

# HARDMath2: A Benchmark for Applied Mathematics Built by Students as Part of a Graduate Class

**James V. Roggeveen, Erik Y. Wang,** David Ettel, Will Flintoft, Peter Donets, Lucy S. Nathwani, Nickholas Gutierrez, Anton Marius Graf, Siddharth Dandavate, Arjun Nageswaran, Raglan Ward, Ava Williamson, Anne Mykland, Kacper K. Migacz, Yijun Wang, Egemen Bostan, Duy Thuc Nguyen, Zhe He, Marc L. Descoteaux, Felix Yeung, Shida Liu, Jorge García Ponce, Luke Zhu, Yuyang Chen, Ekaterina S. Ivshina, Miguel Fernandez, Minjae Kim, Kennan Gumbs, Matthew Scott Tan, Russell Yang, Mai Hoang, David Brown, Isabella A. Silveira, Lavon Sykes, Ahmed Roman, William Fredenberg, Yiming Chen, Lucas Martin, Yixing Tang, Kelly Werker Smith, Hongyu Liao, Logan G. Wilson, Alexander Dazhen Cai, Andrea Elizabeth Biju, Michael P. Brenner  
*School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138*

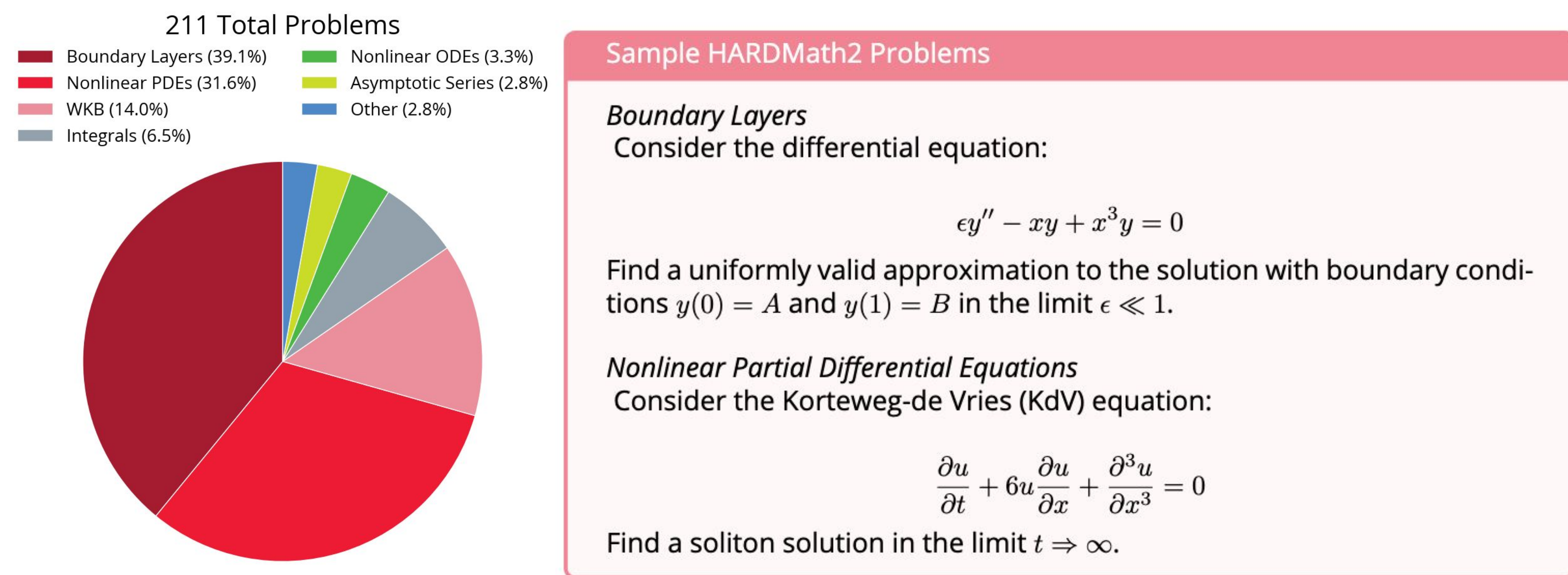
## Why HARDMath2?

