

# An Automated Graphing System for Mathematical Pedagogy

Arya Bulusu, Brandon Man, Ashish Jagmohan, Aditya Vempaty, Jennifer Mari-Wyka



MERLYN MIND™

## Introduction

- Generative AI has significant potential to simplify the tools available to teachers in the classroom
- By creating automated tool-using systems, we allow teachers to more easily control tools such as graphing software
- We present an automated graphing system which converts utterances into mathematical expressions and graphs them with the Desmos interface
- Our design incorporates a mathematical solver into an LLM system, leading to more accurate, consistent responses
- The LLM allows for natural language understanding while the mathematical solver allows for mathematical accuracy

## Dataset

Based on the Common Core standards, we identify a set of learning objectives that can be fulfilled through graphing. These categories form the basis of our two evaluation datasets. The utterance-focused dataset contains simple, single-step commands teachers would use in the classroom to demonstrate intermediate steps. The textbook-focused dataset contains multi-step, complicated problems requiring tool use to solve.

### Example Row from Utterance-Focused Dataset

Natural Language Utterance	Graph Input
Reflect $y$ equals five $x$ minus four across the $y$ -axis	$y = -5x - 4$

## Autoevaluator

Traditional evaluation metrics for text similarity cannot judge mathematical equivalence, so we need a method to precisely compare mathematical statements. We use the computer algebra system SymPy to compare LLM-produced mathematical expressions to our ground-truth expressions, with an LLM as a backup if SymPy cannot parse a given expression. We compare the results of the LLM-only and LLM+SymPy autoevaluators and find that the addition of SymPy significantly increases the accuracy of evaluations.

## Results

We compare the performance of the LLM+Solver system and the LLM-only system, using GPT-4 as the LLM and Wolfram Alpha as the solver for evaluation.

### Accuracy of LLM-only and LLM+Solver Systems

	LLM-only	LLM+Solver
Utterance-Focused Dataset	63%	<b>85%</b>
Textbook-Focused Dataset	55%	<b>75%</b>

The addition of the solver results in a significant performance increase on both datasets. We see the greatest performance increase in categories with difficult, multi-step problems, such as Local Minima and Maxima and Tangents to Circles. The LLM-only system cannot solve these problems as it lacks mathematical reasoning. However, it is easy to write a Wolfram Alpha query to solve these problems, so the LLM+Solver consistently produces accurate answers.

