

Neural theorem proving

Sean Welleck

Carnegie
Mellon
University

W

UNIVERSITY *of* WASHINGTON

Ai2 Allen Institute for AI

Language models and mathematics

ChatGPT

Examples	Capabilities	Limitations
"Explain quantum computing in simple terms" →	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?" →	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?" →	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021

Free Research Preview. ChatGPT may produce inaccurate information about people, places, or facts. [ChatGPT May 24 Version](#)

Question: For every $a, b, b \neq a$ prove that

$$\frac{a^2 + b^2}{2} > \left(\frac{a+b}{2}\right)^2.$$

Model output:

$$\begin{aligned} & \frac{a^2+b^2}{2} > \left(\frac{a+b}{2}\right)^2 \\ \Leftrightarrow & \frac{a^2+b^2}{2} > \frac{a^2+b^2+2ab}{4} \\ \Leftrightarrow & a^2 + b^2 > \frac{a^2+b^2+2ab}{2} \\ \Leftrightarrow & 2a^2 + 2b^2 > a^2 + b^2 + 2ab \\ \Leftrightarrow & a^2 + b^2 > 2ab \\ \Leftrightarrow & a^2 + b^2 - 2ab > 0 \\ \Leftrightarrow & (a - b)^2 > 0 \end{aligned}$$

which is true, because the square of a real number is positive.

Minerva Language Model
2022 Poland National Math Exam
[Lewkowycz et al 2022]

Free-form conversation

...

```
lemma one_plus_one_equals_two:  
  shows "1 + 1 = 2"  
proof -  
  have "1 + 1 = Suc (0 + 1)" by simp  
  also have "1 = 1" by simp  
  also have "1 = 2" by simp  
  finally by simp  
qed
```

Informal

Formal



Show that for any
real number a ,
 $10a \leq 28a^2 + 1$



Informal

Formal

Show that for any
real number a ,
 $10a \leq 28a^2 + 1$



This screenshot shows a chat interface with a dark header bar containing the OpenAI logo and the text "This statement is false for some values of `a`." Below this, the text "For example, if `a = -0.1`, then:" is visible. A large red quote "This statement is false" is overlaid on the interface. Below the quote, the calculation $28a^2 + 1 = 28(-0.1)^2 + 1 = 28 * 0.01 + 1 = 0.29$ is shown, followed by the text "And `-1 > 0.29`, which contradicts the inequality."

← Informal

Formal →

Show that for any
real number a ,
 $10a \leq 28a^2 + 1$



This statement is false for some values of a .

For example, if $a = -0.1$, then:

$$10a = -1$$

$$28a^2 + 1 = 28(-0.1)^2 + 1 = 28 * 0.01 + 1 = 0.29$$

“And $-1 > 0.29$ ”

And $-1 > 0.29$, which contradicts the inequality.

Informal

Formal

Hard to
evaluate

Informal

```
lemma one_plus_one_equals_two:  
  .....  
  shows "1 + 1 = 2"  
proof -  
  have "1 + 1 = Suc (0 + 1)" by simp  
  also have "Suc (0 + 1) = Suc 1" by simp  
  also have "Suc 1 = 2" by simp  
  finally have thesis by simp  
qed
```

Formal

Interactive (formal) theorem proving

$$1 + 1 = 2$$

```
lemma one_plus_one_equals_two:  
  shows "1 + 1 = 2"  
proof -  
  have "1 + 1 = Suc (0 + 1)" by simp  
  also have "... = Suc 1" by simp  
  also have "... = 2" by simp  
  finally show ?thesis by simp  
qed
```



Formal

Interactive theorem proving

$$1 + 1 = 2$$

proof



```
lemma one_plus_one_equals_two:
  shows "1 + 1 = 2"
proof -
  have "1 + 1 = Suc (0 + 1)" by simp
  also have "... = Suc 1" by simp
  also have "... = 2" by simp
  finally show ?thesis by simp
qed
```

