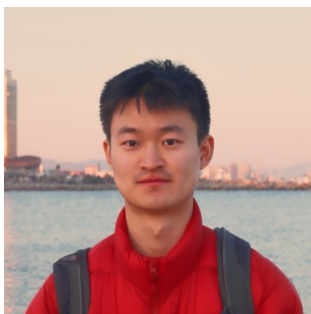


# AbstRaL: Augmenting LLMs' Reasoning by Reinforcing Abstract Thinking



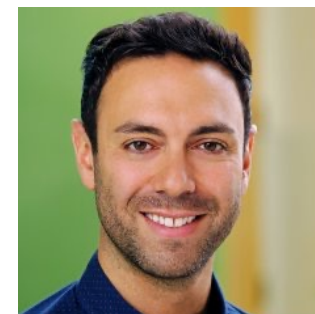
Silin Gao<sup>1,2</sup>



Antoine Bosselut<sup>2</sup>



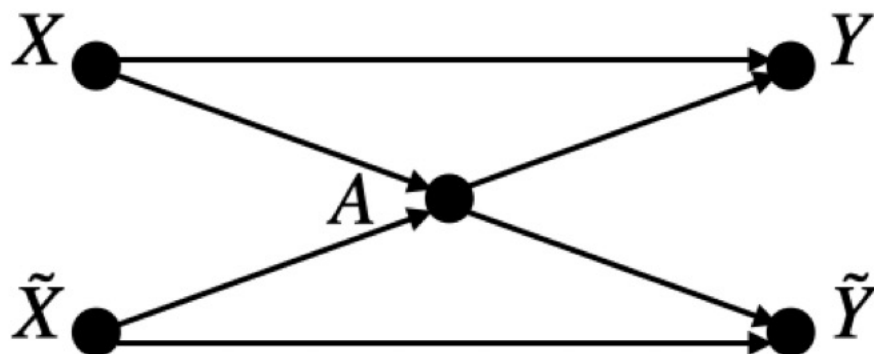
Samy Bengio<sup>1,2</sup>



Emmanuel Abbe<sup>1,2</sup>



# Robustness of Reasoning

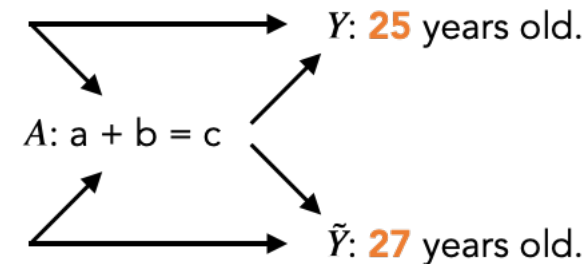


Two problems  $X$  and  $\tilde{X}$ , with solutions  $Y$  and  $\tilde{Y}$ , share the same high-level knowledge or reasoning schema.

Such problems can be handled by a common abstraction  $A$

$X$ : Alice is 20 years old. Bob is 5 years older than Alice. How old is Bob?

$\tilde{X}$ : Tom is 24 years old. Jerry is 3 years older than Tom. How old is Jerry?



$(X, Y) \rightarrow (\tilde{X}, \tilde{Y})$  illustrates a **distribution shift**

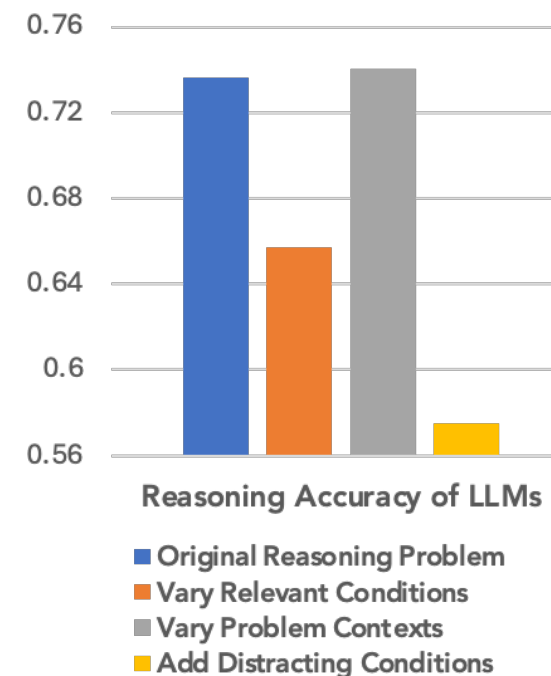
We expect a robust LLM reasoner to understand the underlying abstraction  $A$  (abstract thinking), and therefore achieve:  $p(Y|X) \approx p(\tilde{Y}|\tilde{X})$

