 89998) compared to No Pruning (Setting 1), but at the cost of a accuracy drop (83.4% to 82.4%), indicating that a generic model is insufficient. Standard fine-tuning without EM (Setting 3) improves accuracy to 83.8% (surpassing No Pruning) while improve token efficiency to Setting 2 (89998 to 85451), underscoring the necessity of task-specific training. Critically, our full EquivPruner method with EM-based fine-tuning (Setting 4) achieves both the highest accuracy (84.0%) and the most significant token reduction (106773 to 74194). This highlights that both the fine-tuning process and specifically the EM algorithm are vital for maximizing EquivPruner’s effectiveness in improving accuracy and token efficiency.

5 Conclusion

In this paper, we introduce *EquivPruner*, a simple yet effective approach to address inefficient token usage in LLM reasoning search by identifying and pruning semantically equivalent actions. We also introduce MathEquiv, the first dataset specifically designed for mathematical statement equivalence, which enables the training of an effective lightweight equivalence detector. Extensive experiments demonstrate that *EquivPruner* significantly reduces token consumption—for example, by 48.1% for Qwen2.5-Math-7B-Instruct on GSM8K—while maintaining or often improving reasoning accuracy across various models and tasks. Our findings underscore the substantial benefits of managing semantic redundancy in reasoning search, offering a valuable direction for enhancing the efficiency and effectiveness of LLMs.

Limitations

There are some limitations with our paper, which we reserve for future work. Firstly, due to computational constraints, *EquivPruner* was not evaluated on language models significantly larger than the 7B parameter scale. Secondly, our work focused on *EquivPruner*’s application at inference time, and its potential integration with iterative LLM training or refinement strategies remains an area for future exploration. Lastly, while designed for generalizability, our empirical validation was primarily within mathematical reasoning, influenced by the availability of suitable process reward models

(PRMs). Extending evaluation to other domains such as science tasks or commonsense reasoning, contingent upon broader PRM availability, offers a promising direction for future research.

References

- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. 2024. Alphamath almost zero: Process supervision without process. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. 2024. Generating code world models with large language models guided by monte carlo tree search. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Mehul Damani, Idan Shenfeld, Andi Peng, Andreea Bobu, and Jacob Andreas. 2024. Learning how hard to think: Input-adaptive allocation of lm computation. *arXiv preprint arXiv:2410.04707*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

