Mitigating Misleading Chain-of-Thought Reasoning with Selective Filtering

Yexin Wu, Zhuosheng Zhang, Hai Zhao



SHANGHAI JIAO TONG UNIVERSITY

Accuracy (%) on test split of ECQA ackbone model is UnifiedQA-base. SelF-ner outperforms the pipeline by 3.5%.

Method Vanilla Pipeline SelF-B Accuracy 58.07 54.95 58.48 (+3.5)

Method Vanilla Pipeline SelF-

Accuracy 64.22 76.80 80.06 (+3.26)

Table 4: Accuracy (%) on test split of LastLetter The backbone model is UnifiedQA-base. SelF-

Abstract

- Proposed a novel approach called the selective filtering reasoner (SelF-Reasoner) that assesses the entailment relationship between the question and the candidate reasoning chain. Then, we proceed with CoT reasoning when the reasoning chain demonstrates confidence; otherwise, we opt to predict the answer directly.
- SelF-Reasoner improves the fine-tuned T5 baseline consistently over the ScienceQA, ECQA, and LastLetter tasks.

Motivation

- CoT(Few-shot, zero-shot, chatgpt...) has achieved success with LLM.
- CoT reasoning is considered one of the emergent abilities shown in "Scaling Laws".
- In some situations (Li et al., 2022a; Magister et al., 2022; Ho et al., 2022; Li et al., 2022b; Wang et al., 2023), though small language models (SML) can not do few-shot or zero-shot CoT finetuning, they can still benefit from CoT finetuning.
- But small language models are easier to generate misleading CoTs compared with LLM.
- Ways to mitigate the effect of misleading CoTs

Q: departmental, dome, pressed, fascinating

Predicted rationale: The last letter of the first word 'departmental' is 'I'. The last letter of the second word 'dome' is 'e'. The last letter of the third word 'pressed' is 'd'. The last letter of the forth word 'funny' is 'y' Extracted Answer: The answer is ledy. (Incorrect, the rationale part is wrong) Directly predicted Answer: The answer is ledg. (Correct) Ground Truth Answer: The answer is ledg.

8520

10000

Question: What is th round Truth CoT: Applying for Table 7: Dataset statistics used in our experiments. Generated wrong CoT: P praise wh get praise.

wilt of applying for a job. Be wilt of applying for a job. rrect answer: The approximation Figure 1: An example of an invalid CoT reasoning

from ECQA (Aggarwal et al., 2021; Wang et al., 2023). The errors are highlighted in red. The generated CoT is wrong at the first step, and the error continues to the end. However, when

altering to direct prediction, this one-step reasoning question is solved correctly.

Methodology

Dataset

ECQA

l astl etter

ScienceQA 12726

Datasets

- Three Datasets:
 - ScienceQA(Lu et al., 2022a) ECQA (Aggarwal et al., 2021; Talmor et al., 2019)

Train Validation Test

4241 4241

1221

5000

1221

5000

- LastLetter

Model Architecture

- Three baselines:
- Vanilla finetuning means to generate the answer only.
- Compound generator generates the reasoning chain and the answer in one run. Pipeline first generates the reasoning chain, and then use the question and the reasoning chain to generate the answer.
- SelF-Reasoner will judge the reasoning chain. If the filter thinks the reasoning chain is misleading, then it will adopt vanilla finetuning's method to predict the answer only using the question. Otherwise, it will go as the pipeline using the question and the reasoning chain to generate the answer.
- Our backbone model is T5.



Filter Design

- Training-based Filter.
- Use an insufficiently trained reasoner to generate plausible reasoning chains and correct reasoning chains (correct means it will lead to correct answers).
- T5 encoder is trained to be a filter.
- Load finetuned T5 reasoner's parameter to accelerate training and improve accuracy. Rule-based Filter.
- For the LastLetter task: "The given word should appear in the valid CoT".
- Observation: Many given words in the reasoning chain are different from the ones in the question due to tokenization and sampling, which will cause the incorrect answer extraction



