# DeepSeekMath

Soumyabrata Chaudhuri (Soumya)





# Contributions



- ➤ Introduces **DeepSeekMath-7B**, a SOTA math reasoning model.
- > A meticulous **Iterative data mining pipeline** for crawling the web.
- > Introduces **Group Relative Policy Optimization** (GRPO).

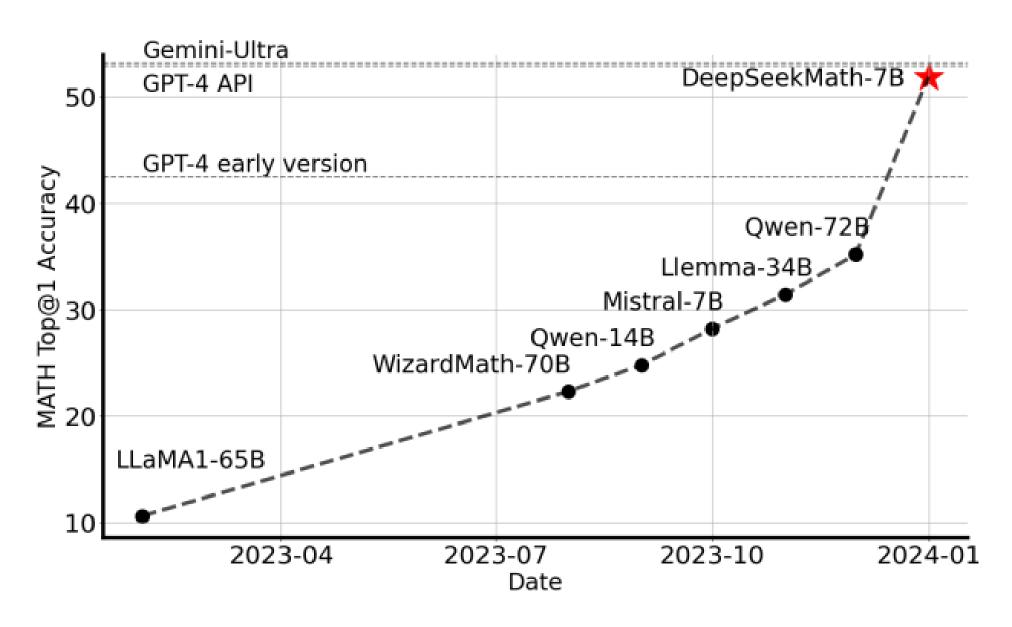


Figure 1 | Top1 accuracy of open-source models on the competition-level MATH benchmark (Hendrycks et al., 2021) without the use of external toolkits and voting techniques.



