

# AbstRaL: Augmenting LLM's Reasoning by Reinforcing Abstract Thinking

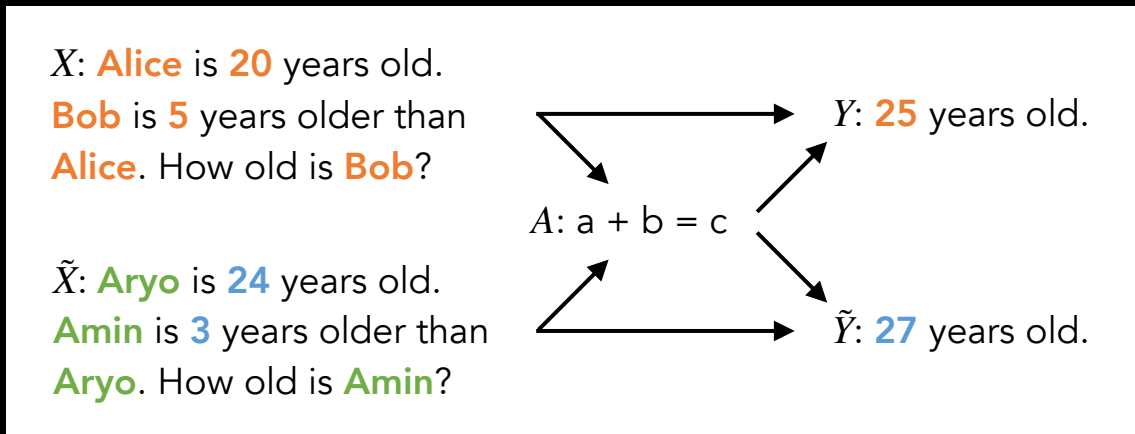
Silin Gao

Antoine Bosselut

Samy Bengio

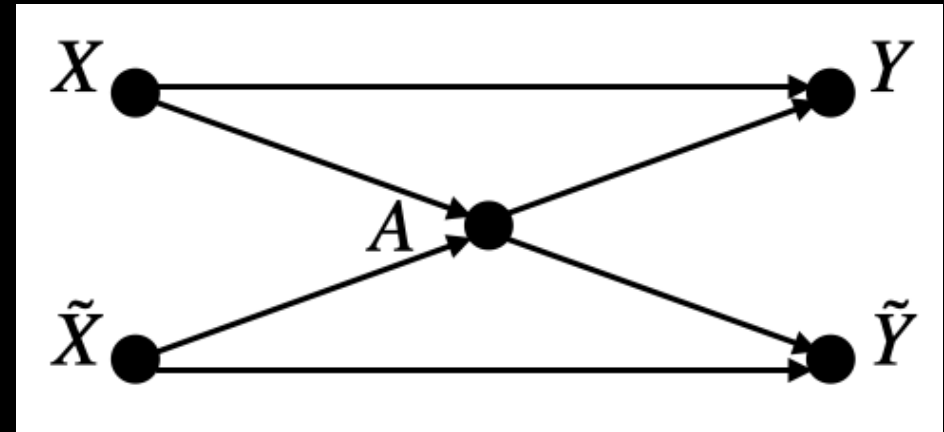
Emmanuel Abbe

# Are LLM truly reasoning? E.g., abstracting?



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

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



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

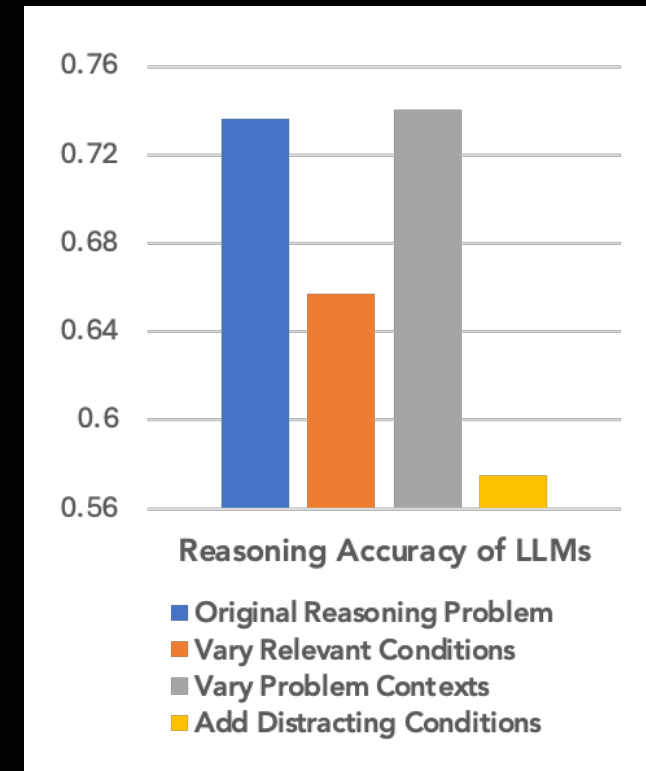
Can problems be handled by the LLM with a common abstraction  $A$ ?

