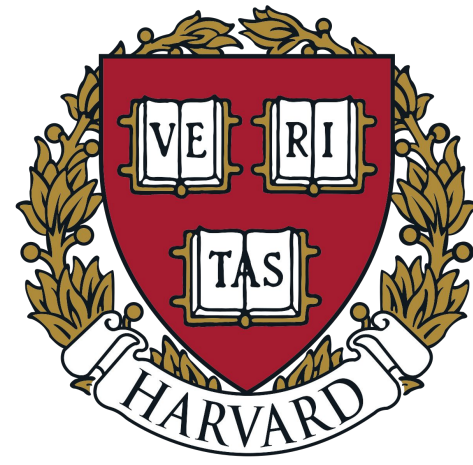


HARDMath: A Benchmark Dataset for Challenging Problems in Applied Mathematics



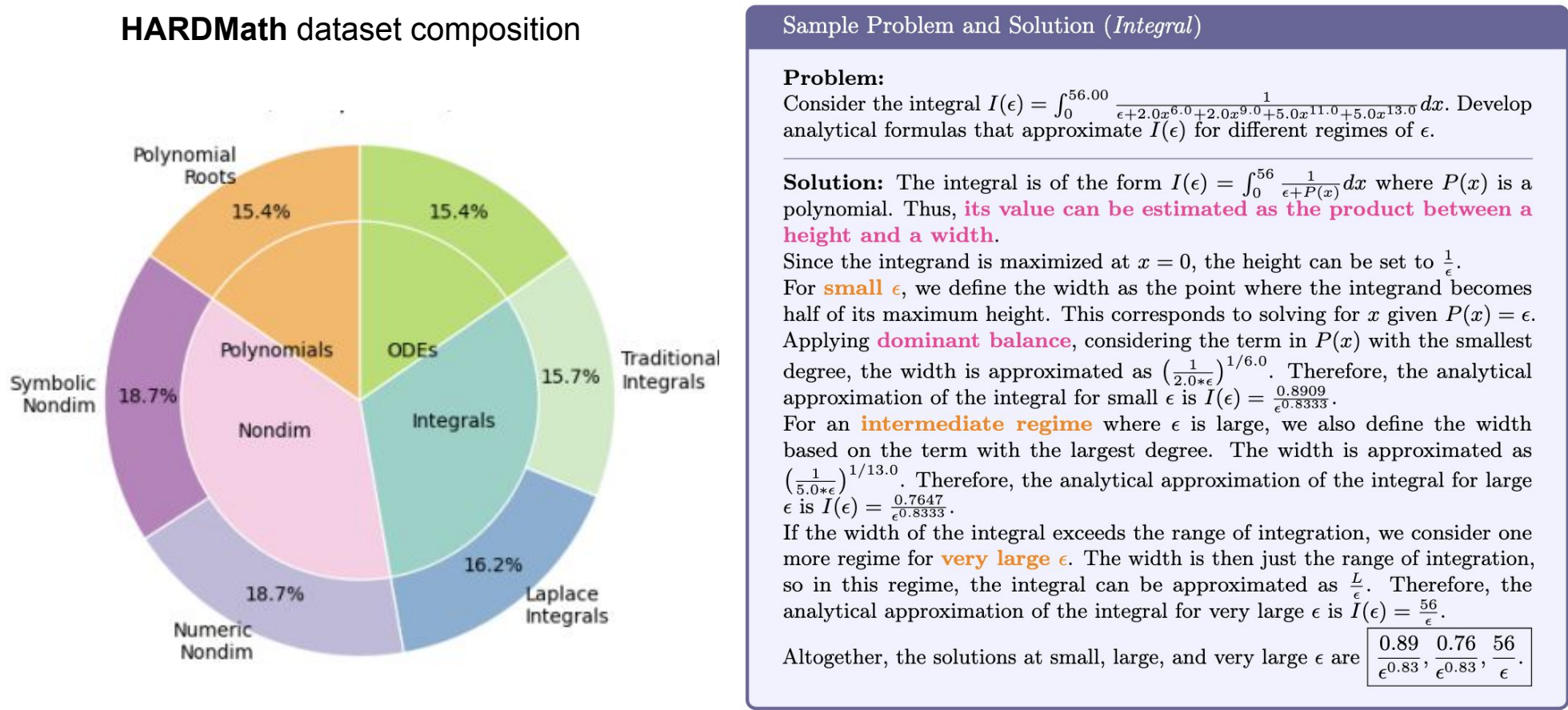
Jingxuan Fan (jfan@g.harvard.edu)*, Sarah Martinson*, Erik Y. Wang*, Kaylie Hausknecht*, Jonah Brenner, Danxian Liu, Nianli Peng, Corey Wang, Michael Brenner
School of Engineering and Applied Sciences, Harvard University

