

Can we scale small LMs to o1 performance? **A Probabilistic Inference Approach to LLM Inference-Time Scaling**



Isha Puri, MIT CSAIL





Results

Model	Method	MATH500	AIME 2024
Closed-Source LLMs			
GPT-40	-	76.2	13.3
o1-preview	-	87.0	40.0
Claude3.5-Sonnet	-	78.3	16.0
Open-Source LLMs			-
Llama-3.1-70B-Instruct	-	65.7	16.6
Qwen 2.5 - Math - 72 B - Instruct	-	82.0	30.0
Open-Source SLMs			
Llama-3.2-1B-Instruct	Pass@1	26.8	0.0
	Ours - PF	59.6	10.0
Llama-3.1-8B-Instruct	Pass@1	49.9	6.6
	Ours - PF	74.4	16.6
phi-4	Pass@1	79.8	16.6
	Ours - PF	83.6	26.6
Mistral-Small-24B-Instruct-2501	Pass@1	69.2	10.0
	Ours - PF	83.4	23.3
Qwen2.5-32B-Instruct	Pass@1	82.8	16.6
	Ours - PF*	89.9	43.3
Open-Source Math SLMs			
Qwen2.5-Math-1.5B-Instruct	Pass@1	70.0	10.0
	Ours - PF	85.4	23.3
Qwen2.5-Math-7B-Instruct	Pass@1	79.6	16.6
	Ours - PF	87.0	23.3

- A 7B model surpasses o1 accuracy in 32 rollouts
- A 1.5B model surpasses GPT4o in only 4 rollouts
- Our method (PF) has a 4-16x better scaling rate than other inference scaling methods / deterministic counterparts
- All of this is done *without any training at all!* Just by intelligently / probabilistically navigating <u>the search space.</u>



- This inference scaling method allows any off the shelf model to punch way above its weight class just by sampling it a few times!

Coauthors: Shiv Sudalairaj, GX Xu, Kai Xu, Akash Srivastava

Scan for more information! \rightarrow