Formal

gpt-f

**Generative Language Modeling for Automated
Theorem Proving**

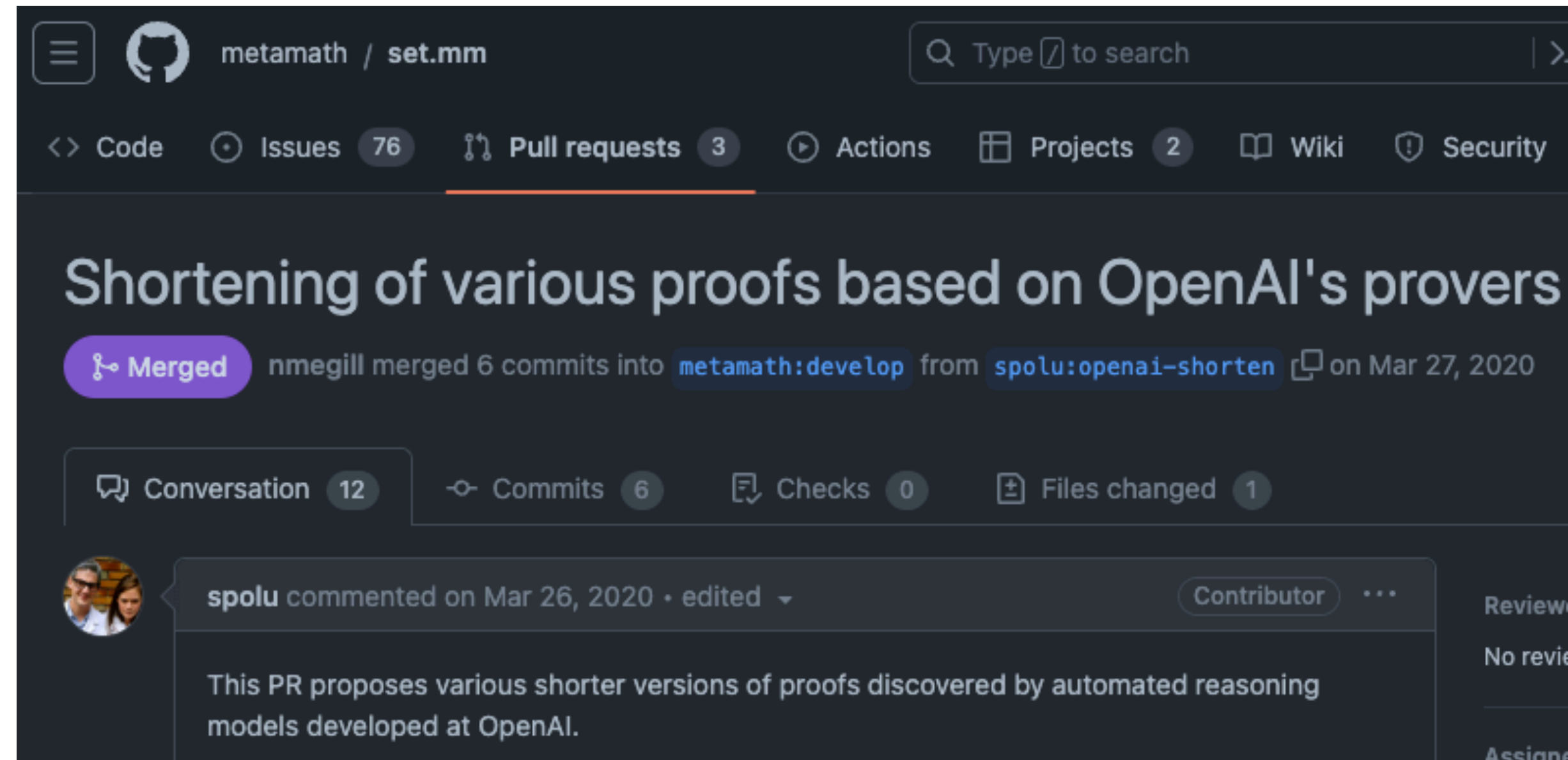
Stanislas Polu
OpenAI
spolu@openai.com

Ilya Sutskever
OpenAI
ilyasu@openai.com



Formal

gpt-f

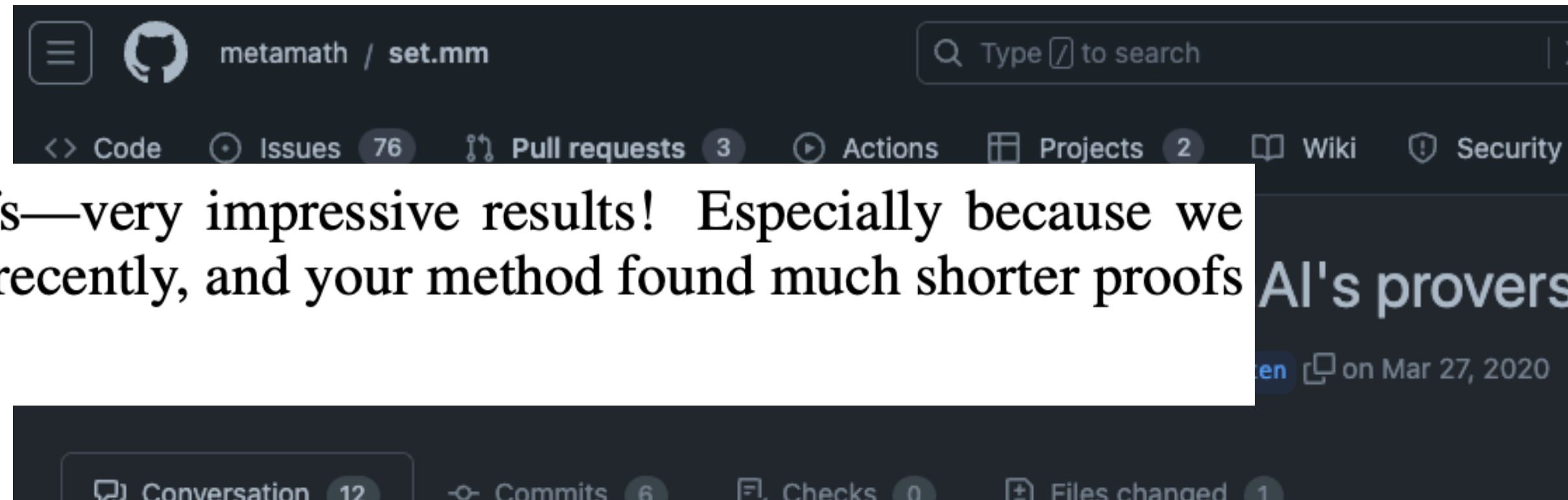


The screenshot shows a GitHub pull request interface. At the top, the repository is identified as 'metamath / set.mm'. The navigation bar includes 'Code', 'Issues (76)', 'Pull requests (3)', 'Actions', 'Projects (2)', 'Wiki', and 'Security'. The main title of the pull request is 'Shortening of various proofs based on OpenAI's provers'. A purple 'Merged' badge is visible, along with the text 'nmegill merged 6 commits into metamath:develop from spolu:openai-shorten on Mar 27, 2020'. Below the title, statistics show 'Conversation (12)', 'Commits (6)', 'Checks (0)', and 'Files changed (1)'. A comment from user 'spolu' dated 'Mar 26, 2020' states: 'This PR proposes various shorter versions of proofs discovered by automated reasoning models developed at OpenAI.'