645	Jinhao Jiang, Zhipeng Chen, Yingqian Min, Jie Chen,	Christian Steinfeldt and Helena Mihaljević. 2024. Eval-	699
646	Xiaoxue Cheng, Jiapeng Wang, Yiru Tang, Hao-	uation and domain adaptation of similarity models	700
647	xiang Sun, Jia Deng, Wayne Xin Zhao, and 1 oth-	for short mathematical texts. In <i>International Con-</i>	701
648	ers. 2024. Technical report: Enhancing llm reason-	<i>ference on Intelligent Computer Mathematics</i> , pages	702
649	ing with reward-guided tree search. <i>arXiv preprint</i>	241–260. Springer.	703
650	<i>arXiv:2411.11694</i> .		
651	Jikun Kang, Xin Zhe Li, Xi Chen, Amirreza Kazemi,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	704
652	Qianyi Sun, Boxing Chen, Dong Li, Xu He, Quan He,	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	705
653	Feng Wen, and 1 others. 2024. Mindstar: Enhancing	and 1 others. 2022. Chain-of-thought prompting elic-	706
654	math reasoning in pre-trained llms at inference time.	its reasoning in large language models. <i>Advances</i>	707
655	<i>arXiv preprint arXiv:2405.16265</i> .	<i>in neural information processing systems</i> , 35:24824–	708
		24837.	709
656	Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen,	Yuxi Xie, Anirudh Goyal, Wenye Zheng, Min-Yen	710
657	Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin,	Kan, Timothy P Lillicrap, Kenji Kawaguchi, and	711
658	Peifeng Wang, Silvio Savarese, and 1 others. 2025.	Michael Shieh. 2024. Monte carlo tree search boosts	712
659	A survey of frontiers in llm reasoning: Inference scal-	reasoning via iterative preference learning. <i>arXiv</i>	713
660	ing, learning to reason, and agentic systems. <i>arXiv</i>	<i>preprint arXiv:2405.00451</i> .	714
661	<i>preprint arXiv:2504.09037</i> .		
662	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	715
663	son Edwards, Bowen Baker, Teddy Lee, Jan Leike,	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	716
664	John Schulman, Ilya Sutskever, and Karl Cobbe.	Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.	717
665	2023. Let’s verify step by step. In <i>The Twelfth Inter-</i>	5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	718
666	<i>national Conference on Learning Representations</i> .		
667	Jieyi Long. 2023. Large language model guided tree-of-	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	719
668	thought. <i>arXiv preprint arXiv:2305.08291</i> .	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.	720
		2023. Tree of thoughts: Deliberate problem solving	721
669	Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat	with large language models. <i>Advances in neural</i>	722
670	Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li, Lei	<i>information processing systems</i> , 36:11809–11822.	723
671	Shu, Yun Zhu, Lei Meng, and 1 others. 2024. Im-		
672	prove mathematical reasoning in language models	Li Yujian and Liu Bo. 2007. A normalized levenshtein	724
673	by automated process supervision. <i>arXiv preprint</i>	distance metric. <i>IEEE transactions on pattern analy-</i>	725
674	<i>arXiv:2406.06592</i> .	<i>sis and machine intelligence</i> , 29(6):1091–1095.	726
675	Goutam Majumder, Partha Pakray, Alexander Gelbukh,	Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue,	727
676	and David Pinto. 2016. Semantic textual similarity	Yuxiao Dong, and Jie Tang. 2024a. Rest-mcts*: Llm	728
677	methods, tools, and applications: A survey. <i>Com-</i>	self-training via process reward guided tree search.	729
678	<i>putación y Sistemas</i> , 20(4):647–665.	<i>Advances in Neural Information Processing Systems</i> ,	730
		37:64735–64772.	731
679	Shuai Peng, Ke Yuan, Liangcai Gao, and Zhi Tang.	Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jia-	732
680	2021. Mathbert: A pre-trained model for math-	tong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang,	733
681	ematical formula understanding. <i>arXiv preprint</i>	Marco Pavone, Yuqiang Li, and 1 others. 2024b.	734
682	<i>arXiv:2105.00377</i> .	Llama-berry: Pairwise optimization for o1-like	735
		olympiad-level mathematical reasoning. <i>arXiv</i>	736
683	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	<i>preprint arXiv:2410.02884</i> .	737
684	Sentence embeddings using siamese bert-networks.	Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu	738
685	In <i>Proceedings of the 2019 Conference on Empirical</i>	Ding, Joshua B Tenenbaum, and Chuang Gan. 2023.	739
686	<i>Methods in Natural Language Processing and the 9th</i>	Planning with large language models for code gener-	740
687	<i>International Joint Conference on Natural Language</i>	ation. In <i>The Eleventh International Conference on</i>	741
688	<i>Processing (EMNLP-IJCNLP)</i> , pages 3982–3992.	<i>Learning Representations</i> .	742
		A Experimental Details	743
689	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	A.1 MathEquiv Dataset Generation Details	744
690	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	The MathEquiv dataset was constructed as follows:	745
691	Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-	Problem Sourcing and Splitting: We selected	746
692	math: Pushing the limits of mathematical reason-	7,500 problems from the MATH training set	747
693	ing in open language models. <i>arXiv preprint</i>	(Hendrycks et al., 2021). These problems were	748
694	<i>arXiv:2402.03300</i> .	divided into three distinct sets for EquivPruner: a	749
695	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-	training set (6,000 problems, 80%), a validation	750
696	mar. 2024. Scaling llm test-time compute optimally	set (750 problems, 10%), and a test set (750 prob-	751
697	can be more effective than scaling model parameters.	lems, 10%). This initial split of problems ensures	752
698	<i>arXiv preprint arXiv:2408.03314</i> .		

no data leakage between the subsequently generated step-pair datasets for EquivPruner. **Step Pair Generation:** For each problem within these three sets, we generated candidate reasoning steps using the Qwen2.5-Math-7B-Instruct model (Yang et al., 2024). This generation was performed via a Step-level Beam Search with the following parameters: beam size (k) = 8, temperature = 0.7, maximum search tree width (tree_max_width) = 10, maximum search tree depth (tree_max_depth) = 50, and maximum new tokens for generation (max_new_tokens) = 1024. **Filtering:** The generated step pairs from each set were then filtered based on their Levenshtein ratio. Only pairs with a ratio between 0.75 and 0.99 (inclusive) were retained. This filtering aimed to capture meaningful variations while excluding nearly identical or overly dissimilar steps. **Sampling and Annotation:** From the filtered pairs of each set, we randomly sampled a large number for annotation: Training set: 80,000 pairs were annotated. Validation set: 10,000 pairs were annotated. Test set: 10,000 pairs were annotated. This process resulted in the final training, validation, and test sets for the MathEquiv dataset, with no overlap in the underlying problems from which the step pairs were derived.

A.2 Implementation Environment and MCTS Parameters

All experiments were conducted using PyTorch version 2.4.0. The GPU infrastructure consisted of eight NVIDIA GeForce RTX 3090 GPUs, each with 24GB, utilizing CUDA version 12.1. The central processing unit was an Intel(R) Xeon(R) Platinum 8255C CPU equipped with 96 cores.

A.2.1 EquivPruner Training

The EquivPruner model, a fine-tuned Longformer-base (Beltagy et al., 2020), was trained using hyperparameters selected via Bayesian optimization. The optimization aimed to maximize the ‘eval/f1’ metric over a maximum of 10 trials. The hyperparameter search spaces are detailed in Table 3.