# **Pretraining Data & Corpus Construction**

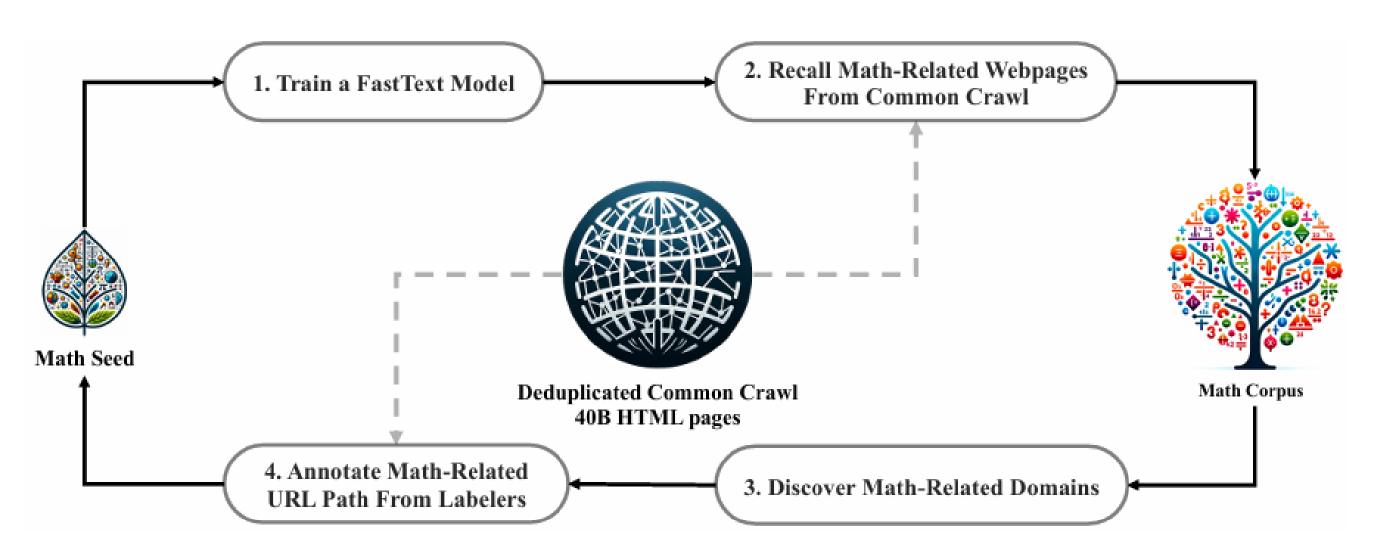


Figure 2 | An iterative pipeline that collects mathematical web pages from Common Crawl.



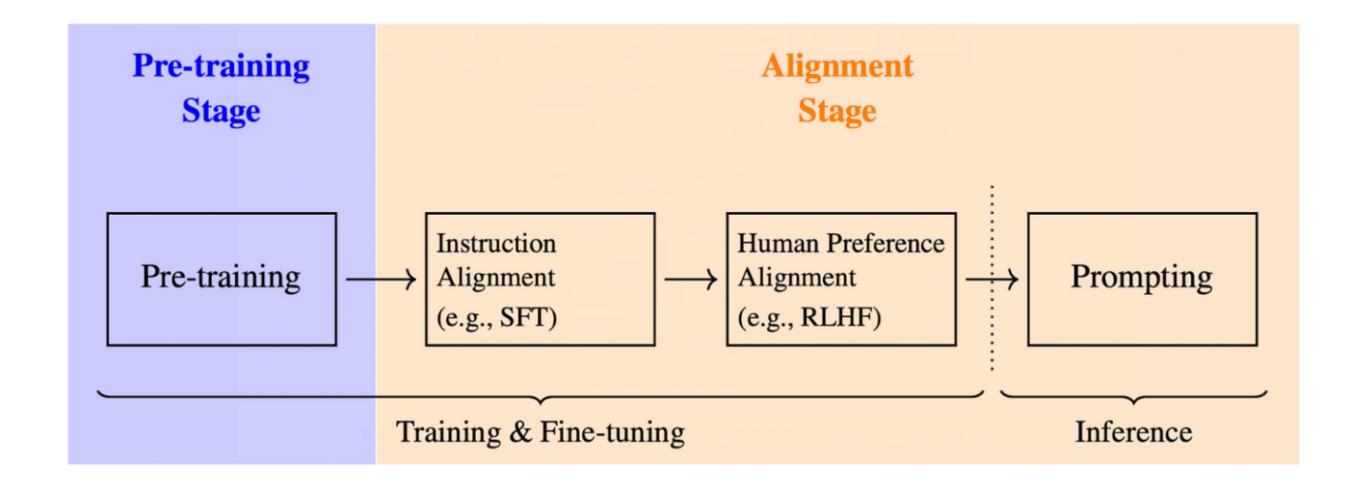
# **Pretraining Data & Corpus Construction**

Math Corpus	Size		English	Bench	marks	Chinese Benchmarks			
		GSM8K	MATH	OCW	SAT	MMLU STEM	СМАТН	Gaokao MathCloze	Gaokao MathQA
No Math Training	N/A	2.9%	3.0%	2.9%	15.6%	19.5%	12.3%	0.8%	17.9%
MathPile	8.9B	2.7%	3.3%	2.2%	12.5%	15.7%	1.2%	0.0%	2.8%
OpenWebMath	13.6B	11.5%	8.9%	3.7%	31.3%	29.6%	16.8%	0.0%	14.2%
Proof-Pile-2	51.9B	14.3%	11.2%	3.7%	43.8%	29.2%	19.9%	5.1%	11.7%
DeepSeekMath Corpus	120.2B	23.8%	13.6%	4.8%	56.3%	33.1%	41.5%	5.9%	23.6%

Table 1 | Performance of DeepSeek-LLM 1.3B trained on different mathematical corpora, evaluated using few-shot chain-of-thought prompting. Corpus sizes are calculated using our tokenizer with a vocabulary size of 100K.