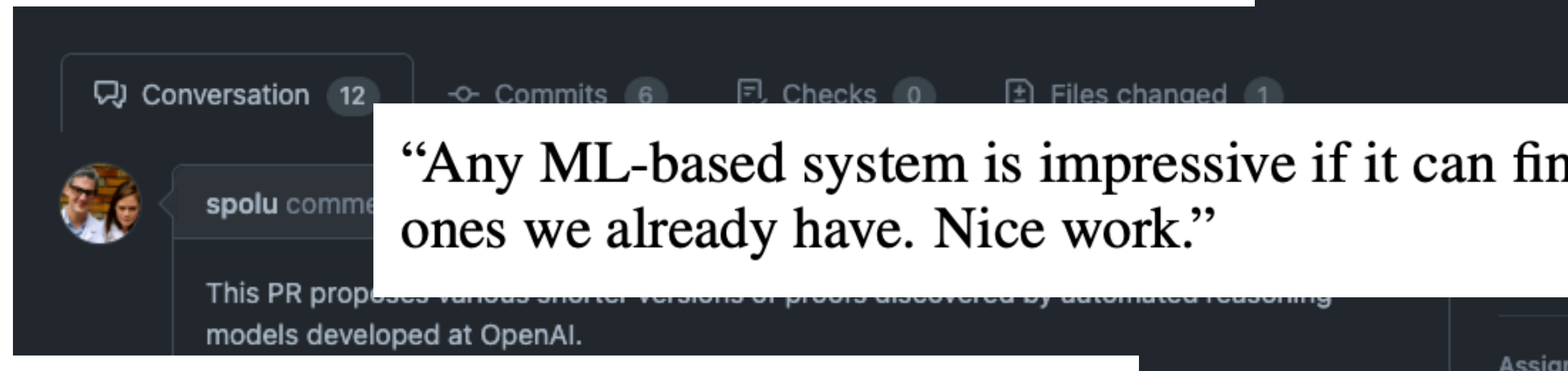


Formal

gpt-f



“I had a look at the proofs—very impressive results! Especially because we had a global minimization recently, and your method found much shorter proofs nevertheless.”



“Any ML-based system is impressive if it can find many shorter proofs than the ones we already have. Nice work.”

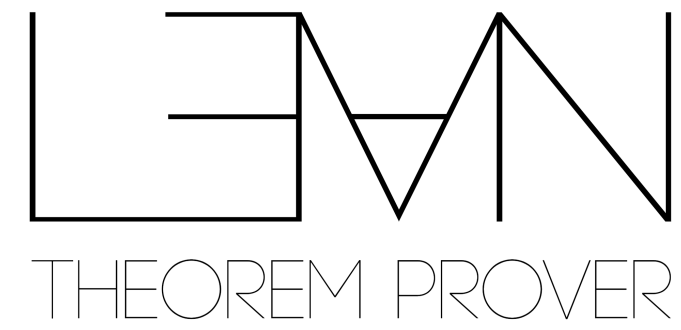
“The shorter proof is easier to translate. It’s more symmetric in that it treats A and B identically. It’s philosophically more concise in that it doesn’t rely on the existence of a universal class of all sets.”



Formal

Demo

If $R \subseteq S$ and $S \subseteq T$ then $R \subseteq T$



Microsoft Research

Application 1: mathematics

- Lean Mathlib
 - 1+ million lines of code
 - > 300 contributors
 - Algebra, Linear Algebra, Topology, Analysis, Probability, Geometry, Combinatorics, ...



