# Bridging Informal and Formal Mathematical Reasoning

Sean Welleck

May 28, 2025

Carnegie Mellon University

Al in expert domains

- Finance
- Science
- Mathematics

### AI in expert domains

- Finance
- Science
- Mathematics
  - Open-ended dialogue
  - Come up with counterexamples
  - Help write proofs
  - ...

Math as raw data (text, images, ...)

- Flexible
- Widely used
- Difficult to check



Language model solution.

#### Math as source code

- Write a specification (e.g., 1+1=2)
- Write a proof
- Automatically checked
  - $\bullet \ \ {\sf Code \ compiles} \equiv {\sf correct \ proof}$



Math as source code.

Math as source code

- Write a specification (e.g., 1+1=2)
- Write a proof
- Automatically checked
  - $\bullet \ \ Code \ \ compiles \equiv correct \ \ proof$

Lean

Isabelle

Coq

Theorem proving languages

# If $R \subseteq S$ and $S \subseteq T$ then $R \subseteq T$ $\lim_{T \to OREM \ PROVER}$

Growing use in mathematics:



Finished formalizing in #Lean4 the proof of an actual new theorem (Theorem 1.3) in my recent paper arxiv.org/abs/2310.05328 :

Terence Tao's Lean formalization project (October 2023)

Growing use in mathematics:



@tao@mathstodon.xyz

Finished formalizing in #Lean4 the proof of an actual new theorem (Theorem 1.3) in my recent paper arxiv.org/abs/2310.05328 :

Terence Tao's Lean formalization project (October 2023)

• New form of **collaboration**: break a big problem into multiple pieces

Formal math for AI

- Tests reasoning
  - From easy: 1+1 = 2
  - To hard: Fermat's Last Theorem

Formal math for AI

- Tests reasoning
  - From easy: 1+1 = 2
  - To hard: Fermat's Last Theorem
- Verifiable
  - Feedback signal for reinforcement learning

Rapid progress in methods based on language models:



7

# $LLMs \cap formal math$

```
theorem imo 1960 p2 (x : R) (h<sub>0</sub> : 0 \le 1 + 2 * x) (h<sub>1</sub> : (1 - Real.sort (1 + 2 *
    x)) 2 \neq 0
    (h_2: 4 * x^2 / (1 - \text{Real.sqrt} (1 + 2 * x))^2 < 2 * x + 9): -(1 / 2)
   \leq x \wedge x < 45 / 8 := bv
 norm num at h0 h1 h2
 have h_3: 0 \le 1 + 2 * x := by linarith
 have h_4: 0 < 1 + Real.sqrt (1 + 2 * x) := by
    nlinarith [Real.sqrt nonneg (1 + 2 * x)]
 have h_5: 4 * x ^ 2 / (1 - Real.sqrt (1 + 2 * x)) ^ 2 < 2 * x + 9 := by
   linarith
 have h_6 : 1 - Real.sort (1 + 2 * x) \neq 0 := by
   intro h
   apply h<sub>1</sub>
    nlinarith
 have h_7: 4 * x ^ 2 / (1 - Real.sart (1 + 2 * x)) ^ 2 = (1 + Real.sart (1 +
   (2 * x))^2 := bv
   field simp [h<sub>6</sub>]
    nlinarith [sq sqrt (show 0 \le 1 + 2 * x by linarith)]
 rw [h7] at h5
 constructor <:> nlinarith [sq sqrt (show 0 \le 1 + 2 * x by linarith)]
```

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



Terence Tao @tao@mathstodon.xyz

Finished formalizing in #Lean4 the proof of an actual new theorem (Theorem 1.3) in my recent paper arxiv.org/abs/2310.05328 :

The ability of Github copilot to correctly anticipate multiple lines of code for various routine verifications, and inferring the direction I want to go in from clues such as the names I am giving the theorems, continues to be uncanny.

Terence Tao's Lean formalization project (October 2023)



Terence Tao @tao@mathstodon.xyz

Finished formalizing in #Lean4 the proof of an actual new theorem (Theorem 1.3) in my recent paper arxiv.org/abs/2310.05328 :

The ability of Github copilot to correctly anticipate multiple lines of code for various routine verifications, and inferring the direction I want to go in from clues such as the names I am giving the theorems, continues to be uncanny.

Terence Tao's Lean formalization project (October 2023)

So...why don't people and LLMs always use formal math?

