

---

# The Effect of Dataset Diversification on Mathematical Problem Solving Performance

---

**Jason Yuan**

Crystal Springs Uplands School  
Hillsborough, CA 94010  
friedepigs@gmail.com

## Abstract

We investigate the impact of dataset diversification on mathematical problem-solving performance and find that diversity control can substantially improve model capabilities. Using Farthest Point Sampling across different diversity levels, we observe performance gains of 12.7 percentage points on GSM8K and 12.4 percentage points on MATH benchmark, with 25% diversity performing optimally across both tasks. Our evaluation reveals that the relationship between diversity level and performance is non-monotonic, with intermediate diversity levels outperforming both random sampling and maximum diversity approaches. Through experiments on NuminaMath, Hendrycks MATH, and MATH-Plus datasets, we demonstrate that these improvements depend on training set size: while 1k examples show minimal diversity benefits, 3k examples exhibit substantial gains. We also find that diverse sampling can harm performance by selecting low-quality examples from noisy datasets, highlighting the importance of quality control in diversification strategies.

## 1 Introduction

The selection of training data influences language model performance, yet many approaches rely on random sampling from available datasets. While random sampling ensures unbiased representation, it may not optimally capture the diversity of problem types necessary for robust generalization, particularly in complex domains like mathematical reasoning. Recent advances in large language models have demonstrated impressive capabilities on mathematical tasks, but their performance remains dependent on the quality and diversity of training examples. Recent work has introduced data diversification strategies [5, 8, 26] and new benchmarks for mathematical reasoning such as HARDMATH [10], MathOdyssey [24], MATH-Vision [27], and advances in chain-of-thought prompting [29]. Mathematical reasoning presents challenges for training data curation. Problems span multiple subjects (algebra, geometry, combinatorics), difficulty levels, and solution approaches, creating a high-dimensional space of reasoning patterns. Random sampling may underrepresent critical problem categories or fail to capture the full spectrum of mathematical concepts needed for comprehensive understanding. This limitation becomes particularly pronounced when training on subsets of large datasets, where random selection might miss important edge cases or problem types that are crucial for generalization.

This paper investigates whether systematic diversity-based sampling can improve mathematical reasoning performance compared to random selection. We employ Farthest Point Sampling to select diverse training examples in embedding space, combined with quality control measures to address potential issues with noisy data. Our approach systematically varies diversity levels from 0% (random sampling) to 100% (maximum diversity) to identify effective configurations for mathematical reasoning tasks.

We find that diversity control achieves substantial performance improvements: 12.7 percentage points on GSM8K and 12.4 percentage points on the MATH benchmark when using 25% diversity levels. However, these benefits depend on training set size—while 1k training examples show minimal

diversity advantages, 3k examples enable substantial gains. We also find that diverse sampling can select low-quality examples from noisy datasets, requiring quality control measures for effective implementation. Our experiments suggest three main patterns: (1) 25% diversity appears to represent an effective level across multiple mathematical reasoning benchmarks, (2) training set size plays a critical role in determining diversity effectiveness, and (3) the relationship between diversity level and performance is non-monotonic, with intermediate diversity often outperforming both random and maximum diversity approaches.

## 2 Related work

Our study is situated within several active areas of research, drawing upon advancements in dataset construction, the application of language models to complex reasoning, and data sampling methodologies.

Benchmarks beyond GSM8K and MATH, such as MiniF2F [32], have pushed formal reasoning evaluation. Recent mathematical reasoning datasets include MuMath [31], MathOdyssey [24], HARDMATH [10], and multimodal extensions such as MATH-Vision [27]. Analyses such as Mehta et al. [16] and recent studies on diversified sampling [5, 8, 1] highlight current directions. Massive supervised data composition analyses [17] explore scaling impacts. Density-based clustering [9] has roots in spatial analysis and its application to data mining.