Experiments

- Datasets:
- ScienceQA, ECQA, LastLetter Metrics: Accuracy
- Result: New SOTA

		Reasoner outperforms the pipeline by 3.26%.				
Method	Method Model		Format	Accuracy		
Lu et al. (2022a)	Human GPT-3 (CoT) UnifiedQA _{Base} UnifiedQA _{Base} UnifiedQA _{Base}	- In-Context Learning fine-tuning fine-tuning fine-tuning	Q-A Q-ALE Q-A Q-AE Q-ALE	88.4 75.17 70.12 73.33 74.11		
Lu et al. (2023)	ChatGPT (CoT) GPT-4 (CoT) Chameleon (ChatGPT) Chameleon (GPT-4)	In-Context Learning In-Context Learning In-Context Learning In-Context Learning	Q-EA Q-EA Q-EA Q-EA	78.31 83.99 79.93 86.54		
Vanilla	UnifiedQA _{Samll} UnifiedQA _{Base} UnifiedQA _{Large}	fine-tuning fine-tuning fine-tuning	Q-A Q-A Q-A	71.54 83.09 86.53		
Compound	UnifiedQA _{Base} UnifiedQA _{Base} UnifiedQA _{Base}	fine-tuning fine-tuning fine-tuning	Q-ALE Q-EA Q-LEA	76.13 77.71 73.97		
Pipeline	UnifiedQA _{Small} UnifiedQA _{Base} UnifiedQA _{Large}	fine-tuning fine-tuning fine-tuning	$\begin{split} Q-E-> QE$-A Q-E-> QE$-$A$ Q-E-> QE$-$A$ Q-E-> QE$-$A$ $\end{subarray}$	66.37 79.32 84.98		
SelF-Reasoner	UnifiedQA _{Small} UnifiedQA _{Base} UnifiedQA _{Large}	fine-tuning fine-tuning fine-tuning	SelF-Reasoner SelF-Reasoner SelF-Reasoner	69.55 83.45 87.24		

Table 1: Accuracy (%) of each baseline on test split. In the format part, Q = Question, A = AnswerE = Explanation, L = Lecture. We list the results from ScienceQA (Lu et al., 2022a), ChatGPT, GP (Lu et al., 2023) for comparison. L and E can be treated as reasoning chain. So LEA/EA and ALE/AE correspond to the standard RA and AR as defined in Section 3.1, respectively. Our SelF-Reasoner (Large) is comparable in accuracy to a human's.

Analysis

Quality of CoT generated by small language models

- Metrics:
 - BLEU-1/4
 - ROUGE
 - Sentence Similarity
- Human evaluation.
- Similar to the ground truth CoT in structure.
- Incorrect in some key parts, causing misleading results.

Model	Split	BLEU-1	BLEU-4	ROUGE-L	Similarity	Complete	Entailment	Correct
Base	Lead to Correct Answer Lead to Incorrect Answer All	0.914 0.789 0.892	0.776 0.660 0.756	0.910 0.797 0.891	0.937 0.860 0.924	1.00 1.00 -	1.00 1.00 -	0.94 0.02 -
Large	Lead to Correct Answer Lead to Incorrect Answer All	0.937 0.775 0.917	0.810 0.642 0.788	0.929 0.784 0.910	0.949 0.847 0.936	1.00 1.00 -	0.98 1.00	0.96 0.02

Automatic metrics (BLEU-1/4, ROUGE-L, Similarity) and human evaluation of generated explanations. We evaluate these metrics on different splits of the produced CoT according to whether they can lead to the correct answer. Details of human evaluation are shown in Appendix A.8

Filter's performance

Mo

Sr

Ba

Conclusion

- Filter can predict whether the CoT will mislead the answer accurately.
- Introducing filter outperforms random assignment.

del	Vanilla	Pipeline	Random	SelF-Reasoner	Generator	Filter	Valid Acc	Invalid Acc	Acc	F1
all	71.54	66.37	68.76	69.55	Base	Base	76.96	76.39	76.84	0.841
se	83.09	79.32	81.61	83.45		Large	81.30	81.64	81.37	0.874
ae	86.53	84.98	86.09	87.24			0.1.00		0.1.01	
9-						Base	74.97	75.03	74.98	0.836
e 5: Ablation on the CoT filter on ScienceOA				Large	Large	80.07	78.17	79.78	0.871	

Table benchmark. Random refers to randomly choosing vanilla fine-tuning and pipeline to produce the answer

Ablation Study: Filter&Pipeline size

- Larger Filter performs better
- The improvement in the larger pipeline is smaller because the reasoner's ability to generate plausible CoT is improved, making it harder for filters to distinguish plausible CoT from correct CoT.



Figure 3: The "scaling law" in the size of the CoT filter and pipeline on ScienceQA benchmark. The dashed line presents the accuracy of the pipeline.

- We proposed a selective filtering reasoner (SelF-Reasoner) to perform CoT only as necessary and mitigate the detrimental effects of erroneous reasoning chains.
- Our SelF-Reasoner outperforms the finetuned CoT/vanilla baseline on ScienceQA, ECQA, and LastLetter datasets, advancing the effectiveness of CoT in small-scale language models.
- We analyze the obstructions of fine-tuning CoT on language models and conclude common types in invalid generated CoT.

- Acknowledgement
- [1] Department of Computer Science and Engineering, Shanghai Jiao Tong University

[2] Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University [3] This paper was partially supported by Joint Research Project of Yangtze River Delta Science and Technology Innovation Community (No. 2022CSJGG1400) and Joint Funds of the National Natural Science Foundation of China (No. U21B2020)

Contact wuyexin_libro_i131[AT]sjtu.edu.cr zhangzs[AT]situ.edu.cn Zhaohai[AT]cs.sjtu.edu.cn

Table 6: Accuracy and F1 score of the CoT filter on classifying the generated reasoning chain on ScienceQA benchmark. Valid/Invalid Acc refers to the filter's accuracy in discriminating valid/invalid reasoning chains. Acc is the overall accuracy

base Filter Size