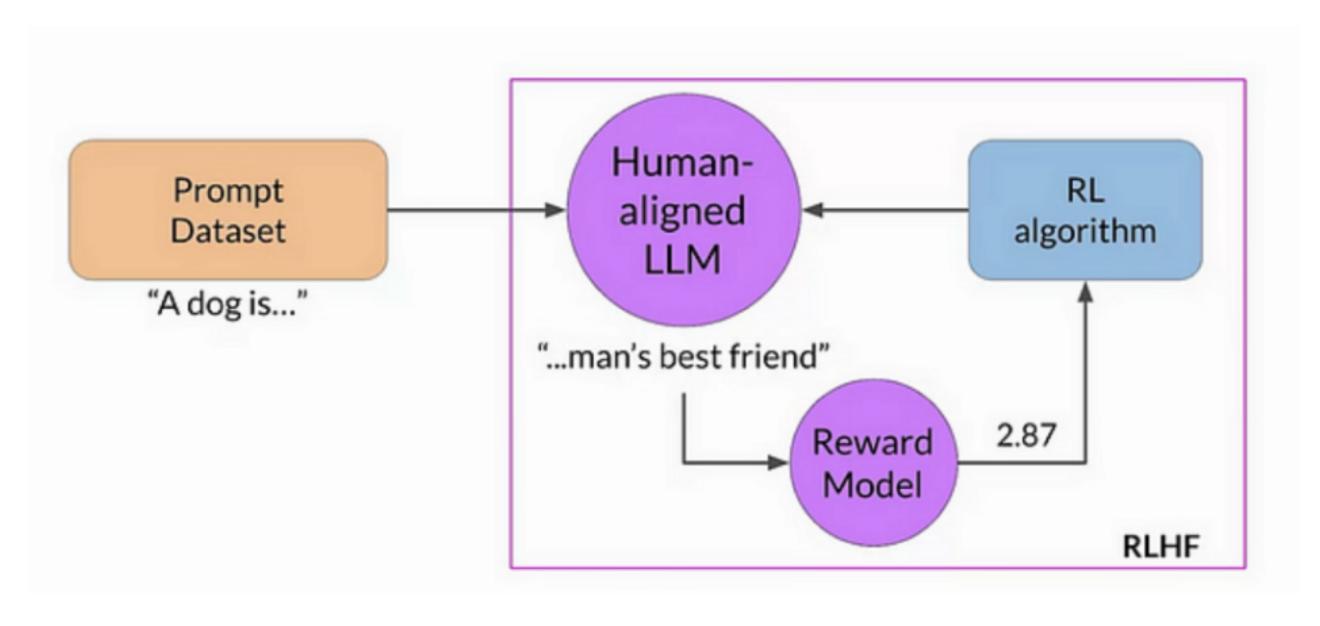


# Three stages of LLM training



The 3 steps of LLM training [1]