**Dataset Diversity and Active Learning Literature:** The concept of diverse data selection has long been a focus in machine learning, particularly within the field of active learning. Active learning aims to reduce the labeling effort by strategically selecting the most informative examples for annotation, often by prioritizing uncertainty or diversity. While active learning typically involves human feedback, the underlying principles of selecting representative or boundary-case examples for improved generalization are relevant to our work. Our research focuses on diverse sampling for *pre-existing* datasets, aiming to optimize the training subset without new annotations. This extends the active learning paradigm to dataset curation in a self-supervised context, where the goal is to capture maximal information or challenge with a limited subset.

**Mathematical Reasoning Datasets and Training Methods:** The development of language models for mathematical reasoning has seen significant progress, marked by the creation of specialized datasets and innovative training paradigms. Benchmarks like GSM8K [7] and MATH [11] have driven advancements in mathematical problem-solving capabilities. Techniques such as Chain-of-Thought (CoT) prompting [29] and fine-tuning on diverse mathematical corpora have enabled models to exhibit more robust reasoning. While much of this work has focused on architectural improvements and scaling, our study examines the dimension of *dataset composition*. We explore how the diversity of a training dataset, beyond its size, may influence a model’s ability to generalize across mathematical problem types.

Clustering algorithms, such as K-means [14, 2] and HDBSCAN [4, 15], group similar data points based on their features or embeddings, and have been used to create balanced or representative subsets [9]. However, their effectiveness depends on the underlying data distribution exhibiting clear cluster structures. Conversely, distance-based sampling techniques like Farthest Point Sampling (FPS) [13] aim to maximize the dissimilarity between selected data points, ensuring broad coverage of the embedding space without assuming predefined clusters. Our work evaluates these methodologies in the context of mathematical problems, finding that while clustering methods often struggle possibly due to the amorphous nature of mathematical problem embeddings, FPS appears effective, particularly when integrated with quality control mechanisms. This comparison of sampling techniques, combined with an emphasis on data quality [22], forms part of our approach to dataset curation. Comprehensive reviews of augmentation techniques [12] and the role of dataset diversity [28] are emerging. Few-shot evaluation strategies [20] and data synthesis approaches [3] are also important for robust benchmarking. Surveys have further analyzed dataset diversification at scale [21, 23, 19].

## 3 Methodology

### 3.1 Datasets

#### 3.1.1 AI-MO/NuminaMath-1.5

NuminaMath 1.5 [25] is a large-scale mathematical reasoning dataset containing approximately 896k competition-level problems sourced from Chinese high school exercises, US and international

mathematics olympiads, and online discussion forums. The dataset contains eight types of problems: Algebra (422,915 problems), Geometry (183,796), Number Theory (91,383), Combinatorics (73,643), Inequalities (41,741), Logic and Puzzles (37,911), Calculus (26,965), and Other (12,645). Problems are formatted with Chain-of-Thought solutions and include metadata for problem type, question type, and answer. The dataset was primarily collected from online exam PDFs and mathematics forums, resulting in varying solution quality due to the web-scraping process [6]. The variation in quality, including incomplete solutions, cross-references to missing problems, and multi-problem entries, makes NuminaMath particularly suitable for investigating the impact of dataset diversification on noisy training data.

### 3.1.2 Hendrycks MATH benchmark

The MATH dataset [11] contains 12,500 high school competition-level mathematics problems designed to evaluate mathematical reasoning capabilities. Problems are sourced from well-established competitions including AMC 10, AMC 12, and AIME, ensuring high solution quality. The dataset contains seven types of problems (Prealgebra, Algebra, Number Theory, Counting and Probability, Geometry, Intermediate Algebra, and Precalculus) with difficulty ratings from 1 to 5. Each problem includes a complete statement, step-by-step solution, final answer, and metadata for subject classification and difficulty. It has a 12,000/500 train/test split and detailed step-by-step solutions for each problem, making it a clean, high-quality benchmark for testing diversification methods without the quality concerns present in web-scraped datasets.

