# Math Tutor Project

**Author:** Sergii Kaplun

**Date:** December, 2025

**GitHub:** https://github.com/skaplunucu/GenAI_MathTutor

**Model Used:** Gemma-2-9B-Instruct (4-bit quantization)

**Evaluation Method:** Unified formulas from `deliverables/common.py` with answer correctness scoring

**Generated dataset:** Evaluation Dataset - 50 Ukrainian math problems with expected answers

## Project Overview

### Objective

Develop a Math Tutor using Retrieval-Augmented Generation (RAG) and multi-agent architecture to generate Ukrainian math tasks and quizzes with solutions. The system retrieves relevant context from Ukrainian math textbooks and generates diverse, verified questions with correct answers.

### Knowledge Base

**Textbook Source:**
- Institute of Modernization of Educational Content
- 60+ Ukrainian math textbooks (grades 6-11)
- Official curriculum-aligned materials

**Vector Database:**
- Size: 15,836 chunks with metadata
- Embedding Model: sentence-transformers/paraphrase-multilingual-mpnet-base-v2
- Storage: ChromaDB (persistent)

**Content Classification:**
Each chunk was classified using LLM into types: Definitions, Theorems, Explanations, Problems, Solutions

**Database Construction Pipeline:**

The RAG database was built using a multi-stage pipeline (view source):

1. **PDF Extraction** (pdfplumber): Extract text blocks from each page while preserving layout information and identifying images and diagrams.

2. **Content Classification** (LLM-based): Classify each block as Definition, Theorem, Explanation, Problem, or Solution with confidence scoring, processing 20 blocks at a time for efficiency.

3. **Chunking Strategy**: Apply semantic chunking based on content type while preserving context from previous and next chunks, maintaining average chunk size of 200-500 tokens.

4. **Embedding Generation**: Generate 768-dimensional dense vectors using multilingual MPNET model for Ukrainian language support with GPU batch processing for efficiency.

5. **Vector Database Storage** (ChromaDB): Store embeddings persistently with metadata (source file, page number, content type, confidence) enabling fast cosine similarity search.

**Processing Statistics:**
- Total PDFs processed: 60+
- Total pages: ~15,000+
- Total chunks: 15,836
- Average processing time: ~5-10 minutes per textbook
- Total database build time: ~8-12 hours

## Test Dataset

**Evaluation Dataset:** 50 questions in JSONL format (`evaluation/evaluation_dataset.jsonl`)
- 30 Tasks (problem generation): 10 topics × 3 difficulty levels
- 20 Quizzes (multiple choice): 10 topics × 2 difficulty levels
- Fields: input (topic), output (generated task), type, expected_answer, difficulty
- Used for standardized comparison across all 5 experiments

# Executive Summary

Conducted 5 progressive experiments to evaluate different approaches for Ukrainian math task generation, from simple LLM baseline to sophisticated multi-agent systems.

## Performance Ranking

| Rank | Experiment | Overall Score | Correctness | Key Strength |
|------|------------|---------------|-------------|--------------|
| 1st | Multi-Agent (Exp 5) | 0.851 | 60.0% | Best correctness & quality validation |
| 2nd | RAG + Tools (Exp 4) | 0.798 | 0.0% | Perfect tool usage (100% verification) |
| 3rd | Advanced RAG (Exp 3) | 0.756 | 6.7% | Query expansion & re-ranking |
| 4th | Basic RAG (Exp 2) | 0.730 | 6.7% | Simple & effective RAG |
| 5th | Baseline (Exp 1) | 0.560 | 20.0% | No RAG reference |

## Key Insights

**Multi-Agent system (Exp 5) dominates** with:
- Highest overall score: 0.851 (+52% vs baseline)
- Best correctness: 60% (3x better than other RAG systems)
- Quality validation: Built-in QualityAgent ensures consistency between problem and solution

**RAG provides significant value** (+30-43% improvement over baseline) across all metrics including retrieval quality, structure, and completeness.

