

# LANGUAGE MODEL SELF-IMPROVEMENT BY REINFORCEMENT LEARNING CONTEMPLATION



Jing-Cheng Pang<sup>\*1,2</sup>, Pengyuan Wang<sup>\*1,2</sup>, Nan Tang<sup>1,2</sup>, Kaiyuan Li<sup>1</sup>, Xionghui Chen<sup>1,2</sup>,  
 Jiacheng Xu<sup>1</sup>, Zongzhang Zhang<sup>1</sup>, Yang Yu<sup>1,2</sup>

<sup>1</sup> National Key Laboratory for Novel Software Technology, Nanjing University

<sup>2</sup> Polixir Technologies



## Background and Motivation

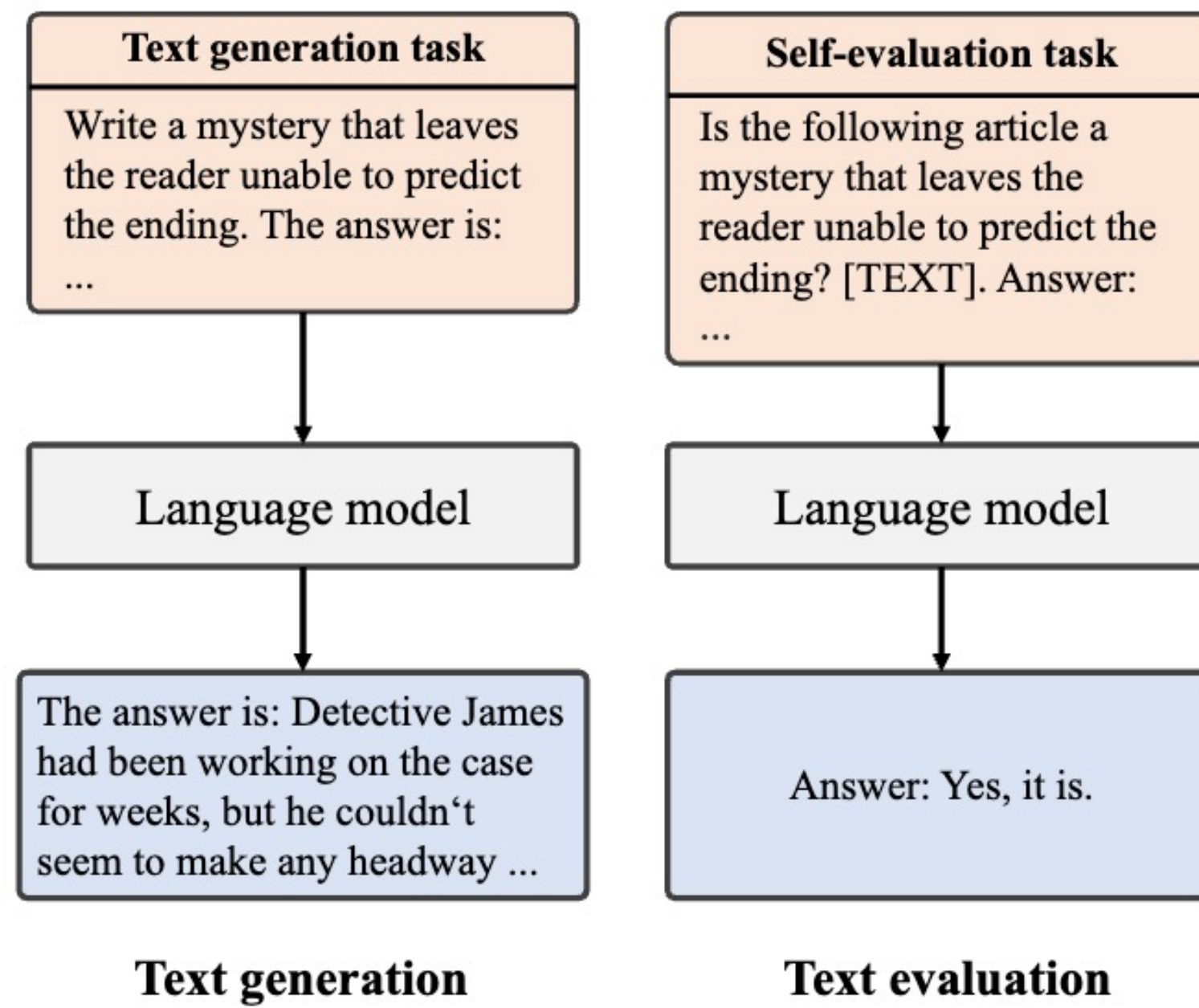
### Background

Reinforcement learning from AI feedback (RLAIF) is a popular technique for language model improvement. RLAIF trains using reinforcement learning (RL), with the preference model as the reward signal which **requires a solid capability**.