The HARDMath Dataset

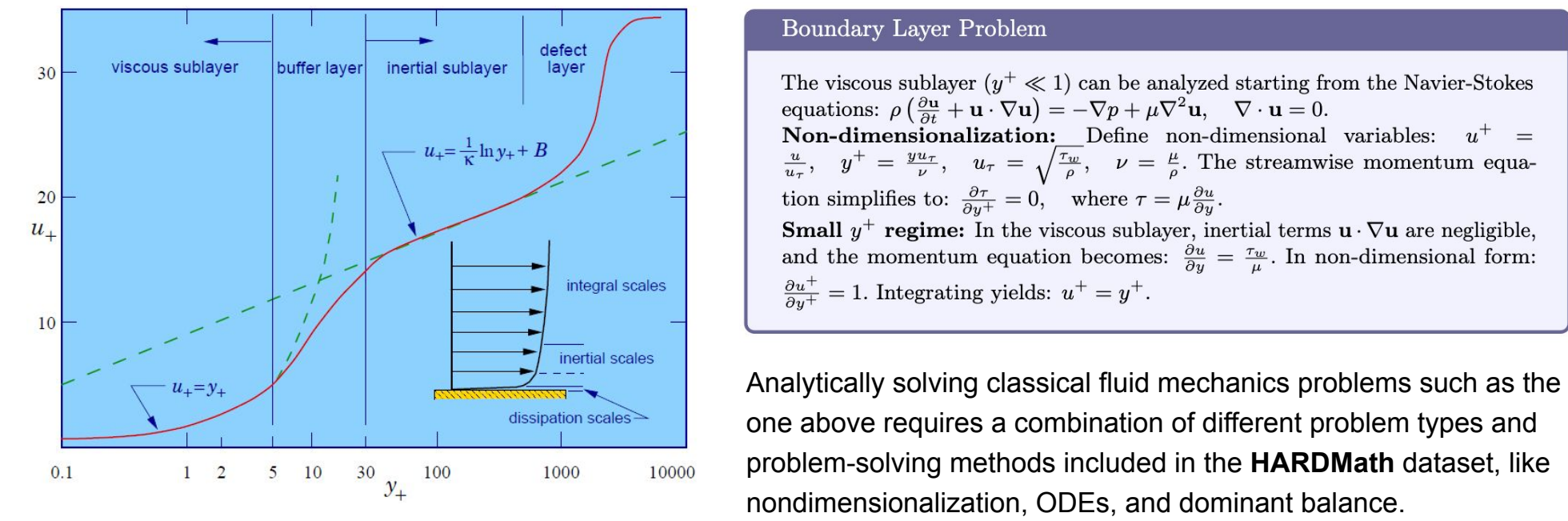
Problem: existing benchmarks lack graduate-level applied math!

We provide:

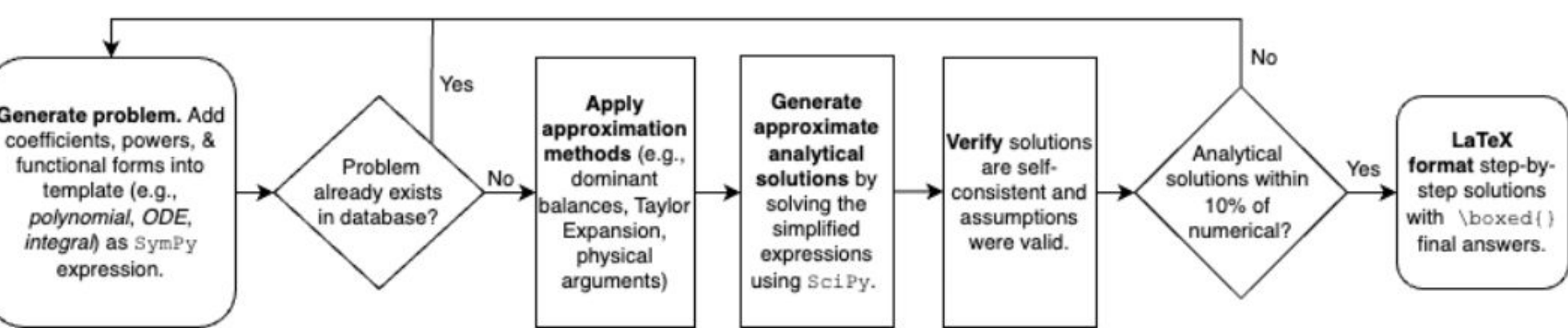
- ➔ 1426 graduate-level applied math problems across 6 domains, targeting **approximation reasoning & asymptotic methods**
- ➔ An **algorithmic framework** that generates questions and solutions with quality guarantees