# 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

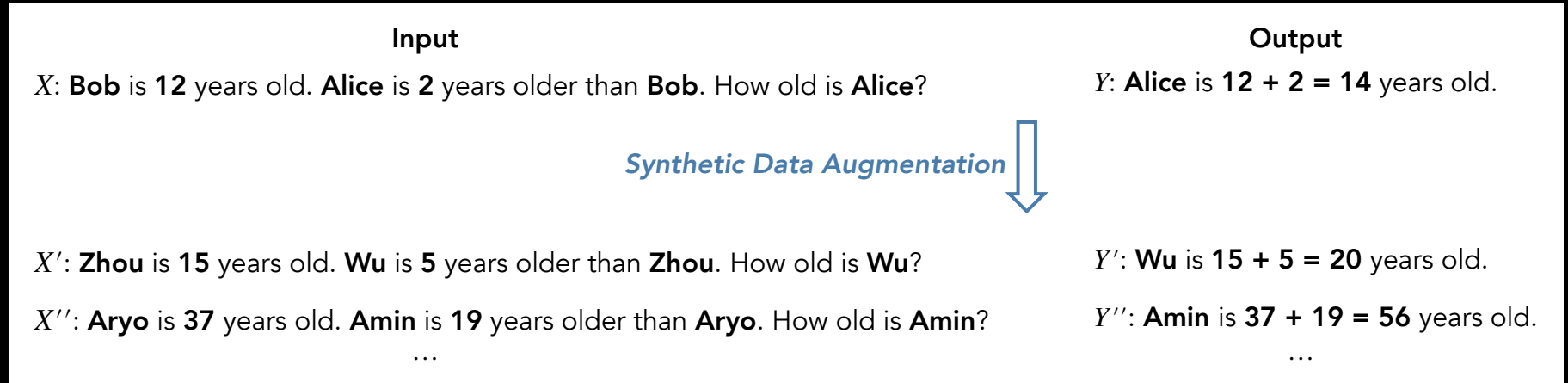
**GSM-Plus<sup>1</sup>:** GSM8K Problems + Distribution Shifts



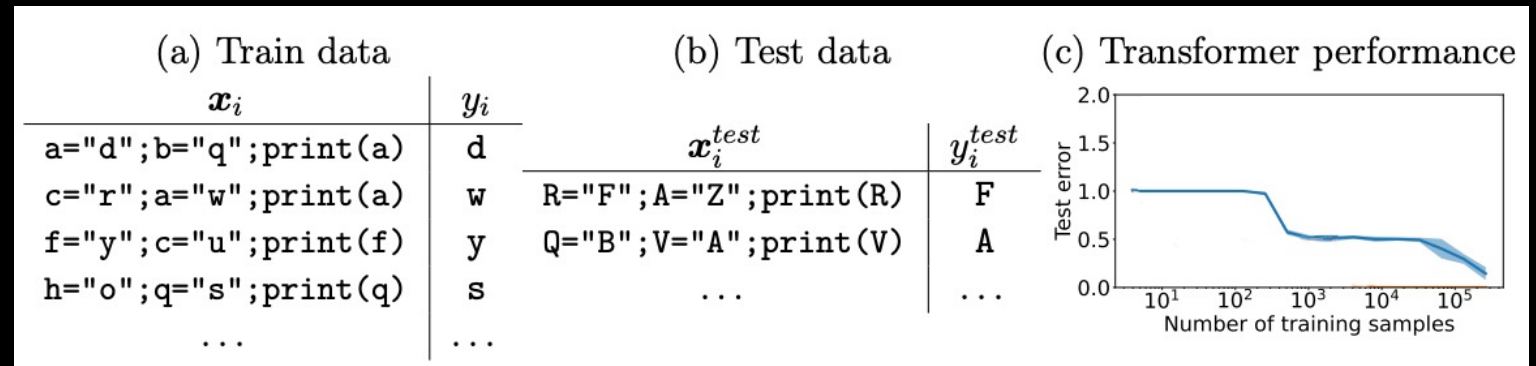
<sup>1</sup>Li et al., 2024. "GSM-Plus: A Comprehensive Benchmark for Evaluating the Robustness of LLMs as Mathematical Problem Solvers."