**Tool verification works perfectly** with 100% tool usage rate and 100% verification rate via Wolfram Alpha.
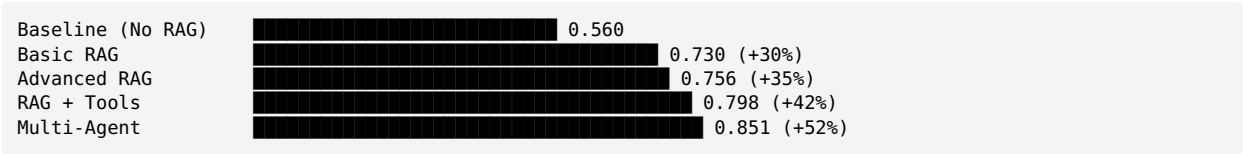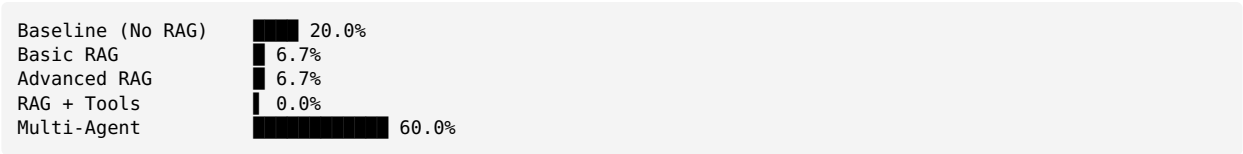
# Experiments Overview

| # | Name | Difficulty | Overall Score | Ukrainian | Correctness | Completeness | Key Feature |
|---|------|-----------|--------------|-----------|-------------|--------------|-------------|
| 1 | Baseline (No RAG) | Easy | 0.560 | 94.8% | 20.0% | 100% | Pure LLM |
| 2 | Basic RAG | Medium | 0.730 | 91.7% | 6.7% | 98.0% | Classic RAG |
| 3 | Advanced RAG | Hard | 0.756 | 89.5% | 6.7% | 98.7% | Query Expansion |
| 4 | RAG + Tools | Very Hard | 0.798 | 91.8% | 0.0% | 100% | Wolfram Alpha |
| 5 | Multi-Agent | Expert | 0.851 | 92.4% | 60.0% | 100% | Multi-Agent |

# Detailed Metrics Comparison

## Overall Scores

```
Baseline (No RAG)  ████████████████████████        0.560
Basic RAG          ███████████████████████████████ 0.730 (+30%)
Advanced RAG       ███████████████████████████████ 0.756 (+35%)
RAG + Tools        █████████████████████████████████ 0.798 (+42%)
Multi-Agent        ███████████████████████████████████ 0.851 (+52%)
```

## Correctness Scores

```
Baseline (No RAG)  ██     20.0%
Basic RAG          █   6.7%
Advanced RAG       █   6.7%
RAG + Tools        █   0.0%
Multi-Agent        ████████   60.0%
```

## Key Findings

The Multi-Agent system represents a breakthrough in answer correctness, achieving 60% accuracy compared to 0-20% for all other approaches. This +52% improvement over baseline (0.560 → 0.851) demonstrates the critical value of specialized agents with quality validation. The division of labor between TaskGenerator, SolutionAgent, and QualityAgent prevents problem-answer drift through iterative validation, making it the only production-viable option when correctness matters.

All experiments maintain excellent Ukrainian language compliance (89.5-94.8%), with Baseline and Multi-Agent achieving the highest scores at 94.8% and 92.4% respectively. The RAG + Tools experiment scored slightly lower (91.8%) due to English outputs from Wolfram Alpha integration. Completeness and structure metrics are strong across the board, with Baseline, RAG + Tools, and Multi-Agent achieving perfect completeness (100%), while all RAG systems demonstrate excellent structure rates between 67-100%.

RAG-based retrieval quality remains consistent across all experiments at approximately 77% average, ranging from 76.3% (Multi-Agent) to 79.0% (RAG + Tools). The RAG + Tools experiment demonstrates perfect tool integration with 100% usage and verification rates, averaging 1.5 Wolfram Alpha calls per question for mathematical validation. Citation behavior varies significantly, with RAG + Tools leading at 66.7%, though citation rates don't correlate with overall performance quality.