# LLMs are Poor at Generalizing to **Distribution Shifts**

**Instantiation Shifts:** paraphrase, varying contexts and/or relevant conditions, etc.

**Interferential Shifts:** adding distracting (topic-relevant but useless) conditions

GSM-Plus<sup>1</sup>: GSM8K Problems + Distribution Shifts



<sup>1</sup>Li et al., 2024. "GSM-Plus: A Comprehensive Benchmark for Evaluating the Robustness of LLMs as Mathematical Problem Solvers."

# A Common Strategy: Robustifying by Instantiation

Learning more instances of the reasoning problem to anticipate potential distribution shifts.

Input

$x$ : Li is 12 years old. Jung is 2 years older than Li. How old is Jung?

Output

$y$ : Jung is  $12 + 2 = 14$  years old.

Synthetic Data Augmentation



$x'$ : Zhou is 15 years old. Wu is 5 years older than Zhou. How old is Wu?

$y'$ : Wu is  $15 + 5 = 20$  years old.

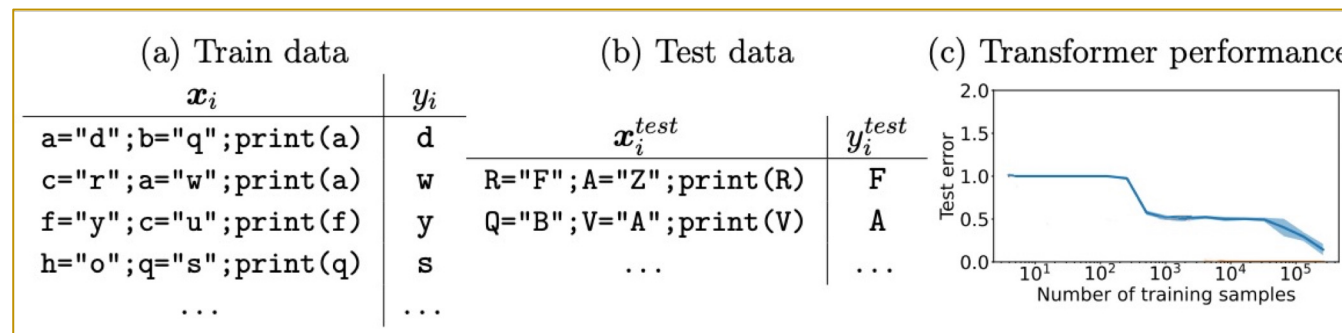
$x''$ : Amy is 37 years old. Sam is 19 years older than Amy. How old is Sam?

$y''$ : Sam is  $37 + 19 = 56$  years old.

...

...