### 3.1.3 TIGER-Lab/MATH-plus

MATH-Plus is an augmented mathematical reasoning dataset containing 893,929 examples created by combining and expanding MetaMath [30], MATH-Orca, and additional MATH problems augmented using GPT-4. The augmentation process involves generating variations of existing problems through techniques such as changing numerical values, rephrasing problem statements, and creating similar problem structures while maintaining mathematical rigor. This results in a larger, high-quality dataset compared to the original component datasets. We split MATH-Plus into 715,143 training examples and 178,786 test examples (80/20 split) for our experiments. The dataset’s large size and consistent quality made it ideal for our diversification experiments, providing sufficient examples to train models effectively while avoiding the quality control issues encountered with web-scraped datasets.

### 3.1.4 openai/gsm8k

GSM8K [7] (Grade School Math 8K) is a dataset of 8,500 high-quality grade school math word problems. Each problem typically requires 2-8 solution steps involving arithmetic operations and concepts no more advanced than Algebra. The dataset was created through freelance contractors and subsequently verified by Surge AI, with solutions provided in natural language. The standard split contains 7,473 training and 1,027 test examples. We use GSM8K as an evaluation benchmark to test whether diversity effects transfer across different mathematical complexity levels.

## 3.2 Diversification Methods

We evaluate three distinct approaches to diverse sampling, each representing different assumptions about the structure of mathematical problem datasets in embedding space.

### 3.2.1 Farthest Point Sampling (FPS)

Our most successful approach uses Farthest Point Sampling to select spaced points in the MathBERT embedding space [13] without making assumptions about cluster structure. The algorithm begins by randomly selecting an initial point, then iteratively selects the point that is farthest from all previously selected points, measured by Euclidean distance in the 768-dimensional embedding space.

We implement a parameterized version of FPS where  $X\%$  diversity means selecting from the  $(100-X)\%$  farthest candidate points at each iteration. Specifically,  $0\%$  diversity corresponds to pure random sampling,  $100\%$  diversity always selects the single farthest point (pure FPS), and  $25\%$  diversity identifies the  $75\%$  farthest points and randomly selects among them. This approach allows fine-grained control over diversity constraints while maintaining the iterative FPS structure and provides a natural interpolation between random and maximally diverse sampling. Alternative strategies such as prismatic synthesis [26], key-point-driven diversification [10], and ARROWS [5] have further advanced reasoning capabilities in new benchmarks.

### 3.2.2 Clustering-Based Methods

We also explored MiniBatchKMeans clustering [15] and HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) [4]. However, these clustering methods were less effective in our experiments, possibly due to the lack of clear structure in mathematical problem embeddings. HDBSCAN consistently identified only 2 meaningful clusters while classifying 42-80% of examples as outliers, suggesting that mathematical problems may not exhibit density-based cluster structure in MathBERT embedding space.

### 3.3 Quality Assessment Framework

To understand the relationship between diverse sampling and data quality, we developed a framework for assessing the trainability of mathematical problems across different sampling strategies.

#### 3.3.1 Pattern-Based Detection

We implement automated pattern-based detection to identify potentially untrainable examples using regular expressions and heuristic rules. Untrainable examples include: (1) problems with missing or incomplete solutions, (2) cross-references to external figures or problems not included in the dataset, (3) formatting errors that make problems incomprehensible, (4) multi-problem entries where a single example contains multiple distinct problems, and (5) solutions that consist only of final answers without reasoning steps.

#### 3.3.2 GPT-4 Evaluation Protocol

To validate and complement our pattern-based analysis, we employ GPT-4 as an expert-level evaluator for a representative subset of examples. For each dataset and sampling method, we randomly select 500 examples for manual evaluation. GPT-4 assesses each example on trainability and difficulty (1-10 scale), providing ground truth for our automated quality metrics.

