

## Can Small Language Models Help Large Language Models Reason Better?: *LM-Guided Chain-of-Thought*

Jooyoung Lee\*, Fan Yang, Thanh Tran, Qian Hu, Emre Barut, Kai-Wei Chang, Chengwei Su The Pennsylvania State University, University Park, PA, USA\* Amazon AGI, Boston, MA, USA Email : jfl5838@psu.edu

Jooyoung Lee, Fan Yang, Thanh Tran, Qian Hu, Emre Barut, Kai-Wei Chang, and Chengwei su. 2024. Can small Language Models Help Large Language Models Reason Better?: LM-Guided Chain-of-Thought. LREC-COLING (2024), May 20-25, 2024, Torino, Italy.



**A**: 8

Chain-of-Thought (CoT) prompting is a technique designed to exploit language models (LMs)' inherent reasoning abilities.

#### **Standard Prompting**

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

V.S.

#### **CoT Prompting**

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there? A : Let's think step by step.



tir up ion) There are 16 balls in total. Half of balls fore golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.



#### **CoT Prompting**

**Q**: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A : Let's think step by step.



(continuation) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.



Key contributions are:



Enhanced end-task performance



#### **CoT Prompting**

**Q**: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A : Let's think step by step.





Enhanced end-task performance



# Supports Interpretability of LMs' output

(continuation) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.



#### **CoT Prompting**

**Q**: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A : Let's think step by step.



Key contributions are:

Supports Interpretability of LMs' output

Enhanced end-task performance

(continuation) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.



No additional training required!



Limitation 1. CoT prompting tends to be only effective when the model is gigantic!





**Limitation 2**. Extremely large LMs may still generate low-quality rationales.





## "How can we help the language model find a optimal CoT in a resource-efficient way?"



# We propose a new framework called **"LM-guided CoT"** that leverages <u>TWO independent LMs</u> for CoT reasoning.

**Q**: Are both directors of films Dear Mr.Prohack and Returning Mickey Stern from the same country?

**Context**: Thornton Freeland was an American film director. Freeland directed the Returning Mickey in 1994. Dear Mr.Prohack is made in 1995 and is first aired in Canada. It is the first movie of John Mason, an American director.

A : Let's think step by step.





# We propose a new framework called **"LM-guided CoT"** that leverages <u>TWO independent LMs</u> for CoT reasoning.

**Q**: Are both directors of films Dear Mr.Prohack and Returning Mickey Stern from the same country?

**Context**: Thornton Freeland was an American film director. Freeland directed the Returning Mickey in 1994. Dear Mr.Prohack is made in 1995 and is first aired in Canada. It is the first movie of John Mason, an American director.

A : Let's think step by step.











- The small LM' inherent reasoning capability is significantly worse than the large LM's.
- We warm-up the small LM by training with a knowledge distillation technique (i.e., *finetune the small LM with rationales generated by the large LM*).
- Prompt used for finetuning

Given a question (Q) and a context, generate a chain of reasoning step-by-step to answer the question. Context: {context} Q: {question} Reasoning: {LLM rationale}



- We attempt to further refine the rationale qualities of the rationale distilled model from Step 1 with reinforcement learning.
- 8 linguistic aspects used for measurement: factuality, relevance, logicality, consistency, coherence, fluency, naturalness, readability

Aspects	Descriptions
Factuality	Percentage (0.0-1.0) measuring if the reasoning is grounded based on the context (Input: $c \& r'$ )
Relevance	Percentage (0.0-1.0) measuring if the reasoning is relevant to the question (Input: $q \& r'$ )
Logicality	Binary (0 or 1) measuring if the reasoning is logical and can reach a final answer
Consistency	Binary (0 or 1) measuring if the reasoning remains consistent and coherent (Input: $q \& r'$ )
Coherence	Binary (0 or 1) measuring if the reasoning is without redundant information (Input: $q \& r'$ )
Fluency	Binary (0 or 1) measuring if the reasoning is well-written and grammatically correct (Input: $r'$ )
Naturalness	Binary (0 or 1) measuring if the reasoning is natural and human-like (Input: $r'$ )
Readability	Binary (0 or 1) measuring if the reasoning is easy to follow and understandable (Input: $r'$ )

Table 1: Descriptions of 8 rationale aspects used for evaluation. q, c, and r' denote a question, context, and a corresponding rationale generated by the small LM, respectively.



