



Our code

# Prover Agent: An Agent-Based Framework for Formal Mathematical Proofs

Kaito Baba<sup>1</sup>, Chaoran Liu<sup>2</sup>, Shuhei Kurita<sup>3,2</sup>, Akiyoshi Sanna<sup>4,5,6,2,7</sup>

<sup>1</sup>The University of Tokyo, Tokyo, Japan <sup>2</sup>Research and Development Center for Large Language Models, National Institute of Informatics, Tokyo, Japan

<sup>3</sup>National Institute of Informatics, Tokyo, Japan <sup>4</sup>Kyoto University, Kyoto, Japan <sup>5</sup>Shiga University, Shiga, Japan

<sup>6</sup>RIKEN Center for Advanced General Intelligence for Science Program, Kobe, Japan <sup>7</sup>National Institute of Science Technology Policy (NISTEP), Tokyo, Japan



Our paper

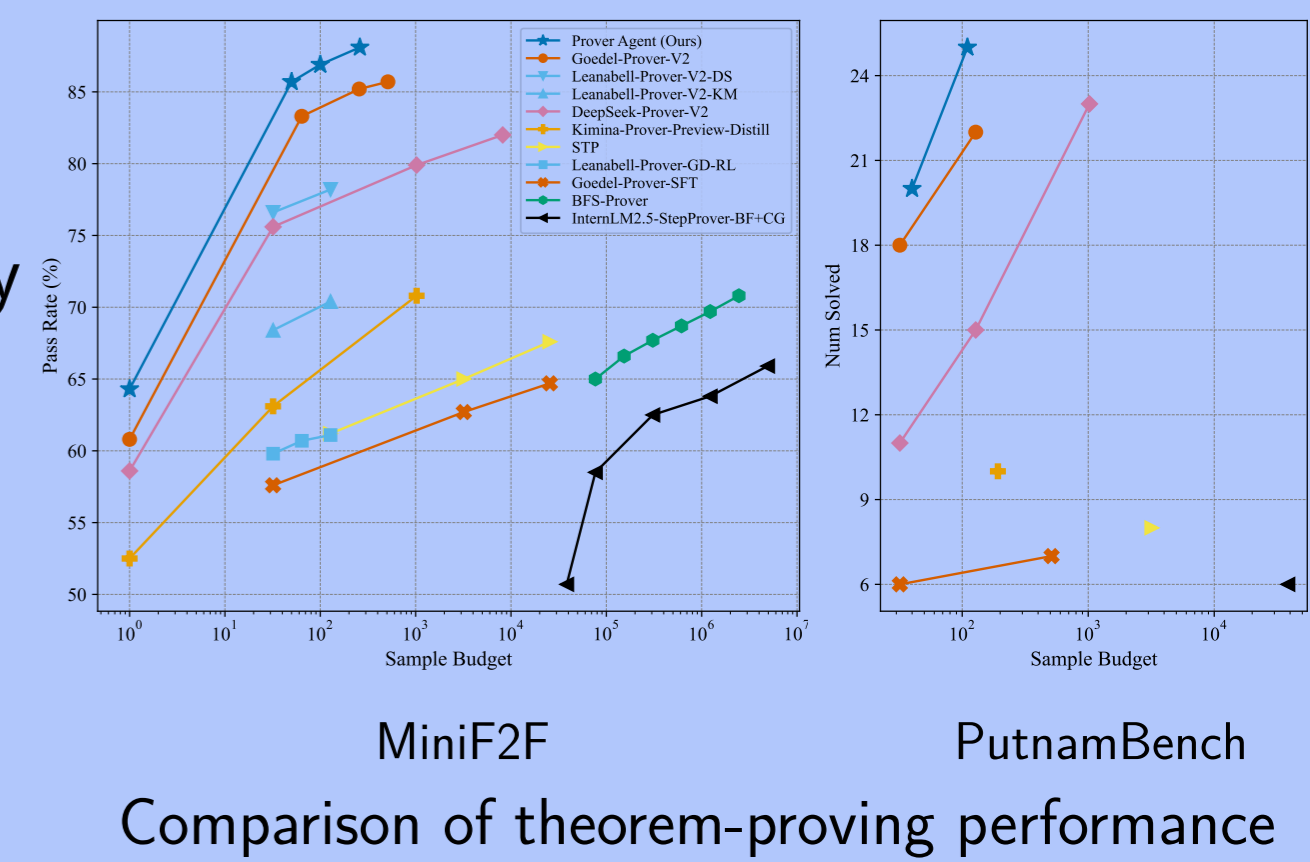
## Motivation

- Large language models (LLMs):
  - Capable of powerful reasoning and generation
  - Prone to errors and hallucinations
- Formal proof assistants (e.g., Lean):
  - Verify mathematical correctness
  - Not generative; requires painstaking meticulous detail
- LLM-based formal proving is gaining attention
- Yet, a large gap remains between informal reasoning and formal proving