Why HARDMath? Scientific Research!



Dataset Generation



- ➔ Parameter combination checks ensure generating unique problems
- ➔ Leveraging SymPy and SciPy to perform the calculations required to derive our approximate, analytical solutions
- ➔ Checking against numerical solver “ground-truth” results ensures correctness of analytical approximation provided in the solution

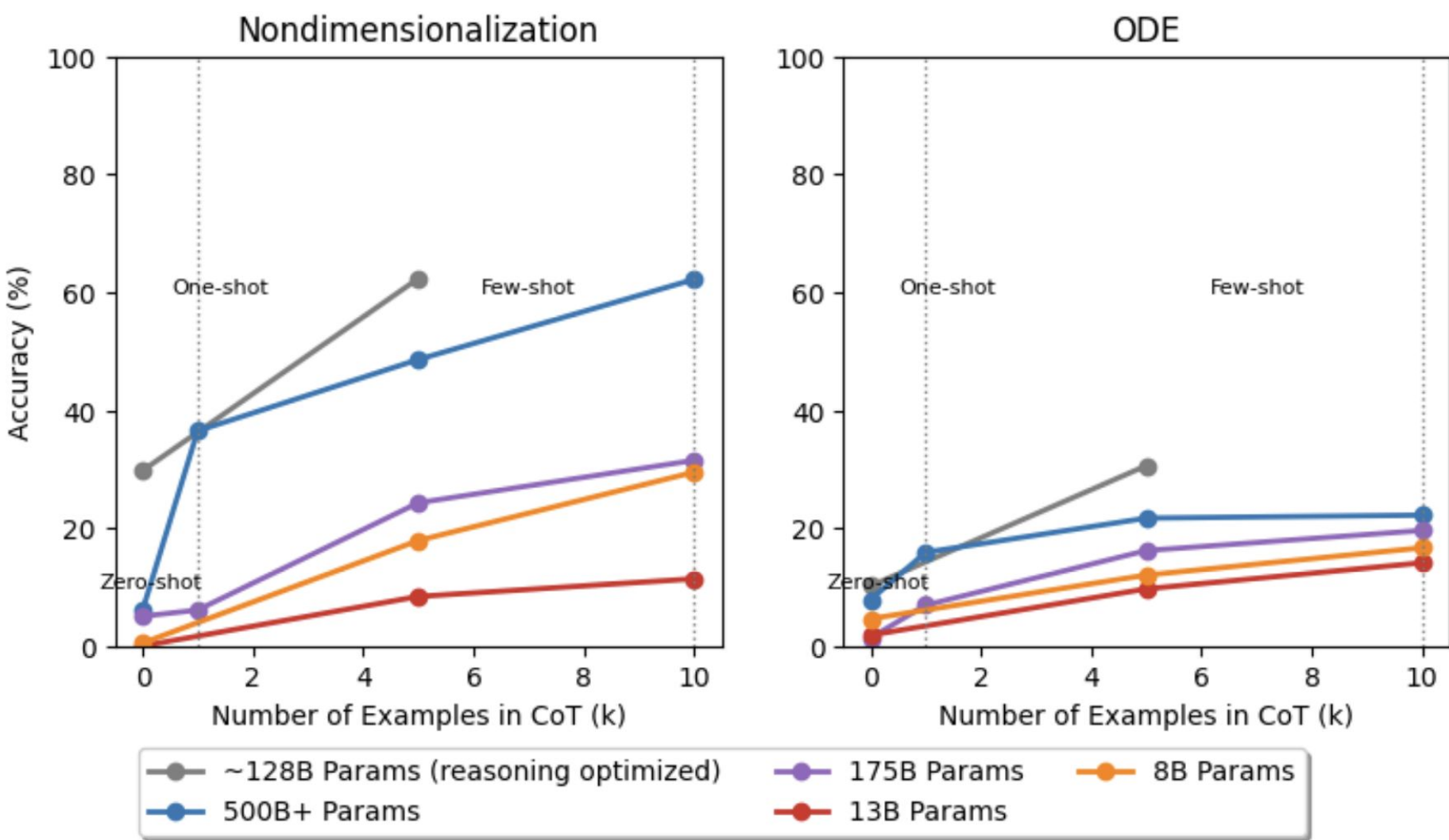
Evaluation

- ➔ Answer + **procedural grading** with **LLM grader and rubrics**
- ➔ Evaluations on frontier models including **STEM-optimized o1**
- ➔ Detailed **few-shot performance** behavior and **error analysis**

Evaluation accuracy (percentage) on HARDMath-mini					
Model	ALL	Nondim	Roots	ODEs	Integrals
Closed-source models					
GPT-3.5 (0 shot)	6.04	5.05	17.2	1.39	3.33
