

# LLMs cannot spot math errors, even when allowed to peek into the solution

University of Artificial Intelligence

KV Aditya Srivatsa, Kaushal Kumar Maurya, Ekaterina Kochmar S<mark>uzhou, China | 中国苏州</mark> Mohamed bin Zayed University of Artificial Intelligence (MBZUAI), Abu Dhabi, UAE

### The space station, Lupus-1, is an enormous spacecraft made up of three identical cylindrical structures that house the living quarters for the crew. The three cylindrical structures are linked together by a series of tunnels that hold them together and allow the crew to move between cylinders. Each cylindrical structure contains 12 bedrooms, 7 bathrooms, and several kitchens. If the entire space station has 72 rooms, how many kitchens are there on the entire space station? << G :: GOLD SOLUTION >> STEP-1: If there are 3 identical cylindrical structures, with a combined total of 72 rooms, then each individual cylindrical structure contains 72/3=24 rooms. STEP-2: In each 24-room cylindrical structure, there are 24-12-7=5 kitchens. STEP-3: Thus, the entire, three-cylinder space station contains a total of 3\*5=15 kitchens. Final Answer: 15 LLM **Poor Step-Level Alignment** Strong Step-Level Alignment << S :: STUDENT SÖLUTION >> << C :: CORRCETED STUDENT SOLUTION >> STEP-1: Each cylindrical structure has 12 + 7 + x = 19 + x rooms, where x is the number of kitchens. STEP-1: Each cylindrical structure has 12 + 7 + x = 19 + x rooms, where x is the number of kitchens. STEP-2:Since there are three cylindrical structures, the total number of rooms is 3(19+x) = 57+3x. STEP-2:Since there are three cylindrical structures, the total number of rooms is 3(19+x) = 57+3x. STEP-3: We know that the total number of rooms is 72, so 57+3x = 72. STEP-3: We know that the total number of rooms is 72, so 57+3x = 72. STEP-4: Subtracting 57 from both sides gives 3x = 15. STEP-4: Subtracting 57 from both sides gives 3x = 15. STEP-5: Dividing both sides by 3 gives x = 5. STEP-5: Dividing both sides by 3 gives x = 5. STEP-6: Since there are 3 identical cylindrical structures, the total number of kitchens is 3x = 3(5) = 15STEP-6: Therefore, there are 5 kitchens on the entire

### **Problem Statement**

Final Answer: 5

Error localization in reasoning chains is a crucial precursor to effective evaluation and automated tutoring.

RQ-1: Can LLMs accurately predict the first error step in an incorrect math solution using the reference solution (w-GS)?

RQ-2: Can localization improve using an intermediate corrected solution (w-Cor)?

### Method

- # Data: LLMs' error localization performance tested on annotated stepwise student solutions from the VtG [1] and PRM800K [2] datasets.
- # Three experimental settings to predict the first error step ID in the student solution:
- (1) w/o-S <Q,S>: Input contains only problem text (Q) and incorrect student solution (S).
- (2) w-GS <Q,S,G>: Input contains Q, S, and a gold annotated reference solution (G).
- (3) w-Cor <Q,S,C>: Input contains Q, S, and a corrected version of S (C) (See below).
- # Corrected Student Solution (C):
- 1. Generated using respective LLMs used for localization later. Prompted to correct the student's solution while sticking to the student's approach.
- 2. Manual evaluation shows most models produce accurate corrections >93% times.

## Results

Final Answer: 15

STEP-7: Therefore, there are 15 kitchens on the entire

- # LLMs cannot localize well:  $\mu$  = 49% & 29%
- **# Even with ref soln: w-GS** > w/o-S for most cases, but **w-GS** remains low (54% & 36%).
- # Intermediate corrections help: (57% & 40%) w-Cor scores highest for most data-LLM pairs.

Model	VtG			PRM800K		
Random		18.32		9.52		
	w/o-S	w-GS	w-Cor	w/o-S	w-GS	w-Cor
Llama3-70B	42.51	49.50	61.28	19.64	24.12	33.03
Llama3.1-70B	49.10	57.98	64.17	24.46	34.23	38.39
Llama3.1-405B	49.90	62.38	64.77	24.12	39.54	47.86
GPT-4o	54.49	63.57	64.57	39.29	43.72	49.40
Qwen2.5-72B-Math	45.01	30.44	19.10	21.86	28.50	21.47
LearnLM-1.5-Pro	54.89	64.07	63.67	42.51	49.69	51.13

Table 1: First error step localization accuracy (in %) on VtG and PRM800K datasets. For each task, within each dataset, the bold value represents the highest accuracy per LLM, whereas the underlined value represents the overall highest accuracy.

## References

[1] Daheim, N., Macina, J., Kapur, M., Gurevych, I., & Sachan, M. (2024). Stepwise verification and remediation of student reasoning errors with large language model tutors. EMNLP 2024. ACL. https://doi.org/10.18653/v1/2024.emnlp-main.478

[2] Lightman H. Kosaraiu V. Burda Y. Edwards H. Baker B. Lee, T., Leike, J., Schulman, J., Sutskever, I., & Cobbe, K. (2023) Let's verify step by step, arXiv preprint arXiv:2305.20050. https://arxiv.org/abs/2305.20050





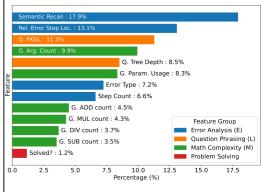


# Analysis

### Q.What impacts error localization the most?

Black-box feature importance study reveals:

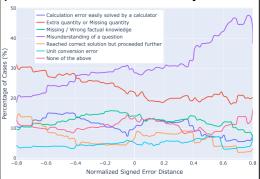
- # Semantic recall b/w the student solution (S) and the reference solution (G or C) has the greatest impact.
- # Questions' phrasing and math complexity rank among the most impactful.
- # LLMs' ability to solve the underlying math problem does not correlate well with their ability to spot errors in an incorrect solution (Chisquared Test: p > 0.01;  $\phi < 0.2$ ).



#### Q. Does the type of error play a role?

Tracking different error types relative to their actual location reveals:

- # Question-independent errors are uniformly distributed, e.g., calculation or unit-conversion mistakes.
- # Errors due to question misunderstanding are predicted much later than they occur.
- # Errors involving missing or extra variables are predicted much before that actually occur.



#### Q.Why does Qwen2.5-72B-Math perform so poorly across settings?

- # Despite being tuned on math problem data, Qwen2.5-Math scores consistently lower than other models in most cases.
- # Manual evaluation reveals that the model struggles to (1) follow instructions to spot first error location and (2) to rectify first error in C.
- # Likely because the model is tuned for problem solving rather than critique or other metareasoning.