It is simpler for language models to **evaluate a sentence than to generate it**, even for small language models.

### Evaluation is Simpler than Generation

Verify the text evaluation ability of LLMs by investigating potential for evaluation to improve LLMs.



**Whether can the evaluation be utilized to improve text generation ?**

## Answer Generation with self-Evaluation on various tasks

	Reasoning about Colored Objects	Logical Deduction (7)	Tracking Shuffled Objects (5)	Object Counting	Tracking Shuffled Objects (3)	Geometric Shapes
w/o SE	30.9%	18.5%	10.1%	34.7%	28.1%	10.7%
w/ SE	<b>31.1%</b>	<b>20.5%</b>	<b>11.1%</b>	<b>34.9%</b>	<b>31.5%</b>	<b>13.5%</b>

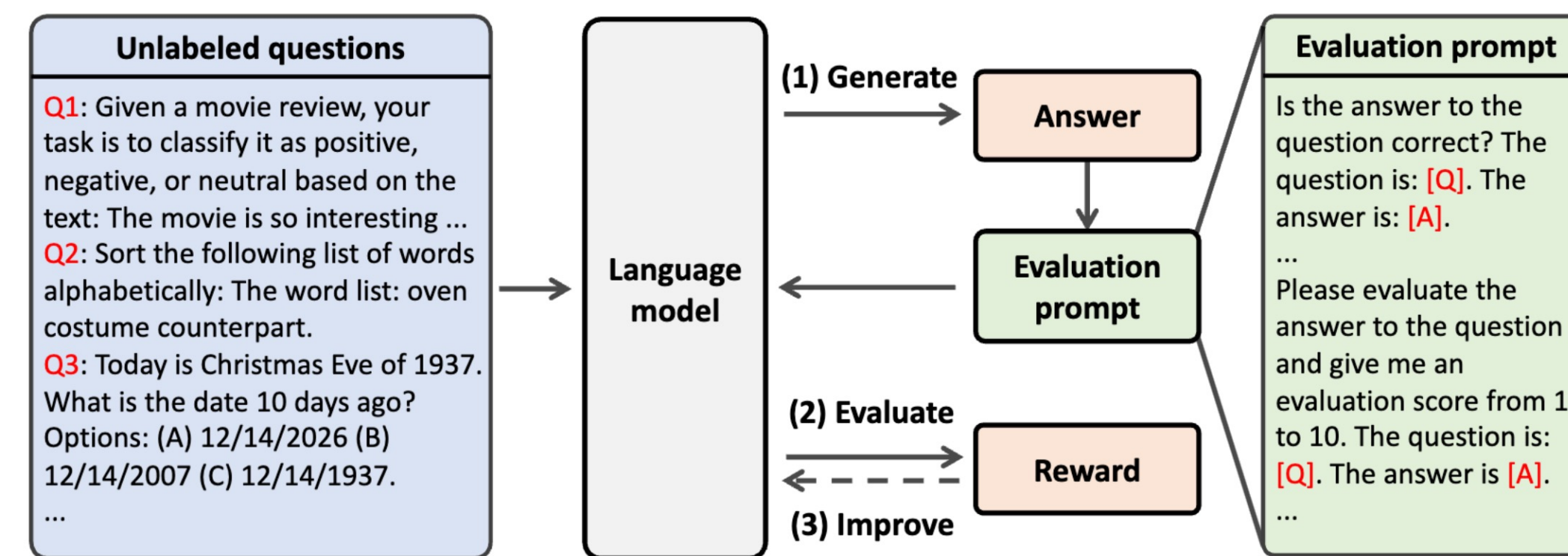
	Web of Lies	Sports Understanding	Logical Deduction (3)	Logical Deduction (5)	Penguins in a Table	Navigate
w/o SE	51.6%	59.7%	34.9%	23.6%	23.5%	47.7%
w/ SE	<b>53.2%</b>	<b>59.7%</b>	<b>38.3%</b>	<b>25.7%</b>	<b>28.8%</b>	<b>50.5%</b>

## Method

### Self-Improvement by Reinforcement Learning Contemplation

Overall training procedure:

- (1) Answer generation to the unlabeled questions.
- (2) Self-evaluation by asking LM using evaluation prompt.
- (3) Update the language model to maximize the reward using reinforcement learning.



Two types of evaluation prompts:

- (1) Correctness Evaluation Prompt (CEP): "Is the answer to the question correct? The question is: [Q]. The answer is: [A]."
- (2) Quality Evaluation Prompt (QEP): "Please evaluate the answer to the question and give me an evaluation score from 1 to 10. The question is: [Q]. The answer is: [A]."