# A Common Strategy: Robustifying by Instantiation

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



But Computational Expensive...



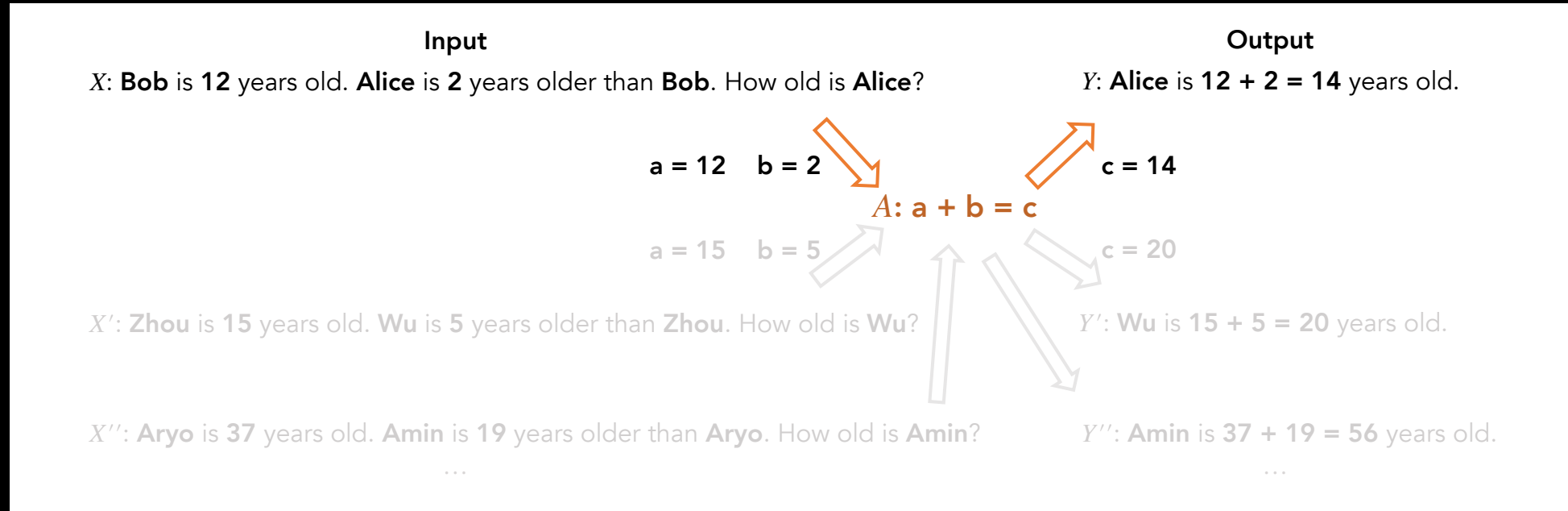
Even on simple token printing task, a large amount of training samples are required to mitigate test-time generalization error in from scratch training.<sup>1</sup>

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

## Ayn Rand

"**Abstract ideas** are conceptual integrations which subsume an incalculable number of concretes—and without abstract ideas you would not be able to deal with concrete, particular, real-life problems. You would be in the position of a new born infant.”

# Our Strategy: Robustifying by **Abstraction**



We train LLMs to directly learn the abstract thinking!

