

SorryDB: Can AI Provers Complete Real-World Lean Theorems?

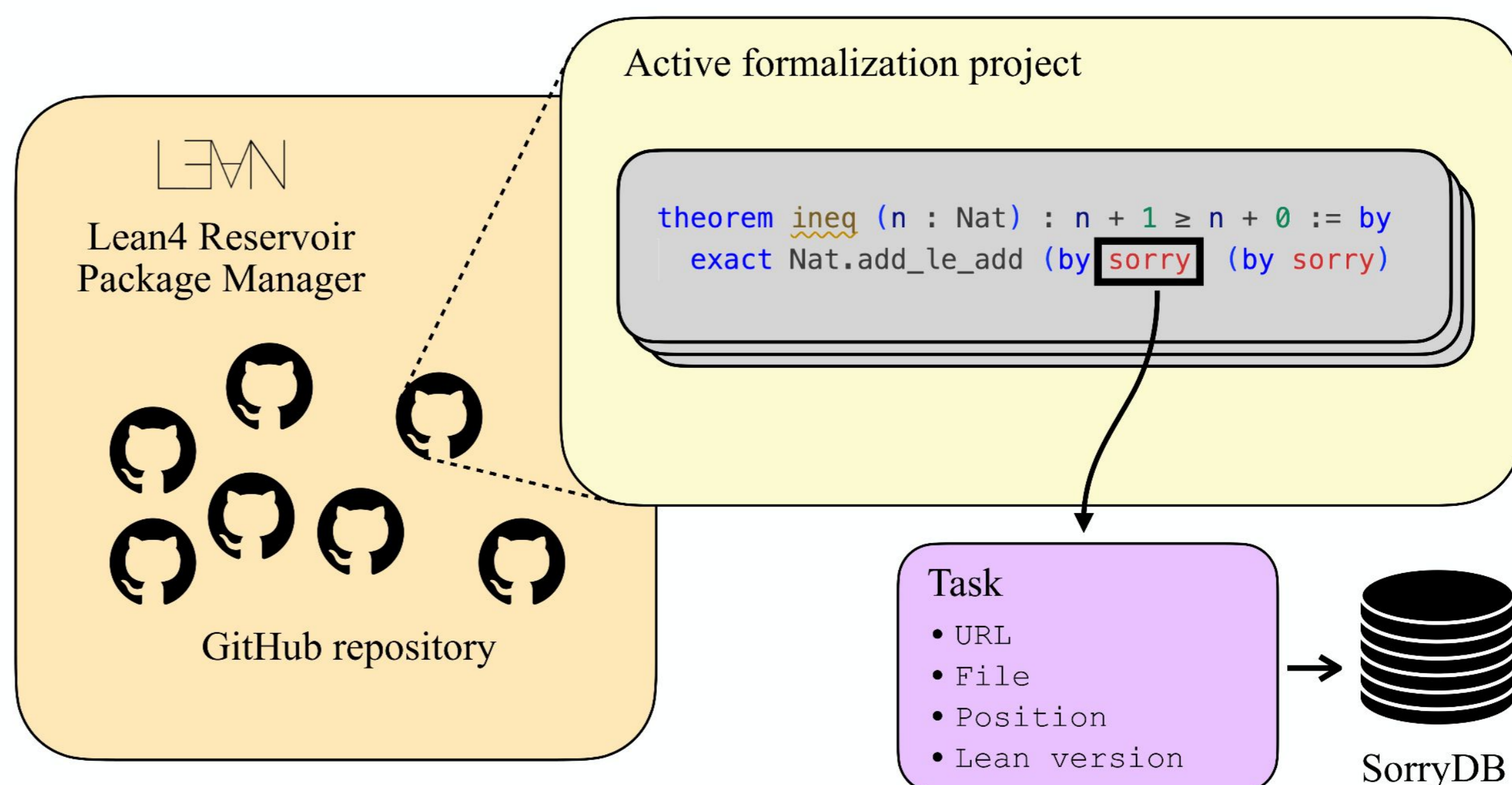
Austin Letson¹, Leopoldo Sarra¹, Auguste Poiroux^{2,3}, Oliver Dressler, Paul Lezeau^{5,6}, Dhyan Aranha^{7,12}, Frederick Pu⁸, Aaron Hill, Miguel Corredera Hidalgo⁹, Julian Berman¹⁰, George Tsoukalas¹¹, Lenny Taelman⁷