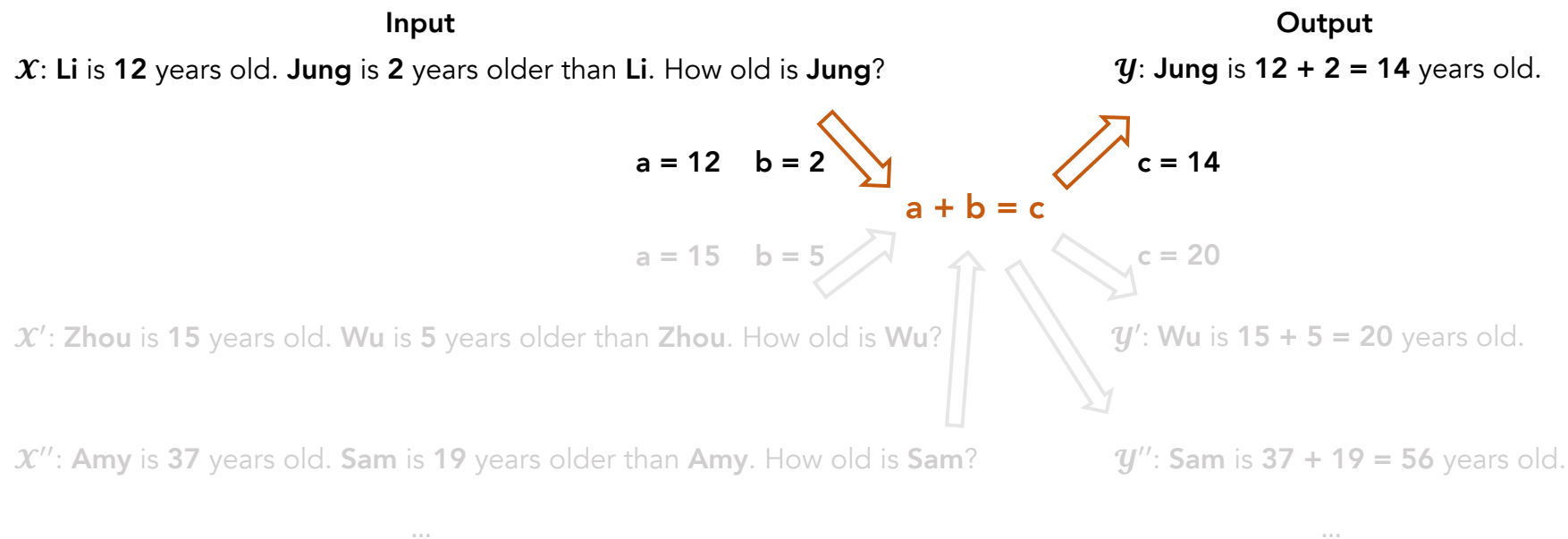
But Computational Expensive...



Even on simple letter printing task, a large amount of training samples were required to mitigate test-time generalization error.<sup>1</sup>

<sup>1</sup>Boix et al., 2024. "When can transformers reason with abstract symbols?"

# Our Strategy: Robustifying by **Abstraction**



**We train LLMs to directly learn the abstract thinking!**

**Modeling more general "abstraction" of reasoning, without scaling up the training data.**

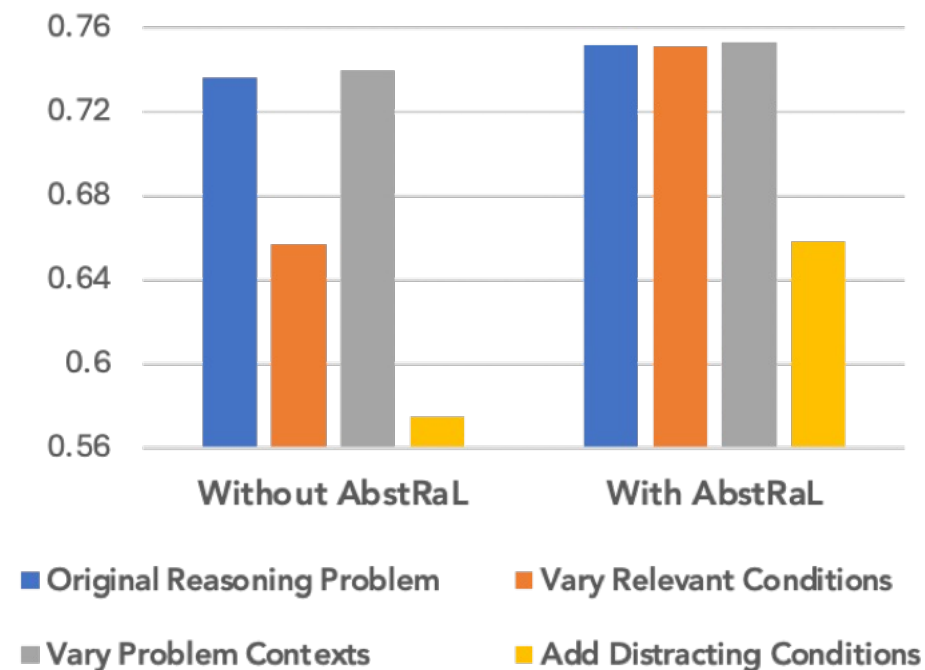
# AbstRaL Improves Robustness to Distribution Shifts