Informal ideas, intuitions, and even proofs are difficult to express formally:



- Each step of reasoning needs to be specified in detail
- Requires a deep knowledge of the formal system



This talk: Bridging Informal and Formal Mathematical Reasoning

- 1. Informal thoughts
- 2. Informal sketches
- 3. Towards research-level mathematics

I: Informal thoughts

### Lean-STaR: Learning to Interleave Thinking and Proving Haohan Lin, Zhiqing Sun, Yiming Yang, Sean Welleck *ICLR 2025* (Spotlight)

#### Neural theorem proving



- Math as checkable code
- Proof: sequence of (state, step)



Can we train a model to "think" before each step of formal reasoning?

# 1. Training models to "think" — Lean-STaR



#### Why?

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

## Lean-STaR (Self-taught reasoner<sup>1</sup>)

Learn to generate thoughts via reinforcement learning

- 1. Initialization
- 2. Reinforcement learning

<sup>&</sup>lt;sup>1</sup>Inspired by STaR: Bootstrapping Reasoning with Reasoning, Zelikman et al 2022

#### 1. Initialization



#### 1. Initialization



#### 2: Reinforcement learning



Algorithm: train on the successful proofs, and repeat:<sup>2</sup>

- Collect (state, thought, tactic) from successful proofs
- Train a new model  $p^1_{\theta}(thought, tactic|state)$
- Generate proofs

• ...

<sup>&</sup>lt;sup>2</sup>I.e. Expert Iteration [Polu et al 2022], Rest-EM [Singh et al 2024]

• miniF2F: competition problems (AMC, AIME, IMO)







# Training models to "think" — Lean-STaR



Increasing the search budget is more effective with thoughts

```
theorem amc12a_2002_p6 (n : N) (h<sub>0</sub> : 0 < n)

: \exists m, m > n \land \exists p, m * p \leq m + p :=

use n + 1, by linarith only [h<sub>0</sub>] -- 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\).
```

#### Example generated thoughts and proof from Lean-STaR

#### Training models to "think" — Lean-STaR

theorem aime\_1983\_p2 (x p : R) (f : R  $\rightarrow$  R) (h<sub>0</sub> : 0 \land p < 15) (h<sub>1</sub> : p < x  $\land$  x < 15)

(h2 : f x = abs (x - p) + abs (x - 15) + abs (x - p - 15)) : 15 < f x :=

cases' h: with h: h: -- Reasoning: Given a function \(f)) defined over the real numbers, where l(f(x)) equals the sum of the absolute values of the differences \(z - p|), \(z - 15|), and \(x - p - 15)), along with conditions on |(p|) and |(a)) that ensure |(a)) lies between |(p|) and |(15)) and |(p)) is positive, we aim to show that \(f(m) \geq 15\). To demonstrate this, we consider all possible cases for the position of \(m)) relative to \(p)) and \(15)), specifically at \(p)) stself, between \(p)) and \(15)), and at \(15)). By evaluating |(f(x)|) for each case, we can establish a lower bound of \(f(x))) that is greater than or equal to \(15)).

simp [h1, abs of nonneg (sub nonneg, 2 h1), abs of nonneg (sub nonneg, 2 h2), \*] -- Reasoning: The objective is to show that |(f(x))) is at least 15 for \(x)) within the bounds | (p \le z |le 15\), where | (p)) is positive but less than 15. The provided expression for \(f(x))) suggests using properties of absolute values and inequalities to simplify and analyze  $\backslash (f(\alpha))$ . By applying simplification rules for absolute values and inequalities, along with the given bounds, we aim to demonstrate the desired inequality \(f(a) | ae 15\). simp only [abs of nonneg] at \* -- Reasoning: The objective is to establish an inequality involuting an expression with absolute values and a sum, where all terms are functions of \(m)). By recognizing that the absolute value functions are always non-negative, we can simplify the given expression by removing the absolute value symbols, thus focusing on the arithmetic manipulation of the terms to demonstrate the inequality.

simp only [abs of nonneg, sub add en add sub] at . -- Reasoning: The objective is to demonstrate that the expression ' $15 \le x - p + |x - 15| + |x - p - 15|$ ' holds given the conditions on 'p', 'z', and the function 'f'. Since the expression in question involves absolute values which inherently involve non-negativity, manipulating these expressions to highlight their non-negativity or equality is a logical initial step. This involves recognising that the absolute value of a difference can be rewritten as the difference of absolute values, which simplifies the manipulation of the expression to prove the inequality.