A.2.2 MCTS Parameters

The Monte Carlo Tree Search (MCTS) based evaluation hyperparameters are detailed in Table 4. These MCTS parameters (temperature, tree_max_width , tree_max_depth , simulations, PUCT values) were kept consistent across base-

Table 3: Hyperparameter search space for EquivPruner using Bayesian optimization.

Hyperparameter	Value or Range
Learning Rate	[1e-6, 5e-5]
Batch Size	4
Training Epochs	Discrete Values {2, 3, 5}
Weight Decay	[0.0, 0.1]

line and EquivPruner-enhanced evaluations unless otherwise specified.

Table 4: Monte Carlo Tree Search (MCTS) hyperparameters.

Hyperparameter	Value
Number of Simulations	20
LLM Generation Temperature	0.7
LLM max_new_tokens	1024
Search Tree Maximum Width	10
Search Tree Maximum Depth	50
PUCT values	1.25

A.2.3 SBS Parameters

The Step-level Beam Search (SBS) based evaluation hyperparameters are detailed in Table 5. These SBS parameters (beam size, temperature, tree_max_width , tree_max_depth) were kept consistent across baseline and EquivPruner-enhanced evaluations unless otherwise specified.

Table 5: Step-level Beam Search (SBS) hyperparameters.

Hyperparameter	Value
Beam Size	3
LLM Generation Temperature	0.7
LLM max_new_tokens	1024
Search Tree Maximum Width	10
Search Tree Maximum Depth	50

Please determine whether the following two sentences are semantically equivalent, and return 0: Not equivalent at all; 1: May not be equivalent; 2: Can't judge; 3: May be equivalent; 4: Exactly equivalent. Please explain the reason, reflect, and provide a more accurate result. Please output in the following Python dictionary format:

{ "reasoning_step": "The reasoning process of the model", "result": "Final result" (int) }

Question:

Sentence1: The inner sum is a geometric series with the first term $\frac{1}{2^{k^2-k+1}}$ and common ratio $\frac{1}{2}$, and it has $2k$ terms. The sum of a geometric series is given by:

$$\sum_{n=a}^{a+b-1} ar^n = a \frac{1-r^b}{1-r}$$

Sentence2: The inner sum $\sum_{n=k^2-k+1}^{k^2+k} \frac{1}{2^n}$ is a geometric series with the first term $\frac{1}{2^{k^2-k+1}}$ and common ratio $\frac{1}{2}$, and it has $2k$ terms. The sum of a geometric series is given by:

$$\sum_{n=a}^{a+r-1} ar^n = a \frac{1-r^r}{1-r}$$

Output:

{ "reasoning_step": "Sentence1 uses the formula $\sum_{n=a}^{a+b-1} ar^n = a \frac{1-r^b}{1-r}$, while Sentence2 uses $\sum_{n=a}^{a+r-1} ar^n = a \frac{1-r^r}{1-r}$. The variables and the structure of the formulas are different, which affects the semantic equivalence. Additionally, the limits of the summation in Sentence2 are explicitly given as $\sum_{n=k^2-k+1}^{k^2+k} \frac{1}{2^n}$, which is not present in Sentence1. These differences indicate that the sentences are not semantically equivalent.", "result": 0 }

Question:

Sentence1: The amplitude is the distance from the midline to the maximum value (or the minimum value). Therefore, we can calculate the amplitude as follows:

$$\text{Amplitude} = \text{Maximumvalue} - \text{Midline} = 3 - 1 = 2$$

Since the amplitude is $|a|$, we have: $|a| = 2$. Since a is a positive constant, we conclude: $a = 2$

Sentence2: The amplitude of the sine wave is half the distance between the maximum value and the minimum value. Therefore, we can calculate the amplitude as follows:

$$\text{Amplitude} = \frac{\text{Maximumvalue} - \text{Minimumvalue}}{2} = \frac{3 - (-1)}{2} = \frac{3 + 1}{2} = \frac{4}{2} = 2$$

Since the amplitude is $|a|$, we have: $|a| = 2$. Since a is a positive constant, we conclude: $a = 2$

Output: { "reasoning_step": "Sentence1 defines the amplitude as the distance from the midline to the maximum value, calculating it as $\text{Amplitude} = \text{Maximumvalue} - \text{Midline} = 3 - 1 = 2$. Sentence2 defines the amplitude as half the distance between the maximum and minimum values, calculating it as $\text{Amplitude} = \frac{\text{Maximumvalue} - \text{Minimumvalue}}{2} = \frac{3 - (-1)}{2} = 2$. While both methods yield the same result ($a = 2$), the definitions and calculations are fundamentally different. This difference in methodology means the sentences are not semantically equivalent.", "result": 0 }

Question:

Sentence1: {sentence1}

Sentence2: {sentence2}

Output:

Figure 4: Complete prompt for labeling.