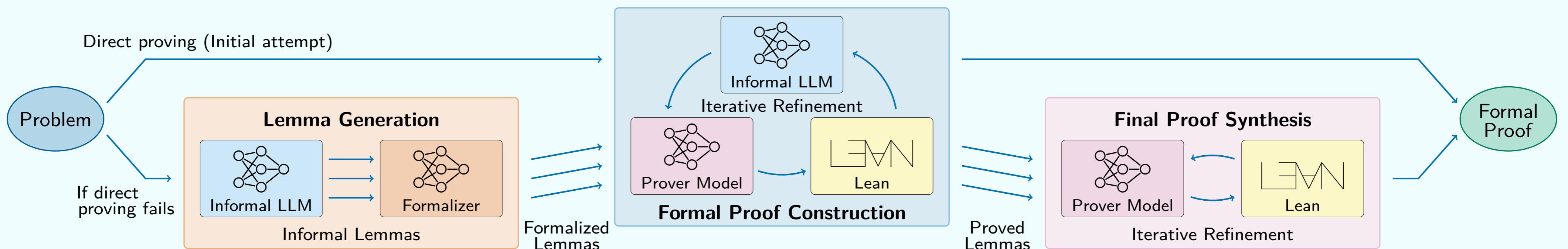
Our Goal: Bridge this gap

## Our Contributions

- Coordination of informal and formal reasoning with Lean feedback
- Auxiliary lemma generation for strategy discovery
  - Helps discover strategies even when the solution path is not apparent at first
- State-of-the-art theorem-proving performance among methods using small language models
- Efficiency in inference-time cost
  - Much smaller sample budget than prior work



## Prover Agent



### Three Key Components of Prover Agent

#### 1 Lemma Generation via Informal Reasoning

- Generate auxiliary lemmas
  - Specific cases
  - Potentially useful intermediate facts
- Not limited to subgoals of predefined proof sketch
  - Key difference from prior approaches
- e.g. Problem: Show that  $n^2 + an$  is even
  - Consider  $n^2 + n$  or  $n^2 + 3n$  ( $n \in \mathbb{N}, a: \text{even}$ )
- Help discover overall proof strategy
- Mirrors how human mathematicians typically work

#### 2 Formal Proof Construction Guided by Informal Reasoning and Iterative Feedback

- Leverage the stronger mathematical ability of the informal LLM
- Construct a formal proof using an informal proof as a guide
- Iteratively refine the proof based on Lean feedback
  - Can be seen as self-correction through in-context learning
  - Akin to how humans improve their understanding based on feedback

#### 3 Final Proof Synthesis Guided by Verified Lemmas and Iterative Feedback

- Consider overall proof using the lemmas
  - Use only the verified lemmas
- Allows bottom-up strategy construction even when the full plan isn't initially clear
  - Prior work: top-down approach requiring the full plan upfront
- Iteratively refine the proof based on Lean feedback

## Experiments

### Experimental Setup

- Informal LLM: DeepSeek-R1-0528-Qwen3-8B
- Formal prover model: Goedel-Prover-V2-8B/DeepSeek-Prover-V2-7B
- Formalizer: Goedel-Formalizer-V2-8B/Kimina-Autoformalizer-7B

