# A Few Open Problems in Neural Theorem Proving

# (in Lean)

Sean Welleck
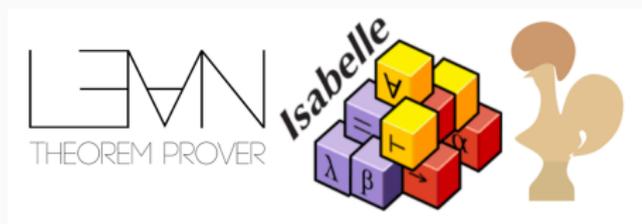
September 5, 2024

Carnegie Mellon University

Use neural networks to:

- Generate proofs in an interactive proof assistant

Rapid progress in methods based on language models:



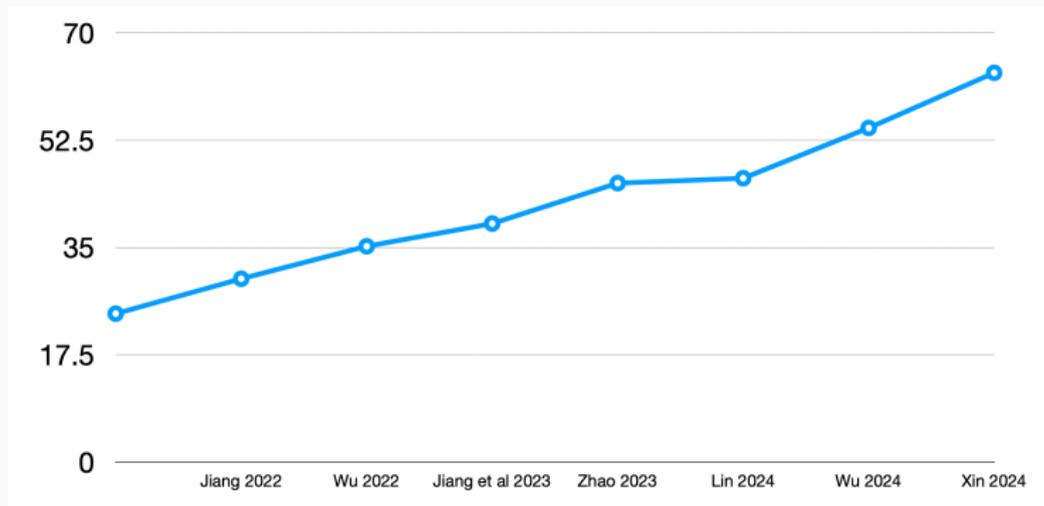**Figure 1:** miniF2F benchmark performance, 2022-2024

```
theorem imo_1960_p2 (x : ℝ) (h₀ : 0 ≤ 1 + 2 * x) (h₁ : (1 - Real.sqrt (1 + 2 *
    x)) ^ 2 ≠ 0)
    (h₂ : 4 * x ^ 2 / (1 - Real.sqrt (1 + 2 * x)) ^ 2 < 2 * x + 9) : -(1 / 2)
    ≤ x ∧ x < 45 / 8 := by
  norm_num at h₀ h₁ h₂
  have h₃ : 0 ≤ 1 + 2 * x := by linarith
  have h₄ : 0 < 1 + Real.sqrt (1 + 2 * x) := by
    nlinarith [Real.sqrt_nonneg (1 + 2 * x)]
  have h₅ : 4 * x ^ 2 / (1 - Real.sqrt (1 + 2 * x)) ^ 2 < 2 * x + 9 := by
    linarith
  have h₆ : 1 - Real.sqrt (1 + 2 * x) ≠ 0 := by
    intro h
    apply h₁
    nlinarith
  have h₇ : 4 * x ^ 2 / (1 - Real.sqrt (1 + 2 * x)) ^ 2 = (1 + Real.sqrt (1 +
    2 * x)) ^ 2 := by
    field_simp [h₆]
    nlinarith [sq_sqrt (show 0 ≤ 1 + 2 * x by linarith)]
  rw [h₇] at h₅
  constructor <;> nlinarith [sq_sqrt (show 0 ≤ 1 + 2 * x by linarith)]
```

**Figure 2:** Generated International Math Olympiad solution in Lean
(DeepSeek Prover-1.5B, Xin et al 2024)

Why talk about Lean?

- Increasing interest from the mathematical community
- Increasing interest from the AI community
- For AI research, the choice of proof assistant matters (not ideal!)