<https://leanprover-community.github.io/>

See also [Adam Topaz, Formal Mathematics and AI, 2023 NASEM AI+Math Workshop](#)

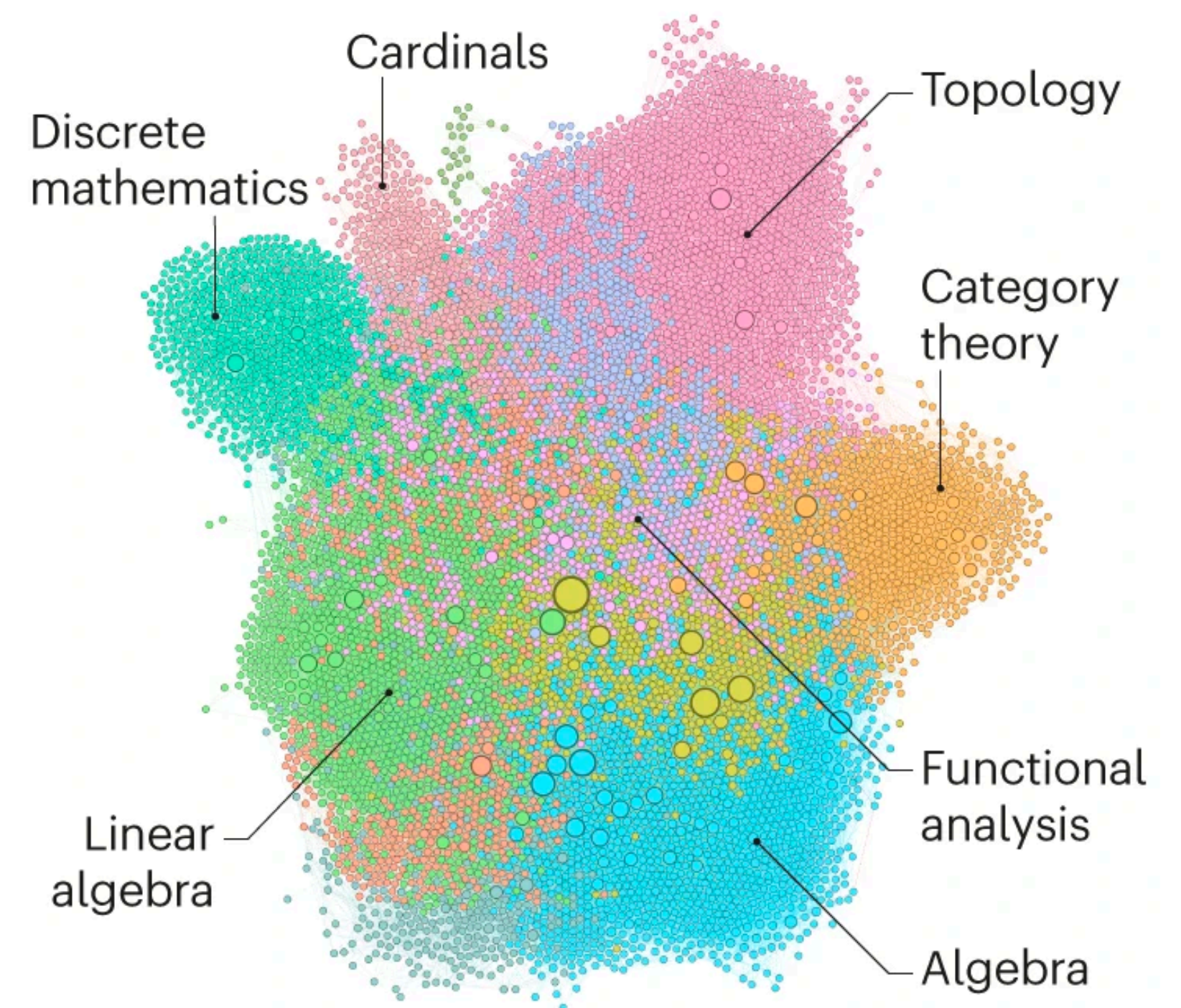
Application 1: mathematics

- Liquid tensor experiment

Mathematicians welcome computer-assisted proof in 'grand unification' theory

Proof-assistant software handles an abstract concept at the cutting edge of research, revealing a bigger role for software in mathematics.

[Davide Castelvechi](#)



Application 1: mathematics

- Education
 - Courses at CMU, Imperial College London, Fordham, JHU, Université Paris-Saclay, ...