Our Reinforced **AbstRaction Learning** framework: **AbstRaL**

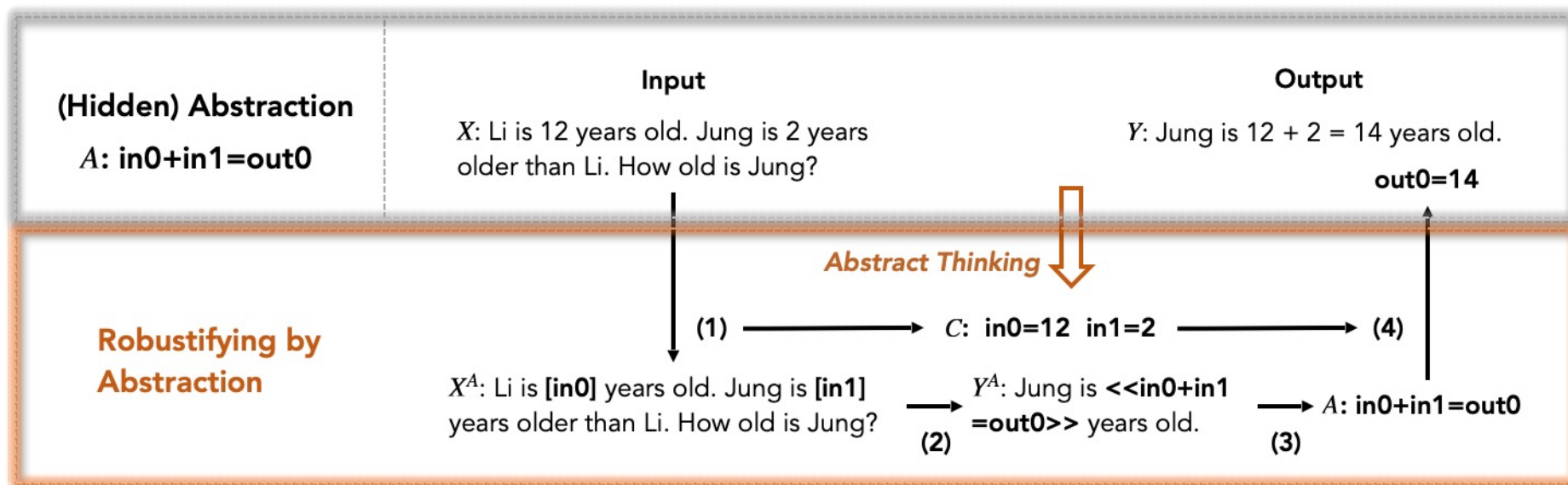
AbstRaL **almost reverts** performance drop caused by varying relevant conditions,

and **significantly mitigates** the interference of distracting conditions.

GSM-Plus: GSM8K Problems + Distribution Shifts



# Overview of **AbstRaL** Framework



(1) Condition Recognition Tool

(4) Symbolic Derivation Tool

(2) Abstract Reasoning (**Core Step**)

(3) Abstraction Retrieval Tool

LLMs are trained on abstract reasoning with **reinforcement learning**, based on our **granularly-decomposed** abstract reasoning data.