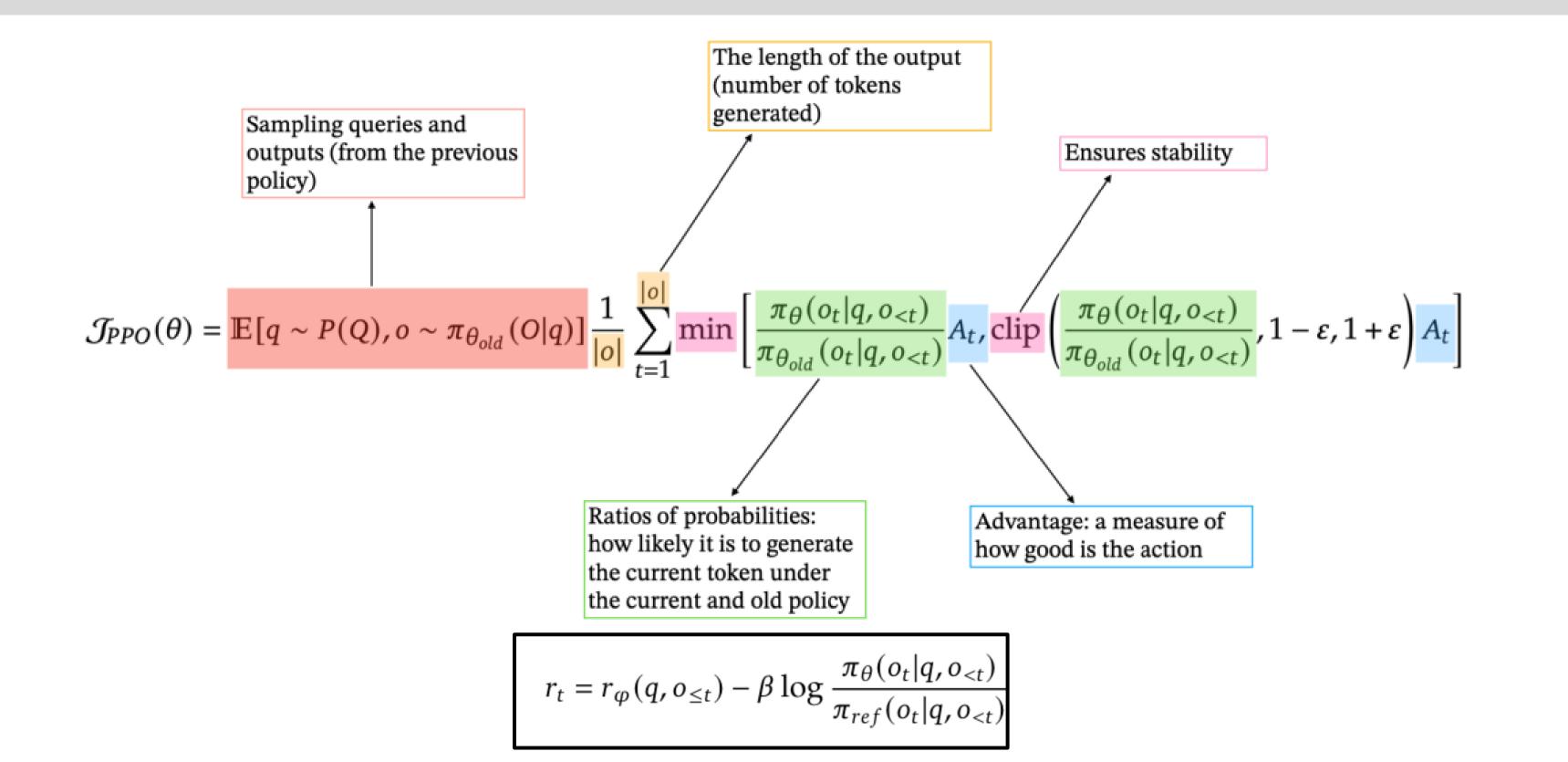
## RL in the context of LLMs



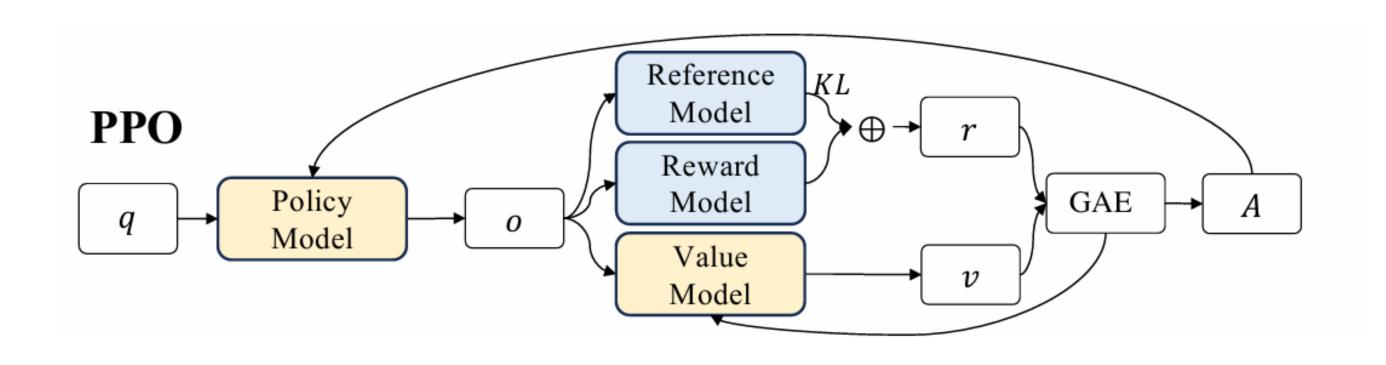
Simplified RLHF Process [3]