- SOTA math benchmarks focus on **exact solutions** or **formal proofs**
- Physics and engineering mathematical problems rely on **approximations**, **asymptotics**, and **perturbation methods** that are underrepresented in current evaluations

### We provide:

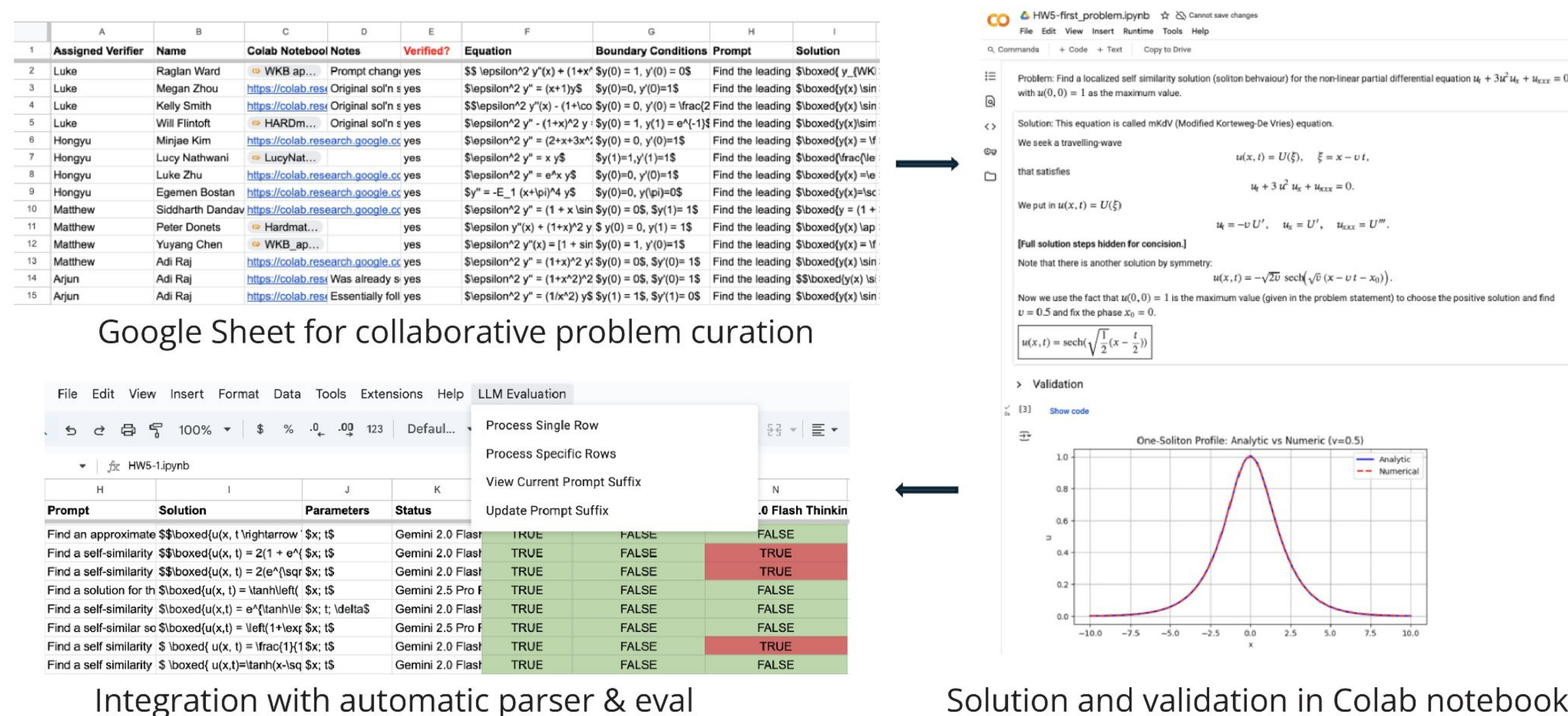
- 211 **original** graduate-level applied math problems across 5+ domains, targeting approximation reasoning & asymptotic methods
- A custom evaluation framework that avoids LLM-as-a-judge
- Unique pedagogical approach combining benchmark creation & education