GPT-3.5 (1 shot CoT)	14.2	6.11	29.3	6.94	18.2
GPT-3.5 (5 shot CoT)	24.6	24.3	35.0	16.2	23.1
GPT-4 (0 shot)	14.0	6.04	33.7	7.87	14.9
GPT-4 (1 shot CoT)	37.6	36.5	52.8	15.9	40.5
GPT-4 (5 shot CoT)	43.8	48.6	57.3	21.7	41.4
o1-mini (0 shot CoT)	29.8	38.1	24.3	10.2	32.5
o1-mini (5 shot CoT)	62.3	84.5	62.1	30.6	46.5
Open-source models					
Llama3-8b (0 shot)	3.67	0.50	11.5	4.63	2.52
Llama3-8b (5 shot CoT)	20.2	17.9	17.1	12.0	28.1
CodeLlama-13b (0 shot)	1.94	0.00	8.73	1.85	0.50
CodeLlama-13b (5 shot CoT)	9.79	8.41	13.1	9.7	9.57

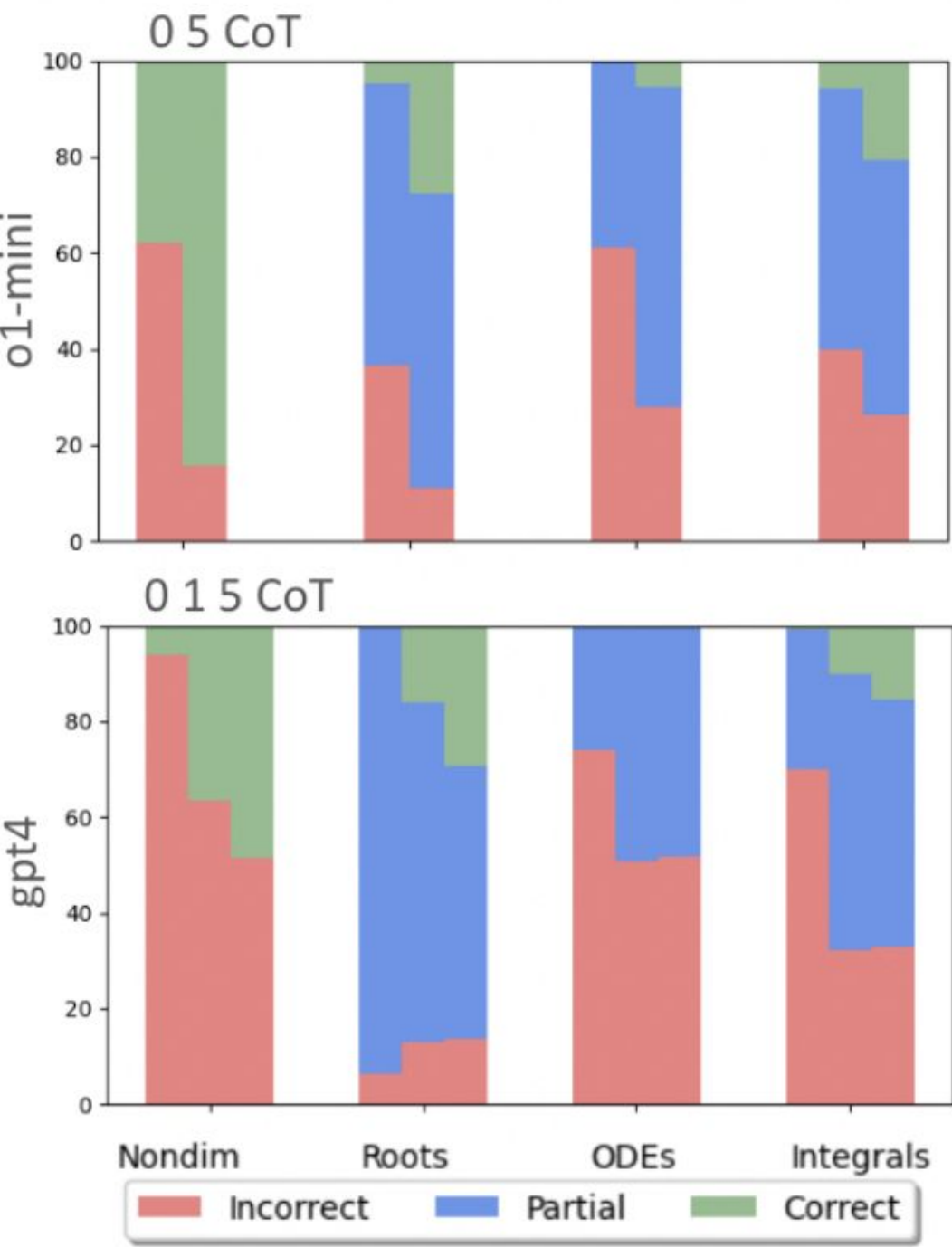
Model	HARDMath	GSM8K	MATH
Closed-source models			
GPT-4 (0 shot)	14.0		73.4 ²
GPT-4 (5 shot CoT)	43.8	92.0 ³	
o1-mini (0 shot)	29.8		90.0 ⁴
o1-mini (5 shot CoT)	62.3		
Open-source models			
Llama3-8b (5 shot CoT)	20.2	79.6 (8-shot) ⁵	30.0 (4-shot) ⁵