### Comparison of Formal Theorem-Proving Performance

Prover System	Method	Model Size	Sample Budget	miniF2F-test
<i>Large Language Models</i>				
DSP+ (Cao et al., 2025)	w/ QwQ, DeepSeek-V3, and BFS-Prover	Informal + Tree search	1	52.5%
			128	74.2%
			1024	79.5%
	w/ DeepSeek-R1, DeepSeek-V3, and BFS-Prover		1024	80.7%
DeepSeek-Prover-V2 (Ren et al., 2025)	Whole-proof	671B	1	61.9%
			8192	88.9%
Delta-Prover (Zhou et al., 2025)	w/ Gemini 2.5 Pro	Agent	16384	95.9%
Seed-Prover (Chen et al., 2025)	Whole-proof	unknown	unknown	99.6%
<i>Medium Language Models</i>				
Kimina-Prover-Preview (Wang et al., 2025)	Whole-proof	72B	1	52.9%
			1024	77.9%
			8192	80.7%
Goedel-Prover-V2 (Lin et al., 2025b)	Whole-proof	32B	32	88.1%
			1024	91.8%
			8192	92.2%
<i>Small Language Models</i>				
DeepSeek-Prover-V1.5-RL + RMaxT5 (Xin et al., 2025a)	Tree search	7B	32 × 16 × 400	63.5%
			256 × 32 × 600	65.9%
			600 × 8 × 400	68.4%
			2048 × 2 × 600	70.8%
HunyuanProver v16 + BFS + DC (Li et al., 2025)	Tree search	7B	128	61.1%
			25600	64.7%
			25600	67.6%
			25600	67.6%
Leanabell-Prover-GD-RL (Zhang et al., 2025)	Whole-proof	7B	1	52.5%
			32	63.1%
			1024	70.8%
			1	58.6%
DeepSeek-Prover-V2 (Ren et al., 2025)	Whole-proof	7B	32	75.6%
			1024	79.9%
			8192	82.0%
			1	60.8%
Leanabell-Prover-V2-KM (Ji et al., 2025)	Whole-proof	7B	32	68.4%
			128	70.4%
			32	76.6%
			128	78.2%
Leanabell-Prover-V2-D5 (Ji et al., 2025)	Whole-proof	7B	1	60.8%
			64	83.3%
			256	85.2%
			512	85.7%
Goedel-Prover-V2 (Lin et al., 2025b)	Whole-proof	7B	1	52.5%
			32	63.1%
			1024	70.8%
			1	58.6%
w/ DeepSeek-Prover-V2	(Direct proving w/o iterative refinement)	Agent	1	61.5%
			50	79.9%
			100	82.0%
			260	87.6%
Prover Agent (Ours)	(Direct proving w/ iterative refinement)	Agent	1	64.3%
			50	84.4%
			100	85.7%
			260	86.5%
w/ Ensemble of Goedel-Prover-V2 and DeepSeek-Prover-V2	(Direct proving w/o iterative refinement)	Agent	1	64.3%
			50	85.7%
			100	86.9%
			260	88.1%

### Performance on Olympiad-Level Problems

Model Size	Sample Budget	Olympiad				MATH			Custom			
		IMO	AIME	AMC	Sum	Algebra	Number Theory	Sum	Algebra	Number Theory	Induction	Sum
8B	1	40.0	53.3	62.2	55.0	71.4	60.0	66.2	55.6	75.0	50.0	58.8
	50	70.0	80.0	82.2	78.8	80.0	88.3	83.8	66.7	75.0	62.5	67.6
	100	70.0	80.0	86.7	81.3	84.3	88.3	86.2	66.7	75.0	62.5	67.6
	260	70.0	80.0	88.9	82.5	84.3	88.3	86.2	66.7	75.0	75.0	70.6
671B	8192	50.0	93.3	77.8	73.8	100.0	96.7	98.5	83.3	87.5	100.0	88.2

- Show strong performance on Olympiad-level problems
  - Suggest that coordination with informal reasoning may be the key
    - Olympiad-level problems require a high degree of mathematical reasoning
- Consistent gap in MATH and Custom
  - Suggests that model size and sample budget may play a more significant role here
    - Prover model also possesses a certain level of mathematical reasoning ability

### Case Study: Success with Lemma-Guided Proofs

#### Problem:

```
1 theorem induction_nfactltnexpnmlngt3
2 (n : ℕ)
3 (h0 : 3 ≤ n) :
4 (n! < n^(n-1)) := by sorry
```

#### Reasoning trace w/ lemmas:

- Consider the specific cases for  $n = 3, 4, 5$ 
  - Clearly identify the use of mathematical induction
- Employ proof techniques used in the lemmas

#### Reasoning trace w/o lemmas:

- Proof strategy is unclear
- The details cannot be worked out sufficiently

#### Generated lemmas:

```
1 theorem base_case_3 : (3 : ℕ) < 3^(3-1) := by
2 have h_main : (3 : ℕ) < 3^(3-1) := by
3 -- Calculate the factorial and the power step-by-step
4 norm_num [Nat.factorial, Nat.pow_succ, Nat.mul_assoc]
5 -- Use 'decide' to confirm the inequality
6 <|> decide
7 exact h_main
```

```
1 theorem exponent_inequality (n : ℕ) (hn : 2 ≤ n) :
2 n^(n-1) < (n+1)^(n-1) := by
3 have h_main : n^(n-1) < (n+1)^(n-1) := by
4 have h0 : n-1+1 = n := by
5 omega
6 have h1 : n ≥ 2 := by omega
7 have h2 : n-1 ≥ 1 := by omega
8 have h3 : n < n+1 := by omega
9 -- Use the fact that if a < b and n > 0, then a^n <
10 exact calc
11 n^(n-1) < (n+1)^(n-1) := by
12 -- Apply the lemma that if a < b and n > 0, then
13 a^n < b^n
14 exact Nat.pow_lt_pow_of_lt_left h3 (by omega)
15 _ = (n+1)^(n-1) := by rfl
16 exact h_main
```

The full reasoning trace can be found in the appendix of our paper.