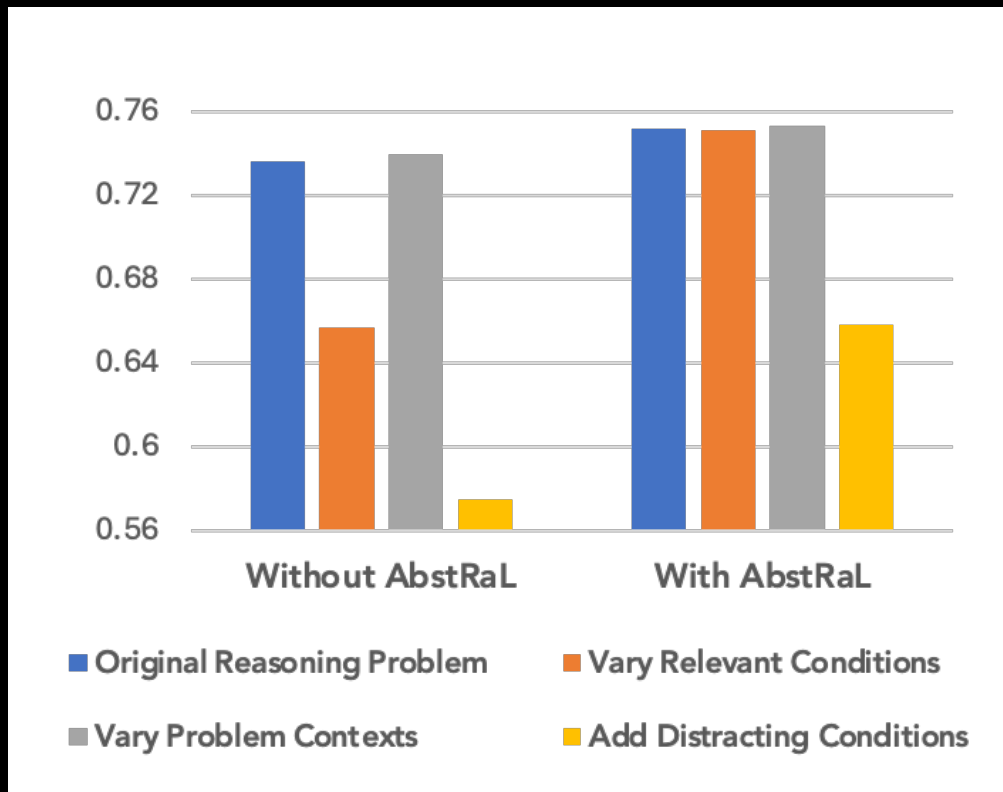
Modelling more general "abstraction" of reasoning, without scaling up the training data.

# AbstRaL Improves Robustness to Distribution Shifts

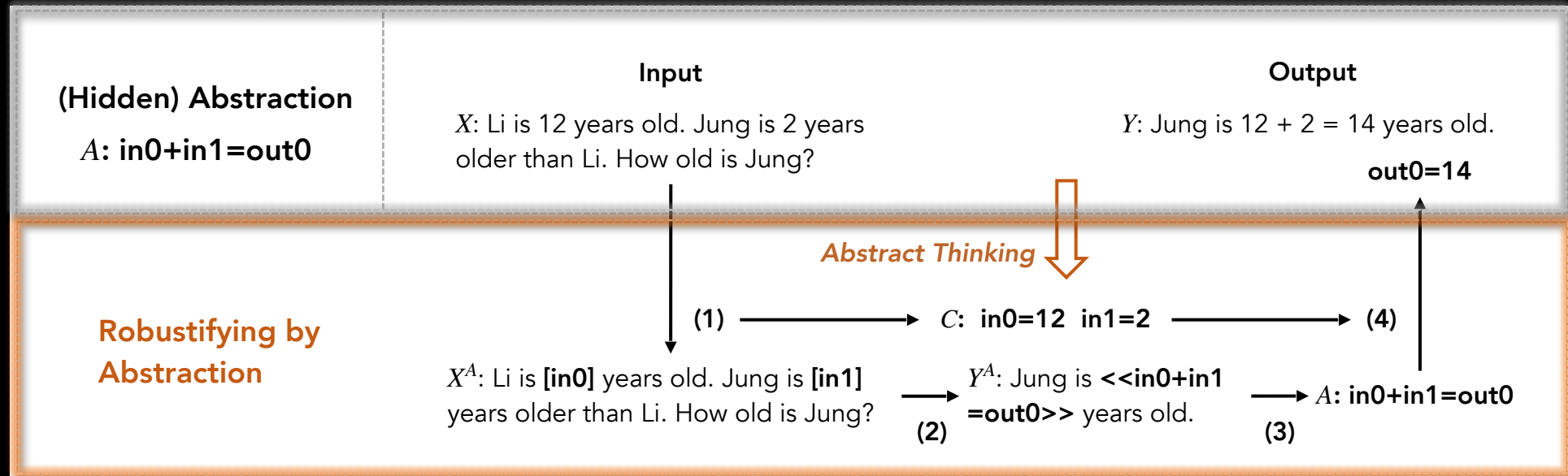
Our Reinforced **AbstRaction** Learning framework: **AbstRaL**

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.