Machine learning potential

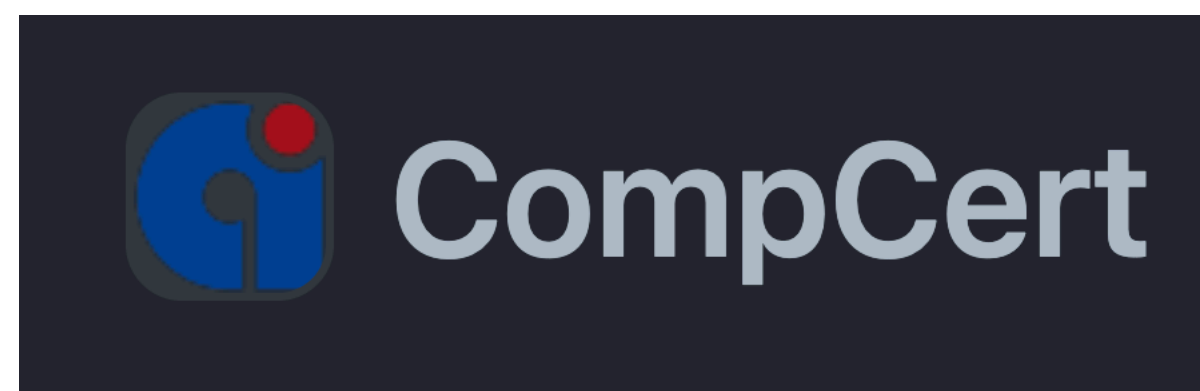
- Mathematician Adam Topaz:
 - “AI Collaborator” ... “It can be useful even if it’s not too smart” [1]
- Automate tedious proofs
- Retrieve definitions/theorems (e.g. $h1.trans\ h2$)

Application 2: software verification

- Specification
 - E.g. “reverse(reverse(list)) == list”
- Proof: code satisfies specification

Application 2: software verification

- Specification
 - E.g. “reverse(reverse(list)) == list”
- Proof: code satisfies specification
- Safety-critical applications
- Certified compilers, low-level systems software



[> Defense Advanced Research Projects Agency](#) [> Our Research](#) [> Proof Engineering, Adaptation, R](#)

Proof Engineering, Adaptation, Repair, and Learning for Software (PEARLS)

See: [QED at Large: A Survey of Engineering of Formally Verified Software](#)
Ringer, Palmskog, Sergey, Gligoric, Tatlock. Foundations and Trends in Programming Languages 5.

Machine learning potential

- Make proof assistants easier to use
- Boilerplate, tedious proofs [1]
- Proof re-use, repair, automation [1,2]

In the end, out of a total of around 550 lemmas, approximately 400 were tedious “infrastructure” lemmas; only the remainder had direct relevance to the meta-theory

Rossberg, Russo, Dreyer, *F-ing Modules*, JFP 2015; from [1]


[1] [Brigitte Pientka, Principles of Programming and Proof Languages, 2023 NASEM AI+Math Workshop](#)

[2] [Talia Ringer, Concrete Problems in Proof Automation, AITP 2022](#)

ML 1: reasoning

- Tests aspects of reasoning

International
Mathematical
Olympiad
2020 Problem 2



Prove that $\forall a, b, c, d \in \mathbb{R}$ with $a \geq b \geq c \geq d > 0$ and $a + b + c + d = 1$,