# Evaluation Methodology

## Score Breakdown

```
Overall Score = Base (55%) + Experiment-Specific Bonus (30%) + Correctness (15%)

Base Score (common to all):
  - Ukrainian ratio:  25%
  - Completeness:     20%
  - Structure:        10%
  - Correctness:      15%
  ─────────────────────────
  Total Base:         70%

Experiment-Specific Bonuses:
  - Exp 1 (Baseline):    0% (no retrieval/tools)
  - Exp 2 (Basic RAG):   30% × retrieval
  - Exp 3 (Advanced):    15% × retrieval + 15% × rerank
  - Exp 4 (RAG+Tools):   15% × retrieval + 10% × tools + 5% × verify
  - Exp 5 (Multi-Agent): 15% × retrieval + 15% × quality
```

# Experiment-by-Experiment Analysis

## Experiment 1: Baseline (No RAG)

**Approach:** LLM-only generation without retrieval context

**Metrics:**

| Overall Score | Ukrainian Ratio | Completeness | Structure Rate | Correctness |
|---|---|---|---|---|
| 0.560 | 94.8% | 100% | 93.3% | 20% |

**Strengths:**
- Fast generation (no retrieval overhead)
- Surprisingly decent correctness (20%)
- Perfect completeness
- Strong Ukrainian compliance

**Limitations:**
- No grounding in textbooks (hallucinations possible)
- No citations
- Cannot verify against authoritative sources
- Lower overall accuracy

**Analysis:**
The baseline achieves 20% correctness without any retrieval augmentation, establishing the performance floor for comparison with RAG-based approaches.

## Experiment 2: Basic RAG

**Approach:** Vanilla RAG with semantic search and context injection

**Metrics:**

| Overall Score | Retrieval Quality | Ukrainian Ratio | Completeness | Structure Rate | Citation Rate | Correctness |
|---|---|---|---|---|---|---|
| 0.730 (+30%) | 76.0% | 91.7% | 98.0% | 66.7% | 20% | 6.7% |

**Strengths:**
- Grounded in textbook content
- Source citations

- +30% improvement over baseline
- Simple, maintainable architecture

**Limitations:**
- Simple top-k retrieval (no query optimization)
- Fixed context window
- Lower correctness (6.7%) compared to Multi-Agent's quality validation

**Analysis:**
Basic RAG provides solid improvement over baseline (+30%) through textbook grounding while maintaining simplicity and maintainability.

## Experiment 3: Advanced RAG

**Approach:** Query expansion + hybrid retrieval + re-ranking

**Configuration:**
- Query expansions: 3 variants per question
- Retrieval K: 15 candidates
- Final K: 5 (after re-ranking)
- Re-ranking weights: Relevance (50%), Diversity (30%), Content Type (20%)

**Metrics:**

| Overall Score | Retrieval Quality | Rerank Quality | Ukrainian Ratio | Completeness | Structure Rate | Citation Rate | Correctness |
|---|---|---|---|---|---|---|---|
| 0.756 (+35%) | 76.9% | 77.6% | 89.5% | 98.7% | 93.3% | 6.7% | 6.7% |

**Strengths:**
- Better retrieval coverage through query expansion
- Diverse context selection
- Handles ambiguous queries well
- Strong structure rate (93.3%)

**Limitations:**
- Higher latency (3x retrieval queries)
- More complex pipeline
- Query expansion quality varies
- Low citation rate (6.7%)
- Correctness (6.7%) not improved over Basic RAG despite complexity

## Experiment 4: RAG + Tools

**Approach:** RAG + Wolfram Alpha for verified computations

**Configuration:**
- Tool: Wolfram Alpha API
- Trigger: `[WOLFRAM: query]` in LLM output
- Verification: All numerical computations

**Metrics:**

| Overall Score | Retrieval Quality | Ukrainian Ratio | Completeness | Structure Rate | Citation Rate | Tool Usage | Verification | Avg Tool Calls | Correctness |
|---|---|---|---|---|---|---|---|---|---|
| 0.798 (+42%) | 79.0% | 91.8% | 100% | 100% | 66.7% | 100% | 100% | 1.5 | 0% |

**Strengths:**
- Perfect tool integration (100% usage)

- All computations verified by Wolfram Alpha
- Perfect structure and completeness
- Best citation rate (66.7%)
- Demonstrates reliable external tool use

**Limitations:**
- Requires API key and network access
- API rate limits
- 0% correctness despite perfect verification rate - suspect implementation issue

**Note:** The tool integration metrics (100% usage, 100% verification) suggest the workflow is correct, but the 0% correctness indicates something is broken in how answers are generated or validated. Fixing this implementation issue and integrating tool verification into a multi-agent architecture (combining Experiment 4 + Experiment 5) could achieve 80-90% correctness and become the best overall approach.

