

Synthetic Training Data Enhances Language Model Performance in Mathematical Reasoning

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How does synthetic training data improve language model performance on mathematical reasoning benchmarks v12. 13 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: TerraGen: A Unified Multi-Task Layout Generation Framework for Remote Sensing Data Augmentation. Research question: How does synthetic training data improve language model performance on mathematical reasoning benchmarks v12.

2 Methodology

Systematic literature search across multiple databases yielded 4 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.0/10.

3 Results

4 papers retrieved. 13 claims extracted; 2 independently verified. Quality review score: 5.0/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The benchmark dataset consists of 1k carefully selected high-quality images representing the full spectrum of challenges	×	0.10
The benchmark images are associated with multiple evaluation scenarios, including generation quality assessment across d	×	0.07
Specialized metrics for spatial accuracy and semantic consistency are introduced, encompassing both pixel-level fidelity	×	0.02
RS Image Quality (RS-IQ) is calculated using an InceptionV3 network fine-tuned on remote sensing datasets.	×	0.08
Content Fidelity is measured using CLIP-T for semantic consistency and DINO-I for visual feature alignment.	×	0.00
Layout Consistency is evaluated using YOLOv8-based models, reporting mAP and AP50 for object detection tasks, and Acc an	×	0.05
The model training adopts a two-stage configuration with a learning rate of 1e-4 using AdamW optimizer for 100k steps wi	×	0.02
The second stage introduces an adaptive mask-weighted loss function to enhance layout consistency and spatial accuracy.	×	0.09
TerraGen is compared against remote sensing-specific conditional generation methods and state-of-the-art natural image g	×	0.10
TerraGen achieves an RS-IQ score of 34.6, Content Fidelity scores of 29.6, and Layout Consistency scores of 64.7 for mAP	×	0.02
TerraGen serves as a universal data-augmentation engine, boosting downstream-task accuracy and exhibiting strong general	✓	0.17
The first large-scale multi-task remote sensing layout generation dataset is constructed, establishing standardized eval	✓	0.25
TerraGen integrates spatial layout information with semantic textual information through geographic spatial-aware condit	×	0.09

References

- <http://arxiv.org/abs/2510.21391v1>
- <http://arxiv.org/abs/2307.00161v1>
- <http://arxiv.org/abs/2406.15126v1>