## Granularly-decomposed Abstract Reasoning (GranularAR)

- LLMs have learned **fine-grained** reasoning strategies at either pre-training<sup>1</sup> or post-training<sup>2</sup> phase, such as Chain-of-Thought (CoT) and Socratic problem decomposition as representatives.
- GranularAR **integrates abstract reasoning** with these pre-learned beneficial strategies.

$X^A$ : Zhang is [in0] times as old as Li. Li is [in1] years old. Zhang's brother Jung is [in2] years older than Zhang. How old is Jung?



$y^A$  (GranularAR):

**(Decomposing and Planning)** Let's think about the sub-questions we need to answer. **Q1**: How old is Zhang? **Q2**: How old is Jung?

**(CoT with Quoting Abstract Symbols)** Let's answer each sub-question one by one.

**Q1**: How old is Zhang? Li is [in1] years old, Zhang is [in0] times as old as Li, so Zhang is  $\ll in0 * in1 = out0 \gg$  years old.

**Q2**: How old is Jung? Zhang is [out0] years old, Jung is [in2] years older than Zhang, so Jung is  $\ll out0 + in2 = out1 \gg$  years old.

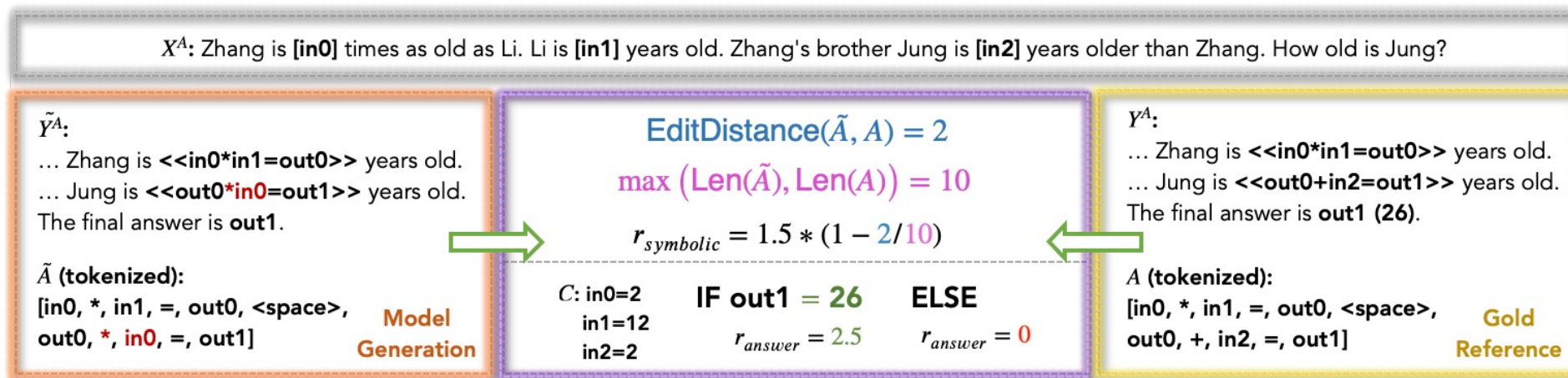
**(Conclusion)** The final answer is [out1].

<sup>1</sup>Yang et al., 2024. "Do large language models latently perform multi-hop reasoning?"

<sup>2</sup>Kumar et al., 2025. "Llm post-training: A deep dive into reasoning large language models."