congr 1 -- Reasoning: To establish the inequality  $15 \le n + |n - 15| + |n - n|$ - 15/ - p', we recognize that the left-hand side of the inequality involves a min of additions and subtractions of real numbers. By introducing a congruence transformation that simplifies operations involving absolute values, we aim to align both sides of the inequality in a form that directly reflects the properties of real numbers under these operations.

simp only [abs\_of\_nonpos (sub\_nonpos\_of\_le h\_), abs\_of\_nonneg (sub\_nonneg\_of\_le h.)] -- Reasoning: The objective is to prove an inequality involving absolute values and subtractions that aligns the given conditions and constraints. Simplification using the properties of absolute values and the provided bounds can directly lead to the desired conclusion by expressing the absolute values in terms of their underlying expressions and applying the given constraints.

simp [abs\_of\_pos, +] -- Reasoning: To prove the given inequality involving absolute values, simplification using algebraic properties and given hupotheses will directly lead us to the desired inequality by recooniging the algebraic manipulation and application of the known bounds 'h1' and 'h2', as well as the non-negativity of 'p'

rw [abs\_of\_nonpos] -- Reasoning: The objective is to demonstrate that '16  $\leq$  15 + /x - p - 15/ - p' under the given constraints on 'p' and 'x'. Given 'p < 15' and 'p ≤ a ≤ 15', it's apparent that 'p - 15' is non-positive. This observation allows us to apply the property that the absolute value of a non-positive number is equal to its negative value, thereby simplifying the expression on the right-hand side of the inequality.

all goals linarith -- Reasoning: The goal involves demonstrating a relationship between expressions involving inequalities and subtraction. The simplification of these eppressions into a form that directly compares their numerical values can lead to a straightforward application of known inequalities and algebraic properties, demonstrating the inequality's validity under the given conditions.

Example generated thoughts and proof from Lean-STaR

# Training models to "think" — After Lean-STaR

After Lean-STaR, incorporating thoughts became a widely-used component of LLM-based theorem proving:



Output:

left-hand side can be rewritten as a perfect square  $(3x \sin x - 2)^2$ . -- We use the fact that  $(3x\sin x - 2)^2$  is non-negative to establish this. nlinarith [sq\_nonneg (3 \* x \* Real.sin x - 2)] ...

generating theorem proof code. Similar to Lean-STaR (Lin et al., 2024), which performs isolated chain-of-thought reasoning (Wei et al., 2022; Feng et al., 2023) before each proof step, our method integrates this reasoning directly as comments within the proof code. We use the DeepSeek-

Informal thoughts in DeepSeek Prover 1.5

# Training models to "think" — After Lean-STaR

After Lean-STaR, incorporating thoughts became a widely-used component of LLM-based theorem proving:



Informal thoughts in Kimina Prover (April 14 2025)
More broadly, *reasoning models* that generate long chains-of-thought have subsequently begun to revolutionize LLM reasoning:

September 12, 2024 Release

#### Learning to reason with LLMs

We are introducing OpenAl of, a new large language model trained with reinforcement learning to perform complex reasoning, of thinks before it answers —it can produce a long internal chain of thought before responding to the user.

OpenAl o1 reasoning model

**W**deepseek

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

#### DeepSeek R1 reasoning model

#### Recap: Lean-STaR

- Training on formal code may be insufficient to learn the underlying thought process needed to produce the code
- Learn to generate thoughts via reinforcement learning

- 1. Informal thoughts
- 2. Informal sketches
- 3. Towards research-level mathematics

II: Informal sketches



How would we write this as a formal proof?





A proof with a high-level sketch and low-level proof steps.

## Low-level provers: Sledgehammer



Sledgehammer [Paulson 2010] calls out to external automated provers.