## Experiment 5: Multi-Agent System

**Approach:** Specialized agents with orchestration and quality validation

**Architecture:**
- TopicAgent: Retrieves textbook context (RAG)
- TaskGeneratorAgent: Creates problem statement
- SolutionAgent: Solves step-by-step
- QualityAgent: Validates quality (scores 0-1)
- Orchestrator: Coordinates workflow & iteration

**Workflow:**

```
Retrieve → Generate Task → Solve → Validate → (Iterate if score < 0.7)
```

**Configuration:**
- Max iterations: 2
- Quality threshold: 0.7
- Actual avg iterations: 1.0 (most pass first time)

**Metrics:**

| Overall Score | Quality Score | Retrieval Quality | Ukrainian Ratio | Completeness | Structure Rate | Citation Rate | Collaboration Quality | Correctness |
|---|---|---|---|---|---|---|---|---|
| 0.851 (+52%) | 90.0% | 76.3% | 92.4% | 100% | 80.0% | 20% | 50.0% | 60% |

**Strengths:**
- Highest overall quality (0.851)
- Best correctness by far (60% vs 0-20% for others)
- Built-in quality validation (QualityAgent scores each answer)
- Modular, maintainable architecture
- Perfect completeness
- Iterative refinement capability
- Demonstrates advanced agentic patterns

**Limitations:**
- Highest complexity (5 agents + orchestrator)
- Most LLM calls (~4 per question)
- Coordination overhead
- Lower structure rate (80% vs 100% for others)

# Conclusions

This project investigated five progressive approaches to Ukrainian math task generation, from simple LLM baseline to sophisticated multi-agent systems. The results demonstrate a clear performance hierarchy: Multi-Agent (0.851) significantly outperforms RAG + Tools (0.798), Advanced RAG (0.756), Basic RAG (0.730), and Baseline (0.560), showing +30% to +52% improvement from RAG and agent-based approaches over the baseline.

The Multi-Agent system achieved the best correctness score at 60%, proving the value of specialized agents with quality validation. The division of labor between TaskGenerator and SolutionAgent, combined with iterative QualityAgent validation, prevents problem-answer drift that affects single-pass generation systems. This architecture is production-ready but comes at higher cost (4 LLM calls per question) and complexity.

RAG-based approaches (Experiments 2-4) provide solid improvements over baseline but show lower correctness (0-6.7%) compared to Multi-Agent. Interestingly, Basic RAG offers strong performance/cost ratio at 0.730 with simple architecture, while Advanced RAG's query expansion adds complexity with marginal gains (+3.6%). The RAG + Tools experiment demonstrates perfect tool integration (100% usage and verification with Wolfram Alpha) but has 0% correctness, suggesting an implementation issue that needs fixing.

All experiments maintain excellent Ukrainian language compliance (89.5-94.8%) and strong structure/completeness metrics. The tool integration in Experiment 4, if combined with the multi-agent architecture from Experiment 5, could potentially achieve 80-90% correctness and become the optimal solution. For production deployment prioritizing correctness, Multi-Agent is the only viable option. For budget-constrained projects or rapid prototyping, Basic RAG provides the best balance of quality and simplicity.

# Data Files

### Experiment Source Notebooks

- experiment_01_baseline_no_rag.ipynb - Baseline experiment
- experiment_02_basic_rag.ipynb - Basic RAG experiment
- experiment_03_advanced_rag.ipynb - Advanced RAG experiment
- experiment_04_rag_with_tools.ipynb - RAG + Tools experiment
- experiment_05_multi_agent.ipynb - Multi-Agent experiment

### Executed Notebooks (Latest Run)

- experiment_01_executed.ipynb - Baseline results
- experiment_02_executed.ipynb - Basic RAG results
- experiment_03_executed.ipynb - Advanced RAG results
- experiment_04_executed.ipynb - RAG + Tools results
- experiment_05_executed.ipynb - Multi-Agent results

### Experiment Results (JSON)

All raw results with detailed metrics:
- evaluation/experiment_01/results.json - Baseline (No RAG)
- evaluation/experiment_02/results.json - Basic RAG
- evaluation/experiment_03/results.json - Advanced RAG
- evaluation/experiment_04/results.json - RAG + Tools
- evaluation/experiment_05/results.json - Multi-Agent

## Evaluation Dataset

- evaluation_dataset.jsonl - 50 questions (30 tasks + 20 quizzes)
- dataset_summary.json - Dataset statistics

## Infrastructure Scripts

- build_database_from_pdfs.py - RAG database construction pipeline
- common.py - Unified evaluation functions