# Probing for Arithmetic Errors in LLMs

### Yucheng Sun\*, Alessandro Stolfo\*, Mrinmaya Sachan

## **E** *H* zürich

#### TL;DR

- 1. Lightweight probes can be trained predict the correctness of arithmetic operations (addition) by LLMs.
- Probes trained on simple arithmetic generalize to 2. chain-of-thought reasoning traces
- Probes can be used as weak oracles for self-correction. 3.

#### **Probing "Pure Arithmetic"**

Question: xxx + yyy =Model Output: <<xxx +yyy = zzz>>

We train multiple types of probes on the hidden states of LLMs. We train these probes to predict:

- the output of the model; 1.
- the ground-truth result; 2.
- the correctness of the model's output 3.



Probes can predict the ground-truth and the correctness with high accuracy.

#### Probing in Structured Chain-of-Thought Traces

**Question:** Sarah is planning to do some baking. She buys 771 pounds of rye flour, 611 [...] How many pounds of flour does she now have?

Model Output: «771+611+505=#987» «1987+758=2745» Error detected!

We then extend our analysis to multi-step arithmetic reasoning in the **GSM8K** dataset.



#### **Self-Correction Using Probes as Weak Oracles**

If the probes detect an error, we **add a follow-up prompt** after the current step. (E.g., That step looks suspicious. Let's redo just this step.)

**Question:** Sarah is planning to do some baking. She buys 771 pounds of rye flour, 611 [...] How many pounds of flour does she now have?

Model Output: «771+611+505=1987»

Re-prompting: That step looks suspicious. Let's re-do just this step:

```
«771+611+505=1887»
«1887+758=2645»
```

#### Using Probes for Self-Correction

Promot Style

TP Correction **FP** Preservation

We prompt the model to format each intermediate computation step as <a+b=c>

Then we apply probes trained on simple arithmetic to hidden states in the complex setting.



Probes tained on the simple arithmetic setting generalize well to the CoT setting.

r tompt Style		
suspicious	11.80%	100%
neutral	11.80%	100%
specific	10.11%	100%
stronger	8.99%	95.45%
detailed	6.18%	100%

Different self-correction prompt lead to varying correction rates, with up to 11.8% of flagged errors corrected and near-perfect preservation of correct answers.

#### Conclusion

Logistic (Separate)

Circular (Separate)

20

25

MLP (loint)

15

MLP (Separate)

Circular (Joint)

Lightweight probing tools offer a practical path toward **extracting** and leveraging this latent knowledge. Probing can be used for self-correction of LLMs in multi-step reasoning.