3 open problems in neural theorem proving in Lean:

- Going beyond human data
- Going beyond competition problems
- Going beyond mathematics

## 1. Going beyond human data

Language model-based proving:

- **Train** a model $p_\theta(y|x)$ on a dataset $\mathcal{D} = \{(x, y)\}$, e.g.,
    - $x$: proof state
    - $y$: next tactic (next "step")
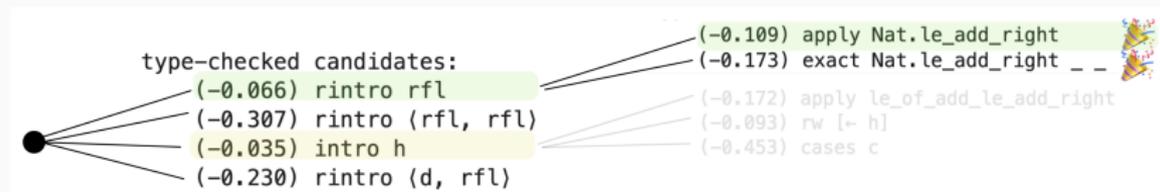    - $\mathcal{D}$: extracted from human-written theorems and proofs

# 1. Going beyond human data

Language model-based proving:

- **Train** a model $p_\theta(y|x)$ on a dataset $\mathcal{D} = \{(x, y)\}$, e.g.,
  - $x$: proof state
  - $y$: next tactic (next "step")
  - $\mathcal{D}$: extracted from human-written theorems and proofs

- **Generate** proofs:



**Figure 3:** Best-first search

## 1. Going beyond human data

- Some models are already trained on $\approx$ all Lean projects!
  - E.g., Lean-GitHub [5]: data from 237 Lean 4 repos

- More human-written data will help, but difficult to scale[1]

---
[1]Please don't stop making more publicly available formal mathematics data!

Open problem I: how do we *synthesize* useful data?

- Proofs
- Theorems
- Augmentations (formal, informal, ...)
- ...

# 1. Going beyond human data

Not a new problem; common methods:

- Statement autoformalization [Wu et al 2022 [4]]
    - Informal theorem $\rightarrow$ formal theorem

- Expert iteration [Polu et al 2022 [3]]
    - Generate proofs with a model, train on successful ones, iterate

# 1. Going beyond human data

Not a new problem; common methods:

- Statement autoformalization [Wu et al 2022 [4]]
  - Informal theorem $\rightarrow$ formal theorem

- Expert iteration [Polu et al 2022 [3]]
  - Generate proofs with a model, train on successful ones, iterate

Used in several state-of-the-art methods, e.g. DeepSeek-Prover 1.5, AlphaProof

**Lean-STaR: Learning to Interleave Thinking and Proving**
Haohan Lin, Zhiqing Sun, Yiming Yang, Sean Welleck
https://arxiv.org/abs/2407.10040

Can we do better by interleaving *informal* steps of reasoning? (right)

Why?

- Plan proof steps
- Diversify search space
- More tokens can give more computational capacity

Why?

- Plan proof steps
- Diversify search space
- More tokens can give more computational capacity

Data doesn't exist! We need to synthesize it.

Lean-STaR (Self-taught reasoner[2])

Step 1: generate an informal "thought" with an off-the-shelf language model retrospectively

- (state, tactic) $\rightarrow$ thought

Train an initial model on a dataset of such examples:

- $p_\theta^0(thought, tactic|state)$

---

[2]Inspired by *STaR: Bootstrapping Reasoning with Reasoning*, Zelikman et al 2022

Step 2: generate proofs with the model



Figure 4: Best-first search: difficult to score (thought, tactic) candidates

Step 2: generate proofs with the model



**Figure 5:** New sampling method

Step 3: train on the successful proofs, and repeat:[3]

- Collect (state, thought, tactic) from successful proofs
- Train a new model $p_\theta^1(thought, tactic|state)$
- Generate proofs
- ...