- **factuality, relevance** : computed by token-level lexical-overlap following Ye and Durrett (2022)
- logicality, consistency, coherence, fluency, naturalness, readability : experimented with two methods
  - Method 1. use a large LM as reference-free NLG evaluators
  - Method 2. train a simple machine learning classifier using human-annotated data

Ye, Xi, and Greg Durrett. "The unreliability of explanations in few-shot prompting for textual reasoning." *Advances in neural information processing systems* 35 (2022): 30378-30392.



- **factuality, relevance** : computed by token-level lexical-overlap following Ye and Durrett (2022)
- logicality, consistency, coherence, fluency, naturalness, readability : experimented with two methods
  - Method 1. use a large LM as reference-free NLG evaluators
  - Method 2. train a simple machine learning classifier using human-annotated data

Answer the question based on the provided information. Question: Can the given reasoning {definition of evaluation metric}? (a) Yes. (b) No. Information: Question: {question} Reasoning: {reasoning} Answer:

Ye, Xi, and Greg Durrett. "The unreliability of explanations in few-shot prompting for textual reasoning." *Advances in neural information processing systems* 35 (2022): 30378-30392.



- **factuality, relevance** : computed by token-level lexical-overlap following Ye and Durrett (2022)
- logicality, consistency, coherence, fluency, naturalness, readability : experimented with two methods
  - Method 1. use a large LM as reference-free NLG evaluators
  - Method 2. train a simple machine learning classifier using human-annotated data

Ye, Xi, and Greg Durrett. "The unreliability of explanations in few-shot prompting for textual reasoning." *Advances in neural information processing systems* 35 (2022): 30378-30392.



- **factuality, relevance** : computed by token-level lexical-overlap following Ye and Durrett (2022)
- logicality, consistency, coherence, fluency, naturalness, readability : experimented with two methods
  - Method 1. use a large LM as reference-free NLG evaluators

Method 2. train a simple machine learning classifier using human-annotated data

Methods	Coherence & Consistency & Logicality			Fluency & Naturalness & Readability				
	Acc	Precision	Recall	F1	Acc	Precision	Recall	F1
Method 1 Method 2	0.62 <b>0.8</b>	0.62 <b>0.79</b>	0.6 <b>0.79</b>	0.6 <b>0.79</b>	0.7 <b>0.9</b>	0.7 <b>0.94</b>	<b>0.81</b> 0.75	0.67 <b>0.9</b>

Ye, Xi, and Greg Durrett. "The unreliability of explanations in few-shot prompting for textual reasoning." *Advances in neural information processing systems* 35 (2022): 30378-30392.





• We further update the knowledge-distilled model from Step 1 with Proximal Policy Optimization (PPO).

• Two type of rewards (aspect-specific reward & task-specific reward)

• incorporate penalties based on the Kullback Leibler (KL) divergence between the learned policy LM and the knowledge-distilled LM.



#### **Experiment Setup**

- <u>Task</u> : Extractive QA
- Model : FLAN-T5 small (80M) for the small LM and FLAN-T5 XXL (11B) for the large LM
- **Dataset** : HotpotQA & 2WikiMultiHopQA

Туре	Description	Template
		Based on the provided context, answer the following question (Q).
Standard	Directly predicting the answer	Context: c
prompting	based on input	Q: <i>q</i>
		A:
CoT prompting		Based on the provided context, answer the following question (Q)
	Predicting the answer after generating the reasoning	by reasoning step-by-step.
		Context: c
		Q: <i>q</i>
		A : Let's think step by step.
LM-guided CoT prompting		Based on the provided context, answer the following question (Q)
	Predicting the answer with conditional generation upon the LM-generated reasoning	by reasoning step-by-step.
		Context: c
(our method)		Q: <i>q</i>
		A : Let's think step by step. $r'$ . Hence, the answer is

Table 3: Descriptions and templates of each prompt used for the answer prediction task. q, c, and r' denote a question, context, and a corresponding rationale generated by the small LM, respectively.