### 3.4 Accuracy Evaluation

We evaluate model accuracy using GPT-4o-mini as an automated judge to determine whether model-generated answers are mathematically equivalent to ground truth solutions. Our evaluation prompt instructs the model to focus on mathematical equivalence rather than exact string matching, accounting for different valid representations of the same result. The evaluation criteria include: (1) mathematical equivalence of final numerical answers, (2) equivalent forms such as fractions vs decimals or simplified vs unsimplified expressions, (3) reasonable rounding differences in decimal approximations, and (4) formatting differences (e.g., "1/2" vs "0.5" vs "50%"). The prompt provides specific examples of equivalent answers and instructs GPT-4o-mini to respond with "YES" for mathematically equivalent answers or "NO" otherwise. This approach aims to ensure consistent and fair evaluation across different answer formats while maintaining mathematical rigor.

## 4 Experimental Setup Sample

### 4.1 Mathematical Reasoning via Fine-tuning

We fine-tune Qwen-1.5B-Instruct on diversely sampled subsets of 3,000 examples from MATH-Plus, evaluating on external benchmarks (GSM8K and MATH) to ensure unbiased assessment. Training uses LoRA fine-tuning with 4-6 epochs, learning rate  $5e-5$ , and standard hyperparameters. We vary diversity levels from 0% (random sampling) to 100% (maximum diversity) using Farthest Point Sampling.

## 5 Results

### 5.1 Optimal Diversity Level: 25% Performs Best

Our primary finding shows that 25% diversity achieves optimal performance across mathematical reasoning benchmarks when using 3k training examples, with substantial improvements over random sampling.

Recent large-scale datasets such as MathOdyssey [24], HARDMATH [10], and multimodal sets such as MATH-Vision [27] further expand the scope of evaluation.

Benchmark	Random (0%)	Optimal (25%)	Improvement
GSM8K	35.8%	48.5%	+12.7 pts
MATH	39.0%	51.4%	+12.4 pts

Table 1: Performance improvements from optimal diversity level on benchmarks using 3k examples

The relationship between diversity level and performance follows a non-monotonic pattern across both benchmarks, with intermediate diversity levels outperforming both random sampling and maximum diversity approaches.

## 5.2 Training Set Size Dependency

An important finding is that diversity benefits depend strongly on training set size. While 3k examples enable substantial improvements, 1k examples show minimal to no diversity advantages.

Dataset	Training Size	Random (0%)	Best Diversity	Best Accuracy	Advantage
GSM8K	1k	35.6%	0%	35.6%	0.0 pts
GSM8K	3k	35.8%	25%	48.5%	+12.7 pts
MATH	1k	42.6%	50%	44.8%	+2.2 pts
MATH	3k	39.0%	25%	51.4%	+12.4 pts

Table 2: Training set size dependency showing diversity benefits emerge with sufficient data

This suggests a threshold effect: diversity-based sampling appears to require approximately 2-3k training examples to be effective, with benefits potentially scaling beyond this threshold.

## 5.3 Quality Control Discovery

Initial experiments on NuminaMath revealed that diverse sampling can systematically select low-quality examples, leading to performance degradation rather than improvement.

Sampling Method	Trainable Examples	Performance
Random	91.5%	20.0%
Diverse (no quality control)	80.6%	12.5%

Table 3: Diverse sampling selects more untrainable examples from noisy datasets

This finding motivated the development of quality assessment frameworks and suggested that diversity benefits emerge more readily when applied to high-quality datasets.

# 6 Discussion

## 6.1 Key Findings

Our experiments suggest several insights for applying diversification methods to mathematical reasoning tasks:

- Quality Control Considerations:** Diverse sampling can select low-quality examples when applied to noisy datasets, potentially leading to performance degradation rather than improvement.
- Optimal Diversity Level:** Both GSM8K and MATH benchmarks show optimal performance at 25% diversity level in our experiments, suggesting this may be an effective configuration across different mathematical reasoning tasks.
- Training Set Size Dependency:** Diversity benefits depend strongly on training set size, with 3k examples enabling substantial improvements (+12 percentage points) while 1k examples show minimal benefits.
- Non-Monotonic Relationship:** The relationship between diversity level and performance is non-monotonic, with intermediate diversity levels outperforming both random sampling and maximum diversity approaches.