# RL with Abstraction Rewards

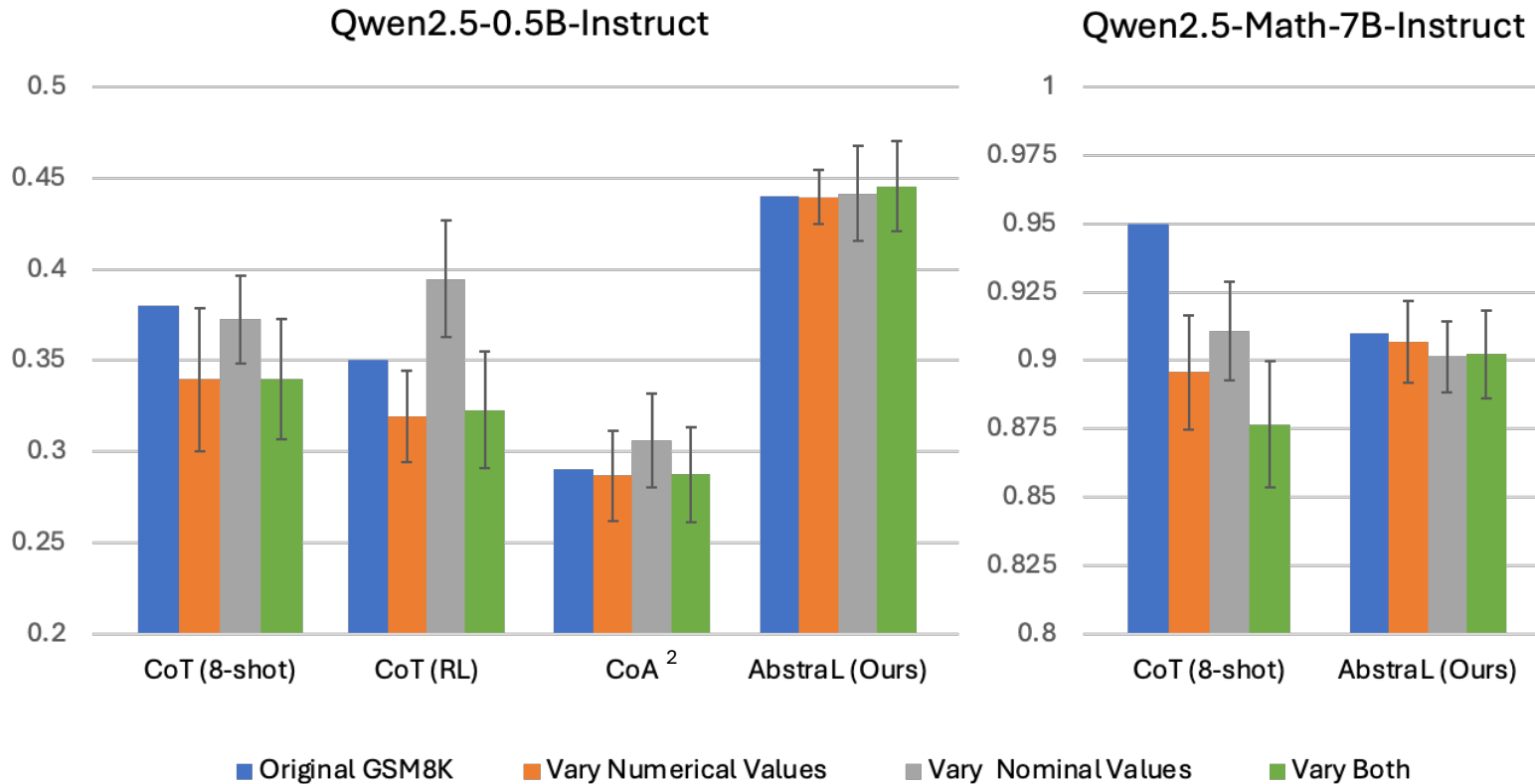
- LLMs are **poor at following in-context** demonstrations to reason in abstract manner<sup>1</sup>.
- Limitation of plain supervised fine-tuning (SFT)**: auto-regressive training objective forces LLMs also modeling the specific contexts.



- Symbolic Distance Reward** ( $r_{\text{symbolic}}$ ) granularly measures how the generated abstraction is aligned with (or close to) the expected abstraction, serving as a **milestone-style** reward that **more closely monitors** the progress of learning.
- Answer Correctness Reward** ( $r_{\text{answer}}$ ) checks whether the generated abstraction can derive the **correct final answer** given the gold input conditions.

<sup>1</sup>Gao et al., 2025. "Efficient Tool Use with Chain-of-Abstraction Reasoning."