• First-order logic, higher-order logic, SMT

## Low-level provers: Sledgehammer



#### Struggles due to the large search space of possible proofs

## Idea: combine high-level and low-level proving



Idea: combine high-level (human, LLM) and low-level proving

# Draft, Sketch, Prove: Guiding Formal Theorem Provers with Informal Proofs

Albert Q. Jiang, Sean Welleck, Jin Peng Zhou Jiacheng Liu, Wenda Li, Mateja Jamnik Guillaume Lample, Yuhuai Wu *ICLR 2023* (Oral) Given informal theorem  $x_I$ , formal theorem  $x_F$ 

1. Draft  $y_l \sim p(\cdot|x_l)$ 



Human-written or LLM-generated draft

Given informal theorem  $x_I$ , formal theorem  $x_F$ 

- 1. Draft  $y_l \sim p(\cdot|x_l)$
- 2. Sketch  $z_F \sim p(\cdot|x_F, x_I, y_I)$



LLM-generated sketch

Given informal theorem  $x_I$ , formal theorem  $x_F$ 

- 1. Draft  $y_l \sim p(\cdot|x_l)$
- 2. Sketch  $z_F \sim p(\cdot | x_F, x_I, y_I)$
- 3. **Prove**  $y_F = f(x_F, z_F)$



Low-level prover (Sledgehammer) fills in the gaps

# Proof search



## Draft-sketch-prove



Inference-time proof search scaling

## Draft-sketch-prove







<b>Statement:</b> Prove that the fraction $\frac{21n+4}{2}$ is irreducible for every natural number $n$
biatement in the fille interior $14n+3$ is interactive for every initial infinite in $\frac{1}{14n+3}$
Informal Proof (Minerva 62B):
We must show that $gcd(21n + 4, 14n + 3) = 1$ .
The Euclidean algorithm gives
$21n + 4 = 1 \cdot (14n + 3) + 7n + 1$
$14n + 3 = 2 \cdot (7n + 1) + 1$
Since $gcd(7n + 1, 1) = 1$ , we have $gcd(21n + 4, 14n + 3) = 1$ .
Formal Proof:
theorem imo 1959 pl:
fixes n :: nat
shows " $acd (21 * n + 4) (14 * n + 3) = 1$ "
proof -
(* The Euclidean algorithm gives
21n+4=1\cdot (14n+3)+7n+1
14n+3=2\cdot(7n+1)+1. *)
have c0: "21*n + 4 = 1*(14*n + 3) + 7*n + 1"
by auto [ATP]
(have c1: "14*n + 3 = 2*(7*n + 1) + 1" using c0
by auto [ATP]
$(* Since \d(7n+1,1)=1, we have \d(21n+4,14n+3)=1, *)$
then have "acd $(7*n + 1)$ 1 = 1"
using cl
by auto [ATP]
then have "gcd (21*n + 4) (14*n + 3) = 1"
using cl
by (smt (z3) BitM_plus_one ab_semigroup_add_class.add_ac(1)
add.assoc c0 gcd.commute gcd_add2 gcd_add_mult mult_numeral_1
<pre>numeral_One numeral_eq_Suc numerals(1) semiring_norm(3)) [ATP]</pre>
then show ?thesis
using cl
by blast [AIP]
dea

International Math Olympiad problem

#### Recap:

- Draft-Sketch-Prove: generate high-level sketches and fill in gaps
- Isabelle's Sledgehammer calls out to external provers to fill in gaps

#### Recap:

- Draft-Sketch-Prove: generate high-level sketches and fill in gaps
- Isabelle's Sledgehammer calls out to external provers to fill in gaps

Next: can we build a Sledgehammer for Lean?

## Premise Selection for a Lean Hammer

Thomas Zhu, Joshua Clune Jeremy Avigad, Albert Q. Jiang, Sean Welleck *Under Review 2025* 

# A hammer pipeline

#### A standard hammer pipeline:



# A hammer pipeline

### A standard hammer pipeline:



Existing components:

- Translation: LeanAuto [Qian et al 2025]
- ATP: Zipperposition [Cruanes et al 2015]
- Reconstruction: Duper [Clune et al 2024]

# A hammer pipeline

### A standard hammer pipeline:



Our challenge:

- Premise selection
- Combine pieces to create LeanHammer

Idea: frame premise selection as retrieval with a neural language model



Idea: frame premise selection as retrieval with a neural language model



• Transformer encoder embeds the state and candidate premises

Idea: frame premise selection as retrieval with a neural language model



- Transformer encoder embeds the state and candidate premises
- Contrastive loss on (state, {premise<sup>+</sup>}, {premise<sup>-</sup>}) examples
  - Nuance in how to collect and format examples

## LeanHammer — Putting it all together

Idea: combine the premise selector and ATP with a tree search



Tree search: Aesop [Limperg & From 2023]

## LeanHammer — Putting it all together

Idea: combine the premise selector and ATP with a tree search



Tree search: Aesop [Limperg & From 2023]

- 1. Queries the automated theorem prover using the premises
- 2. Applies tactics (e.g. apply, simp\_all) using the premises

## LeanHammer — Putting it all together

Idea: combine the premise selector and ATP with a tree search



Tree search: Aesop [Limperg & From 2023]

- 1. Queries the automated theorem prover using the premises
- 2. Applies tactics (e.g. apply, simp\_all) using the premises

Goes beyond the standard hammer pipeline!