5. **Method Effectiveness:** In our experiments, Farthest Point Sampling outperformed clustering-based approaches, possibly due to the lack of clear cluster structure in mathematical problem embeddings.

## 6.2 Future Work

Our findings suggest several concrete directions for extending this research:

Comprehensive reviews have analyzed dataset diversification at scale [21], instruction-tuned programmatic reasoning [1], generalizability perspectives [23], and sampling advancements [19], and new benchmarks continue to emerge.

1. **Training Set Size Scaling:** Our experiments reveal a threshold between 1k and 3k examples where diversity benefits emerge. Systematic evaluation across intermediate sizes (1.5k, 2k, 2.5k) would help characterize this transition and establish practical guidelines for minimum dataset requirements.
2. **Larger Scale Validation:** Testing whether the 25% optimal diversity level holds with larger training sets (5k, 10k+ examples) would determine if this represents a fundamental property or if optimal levels shift with scale.
3. **Mechanistic Understanding:** The non-monotonic relationship between diversity and performance warrants investigation. Analyzing which types of problems are selected at different diversity levels could reveal why intermediate diversity outperforms maximum diversity.
4. **Domain Generalization:** Our 25% optimal finding should be tested on other reasoning domains (code generation, logical reasoning, scientific problem-solving) to determine the generalizability of this configuration beyond mathematical reasoning.

## 7 Limitations

This work has several limitations that should be addressed in future research:

1. **Limited Domain Scope:** Experiments focused specifically on mathematical reasoning tasks; generalizability to other domains requires further investigation.
2. **Model Size Constraints:** Computational limitations restricted experiments to smaller models (1.5B-7B parameters); larger models might show different diversification patterns.
3. **Embedding Quality:** The effectiveness of diversification methods depends on MathBERT embeddings [18], which may not capture all relevant mathematical reasoning features.

## 8 Conclusion

This work investigates dataset diversification methods for mathematical reasoning tasks, examining how different diversity levels affect model performance. Our experiments suggest three main findings: (1) 25% diversity performed optimally in our experiments across GSM8K and MATH benchmarks, yielding improvements of 12+ percentage points over random sampling, (2) diversity benefits appear to depend on training set size, with minimal effects observed at 1k examples but substantial improvements at 3k examples, and (3) diverse sampling may select lower-quality examples from noisy datasets, suggesting the importance of data quality considerations. These results indicate that diversification methods may offer benefits for mathematical reasoning tasks when applied appropriately. The consistency of the 25% diversity result across both benchmarks suggests this configuration may be worth exploring in similar contexts, though further validation would be valuable. Our findings suggest practical considerations for similar work: diversity-based sampling may be most beneficial for training sets above 2k examples, and intermediate diversity levels may offer advantages over both random sampling and maximum diversity approaches. However, these patterns should be validated across different models, datasets, and domains before drawing broader conclusions.

## References

- [1] Anonymous. Sprint: Scaling programmatic reasoning for instruction tuning in math. *Neuro-computing*, 2025.
- [2] David Arthur and Sergei Vassilvitskii. K-means++: The advantages of careful seeding. In *ACM-SIAM SODA*, 2007.
- [3] K. Brown and et al. Synthesizing mathematical datasets for enhanced llm performance: A case study. In *International Conference on Machine Learning*, 2025.