_____

[3]I.e. Expert Iteration [Polu et al 2022 [3]]

- **miniF2F** [7]: competition problems (AMC, AIME, IMO)



**Problem**    1959 IMO Problems/Problem 1

Prove that the fraction $\dfrac{21n+4}{14n+3}$ is irreducible for every natural number $n$.

```
theorem imo_1959_p1
  (n : ℕ)
  (h₀ : 0 < n) :
  Nat.gcd (21*n + 4) (14*n + 3) = 1 := by sorry
```

- train
- test
- mathlib

16

Figure 6: MiniF2F test

Figure 7: MiniF2F test

```
theorem amc12a_2002_p6 (n : ℕ) (h₀ : 0 < n)
  : ∃ m, m > n ∧ ∃ p, m * p ≤ m + p :=
  use n + 1, by linarith only [h₀] -- Reasoning: To demonstrate the existence of
  such an `m` and `p` that satisfy the given inequalities, we leverage the
  property of natural numbers that ensures an increase in `n` from `0` leads to
  the possibility of finding such `m` and `p`. This approach is supported by the
  fundamental property of natural numbers and their arithmetic operations,
  suggesting that for a given `n` greater than `0`, there exists a successor `m`
  and a number `p` that satisfy the conditions outlined.
  exact ⟨0, by simp ⟩ -- Reasoning: To demonstrate that there exists an \(n \in ℕ
  \) such that \((n + 1) * p \leq n + 1 + p\) holds for some \(p \in ℕ\), we
  recognize that setting \(p = 0\) satisfies the inequality by the properties of
  zero in arithmetic. This is because \(n + 1 > 0\) by assumption, ensuring the
  left-hand side is non-negative and the right-hand side simplifies to \(n +
  1\), making the inequality trivially true.
```

Figure 8: Example generated thoughts and proof from Lean-STaR

**Figure 9:** Example generated thoughts and proof from Lean-STaR

**Figure 10:** Increasing the search budget is more effective with thoughts

3 open problems in neural theorem proving in Lean:

- Going beyond human data
  - Synthesizing data: problems, proofs, plans, …
- **Going beyond competition problems**
- Going beyond mathematics

## 2. Going beyond competition problems

Lots of exciting progress! Some methods can solve IMO problems!

However, not much impact on proving in practice.

Accessibility gap:

- Some methods are hard to integrate into tools
  - Not open-source (AlphaProof, ...)
  - Expensive to run (MCTS, ...)

## 2. Going beyond competition problems

Accessibility gap:

- Some methods are hard to integrate into tools
  - Not open-source (AlphaProof, ...)
  - Expensive to run (MCTS, ...)

However, there are model-agnostic tools available to plug into!

Figure 11: https://github.com/cmu-l3/llmlean

Figure 12: https://github.com/cmu-l3/llmlean

**Figure 13:** Example on Polynomial Freiman Rusza Conjecture project
https://github.com/cmu-l3/llmlean

# 2. Going beyond competition problems

Benchmarking gap:

- Benchmark improvements (e.g., on competition problems) do not measure improvement in real-world proving conditions

Python
Interview Question

Code in a real repository

**Figure 14:** Interview questions $\neq$ real code development

**Figure 15:** Competition problems ≠ real proof development

Real-world proving is **context-dependent**:

- (context, theorem) $\rightarrow$ proof
  - Context: repository of code, new definitions, auxiliary lemmas

Generalization to new contexts is studied in other proof assistants, e.g., online setting[4], testing on held-out repositories[5]

Not a focus for state-of-the-art models/benchmarks in Lean!

---

[4]Tactician [2], Graph2Tac [1]
[5]CoqGym [6]

miniCTX: Neural Theorem Proving with (Long-)Contexts
Jiewen Hu, Thomas Zhu, Sean Welleck
https://www.arxiv.org/abs/2408.03350

miniCTX:

Collect (context, theorem) examples from real Lean projects:[6]

- "Future mathlib": theorems added after a time cutoff
- Recent projects: PFR, PrimeNumberTheorem
- Textbook exercises: How To Prove It, Math 2001

---

[6] + tools for easily adding new projects: https://github.com/cmu-l3/ntp-toolkit

miniCTX:

Collect (context, theorem) examples from real Lean projects:[6]

- "Future mathlib": theorems added after a time cutoff
- Recent projects: PFR, PrimeNumberTheorem
- Textbook exercises: How To Prove It, Math 2001

Goal: generalize to new theorems/contexts/repositories

---

[6]+ tools for easily adding new projects: https://github.com/cmu-l3/ntp-toolkit

Context:

- Preceding code in the file
- All accessible premises
- Repository metadata (to recover any other code)

Does context actually matter? A simple experiment.