## Experiments

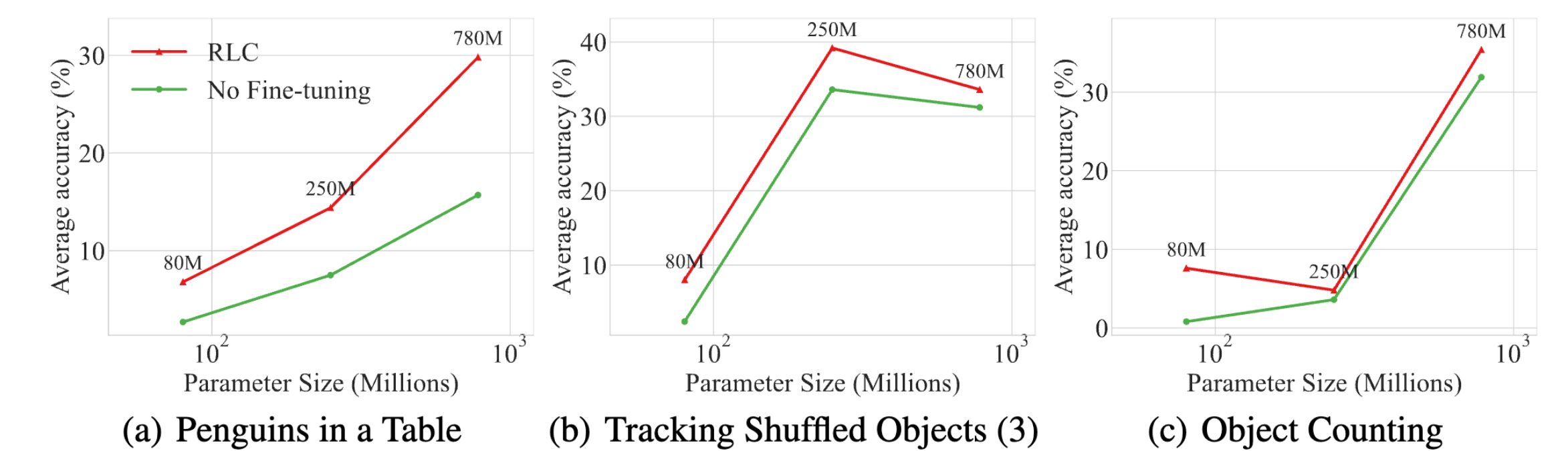
### Main Results on Reasoning Tasks

	Reasoning about Colored Objects	Logical Deduction (7)	Tracking Shuffled Objects (5)	Object Counting	Tracking Shuffled Objects (3)	Geometric Shapes
RLFT	32.1%	45.7%	12.4%	42.6%	33.6%	18.9%
DG	32.0%	35.2%	12.4%	31.9%	31.2%	5.2%
SC	<b>39.6%</b>	27.6%	12.4%	24.0%	33.6%	15.6%
Self-train	19.5%	13.1%	<b>15.5%</b>	11.7%	33.1%	12.4%
Self-refine	25.2%	13.2%	8.0%	18.0%	25.2%	10.0%
Best-of-N	26.8%	12.8%	12.1%	14.0%	30.0%	8.4%
RLAIF	30.4%	36.9%	11.4%	32.5%	32.8%	14.0%
RLC	35.0%	<b>39.2%</b>	12.2%	<b>35.4%</b>	<b>33.6%</b>	<b>17.8%</b>

	Web of Lies	Sports Understanding	Logical Deduction (3)	Logical Deduction (5)	Penguins in a Table	Navigate
RLFT	72.2%	68.8%	58.6%	41.9%	44.2%	55.6%
DG	43.6%	53.2%	39.6%	28.4%	15.7%	46.4%
SC	48.4%	53.6%	42.8%	30.8%	<b>35.2%</b>	<b>62.8%</b>
Self-train	51.1%	51.1%	34.0%	18.4%	19.7%	48.7%
Self-refine	47.2%	50.0%	28.4%	17.2%	17.8%	46.0%
Best-of-N	50.0%	<b>59.2%</b>	42.0%	22.0%	17.8%	45.2%
RLAIF	52.1%	56.1%	22.0%	33.7%	19.8%	48.8%
RLC	<b>52.9%</b>	53.5%	<b>44.0%</b>	<b>34.6%</b>	29.8%	57.1%

### Performance on Different Sizes of Language Model



### Answer Accuracy on Unseen Datasets

	Logical Deduction (7)	Object Counting	Penguins in a Table	Sports Understanding	Tracking Shuffled Objects (5)	Average
Acc.	36.7 (+1.5)	32.7 (+0.7)	18 (+2.2)	52.8 (-0.4)	12.3 (-0.1)	<b>30.5 (+0.8)</b>