#### SC : self-consistency

#### 2WikiMultiHopQA Rationale HotpotQA Prompt Provision? Answer Answer EM F1 EM F1 Inclusion Inclusion standard prompting X 0.5 0.714 0.583 0.5 0.625 0.647 CoT prompting 0.483 0.686 0.611 0.4 0.532 0.561 **Baselines** 0.624 0.625 CoT prompting + SC X 0.503 0.70 0.471 0.603 LM-guided CoT prompting (KD) 0.507 0.702 0.625 0.506 0.626 0.661 1 LM-guided CoT prompting (KD + SC) 0.513 0.714 0.635 0.524 0.644 0.679 X LM-guided CoT prompting 0.503 0.698 0.625 0.507 0.631 0.665 1 $(\mathsf{KD} + R_{aspect})$ LM-guided CoT prompting 0.508 0.704 0.627 0.503 0.622 0.657 1 $(\mathsf{KD} + R_{aspect} + R_{taskAcc})$ LM-guided CoT prompting 1 0.5 0.698 0.623 0.501 0.619 0.653 $(KD + R_{aspect} + ranking)$

Results

Table 2: Answer prediction performance results of baselines and our approach. We regard SC decoding as a non-rationale provision because this method can result in multiple variations of rationales, rather than a single one. Values in bold represent the highest scores and underlined values are the second highest scores.



#### SC : self-consistency

#### Rationale HotpotQA 2WikiMultiHopQA Prompt Provision? Answer Answer EM F1 EM F1 Inclusion Inclusion standard prompting 0.5 0.714 0.583 0.5 0.625 0.647 X CoT prompting 0.483 0.686 0.611 0.4 0.532 0.561 CoT prompting + SC 0.624 0.503 0.70 0.471 0.603 0.625 X LM-guided CoT prompting (KD) 0.507 0.702 0.625 0.506 0.626 0.661 1 LM-guided CoT prompting (KD + SC) X 0.513 0.714 0.635 0.524 0.644 0.679 LM-guided CoT prompting 1 0.503 0.698 0.625 0.507 0.631 0.665 Ours $(\mathsf{KD} + R_{aspect})$ LM-guided CoT prompting 0.508 0.704 0.627 0.503 0.622 0.657 1 $(\mathsf{KD} + R_{aspect} + R_{taskAcc})$ LM-guided CoT prompting 1 0.5 0.698 0.623 0.501 0.619 0.653 $(KD + R_{aspect} + ranking)$

Results

Table 2: Answer prediction performance results of baselines and our approach. We regard SC decoding as a non-rationale provision because this method can result in multiple variations of rationales, rather than a single one. Values in bold represent the highest scores and underlined values are the second highest scores.



#### Results



Table 2: Answer as a non-rationa than a single or highest scores. Figure 2: Average answer prediction performance (across three evaluation metrics) and average rationale quality scores (*i.e.*,  $R_{aspect}$ ) for HotpotQA (left) and 2WikiMultiHopQA (right). The right y-axis represents the mean answer prediction scores, and the left y-axis represents the mean rationale quality scores.

egard SC decoding f rationales, rather es are the second



#### Conclusion

- CoT prompting does not always outperform standard prompting.
  - Especially when the context gets longer like 2WikiMultiHopQA, models tend to make errors.
- Smaller models trained for rationale generation can help black-box larger models with reasoning.
  - Our approach outperforms the original CoT prompting, where the large LM is asked to reason first and generate the final answer.
  - With self-consistency methods, its performance gain enhances.
- In-depth analyses on how to balance the performance and rationale quality are required.
  - Our final approach (rationale ranking with reinforced policy models) significantly improves the rationale quality, but at the same time, shows a minor drop in answer prediction.



# PennState LREC-COLING 2024



#### Check out our paper!

# Thank you:)

Jooyoung Lee Email : j<u>fl5838@psu.edu</u>



SCAN ME