# GranularAR Schema

## 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.

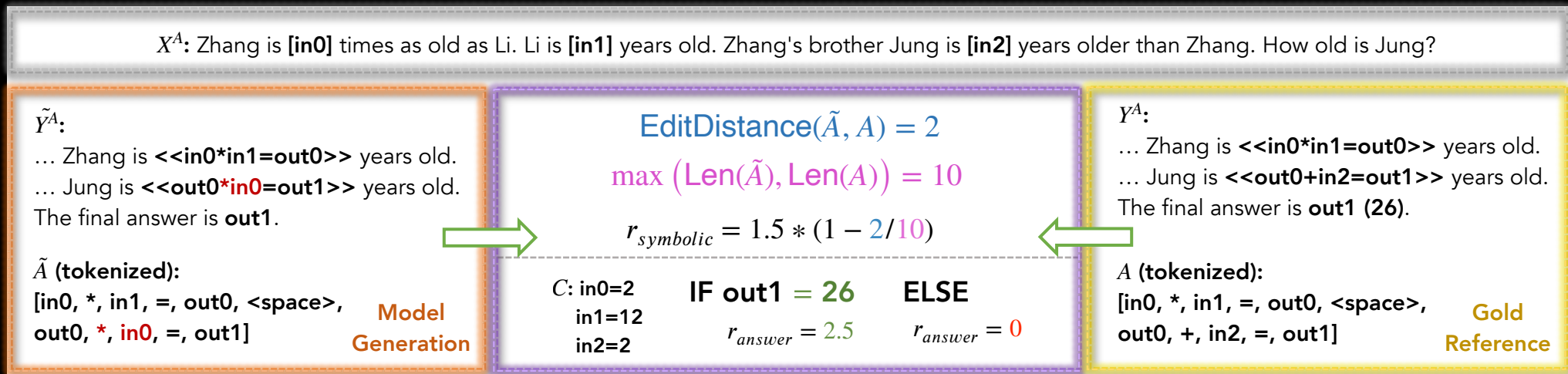
**(Conclusion)** The final answer is [out1].

