

Graph-Based Arithmetic Word Problem Generation for LLMs

Graph-based approaches for generating arithmetic word problem (AWP) transfer sequences with controllable complexity represent a sophisticated fusion of graph theory, natural language processing, and educational technology that enables systematic creation of high-quality training datasets for large language models. **Recent research demonstrates that leveraging graph structures to model mathematical relationships and control problem complexity can improve LLM performance by up to 62.4% on challenging mathematical reasoning benchmarks like MATH,** (MDPI) matching GPT-4's capabilities while providing systematic curriculum learning progression. (OpenReview)

The significance extends beyond mere performance gains. Graph-based generation enables precise control over problem difficulty through mathematical properties like graph diameter, centrality measures, and structural complexity, (TutorialsPoint) allowing researchers to create targeted training sequences that address specific reasoning deficits. This approach addresses the critical challenge of creating diverse, solvable mathematical problems at scale while maintaining educational validity and logical consistency.

Graph architectures unlock precise mathematical reasoning control

The mathematical foundation rests on five primary graph structures, each optimized for different aspects of transfer sequence modeling. **Directed Acyclic Graphs (DAGs) excel at representing multi-step transfer chains where nodes represent entities and edges capture transfer relationships while preserving temporal dependencies.** (ACL Anthology) (aclanthology) The GraphMR framework from EMNLP 2021 demonstrates this approach, using DAGs to model mathematical expressions where interior nodes represent operators and exterior nodes represent variables, achieving significant improvements over sequence-to-sequence baselines. (ACL Anthology)

Bipartite graphs naturally model agent-object relationships in transfer problems, where one node set represents agents (people, entities) and another represents objects (money, items). (Brilliant) (ScienceDirect) This structure proves particularly effective for assignment problems and resource allocation scenarios common in elementary mathematics. (Study.com) **Temporal graphs extend this capability by incorporating time-sequenced transfers through edges with timestamps, enabling modeling of complex multi-step problems where order matters.** (arXiv)

Tree structures provide hierarchical organization of transfer patterns, with the Graph2Tree architecture demonstrating superior performance on Math23K and similar benchmarks. (ACL Anthology) (ACL Anthology) Flow networks complete the architectural spectrum by modeling complex multi-agent transfers with capacity constraints, using algorithms like Ford-Fulkerson for maximum flow computation and minimum cost flow for optimization problems. (Study.com +3)

Complexity control through mathematical graph properties

The breakthrough in controllable complexity generation lies in establishing quantitative relationships between graph-theoretic properties and problem difficulty. **Research from MDPI Mathematics (2020)**

reveals that graph complexity measures correlate directly with reading comprehension difficulty, with median reading times increasing from 5.12 to 6.68 seconds per word when transitioning from simple to complex graph structures. [MDPI](#) [Journal of Cloud Computing](#)

Graph diameter emerges as the most critical complexity indicator, with diameter values exceeding 3 leading to exponential increases in solving difficulty. [Springer](#) The mathematical relationship follows the formula: $\text{ComplexityScore}(G) = \alpha \cdot \text{Diameter}(G) + \beta \cdot \text{Density}(G) + \gamma \cdot \text{AvgBranching}(G) + \delta \cdot \text{CycleCount}(G)$, where typical weight distributions are $\alpha = 0.3$, $\beta = 0.2$, $\gamma = 0.3$, $\delta = 0.2$.

Centrality measures provide additional granular control. Degree centrality quantifies direct connections per node, while betweenness centrality measures how frequently a node appears on shortest paths between other nodes. [TutorialsPoint](#) **High betweenness centrality correlates with increased multi-step reasoning complexity,** [Wikipedia](#) [ResearchGate](#) enabling systematic generation of problems requiring specific reasoning depths.

Graph density and clustering coefficients offer complementary complexity dimensions. Optimal clustering coefficients for educational problems range from 0.3-0.5, balancing structural coherence with cognitive tractability. Dense graphs with clustering coefficients exceeding 0.6 create cognitive overload, while sparse graphs below 0.3 may lack sufficient structural complexity for meaningful learning.

Systematic generation methodologies ensure dataset quality

Large-scale dataset generation requires sophisticated methodologies balancing diversity, quality, and educational validity. The **TemplateGSM framework demonstrates industrial-scale generation capability, creating over 7 million problems from 7,473 unique meta-templates using GPT-4**, with systematic parameter variation ensuring comprehensive coverage of mathematical scenarios.

The generation process follows a multi-stage pipeline: template creation, systematic parameter sampling, solvability verification, and quality assessment. Template-based approaches provide controllable diversity while maintaining mathematical consistency. [arXiv](#) The "answer-first" methodology proves particularly effective, sampling target answers before generating questions to guarantee solvability.