# **Proximal Policy Optimization (PPO)**

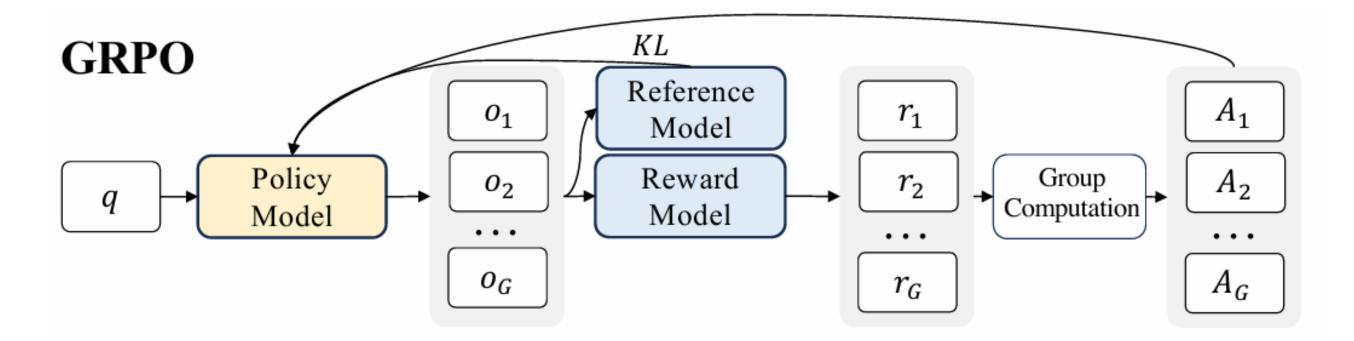


## **PPO to GRPO**



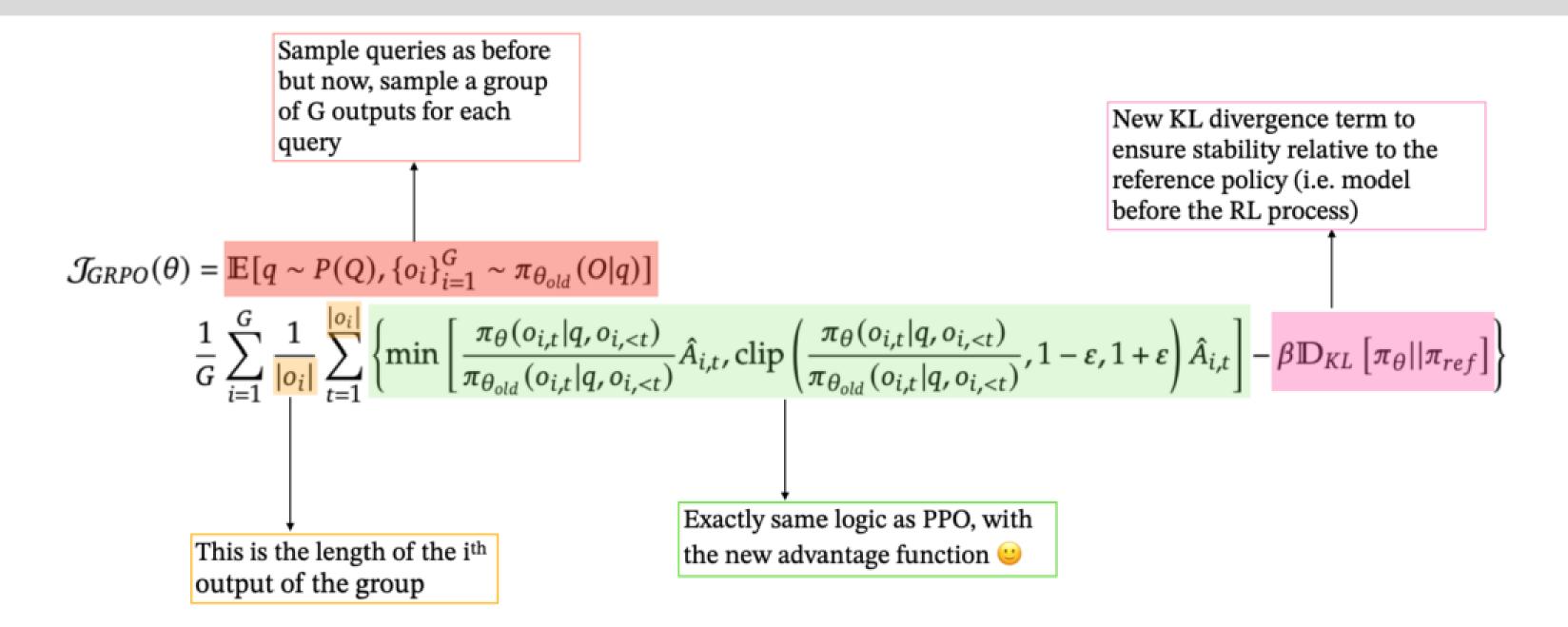
Trained Models

Frozen Models





# **Group Relative Policy Optimization (GRPO)**



$$\mathbb{D}_{KL}\left[\pi_{\theta}||\pi_{ref}\right] = \frac{\pi_{ref}(o_{i,t}|q,o_{i,< t})}{\pi_{\theta}(o_{i,t}|q,o_{i,< t})} - \log \frac{\pi_{ref}(o_{i,t}|q,o_{i,< t})}{\pi_{\theta}(o_{i,t}|q,o_{i,< t})} - 1$$



#### **GRPO Variants**

Outcome Supervision RL:

$$\hat{A}_{i,t} = \widetilde{r}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}$$

Process Supervision RL:

$$\mathbf{R} = \{\{r_1^{index(1)}, \cdots, r_1^{index(K_1)}\}, \cdots, \{r_G^{index(1)}, \cdots, r_G^{index(K_G)}\}\}$$

$$\tilde{r}_i^{index(j)} = \frac{r_i^{index(j)} - \text{mean}(\mathbf{R})}{\text{std}(\mathbf{R})}$$

$$\hat{A}_{i,t} = \sum_{index(j) > t} \tilde{r}_i^{index(j)}$$



#### **GRPO Variants**

#### **Algorithm 1** Iterative Group Relative Policy Optimization

```
Input initial policy model \pi_{\theta_{\text{init}}}; reward models r_{\varphi}; task prompts \mathcal{D}; hyperparameters \varepsilon, \beta, \mu
```

```
1: policy model \pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}
2: for iteration = 1, ..., I do
         reference model \pi_{ref} \leftarrow \pi_{\theta}
         for step = 1, \ldots, M do
             Sample a batch \mathcal{D}_b from \mathcal{D}
5:
             Update the old policy model \pi_{\theta_{old}} \leftarrow \pi_{\theta}
6:
             Sample G outputs \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot \mid q) for each question q \in \mathcal{D}_b
             Compute rewards \{r_i\}_{i=1}^G for each sampled output o_i by running r_{\varphi}
8:
             Compute \hat{A}_{i,t} for the t-th token of o_i through group relative advantage estimation.
9:
             for GRPO iteration = 1, ..., \mu do
```

10:

Update the policy model  $\pi_{\theta}$  by maximizing the GRPO objective (Equation 21) 11:

Update  $r_{\varphi}$  through continuous training using a replay mechanism. 12:

#### Output $\pi_{\theta}$

# Discussion: Code Training Benefits Mathematical Reasoning

Training Setting	Training Tokens			w/o Tool Use			w/ Tool Use		
Training Setting	General	Code	Math	GSM8K	MATH	CMATH	GSM8K+Python	MATH+Python	
No Continual Training	_	_	_	2.9%	3.0%	12.3%	2.7%	2.3%	
			Two-	Stage Tra	ining				
Stage 1: General Training	400B	_	_	2.9%	3.2%	14.8%	3.3%	2.3%	
Stage 2: Math Training	-	_	150B	19.1%	14.4%	37.2%	14.3%	6.7%	
Stage 1: Code Training	_	400B	_	5.9%	3.6%	19.9%	12.4%	10.0%	
Stage 2: Math Training	_	_	150B	21.9%	15.3%	39.7%	17.4%	9.4%	
One-Stage Training									
Math Training	_	_	150B	20.5%	13.1%	37.6%	11.4%	6.5%	
Code & Math Mixed Training	; –	400B	150B	17.6%	12.1%	36.3%	19.7%	13.5%	

Table 6 | Investigation of how code affects mathematical reasoning under different training settings. We experiment with DeepSeek-LLM 1.3B, and evaluate its mathematical reasoning performance without and with tool use via few-shot chain-of-thought prompting and few-shot program-of-thought prompting, respectively.



