
Prompt-Driven Verification of LLM-Generated Mathematical Solutions

Hengzhi Zhang
hz2663@columbia.edu

Abstract

We propose a self-verification pipeline that transforms an LLM’s chain-of-thought into an explicit proof object, enabling the model to audit its own reasoning with only lightweight calls to itself (or an even smaller variant)—no external solver or larger teacher is required. Our method decomposes multi-step mathematical solutions into five modular stages: (1) solution generation via Vertex AI’s Gemini backends, (2) sentence-level segmentation to isolate individual expressions, (3) rule attribution, where each derivation step is labeled with the specific theorem or identity applied, (4) application checking, in which the model verifies each step’s logical validity against its stated premises, and (5) error localization via majority voting over multiple trials to pinpoint the earliest incorrect sentence. On an initial sample of ten hard, level-5 MATH problems, our pipeline localizes the first error with high accuracy while providing compact, JSON-structured justifications for each ruling. This traceable, data-rich output not only aids human reviewers but also yields labeled triples (premises, rule, conclusion) that can supervise future theorem-rationale learning. Our results demonstrate that modeling mathematical reasoning as explicit rule applications empowers LLMs to self-inspect and correct their own proofs, laying groundwork for more reliable, reflective reasoning agents.

1 Introduction

Large language models (LLMs) now routinely derive multi-step solutions to competition-level mathematics problems, yet they still hallucinate steps, mis-apply algebraic identities, or lose track of definitions. We present a *self-verification pipeline* that turns the model’s own chain-of-thought into an *object of reasoning*: every generated sentence is treated as a candidate statement that must be justified by explicit applications of previously accepted sentences and general mathematical rules. The pipeline shows that an LLM can *audit itself* with only small, inexpensive calls to the same (or an even smaller) model—no external symbolic solver or larger teacher model is required.

2 Related Work

Chain-of-Thought Prompting Chain-of-Thought (CoT) prompting was first introduced by Wei et al. [?], who showed that guiding large language models to generate intermediate reasoning steps markedly improves performance on multi-step arithmetic and commonsense problems. Wang et al. [?] later proposed the *self-consistency* strategy, in which multiple CoT samples are generated and aggregated by majority-vote to further boost answer accuracy.

Self-Verification and Error Detection A line of work has explored using LLMs to verify their own reasoning. Weng et al. [?] introduced backward verification of CoT outputs, assigning interpretable validation scores to each reasoning chain and demonstrating gains on arithmetic, logical, and commonsense benchmarks. Building on this, Miao et al. [?] proposed *SelfCheck*, a zero-shot

schema in which an LLM flags errors in its step-by-step reasoning and combines multiple checks via weighted voting to boost final accuracy on GSM8K, MathQA, and the MATH dataset [?]. More recently, Chowdhury and Caragea [?] designed COT-STEP prompts to decompose reasoning, plus a zero-shot verifier that classifies the correctness of each step without any fine-tuning.

Deductive and Dataset-Driven Verification Ling et al. [?] introduced the *Natural Program* formalism for rigorous, stepwise deductive checking of CoT outputs, providing a template for exact logical verification. Hong et al. [?] offered an in-depth empirical study of LLMs’ self-verification skills, constructing the *Fallacies* dataset to benchmark models’ abilities to detect logical missteps in their own reasoning chains. Beyond algorithmic advances, Chen et al. [?] present a broad survey of CoT methods and long-form reasoning strategies, situating verification-guided pipelines within the evolving *Reasoning Era* of large language models.

Theorem Rationale Learning Sheng et al. [?] introduce *Theorem Rationale* (TR), a framework that explicitly guides LLMs to select and apply the most pertinent mathematical theorems when solving complex problems. They construct a dedicated dataset of problem–theorem–solution triples and train models to generate a concise theorem rationale before each derivation step, yielding significant gains on high-school and collegiate mathematics benchmarks in AAAI-25.

3 Pipeline Overview

The code in supplementary material decomposes the task “*point to the first incorrect step in a large-language-generated solution to a MATH problem*” into five well-isolated stages, each implemented by an explicit module (Fig. 1). These stages are executed sequentially by the driver class `VerifyCotTheorems`, though any individual stage can be run in isolation for ablation studies or debugging. More detail could be found in ‘<https://github.com/henryzhang11/verify-cota>’.

Solution generation. The `model.py` wrapper supports two Vertex AI backends—`gemini-2.0-flash-001` by default, and an optional higher quality leader model `gemini-2.5-flash-preview-04-17`. Calls to `VertexAI.generate` use a near-deterministic sampling policy (temperature 0.05, $p = 0.95$, $k = 20$) with an exponential back-off retry mechanism of up to six attempts before failing gracefully. Problems are drawn from the Hendrycks MATH dataset via helpers in `dataloader.py`, which expose methods to load the full training split, a hard (non-geometry level-5) subset, and toy instances for unit tests. The script `find_incorrect_solution.py` streams each problem to the model, extracts the `\boxed{...}` final answer, and records any mismatches against the official key in a JSONL file for downstream verification.

Sentence segmentation. Two prompts—`VerifyCotTheorems.parse_text` for the problem statement and `cleanup_answer` for the generated solution—ask the Gemini model to rewrite text so that every mathematical expression appears within a complete English sentence and each sentence sits on its own line. The model’s response is delimited by back-ticks, split on newlines, and filtered for empties, yielding ordered lists `problem_sentences` and `solution_sentences`. These are concatenated into `all_sentences`, which serves as the indexed proof state for later stages.