¹Axiomatic AI, ²Math, Inc., ³EPFL, ⁵The London School of Geometry and Number Theory, ⁶Imperial College, ⁷University of Amsterdam, ⁸University of Toronto, ⁹ENSEIRB-MATMECA, INP-Bordeaux, ¹⁰Columbia University, ¹¹The University of Texas at Austin, ¹²Côte d'Azur University,

The SorryDB Dataset

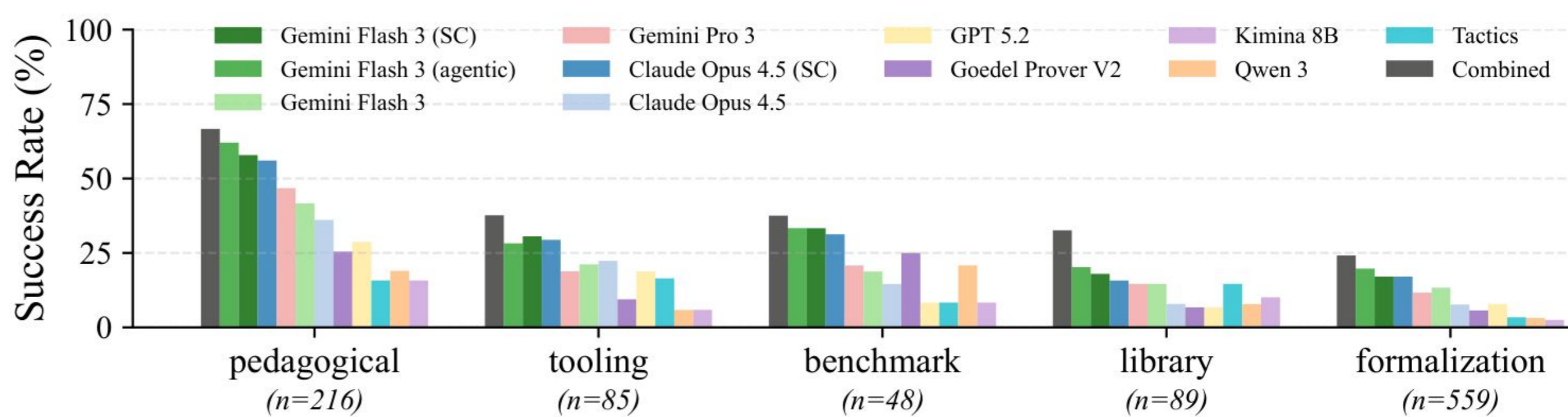
SorryDB is a dynamically-updating benchmark of open Lean tasks drawn from real world formalization projects on GitHub.

Hillclimbing the SorryDB benchmark will yield tools that are aligned to the community needs, more usable by mathematicians, and more capable of understanding complex dependencies.