# Results on GSM-Symbolic<sup>1</sup>

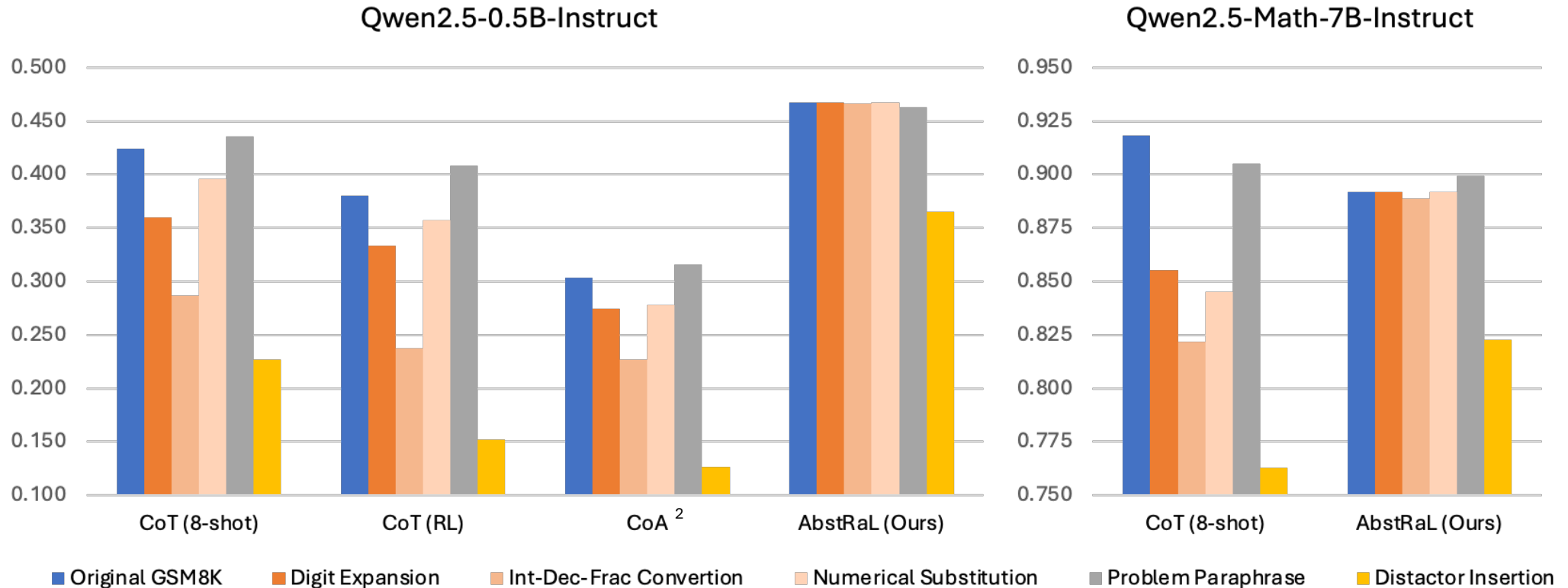


- Compared to baseline learning schemes (CoT with RL and CoA), AbstraL is **more reliable** to augment reasoning.
- AbstraL may **mitigate LLMs' overfitting** to the existing input conditions (or numbers), caused by potential **data contamination** at the pre-training or post-training stage.

<sup>1</sup>Mirzadeh et al., 2024. "GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models."

<sup>2</sup>Gao et al., 2025. "Efficient Tool Use with Chain-of-Abstraction Reasoning."

