

# **AI MATH TUTOR**

## **ASSIGNMENT 3 – PROMPT EVALUATION**

---

Jóhannes Reykdal Einarsson · Sævar Breki Snorrason · Sölvi Santos

Course: REI603M — The AI Lifecycle, Tablet-first AI Math Tutor (Icelandic feedback), February 2026,

- Students often get stuck **mid-solution** and need step-level guidance, not final answers.
- Target setting:
  - secondary school algebra & calculus,
  - handwritten work on tablet + stylus.
- Tutor behavior is **explicitly mode-controlled**:
  - **Check Solution**: verify correctness and point out the first error
  - **Hint**: Socratic guidance without revealing the next step
  - **Reveal**: full worked solution only on explicit request
- Assignment goal:
  - build a structured evaluation dataset,
  - iterate prompt designs systematically,
  - analyze failures and define next evaluation steps.

- Dataset design principles:
  - cases grouped by tutoring intent and failure risk,
  - includes normal, incorrect, ambiguous, and leakage-attempt scenarios,
  - each case annotated with expected **verdict** and **allowed behavior**.
- Size: 50 image-based test cases

## Distribution (by mode)

- Hint mode: 23
- Check Solution mode: 22
- Reveal mode: 5

## Expected verdicts

- fully\_solved: 17
- incorrect: 17
- correct\_so\_far: 15
- unclear: 1

- What the dataset explicitly tests:
  - pointing out errors **without** advancing the solution,
  - refusing to continue tutoring after a problem is already solved,
  - handling unreadable or ambiguous handwriting,
  - resisting prompt-injection or answer-leakage attempts.
- Annotation enables deterministic scoring:
  - verdict correctness vs ground truth,
  - feasibility check (response type must match verdict),
  - spot checks for Icelandic language, concision, and pedagogical tone.

- Four prompt versions evaluated:
  - **v1**: baseline zero-shot with role, mode rules, JSON schema
  - **v2**: expanded instruction prompt with stricter reasoning order
  - **v3**: few-shot with curated examples per mode
  - **v4**: refined few-shot + stricter policy wording
- All versions enforced a structured output contract:
  - verdict
  - response\_type
  - message\_is
- Model:
  - Gemini 3 Flash (multimodal preview)
  - schema validated with Pydantic

- For each prompt version (v1–v4):
  - run all 50 dataset images,
  - input = prompt + mode + image,
  - output parsed as strict JSON.
- Tooling:
  - Python scripts (agentic.py, review.py)
  - retries with exponential backoff for transient failures
- Metrics:
  - **Verdict accuracy**: predicted == expected
  - **Non-feasible ratio**: response type violates policy map
- Controlled setup:
  - identical dataset and scoring logic across all runs.

## Quantitative summary

- v1: 96% accuracy, 2% non-feasible
- v2: 94% accuracy, 12% non-feasible
- v3: 94% accuracy, 0% non-feasible
- v4: 96% accuracy, 0% non-feasible

## Key observation

- v4 is the best overall tradeoff (high accuracy + policy consistency).

## Qualitative analysis (LLM-judge, guideline v2)

Prompt	Correctness %	Hint Usefulness	Clarity
v1	90.0	3.87	4.64
v2	78.0	3.57	4.66
v3	94.0	4.30	4.80
v4	96.0	4.35	4.62

## Persistent errors (all versions)

- Ambiguous handwriting misclassified as correct\_so\_far instead of unclear
  - ambiguity detection weaker than algebra checking
- Dense symbolic solutions misclassified as incorrect
  - false negatives on fully solved expressions

## Version-specific failure

- v2 produced incorrect verdicts with hint-style responses
  - direct policy violations in 6 cases

## Lessons

- Instruction clarity beats instruction length
- Few-shot examples anchor behavior better than prose
- Ambiguity needs explicit triggers and more data.

- What worked:
  - structured outputs + explicit response mapping,
  - few-shot prompting (v3/v4).
- What surprised us:
  - longer prompts reduced policy compliance (v2).
- Evaluation gaps:
  - only one explicit unclear case,
  - limited handwriting variability.
- Next iteration:
  - expand ambiguous and OCR-noisy cases,
  - add harder fully-solved symbolic problems,
  - normalize response\_type strings in evaluator,
  - compare against at least one additional model.

- Goal: measure tutoring effectiveness in real sessions.
- Core events:
  - problem\_started, step\_submitted, mode\_requested, llm\_response\_generated, user\_retry, reveal\_clicked, problem\_completed, session\_abandoned.
- Quality signals:
  - valid next step within two actions after a hint,
  - correction after first-error feedback,
  - user ratings (thumbs up/down).
- Failure indicators:
  - repeated retries,
  - frequent clarifications,
  - abandonment after check or hint.
- Privacy-first:
  - hashed user IDs,
  - minimal image retention,
  - metadata-only analytics.

## Failed case 1 (expected unclear, predicted correct\_so\_far)

Leyfði jöfnunum  $(x-2)(x+3) = x^2 - 1$

Líðum vinstrum hljóðum:

$$x^2 + 3x - 2x - 6 = x^2 + x - 6$$

Við er  $x^2 + x - 6 = x^2 - 1$

$$\Rightarrow x - 6 = \cancel{x}$$

## Failed case 2 (expected fully\_solved, predicted incorrect)

b) Skrifði til meðalna  $e^{-\frac{3}{2}i - 2i \ln(2)} e^{\frac{2i(1 - \frac{3}{2})}{i} - 2i \ln(2)}$  á formuna  $x + iy$

Við hafum  $e^{-\frac{3}{2}i - 2i \ln(2)}, \frac{2i(1 - \frac{3}{2})}{i} - 2i \ln(2) \in \mathbb{C}$

$$= e^{-\frac{3}{2}i - 2i \ln(2)}, \frac{2i(-\frac{1}{2})}{i} - 2i \ln(2)$$

$$= e^{-\frac{3}{2}i - 2i \ln(2)}, i(-\frac{3}{2} - \frac{2i(1 - \frac{3}{2})}{i})$$

$$= e^{-\frac{3}{2}i - 2i \ln(2)} \cdot i(-\frac{3}{2} - \frac{2i(-\frac{1}{2})}{i})$$

$$= e^{-\frac{3}{2}i - 2i \ln(2)} \cdot i(-\frac{3}{2} + \frac{2}{i})$$

$$= e^{(\ln(2)^2)} (\cos(-\frac{3}{2}) + i \sin(-\frac{3}{2}))$$

$$= 2^{-3}(0 - i)$$

$$= -i \cdot 2^{-3}$$

## Good case (v4 corrected to incorrect + fix\_first)

17. Finnið afleiðan  $f(x) = -4 \log_2(\cos(11x))$ .

$$\frac{df}{dx} = -4 \frac{1}{\log_2(\cos(11x))} \cdot (-\sin(11x)) \cdot 11$$

**Deployment takeaway:** default to clarification under uncertainty and keep a targeted regression set for ambiguity + dense symbolic algebra.