## The HARDMath2 Dataset

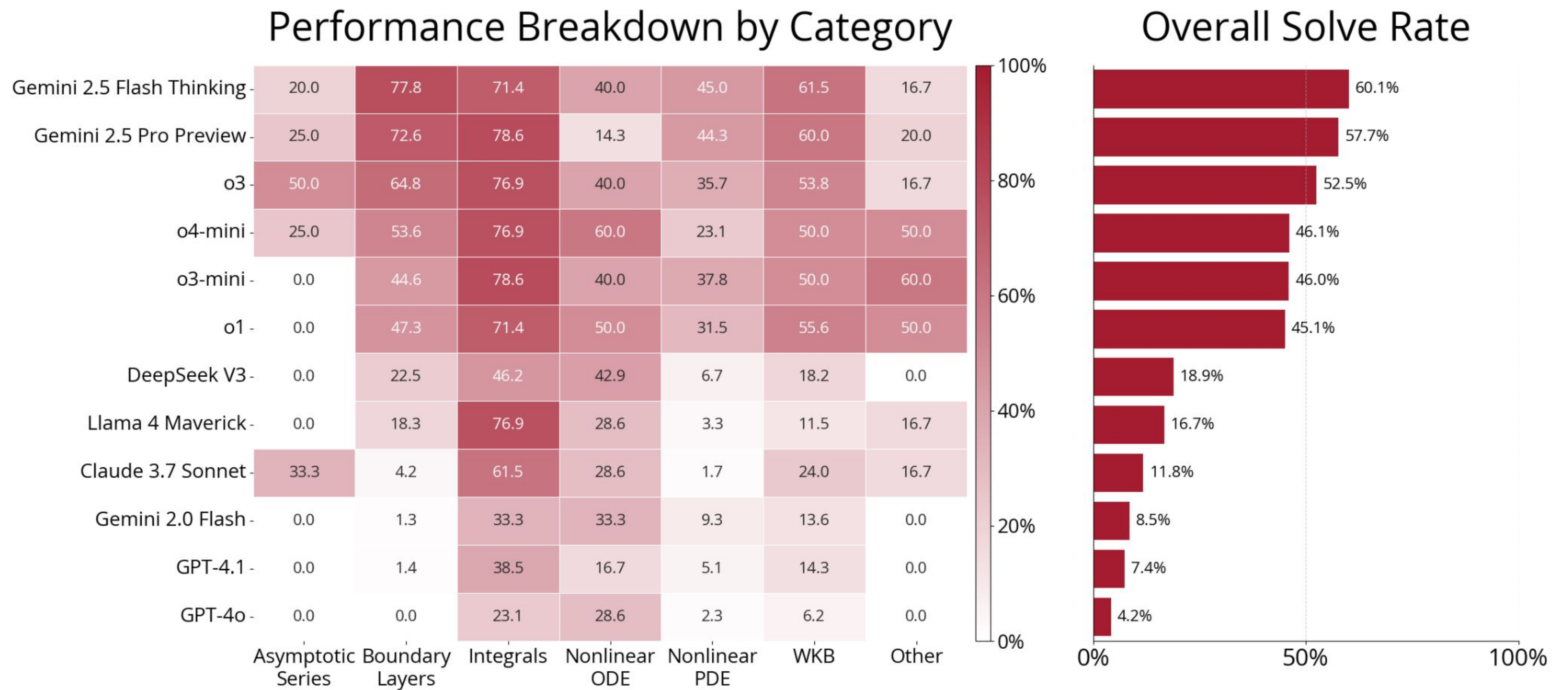


## Pedagogical Approach

- Built by students as part of a Harvard graduate applied math class
- Homework involved creating problems that SOTA models **could not solve**
- Students' problems were evaluated by custom Google Sheet integration, giving **real-time feedback** on problem difficulty
- Peer verification for correctness of solutions
- By competing against models, students identified **critical failure modes**



## Model Performance



## Evaluation Scheme

- HARDMath2 implements a **custom LaTeX to SymPy parser** to avoid LLM-as-a-judge
- Numerical evaluation accommodates equivalent forms of expressions without needing symbolic simplification
- Model and reference solutions evaluated numerically and compared to determine equivalence up to some tolerance
- Problems judged based on equivalence to ground-truth solution **without** subjective partial credit scores

## Conclusions

- **Overall Performance:** reasoning models struggle, hovering around 50-60% accuracy.
- Non-reasoning models perform significantly worse (<10% on many categories).
- **Specific Weaknesses:** Models are weakest on Nonlinear PDEs (requires physical ansatz selection) and Boundary Layers (requires matching inner/outer solutions).

## Common Failure Modes

- Incorrect dominant balances
- Solving to a different asymptotic order than specified
- Not fully solving the problem by leaving off key matching steps to get final answer (see below sample)
- Failure to follow formatting guidelines

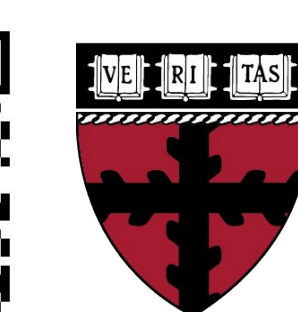
### Sample Gemini 2.5 Pro Output for Boundary Layer

The leading order behavior of  $y(x)$  as  $x \rightarrow 0^+$  is given by  $y(x) \sim Ax^{4/3} \exp(-\frac{3}{2x^{1/3}}) \cos(\frac{3\sqrt{3}}{2x^{1/3}} - \phi)$  where  $A$  is a positive constant, and  $\phi$  is a phase constant. Since the constants  $A$  and  $\phi$  are fixed by the boundary conditions and are non-zero, this gives the leading behavior.

Final Answer: The final answer is  $y(x) \sim Ax^{4/3} e^{3/(2x^{1/3})} \cos\left(\frac{3\sqrt{3}}{2x^{1/3}} - \phi\right)$

**Note:** The LLM makes some attempt to narrow down what the constants are, but it does not solve for them completely. It mentions that the initial conditions fix the constants, but it does not list the equations it would use to solve for these constants, nor does it solve them.

Paper link:



**Harvard John A. Paulson School of Engineering and Applied Sciences**