Figure 16: "File tuning": train on (preceding code, state, next-tactic) examples

Two methods can have similar performance on competition problems, but vastly difference performance on actual projects:

| Models | MiniF2F Test | MiniCTX Prime | PFR | Mathlib | HTPI | Avg. |
|---|---|---|---|---|---|---|
| GPT-4o (full proof) | - | 1.15% | 5.56% | 2.00% | 9.73% | 5.59% |
| GPT-4o (+ context) | - | 13.79% | 1.85% | 18.00% | 31.89% | 22.07% |
| State-tactic prompting | 28.28% | 19.54% | 5.56% | 16.00% | 19.15% | 20.61% |
| State-tactic tuning | 32.79% | 11.49% | 5.56% | 22.00% | 5.95% | 9.31% |
| File tuning | **33.61%** | **32.18%** | **5.56%** | **34.00%** | **38.38%** | **31.65%** |

File-tuned model is deployed in LLMLean:



```
LLM on your laptop:
1. Install ollama.

2. Pull a language model:

ollama pull wellecks/ntpctx-llama3-8b
```

Figure 17: https://github.com/cmu-l3/llmlean

Several open-source artifacts:

- Data/models: `https://huggingface.co/l3lab`
- Data extraction: `https://github.com/cmu-l3/ntp-toolkit`
- Evaluation: `https://github.com/cmu-l3/minictx-eval`

Many approaches to explore in the future:

- "File tuning": context is preceding code
- Premise selection: context is a set of definitions and theorems
- Full repo: context is all other code in the repository
- …

Many other potential tools beyond proof completion!

3 open problems in neural theorem proving in Lean:

- Going beyond human data
    - Synthesizing data
- Going beyond competition problems
    - Have actual tools as a goal
- **Going beyond mathematics**

**miniCodeProps: a Minimal Benchmark for Proving Code Properties**
Evan Lohn, Sean Welleck

https://arxiv.org/abs/2406.11915

Interactive theorem provers

- Mathematics:
  - Math as code
  - Guarantees on proof correctness
- Code:
  - Prove properties of code

Formally verified code

**Figure 18:** https://aws.amazon.com/blogs/opensource/lean-into-verified-software-development/

AI/neural theorem proving for program verification is actively studied in other proof assistants, such as Coq and Isabelle.

Not in Lean!

## 3. Going beyond mathematics

Our question:

- What is the simplest program verification scenario that:
    - Is a subproblem of the full 'verification problem'
    - Breaks current neural theorem proving methods

Our question:

- What is the simplest program verification scenario that:
    - Is a subproblem of the full 'verification problem'
    - Breaks current neural theorem proving methods

*"Simple"*:

- Self-contained, no complex dependencies
- Relatively small (fast, cheap evaluation)

Formally verified code

**Code**

```
inductive MyTree (α: Type) where
| leaf : MyTree α
| node : MyTree α → α → MyTree α → MyTree α

def tree_size : MyTree α → ℕ
  | .leaf => 1
  | .node l _x r => 1 + (tree_size l) + (tree_size r)

def balanced : MyTree α → Prop
  | .leaf => true
  | .node l _x r => ((tree_size l) =
      (tree_size r)) ∧ (balanced l) ∧ (balanced r)
```

Tree implementation

**Property**

```
-- The size of a balanced tree is odd
theorem balanced_tree_size_odd
  (t: MyTree α) (hb: balanced t): Odd (tree_size t) :=
```

"The size of a balanced tree is odd"

**Proof**

```
by
cases t with
  | leaf => simp [tree_size]
  | node p x q =>
      unfold tree_size
      unfold balanced at hb
      simp [hb.1]
```

Subproblem: theorem proving! Given (code, property), generate proof

47

Code blocks and 201 properties from *Tons of Inductive Problems*[7], translated from Haskell to Lean.



---

[7]https://tip-org.github.io/, Claessen et al 2015

MiniCodeProps

- Implementation + properties about lists, trees, and heaps
- Classified into difficulties:
    - Easy: Data structure properties
    - Medium: Termination properties
    - Hard: Sorting algorithm properties

Evaluation:

- Given property and all dependent code, generate a proof

Models:

- **GPT-4o**: generate full proof, 32 attempts + 1 round of refinement
- **ntp-ctx**: generate a proof via best-first search

      https://github.com/cmu-l3/minicodeprops-eval

| Model | Easy | Medium & Hard | Overall |
|---|---|---|---|
| GPT-4o (32 samples) | 75.6% (65/86) | 4.34% (5/115) | 34.8% (70/201) |
| + refinement | 77.9% (67/86) | 6.96% (8/115) | 37.3% (75/201) |
| ntp-context-1.3B | 72.1% (62/86) | 8.69% (10/115) | 35.8% (72/201) |