- ➔ Performance on **HARDMath** is much worse than widely used math benchmarks under comparable prompting conditions, **even with partial credit**
- ➔ Offers new opportunities for evaluation and model development



Performance gain from few-shot CoT prompting is model size and problem type dependent: the more challenging type ODE shows minimal performance increase after few-shot CoT

Error Analysis



Example error composition analysis for *Polynomial Roots* problems

Polynomial Roots (error mode: incorrect dominant balance terms)

Question: Consider the polynomial $P(x) = \epsilon x^6 + x^4 - 1$. Find approximate expressions for all roots of the polynomials in the limit of small positive ϵ and large positive ϵ (single term approximation).

Model Response: For large ϵ When ϵ is large, the term ϵx^6 dominates the polynomial for most x . However, focusing on the dominant term gives: $P(x) \approx \epsilon x^6$. Setting this equal to zero for an approximation suggests that we don't get meaningful roots from this term alone since $x^6 = 0$ doesn't provide roots that satisfy the original polynomial's requirement for non-zero roots.

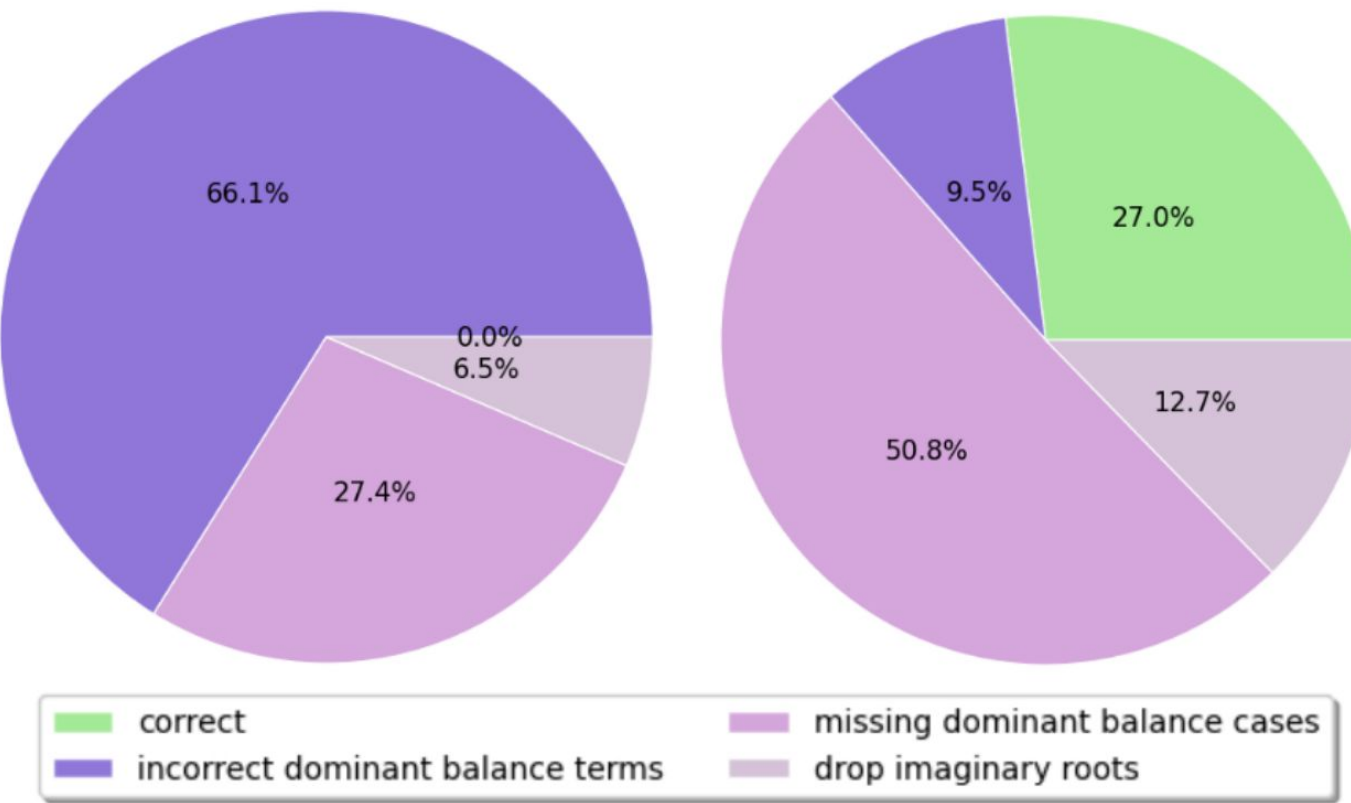
Grading: The response wrongly sets up the dominant balance as $A = 0$ and does not provide a concrete list of roots for large ϵ . Therefore, score for large positive ϵ is $[0.0]$

Polynomial Roots (error mode: missing dominant balance cases)