Evaluation

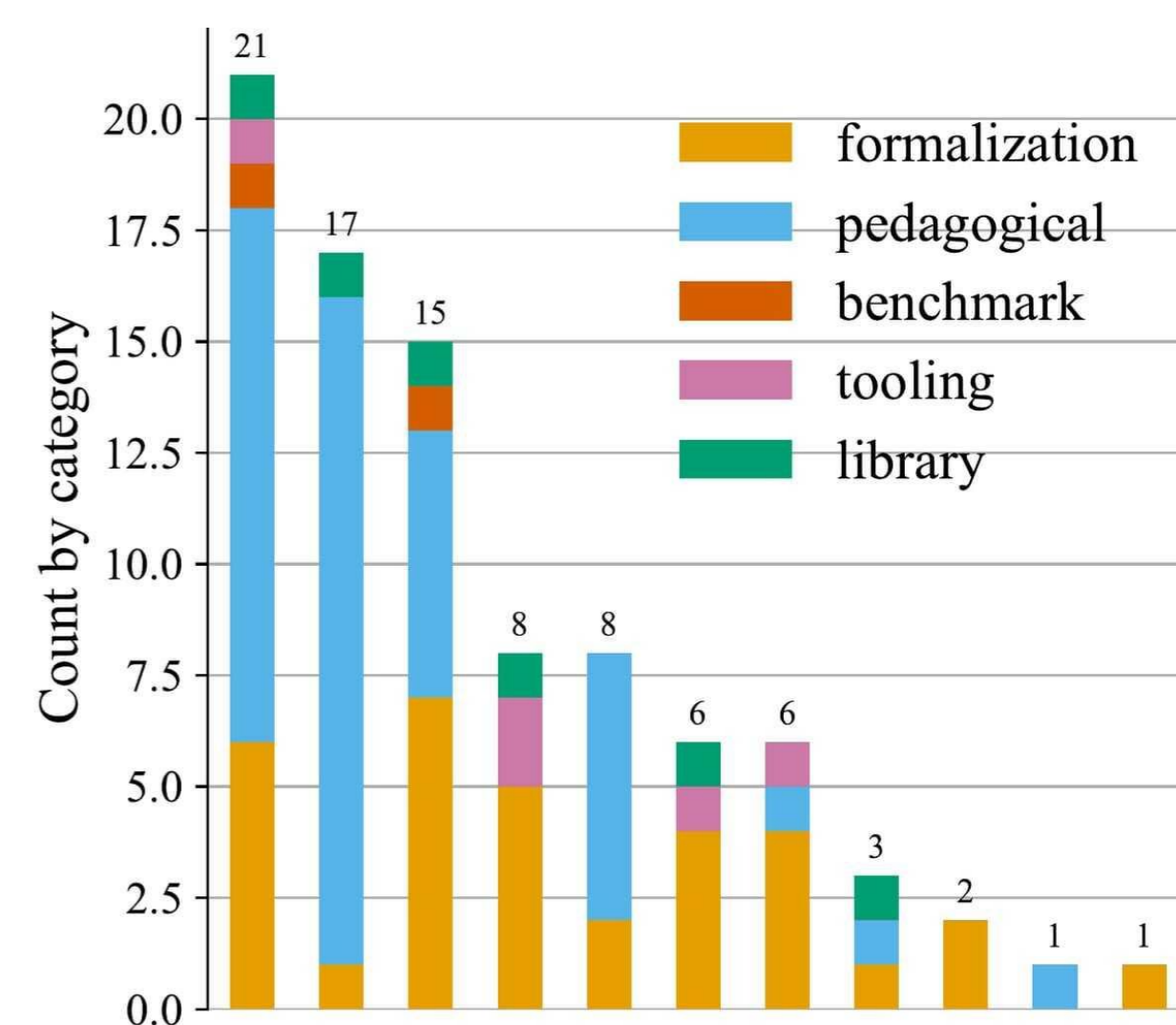
APPROACH	PASS@1	PASS@32
<i>Deterministic</i>		
TRIVIAL	2.1 %	N/A
TACTICS	8.4 %	N/A
<i>General-purpose LLM</i>		
GPT 5.2	6.2 %	13.2 %
CLAUDE OPUS 4.5	7.8 %	15.4 %
GEMINI FLASH 3	10.8 %	20.5 %
GEMINI PRO 3	11.0 %	20.5 %
QWEN 3	5.0 %	8.1 %
<i>Specialized LLM</i>		
KIMINA PROVER 8B	1.0 %	6.6 %
GOEDEL PROVER V2 32B	2.7 %	11.3 %
<i>Iterative</i>		
CLAUDE OPUS 4.5 (SC)	27.1 %	N/A
GEMINI FLASH 3 (SC)	27.9 %	N/A
GEMINI FLASH 3 (AGENTIC)	30.3 %	N/A
COMBINED	35.7 %	



Evaluation by repo category

Pedagogical and benchmark repos are easier than larger formalizations

Try it here!



Evaluation results

- Deterministic strategies already show baseline performance
- Iterative strategies outperform pass@k
- General purpose LLMs outperform specialized provers

Provers are complementary

The most powerful prover does not cover all solved tasks.

Upset plot of sorries solved by each model. Each column represents the number of theorems solved by all the provers marked in the table with a black circle, and not solved by the others. Deterministic tactics solve some sorries that none of the LLM-based approaches could solve.