Rule attribution. A lightweight classifier in `categorize` tags each solution sentence as either a *derivation* or *other*, using a majority vote over up to seven LLM passes (typically converging in three). For each derivation, the `name_theorem` module re-prompts the model to emit a single rigid JSON quadruplet containing: a concise rule description (e.g. “quadratic formula,” “AM–GM inequality”), the indices of premises and any earlier conclusions reused, and the natural-language form of the claimed conclusion. These quadruplets are cached in `theorems_applied[k]` for sentence k .

Application checking. The `check_application` module prompts the Gemini model with each rule application’s premises, rule name, and claimed conclusion. It returns for each application a Boolean `verdict_i` asserting whether the step is valid under first-order logic and a global “`relation`” label (RESTATE, CONTRADICT, or NEITHER) relative to the original sentence. A strict regular expression extracts the fenced JSON (‘`json ...`’), with malformed or incomplete replies triggering up to five retries before aborting. Parsed results populate `application_correctness[k]` and `application_relevance[k]`.

Error localisation and majority voting. A sentence is deemed sound only if all its rule applications are correct and its final conclusion restates the original sentence. The earliest sentence failing this criterion is reported as the first mistake; if none fail, the proof is certified correct (index -1). To mitigate stochasticity in LLM grading, the entire pipeline (`find_first_mistake`) runs three times, pooling the reported indices and taking a majority vote—ties are broken by choosing the smallest index. The harness in `find_first_mistake.py` processes batches from `clean_sentences_2.txt` and logs aggregate statistics.

Implementation notes Robust logging is provided by `logger_setup.py`, which captures every prompt, response, and intermediate data structure at the DEBUG level while streaming progress at INFO, with automatic rotation above 100000 characters. Numerical answers undergo symbolic normalization in `math_equivalence.py`, ensuring canonical representation of fractions, radicals, and units before string comparison. For oracle baselines, `prepare_baseline_prompt_2.py` generates prompts for feeding into fresh LLMs to locate the first mistake with or without the official answer, quantifying the benefit of our rule-checking decomposition.

4 Reasoning as Rule Applications

A “rule” is any conditional map

$$\text{premises} \longrightarrow \text{conclusion},$$

such as “if $p \rightarrow q$ and p then q ” or “if $a + b = c$ and $a, b \geq 0$ then $c \geq a$ ”.

As a side note, the model frequently do not follow the instruction ‘let’s reason step by step’ if a step is defined by the application of a single theorem. Instead, large portions of the sentences generated by Gemini jumps several steps from the previous sentence, so it’s necessary to ask the checking model to generate a proof instead of pick a single theorem to show the next sentence based on previous sentences.

Forcing the model to name its rules yields:

- **Traceability** – Human reviewers can skim compact JSON objects instead of deciphering free-text proofs.
- **Data synthesis** – Each verified (or rejected) application becomes a labelled triple $\langle \text{premises}, \text{rule}, \text{conclusion}, \text{verdict} \rangle$ that favours future models which *inspect*, not skip, intermediate steps.
- **Modularity** – Additional checkers (e.g. for theorem correctness) can be layered without altering generation.

5 Experimental Setup

We load 850 *level-5* non-geometry problems from the public MATH train split via `load_MATH_hard` (geometry problems are not tested due to lack of proper diagrams as testing data and price of model processing images). After a test generation sweep, the naive expression equivalence checking function provided by the MATH paper flagged 60 problems to be incorrect. The first ten of these problems are used to conduct an initial test of our above pipeline. As a baseline we use **ChatGPT o4-mini** to decide the first faulty sentence followed by *direct human inspection* to ascertain the existence of the mistake.

6 Results

30 problems that have not been used to tune the algorithm are used to test the algorithm (the author tuned the algorithm further on a batch of 20 problems separate from the initial batch of 10 problems mentioned in the paper and updated the pipeline). Among the 30 problems, 28 are non-graph related. Among them the model correctly identifies the first mistake (or lack of any mistake) in 20 problems. The model correctly flags two other solutions as incorrect, missing the first mistake but catching subsequent mistakes. Counting these 2 problems, the model has an accuracy rate of 78%.

The model made serious mistakes when verifying 6 out of 28 problems. In 3 of the problems, the model incorrectly flags correct sentences as incorrect. In 2 others, the model failed to flag incorrect sentences. In the remaining problem, the parsing function malfunctioned, filling in extra repeating sentences.

The pipeline uses a single pass with multiple reattempts at each sentence if the provided validation proof contains misapplications of theorems or doesn't prove/disprove the validated sentence.

7 Discussion

The pipeline could potentially serve as real-time guardrail on problems drawn from domains the model hasn't seen—e.g. novel theorem families, atypical presentation styles, or entirely new subject areas—to test whether self-verification prevents spurious leaps outside the training distribution.

The preliminary sample is encouraging but too small for definitive claims.

8 Conclusion

Modeling mathematical reasoning as an explicit sequence of rule applications lets an LLM audit itself, surface the earliest slip, and provide human-legible justifications. Even on a small public sample, the pipeline reaches an effective 70 % localisation accuracy—without consulting a larger teacher model and while remaining text-only.

References