Figure 19: Baselines perform well on easy properties

| Model | Easy | Medium & Hard | Overall |
|---|---|---|---|
| GPT-4o (32 samples) | 75.6% (65/86) | 4.34% (5/115) | 34.8% (70/201) |
| + refinement | 77.9% (67/86) | 6.96% (8/115) | 37.3% (75/201) |
| ntp-context-1.3B | 72.1% (62/86) | 8.69% (10/115) | 35.8% (72/201) |

Figure 20: Poor performance on medium/hard properties

```
import Mathlib

def butlast : List α → List α
  | [] => []
  | [_x] => []
  | x::xs => x::(butlast xs)

def butlastConcat : List α → List α → List α
  | xs, [] => butlast xs
  | xs, ys => xs ++ butlast ys

theorem prop_49 (xs: List Nat) (ys: List Nat) :
  (butlast (xs ++ ys) = butlastConcat xs ys) := by
  induction ys generalizing xs
  case nil =>
    simp [butlast, butlastConcat]
  case cons y ys ih =>
    simp [butlast, butlastConcat]
    induction xs
    case nil =>
      simp [butlast, butlastConcat]
    case cons x xs ih' =>
      simp [butlast, butlastConcat, List.cons_append, ih']
```

Figure 21: Successful proof (GPT-4o)

```
import Mathlib

def filter : List Nat → (Nat → Bool) → List Nat
| [], _f => []
| x::xs, f => if f x then x::(filter xs f) else (filter xs f)

theorem qsort_term2 (x:Nat) (xs: List Nat) :
List.length (filter xs fun y => decide (y > x)) < Nat.succ (List.length xs) := by
  induction xs with
  | nil =>
    simp [filter, Nat.zero_lt_succ]
  | cons y ys ih =>
    simp only [filter]
    split_ifs with h
    · simp only [List.length]
      exact Nat.succ_lt_succ ih
    · simp only [List.length]
      exact Nat.lt_succ_of_lt ih
```

Figure 22: Successful proof (GPT-4o)

52

Figure 23: Human-written proof showing potential length of proofs

52

## This talk

3 open problems in neural theorem proving in Lean:

- Going beyond human data
  - Synthesizing data: problems, proofs, plans, ...
- Going beyond competition problems
  - Have actual tools as a goal
- Going beyond mathematics
  - Program verification

# Thank you!

Haohan Lin (Tsinghua)
Evan Lohn (CMU)
Jiewen Hu (CMU)
Zhiqing Sun (CMU)
Yiming Yang (CMU)
Thomas Zhu (CMU)

*Lean-STaR: Learning to Interleave Thinking and Proving.*
Haohan Lin, Zhiqing Sun, Yiming Yang, Sean Welleck, 2024.

*miniCTX: Neural Theorem Proving with (Long-)Contexts.*
Jiewen Hu, Thomas Zhu, Sean Welleck, 2024.

*miniCodeProps: a Minimal Benchmark for Proving Code Properties.*
Evan Lohn, Sean Welleck, 2024.

Sean Welleck (CMU)
Learning, Language, and Logic (L3) Lab

L. Blaauwbroek, M. Olšák, J. Rute, F. I. S. Massolo, J. Piepenbrock, and V. Pestun.
**Graph2tac: Online representation learning of formal math concepts.**
In *Forty-first International Conference on Machine Learning*, 2024.

L. Blaauwbroek, J. Urban, and H. Geuvers.
*The Tactician: A Seamless, Interactive Tactic Learner and Prover for Coq*, page 271–277.
Springer International Publishing, 2020.

S. Polu, J. M. Han, K. Zheng, M. Baksys, I. Babuschkin, and I. Sutskever.
Formal mathematics statement curriculum learning.
In *The Eleventh International Conference on Learning Representations*, 2023.

Y. Wu, A. Q. Jiang, W. Li, M. N. Rabe, C. E. Staats, M. Jamnik, and C. Szegedy.
Autoformalization with large language models.
In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

Z. Wu, J. Wang, D. Lin, and K. Chen.
Lean-github: Compiling github lean repositories for a versatile lean prover, 2024.

📄 K. Yang and J. Deng.
Learning to prove theorems via interacting with proof
assistants.
*ArXiv*, abs/1905.09381, 2019.

📄 K. Zheng, J. M. Han, and S. Polu.
minif2f: a cross-system benchmark for formal olympiad-level
mathematics.
In *International Conference on Learning Representations*, 2022.