# Results on GSM-Plus<sup>1</sup>

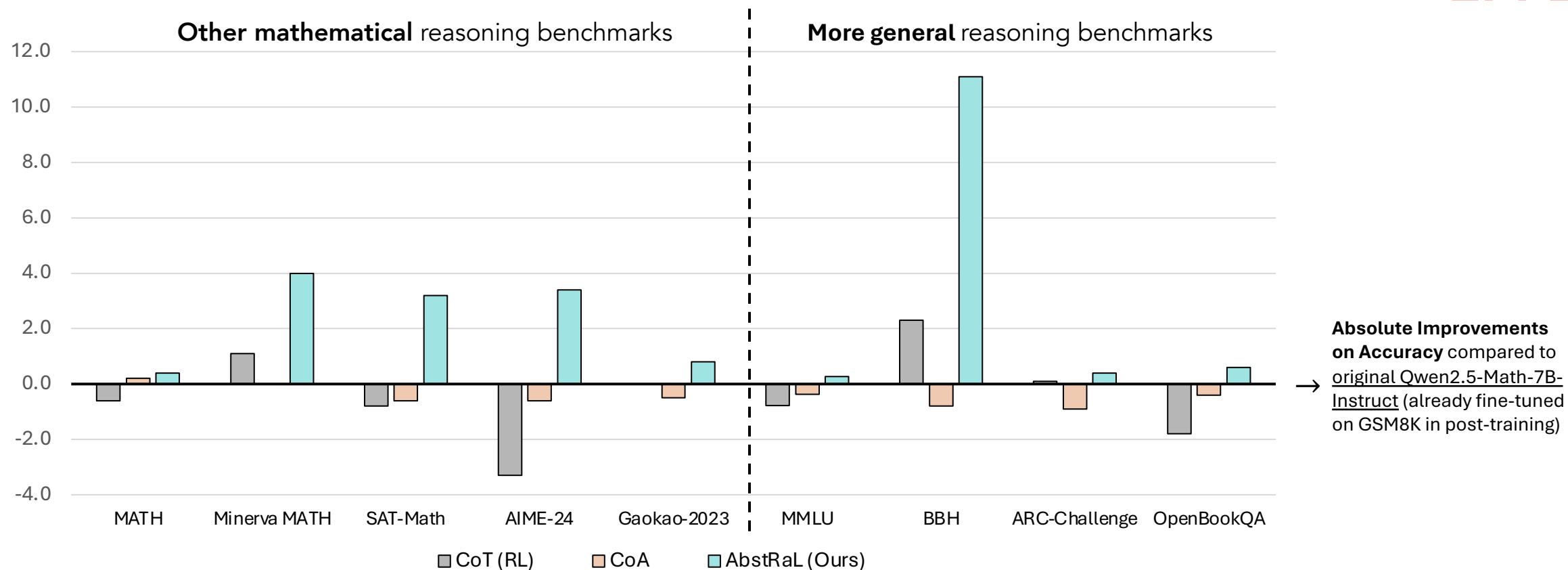


- AbstraL **almost reverts** the performance degradation caused by variations of input numbers (**orange**).
- AbstraL **mitigates the interference** of distracting conditions (**yellow**).

<sup>1</sup>Li et al., 2024. "GSM-Plus: A Comprehensive Benchmark for Evaluating the Robustness of LLMs as Mathematical Problem Solvers."

<sup>2</sup>Gao et al., 2025. "Efficient Tool Use with Chain-of-Abstraction Reasoning."

# Generalization to OOD Tasks



- AbstRaL **consistently outperforms** baseline methods, even though the learned GSM reasoning scheme may not be directly applicable to solving the OOD tasks.
- Learning the abstract thinking underlying GSM can already **implicitly benefit** other mathematical reasoning and more general reasoning of LLMs.

# Conclusion

- LLMs, particularly **small ones**, were shown to be **non-robust** at reasoning, representatively on the task of grade school mathematics (GSM)
  - **AbstRaL** effectively **offsets** the performance drop caused by distribution shifts!
- Our **RL with abstraction rewards** and **GranularAR** reasoning schema appear to be the key to reasoning robustness; while other baseline methods (CoT with RL, CoA, etc.) often fail.
- Improving GSM robustness via AbstRaL also implicitly benefits LLMs' capabilities on OOD mathematical and general reasoning tasks, indicating that **abstract thinking broadly enables better generalizability**.