<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

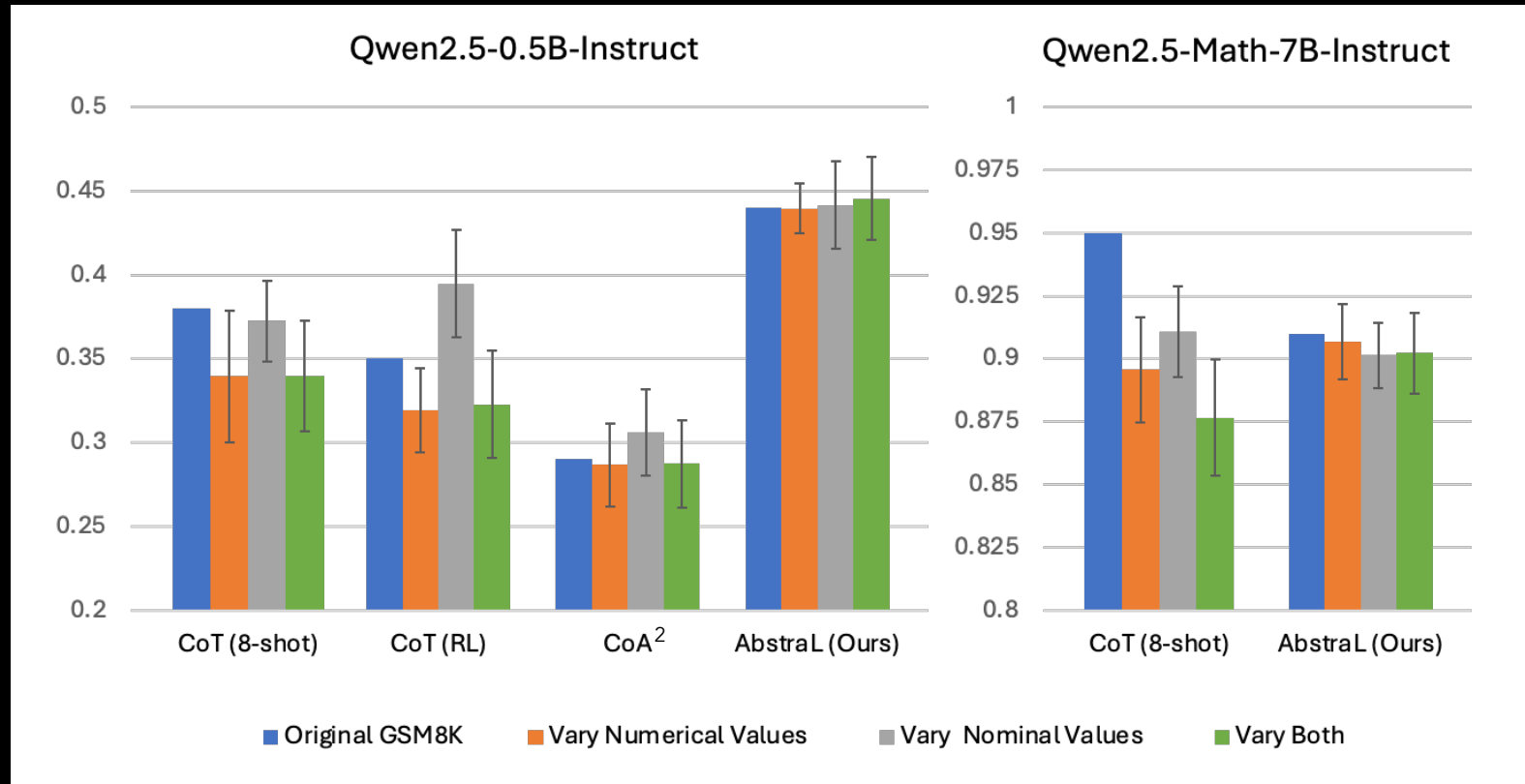
- LLMs are **poor at following in-context** demonstrations to reason in abstract manner<sup>1</sup>.
- Limitation of plain supervised fine-tuning (SFT)**: auto-regressive training objective forces LLMs also modeling the specific contexts.



- Symbolic Distance Reward** ( $r_{symbolic}$ ) 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** ( $r_{answer}$ ) 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 Grade School Mathematics (GSM-Symbolic<sup>1</sup>)