- [4] Ricardo J. G. B. Campello, Davoud Moulavi, and Jörg Sander. Density-based clustering based on hierarchical density estimates. In *Advances in Knowledge Discovery and Data Mining*, pages 160–172, 2013.
- [5] Sirui Chen, Wenhao Zheng, Ziyu Jiang, and et al. Arrows of math reasoning data synthesis for large language models: Diversity, complexity and correctness. *arXiv:2508.18824*, 2025.
- [6] Web Scraping Club. Ensuring data quality in web scraping projects. *Web Scraping Club Blog*, 2023.
- [7] Karl Cobbe, Vineet Kosaraju, Markus Bavarian, and et al. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- [8] Berkan Dokmeci, Qing Wu, Ben Athiwaratkun, and et al. Data diversification methods in alignment enhance math performance in llms. *arXiv:2507.02173*, 2025.
- [9] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996.
- [10] Jingxuan Fan, Tianyu Liu, Yuxuan Wang, and et al. Hardmath: A benchmark dataset for challenging problems in applied mathematics. In *NeurIPS Datasets and Benchmarks*, 2024.
- [11] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, and et al. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- [12] Y. Li, Y. Liang, Y. Wu, and et al. Data augmentation for mathematical reasoning: A comprehensive review. *arXiv:2405.xxxx*, 2024.
- [13] Yanjie Liu and Xinyuan Yu. Farthest point sampling in property designated chemical feature space as a general strategy for enhancing machine learning model performance. *arXiv:2404.11348*, 2024.
- [14] Stuart Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 1982.
- [15] L. McInnes and J. Healy. Accelerated hierarchical density based clustering. In *IEEE ICDMW*, pages 33–42, 2017.
- [16] Yash Mehta and Karthik Seetharaman. Mathematical reasoning through llm finetuning. *Stanford CS224N*, 2024.
- [17] OpenAI. Massive supervised fine-tuning experiments reveal how data composition affects model performance. *arXiv:2506.14681*, 2024.
- [18] Shuai Peng, Ke Yuan, Liqiang Gao, and Zichao Tang. Mathbert: A pre-trained model for mathematical formula understanding. *arXiv:2105.00377*, 2021.
- [19] Konstantin Rusch, David Dao, and Yizhe Zhang. How ai is improving simulations with smarter sampling techniques. *MIT News*, 2024.
- [20] J. Smith and et al. Evaluating the impact of dataset size and diversity on few-shot mathematical problem solving. In *Conference on Artificial Intelligence and Mathematical Reasoning*, 2024.
- [21] Abaka AI Team. Best datasets for math in 2025. *Abaka AI Blog*, 2024.
- [22] Anomalo Team. Data quality in machine learning: Best practices and techniques. *Anomalo Blog*, 2024.
- [23] Curate ND Team. Exploring human-like mathematical reasoning: Perspectives on generalizability and efficiency. *Curate ND*, 2024.
- [24] MathOdyssey Team. Mathodyssey: Benchmarking mathematical problem-solving skills. *Nature Scientific Data*, 2025.
- [25] NuminaMath Team. Numinamath: The largest public dataset in ai4maths with 860k pairs. *AI4Maths Dataset*, 2024.

- [26] Nvlabs Research Team. Prismatic synthesis: Gradient-based data diversification for reasoning models. Nvlabs Research Blog, 2024.
- [27] Corey Wang, Jiajie Zhang, Yue Wu, and et al. Measuring multimodal mathematical reasoning with math-vision dataset. In *NeurIPS*, 2024.
- [28] L. Wang and Y. Zhang. The role of data diversity in training effective mathematical reasoning models. *Journal of AI in Mathematics*, 10:115–130, 2025.
- [29] Jason Wei, Xuezhi Wang, Dale Schuurmans, and et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 2022.
- [30] Lei Yu, Wanjun Jiang, Han Shi, and et al. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv:2309.12284*, 2023.
- [31] Xiaoyu Yue, Xin Qu, Guangsheng Zhang, and et al. Mumath: Multi-perspective data augmentation for mathematical reasoning. *NAACL Findings*, 2024.
- [32] Kaiyu Zheng, J. M. Han, and Stanislas Polu. Minif2f: A cross-system benchmark for formal olympiad-level mathematics. *ICLR*, 2021.