## Discussion: Arxiv Papers seem Ineffective in improving Math Reasoning

				English	Bench	marks	Chinese Benchmarks			
Model	Size	ArXiv Corpus	GSM8K	MATH	OCW	SAT	MMLU STEM	СМАТН	Gaokao MathCloze	Gaokao MathQA
DeepSeek-LLM 1		No Math Training	2.9%	3.0%	2.9%	15.6%	19.5%	12.3%	0.8%	17.9%
	1.3B	MathPile ArXiv-RedPajama	2.7% 3.3%	3.3% 3.4%		12.5% 9.4%	15.7% 9.0%	1.2% 7.4%	0.0% 0.8%	2.8% 2.3%
		No Math Training	29.0%	12.5%	6.6%	40.6%	38.1%	45.9%	5.9%	21.1%
DeepSeek-Coder-Base-v1.5	7B	MathPile ArXiv-RedPajama	23.6% 28.1%	11.5% 11.1%		46.9% 50.0%		37.9% 42.6%	4.2% 7.6%	25.6% 24.8%

Table 8 | Effect of math training on different arXiv datasets. Model performance is evaluated with few-shot chain-of-thought prompting.

ArXiv Corpus	miniF2F-valid	miniF2F-test
No Math Training	20.1%	21.7%
MathPile ArXiv-RedPajama	16.8% 14.8%	16.4% 11.9%

Table 9 | Effect of math training on different arXiv corpora, the base model being DeepSeek-Coder-Base-v1.5 7B. We evaluate informal-to-formal proving in Isabelle.



# Discussion: Why RL works?

- ➤ RL enhances Maj@K's performance but not Pass@K.
- It seems that the improvement is attributed to boosting the correct response from TopK rather than the enhancement of fundamental capabilities.

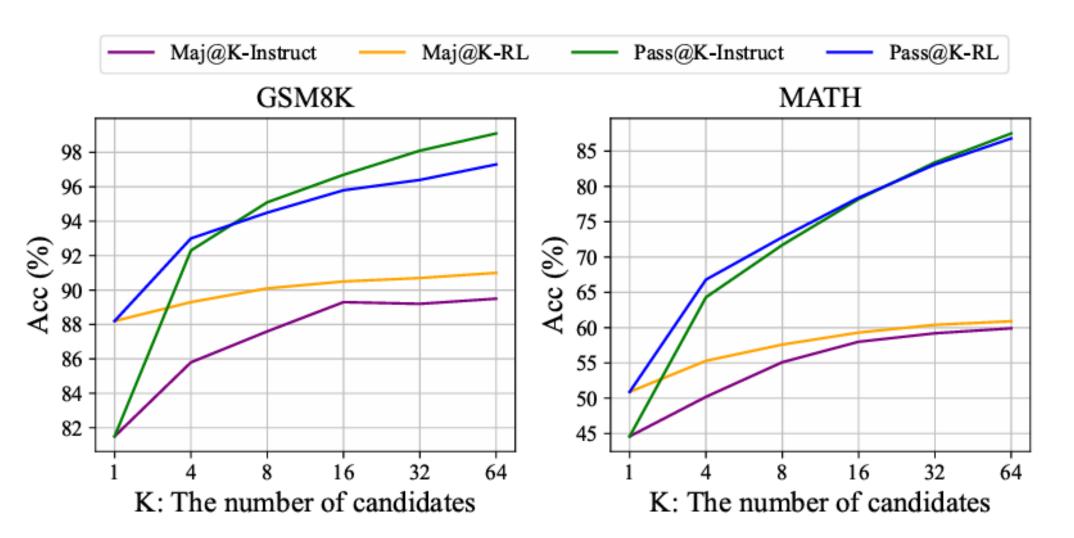


Figure 7 | The Maj@K and Pass@K of SFT and RL DeepSeekMath 7B on GSM8K and MATH (temperature 0.7). It was noted that RL enhances Maj@K but not Pass@K.



### Conclusion

- ➤ **Data Curation:** DeepSeekMath surpasses all open-source models on the MATH benchmark and nears closed-source model performance through large-scale training with rich mathematical web data.
- ➤ Algorithm: They introduce Group Relative Policy Optimization (GRPO) which effectively polishes its response distribution with lower memory usage compared to PPO.
- Limitations & Future Work: The model struggles with geometry and theorem-proof tasks and lacks strong few-shot learning ability. Their future work aims to refine data selection and reinforcement learning approaches to address these gaps.