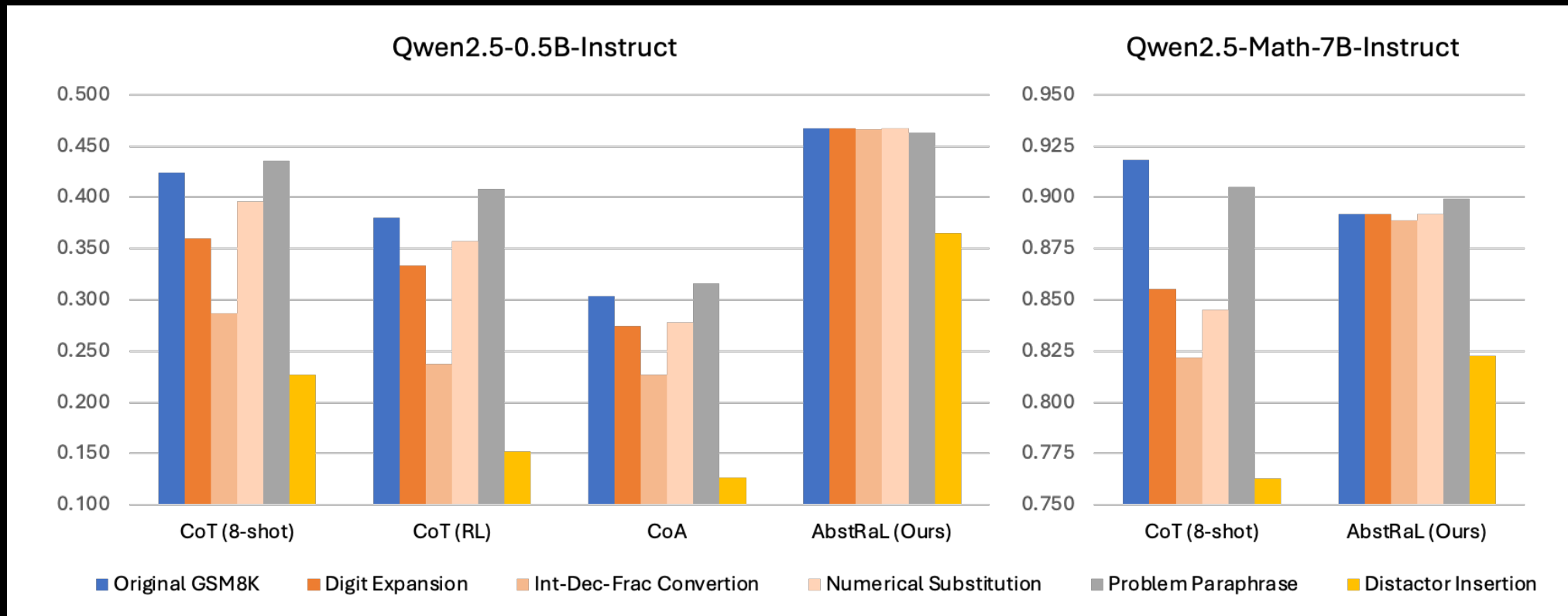


- Compared to baseline learning schemes (CoT with RL and CoA), AbstraL is **more reliable** to augment reasoning.
- AbstraL may **mitigate LLMs' overfitting** to the existing input conditions (or numbers), caused by potential **data contamination** at the pre-training or post-training stage.

<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 Grade School Mathematics (GSM-Plus)



- 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.

# OOD performance (trained with GSM and tested beyond GSM)

Table 3: Zero-shot evaluation results on OOD datasets. Best results on each model are **bold**.

Model	GSM Train	MATH	Minerva MATH	SAT- Math	AIME24	MMLU				BBH
	Method					STEM	Social	Humanities	Other	
Qwen2.5-0.5B-Instruct	Ori-SFT	30.0	5.5	46.9	0.0	37.5	51.5	41.6	49.7	21.0
	CoT-RL	19.2	5.5	43.8	0.0	38.5	50.8	41.2	49.1	23.2
	CoA	31.6	5.5	43.8	0.0	38.5	50.8	41.9	49.3	21.8
	AbstRaL	<b>34.7</b>	<b>7.4</b>	<b>62.5</b>	0.0	<b>38.8</b>	<b>53.5</b>	<b>42.2</b>	<b>50.7</b>	<b>26.3</b>
Qwen2.5-Math-7B-Instruct	Ori-SFT	83.5	34.6	90.6	13.3	68.4	82.5	63.2	76.4	43.7
	CoT-RL	82.9	35.7	89.8	10.0	68.5	81.5	63.0	74.4	46.0
	CoA	83.7	34.6	90.0	12.7	68.5	81.7	63.3	75.5	42.9
	AbstRaL	<b>83.9</b>	<b>38.6</b>	<b>93.8</b>	<b>16.7</b>	<b>68.6</b>	<b>82.7</b>	<b>63.8</b>	<b>76.5</b>	<b>54.8</b>

# Conclusion

- LLMs, particularly **small ones**, were shown to be **non-robust** at reasoning, such as on **Apple's GSM-Symbolic**  
→ **AbstRaL** allows to essentially **offset** the performance drop caused by distribution shifts!
- Our **RL with abstraction rewards** and **Granular** reasoning schema appear to be the key to reasoning robustness; while other baseline methods (CoT with RL, CoA, etc.) often fail.
- **Next:** expand our AbstRaL framework in a **domain adaptive** manner, to augment more **general purpose** LLMs.

## More general conclusion

LLMs can have difficulties solving challenging reasoning tasks on their own (barrier phenomena).

Rather than teaching LLMs how to produce good solutions directly, teach LLMs how to produce good explanations.

This goes beyond just training on thinking traces, it calls upon principles of reasoning (induction, abstraction and more).