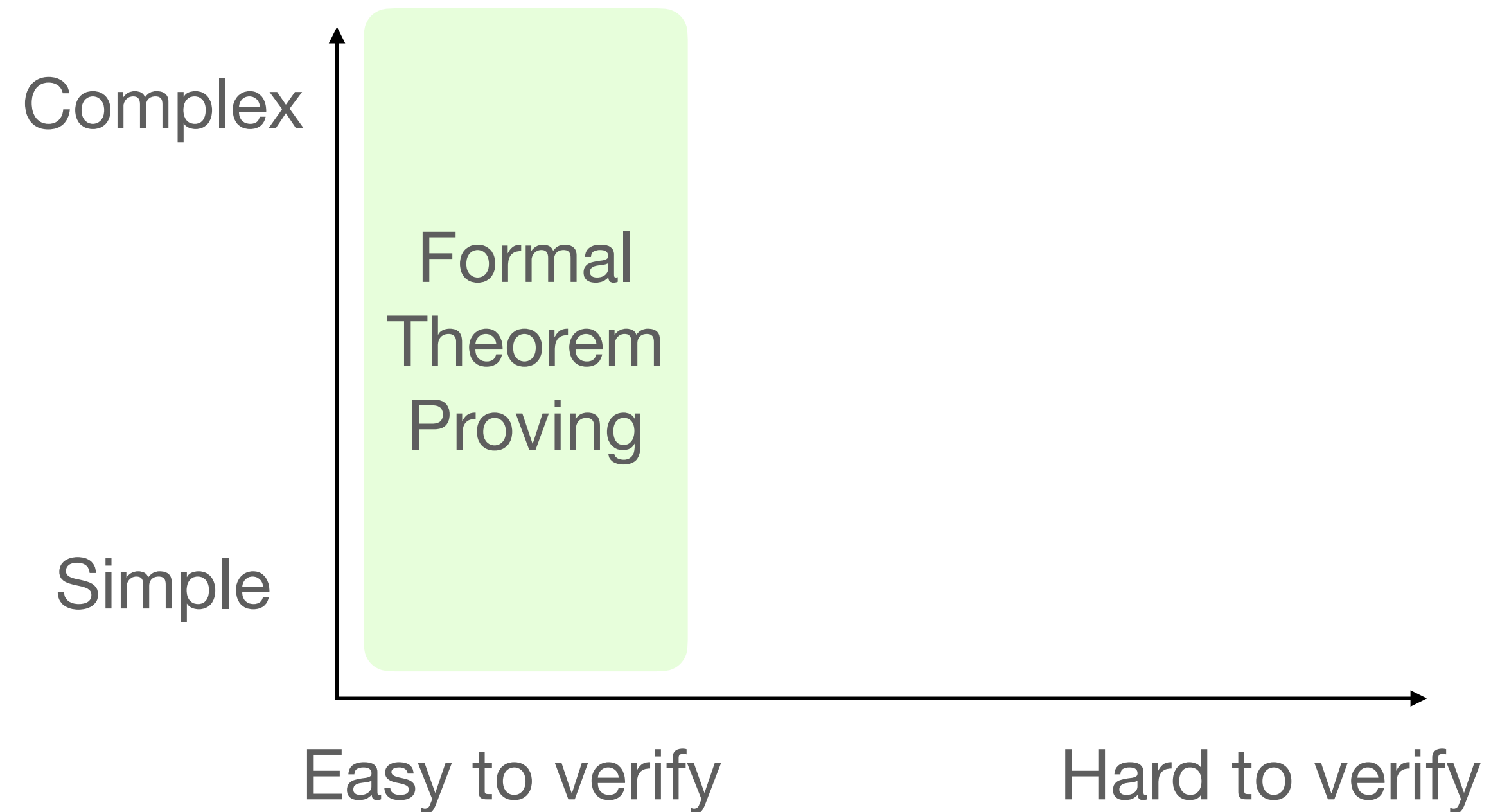
$$(a + 2b + 3c + 4d)a^a b^b c^c d^d < 1.$$

→

Planning
Pattern recognition
Background knowledge
Backtracking
Creativity
Computation
...