[Seriousgamessociety +3](#)

Validation pipelines incorporate multiple verification layers. Code execution ensures mathematical correctness, while constraint checking validates logical consistency. [arXiv](#) [MDPI](#) **Human expert evaluation provides final quality assurance, with graduate-level mathematicians reviewing generated problems for educational appropriateness and mathematical rigor.** [TU Delft Research Portal](#)

[CoLab](#)

Reject sampling maintains quality standards by discarding invalid or inappropriate problem instances. [arXiv](#) [MDPI](#) The MaKE framework uses conditional variational autoencoders for complexity control, enabling precise targeting of specific difficulty levels through latent space manipulation. [ACL Anthology +2](#)

Training applications demonstrate significant performance gains

Graph-generated problems enable sophisticated training strategies that leverage structural information for enhanced learning. **The WISDOM curriculum learning approach achieves remarkable results, with WISDOM-7B reaching 62.4% accuracy on the MATH dataset, matching GPT-4's performance through progressive exposure to increasingly challenging graph-structured problems.** [OpenReview](#)

[MDPI](#)

Multi-task learning across different graph types provides additional benefits. The Template-based Multi-Task Deep Neural Networks (T-MTDNN) achieve 85.2% accuracy on MAWPS and 85.3% on Math23k by training across diverse graph structures simultaneously. [ACL Anthology +2](#) Graph-Aware Language Model pre-training (GaLM) demonstrates how single models can serve multiple mathematical reasoning applications through unified graph-text representations.

Robustness testing reveals critical limitations in current approaches. **The Math-RoB benchmark identifies four key robustness issues: positional bias (48% accuracy drop in smaller models), instruction sensitivity (5.0-7.5% performance decline), numerical fragility (GPT-4o dropping from 97.5% to 82.5% accuracy with symbol substitution), and memory dependence (over 50% inappropriate information completion rates).** [MDPI](#)

These findings highlight the continued reliance on pattern matching rather than genuine logical reasoning, even in state-of-the-art models. Graph-based generation helps identify these failure modes through systematic perturbation testing and edge case construction.

Implementation considerations and practical deployment

Successful implementation requires careful attention to computational efficiency and scalability. [Medium](#) **NetworkX provides comprehensive Python-based graph analysis capabilities, while distributed generation systems enable creation of millions of problems with consistent quality.** [github](#) [Journal of Cloud Computing](#) The Graph2Tree implementation demonstrates practical deployment with significant improvements over baseline approaches. [ACL Anthology +2](#)

Quality control mechanisms remain paramount. Expert validation ensures educational appropriateness, while automated verification systems check mathematical consistency. The integration of multiple validation stages—generation, verification, quality assessment, and filtering—ensures high-quality outputs suitable for LLM training. [arXiv +2](#)

Storage and processing considerations become critical at scale. Compressed template storage and streaming generation enable efficient resource utilization, while cached computations reduce redundant processing across similar problem structures.

Future directions and emerging opportunities

Graph-based AWP generation stands poised for significant advancement through several emerging

research directions. **Multi-modal integration combining textual, visual, and symbolic mathematical representations promises more comprehensive reasoning capabilities**, (ACL Anthology +2) while cross-cultural adaptation could enable globally relevant mathematical training datasets.

Theoretical understanding of why graph-based approaches prove effective requires deeper investigation. Current success appears linked to the alignment between graph structures and human cognitive processing of mathematical relationships, but formal characterization remains incomplete.

The integration of symbolic reasoning with neural methods represents another promising frontier. Hybrid approaches combining graph-theoretic precision with neural flexibility could address current robustness limitations while maintaining scalability advantages. (arXiv)

Conclusions and strategic implications

Graph-based approaches for generating arithmetic word problem transfer sequences represent a mature yet evolving field with demonstrated practical value for LLM training. **The convergence of graph theory, complexity control mechanisms, and scalable generation methodologies enables systematic creation of high-quality mathematical reasoning datasets that significantly improve model performance while providing unprecedented control over learning progression.** (Medium +4)

The research reveals that successful implementation requires attention to multiple dimensions: graph structure selection, complexity control through mathematical properties, systematic generation methodologies, and robust evaluation frameworks. (Number Analytics) Current approaches demonstrate clear benefits but also expose fundamental limitations in model reasoning capabilities that future research must address.

The strategic implications extend beyond mathematical reasoning to any domain requiring structured problem generation with controllable complexity. The methodological frameworks developed for AWP generation provide templates for systematic dataset creation across diverse reasoning domains, positioning graph-based approaches as a foundational technology for advanced AI training systems.