Question: Consider the polynomial $P(x) = \epsilon x^6 + x^4 - 1$. Find approximate expressions for all roots of the polynomials in the limit of small positive ϵ and large positive ϵ (single term approximation).

Model Response: When ϵ is small, the term ϵx^6 becomes negligible, so we balance term B and C and the polynomial simplifies to: $P(x) \approx x^4 - 1$. Solving the equation gives roots $[1, -1, i, -i]$.

Grading: The response only includes the roots from the balance $B + C = 0$ and completely misses the roots from the balance $A + B = 0$. Therefore, score for small positive ϵ is $[0.5]$



GPT-4 0 shot CoT error composition

GPT-4 5 shot CoT error composition

- ➔ Few-shot CoT prompting increases model performance by changing correct-partial compositions
- ➔ Highly dependent on model and problem types

References

- LearnCAX (n.d.). “Basics of y-plus, boundary layer, and wall function in turbulent flows.” In: LearnCAX Knowledge Base. Available at: <https://learncax.com>
- OpenAI (2024). “Simple Evals.” In: GitHub Repository. Available at: <https://github.com/openai/simple-evals>
- Achiam, Josh, Adler, Steven, Agarwal, Sandhini, Ahmad, Lama, Akkaya, Ilge, Aleman, Florencia Leoni, Almeida, Diogo, Altschmidt, Janko, Altman, Sam, Anadkat, Shyamal, et al. (2023). “GPT-4 Technical Report.” In: arXiv preprint arXiv:2303.08774.
- OpenAI (2024). “OpenAI o1-mini.” In: Advancing Cost-Efficient Reasoning. Available at: <https://openai.com/index/openai-o1-mini-advancing-cost-efficient-reasoning/>
- Meta AI (2024). “Meta Llama 3.” In: Meta AI Blog. Available at: <https://ai.meta.com/blog/meta-llama-3/>

Arxiv



Github

