## Improving Language Model Reasoning with Self-motivated Learning

Yunlong Feng, Yang Xu, Libo Qin, Yasheng Wang, Wanxiang Che Harbin Institute of Technology, China Huawei Noah's Ark Lab

### **Motivation & Background**

- Models gain reasoning capability
  - After trained with data that has rationales
- Lacks of datasets with high-quality rationales
  - high annotation cost
- How to use existing datasets without rationales?
  - Utilize the correctness



• Correct is better than Wrong





#### Method

Question Answer / Answer * Rationale ① Generation 1. Rationale Generation			Questio SFT	n, <i>Rationale, A</i> 1 Tinetuning Model	Inswer PPO Model	Rationale Generation [Instruction and Question] $\langle q_i \rangle$ [Answer] $\langle a' \rangle$ [Pationalo]	Final Answer Generation[Instruction and Question] $\langle q_i \rangle$ [Rationale] $\langle \hat{r}_i \rangle$ [Answer] $\langle \hat{a}_i \rangle$	
Condition Rank	ondition $Answer = Answer'$ $Answer \neq Answer'$ $\checkmark$ Filter $\checkmark$ FilterRank $Rationale$ $\checkmark$ Rationale $\overset{\star}{\sim}$		Answer ✓ Filter	Reward Model   2 Finetuning   Question, Rationale   Question, Rationale				
Dataset	2. Ration Choices	Training Samples	Test Samples	Data Split	3. IVIO License	References	[Instruction and Question] Would a pear sink in water? [Answer]	[Instruction and Question] Would a pear sink in water? [Answer]
SingleEq AddSub MultiArith SVAMP	- - - -	356 276 420 700	152 119 180 300	70:30 70:30 70:30 70:30	None Unspecified Unspecified MIT	Koncel-Kedziorski et al. (2015a) Hosseini et al. (2014) Roy and Roth (2016) Patel et al. (2021a)	No [Rationale] The density of a pear is about 0.6a/cm3, which is less than	Yes [Rationale] The density of a pear is about 0.60/cm3, which is less than

• • • • • • • • • • • • • • • • • • • •						()	U 60/CM3 Which is less than	I 60/cm3 Which is less than
GSM8K	-	7473	1319	Original	MIT	Cobbe et al. (2021)		
Date Understand	ling 5-6	258	111	70.30	Anache-20	Srivastava et al. (2022)	water. Objects less dense than	water. Objects less dense than
		200	1001			$\frac{1}{2}$	water float Thus a pear would	water float Thus a near would
CommonSenseC	QA 5	9741	1221	Original	Unspecified	Talmor et al. (2018)	Mater neut. Thae, a pear weard	
StrategyQA	2	1603	687	70:30	Apache2.0	Geva et al. (2021a)	float.	SINK.

#### Experiments

Method	Param	Single Eq	Add Sub	Multi Arith	SVAMP	GSM8K	Date Understanding	Common SenseQA	Strategy QA
			Close-	Source	Models				
text-davinci-003	175B	86.4	81.3	83.7	73.6	59.5	77.0	70.0	61.1
text-davinci-002	175B	82.24	78.99	78.89	64.67	40.26	73.87	61.75	53.57
			Open-	Source I	Models				
StableVicuna	13B	62.50	57.14	43.33	46.67	40.26	45.95	58.64	41.34
LLama2-Chat	7B	73.03	68.91	67.22	53.67	28.35	35.14	56.67	38.14
			Method	s on Lla	ma2 7B				
Few-shot-CoT	7B	63.82	54.62	35.00	39.00	14.60	53.15	50.61	61.28
Few-shot-CoT <sup>SC=8</sup>	7B	67.76	67.23	55.56	44.67	15.09	35.13	48.40	62.45
Fine-tune	7B	71.05	63.87	11.67	45.67	12.58	64.87	76.58	65.21
Fine-tune-CoT (text-davinci-002)	7B	70.39	72.27	76.67	47.33	_	73.88	_	58.95
Fine-tune-CoT (STaR)	7B	75.66	67.23	72.78	44.33	17.29	81.98	63.63	64.63
Fine-tune-CoT (Llama2)	7B	71.05	65.55	53.33	40.67	13.72	83.78	69.53	60.84
Self-motivated Learning	7B	76.32	76.47	80.00	55.33	18.88	87.39	77.97	66.08

Method	SingleEq	AddSub	MultiArith
Fine-tune-CoT (text-davinci-002)	70.39	72.27	76.67
+RL	71.71	78.16	80.56
Increase	+1.32	+5.89	+3.89

The reward model trained with rank information  $\bullet$ exhibits a certain degree of generalization.



• Accuracy (%) in 8 tasks under our different models and methods. Note that the methods based on the LLama2 7B are trained in different datasets separately.

The Self-motivated Learning applies PPO for reinforcement learning in the training dataset  $\bullet$ using the Fine-tune-CoT model, resulting in an average increase of 10.68%.

#### Conclusion

• We propose "Self-motivated Learning", a framework is grounded in the idea that a rationale leading to the correct answer is superior to one leading to an incorrect answer. • We conducted experiments across different datasets encompassing three categories of complex reasoning, demonstrating that our method can significantly enhance model performance without external annotation.

#### The reward of model can reflect the performance. $\bullet$

The performance is increasing during training.

#### Acknowledgement

We gratefully acknowledge the support of the National Natural Science Foundation of China (NSFC)via grant 62236004 and 62206078.



# LREC-COLING 2024