As a user, simply issue hammer at any step of a proof:



LeanHammer in action

## LeanHammer — Example

Demo: start with human-written proof sketch (from Mathematics in Lean)

```
/-- Theorem taken from Mathematics in Lean -/
theorem irrational_sqrt_two {m n : N} (coprime_mn : m.Coprime n) :
    m^{2} \neq 2 \neq n^{2} := by
 intro sqr eq
 have : 2 \mid m := by
  sorry
  obtain (k, meq) := dvd_iff_exists_eq_mul_left.mp this
  have : 2 * (2 * k^2) = 2 * n^2 := by
  sorrv
  have : 2 * k^2 = n^2 := by
   sorrv
 have : 2 \mid n := by
   sorry
 have : 2 \mid m_{n} \neq n := by
   sorry
 have : 2 | 1 := by
   sorry
  sorry
```

## LeanHammer — Example

## Demo: fill in the gaps (sorrys) with LeanHammer



Varying the premise selector within LEANHAMMER:



## Sketching proofs and filling in the gaps

## Sketching proofs and filling in the gaps

• Draft-Sketch-Prove (DSP)

## Sketching proofs and filling in the gaps

- Draft-Sketch-Prove (DSP)
- LeanHammer

- 1. Informal thoughts
- 2. Informal provers
- 3. Towards research-level mathematics
**III: Research-level mathematics** 



<sup>&</sup>lt;sup>3</sup> Formalizing the proof of PFR in Lean4 using Blueprint: a short tour by Terence Tao



<sup>&</sup>lt;sup>3</sup> Formalizing the proof of PFR in Lean4 using Blueprint: a short tour by Terence Tao

As a start, can AI help with filling in small parts of the blueprint?







LLMLean: https://github.com/cmu-l3/llmlean

#### Where can AI help? — Existing tools



LLMLean: https://github.com/cmu-l3/llmlean



#### LLMLean example on Polynomial Freiman Rusza Conjecture project

#### Math competition problems

1		
1		
1	INTERNATIONAL MATHEMATICAL OLYMPIAD	
	INO 2824 60 INO 2025	
	Theorem to prove	
theorem sorry	imo_1964_p1_2 (n : N) : ¬7   2 ^ n	+ 1 := by

- Self-contained
- Uses standard results

## Where can AI help? — Benchmarking gap



- Self-contained
- Uses standard results

- Part of a project
- Uses new definitions and lemmas

# Where can AI help? — Benchmarking gap



- Self-contained
- Uses standard results

- Part of a project
- Uses new definitions and lemmas
- New Benchmark: miniCTX

miniCTX: Neural Theorem Proving with (Long-)Contexts Jiewen Hu, Thomas Zhu, Sean Welleck. ICLR 2025 (Oral).

- 1. Informal thoughts
  - Training models to think informally
    - Lean-STaR
- 2. Informal provers
  - Sketching proofs and filling in the gaps
    - Draft, Sketch, Prove
    - LeanHammer
- 3. Towards research-level mathematics
  - Assisting in research-level projects
  - Practical tools
  - MiniCTX

### Thank you!

Collaborators on works in this talk (alphabetical by last name):

- Jeremy Avigad (CMU)
- Joshua Clune (CMU)
- Jiewen Hu (CMU)
- Mateja Jamnik (Cambridge)
- Albert Q. Jiang (Cambridge, Mistral)
- Timothee Lacroix (Meta, Mistral)
- Guillaume Lample (Meta, Mistral)
- Haohan Lin (Tsinghua)

- Wenda Li (Edinburgh)
- Jiacheng Liu (Washington)
- Zhiqing Sun (CMU, OpenAI)
- Yuhuai (Tony) Wu (Google, X.ai)
- Yiming Yang (CMU)
- Jin Peng Zhou (Cornell)
- Thomas Zhu (CMU)

Sean Welleck CMU School of Computer Science Learning, Language, and Logic (L3) Lab www.wellecks.com wellecks@cmu.edu