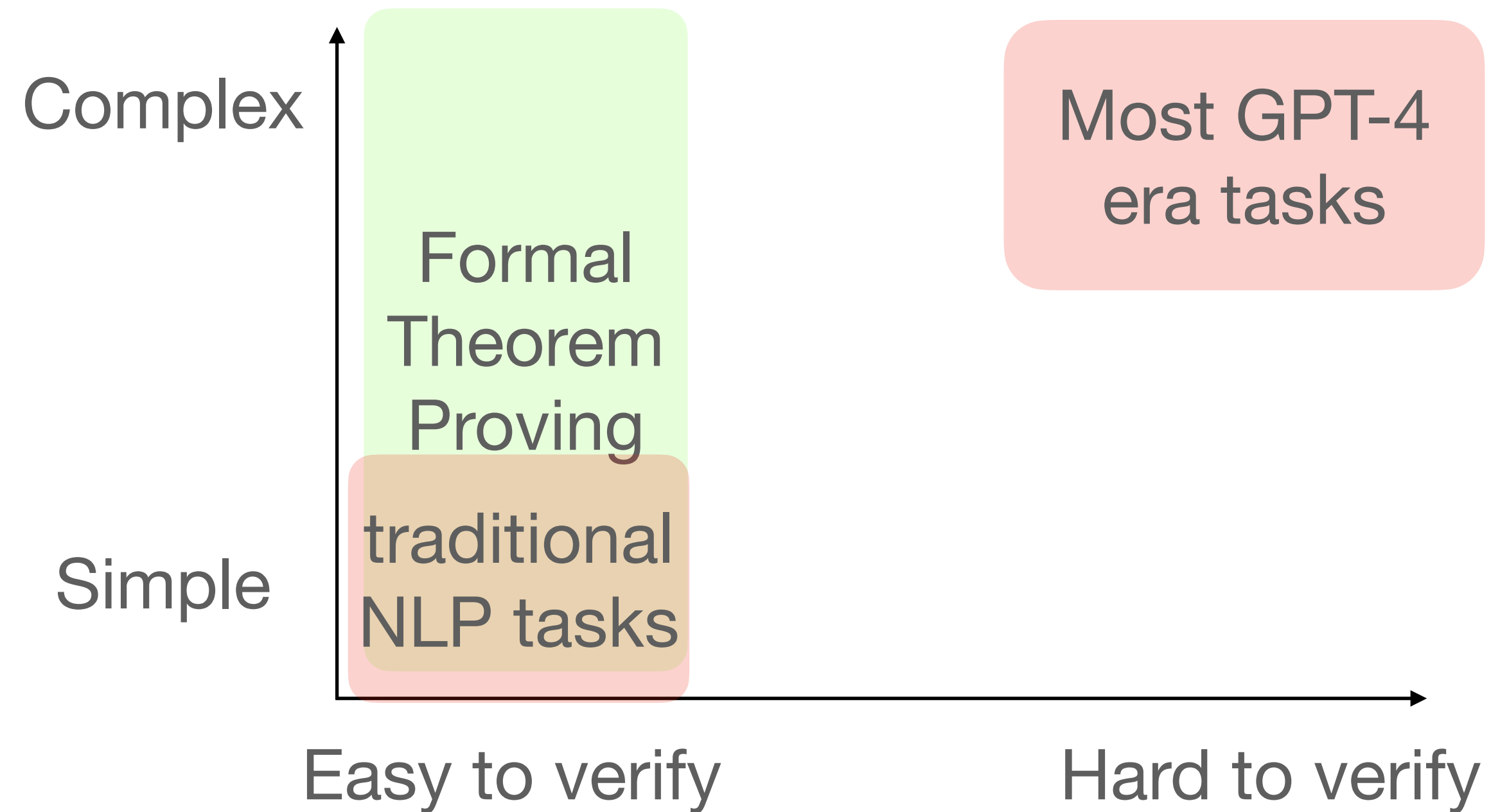
ML 1: reasoning

- Tests aspects of reasoning
- **Verifiable**, yet complex



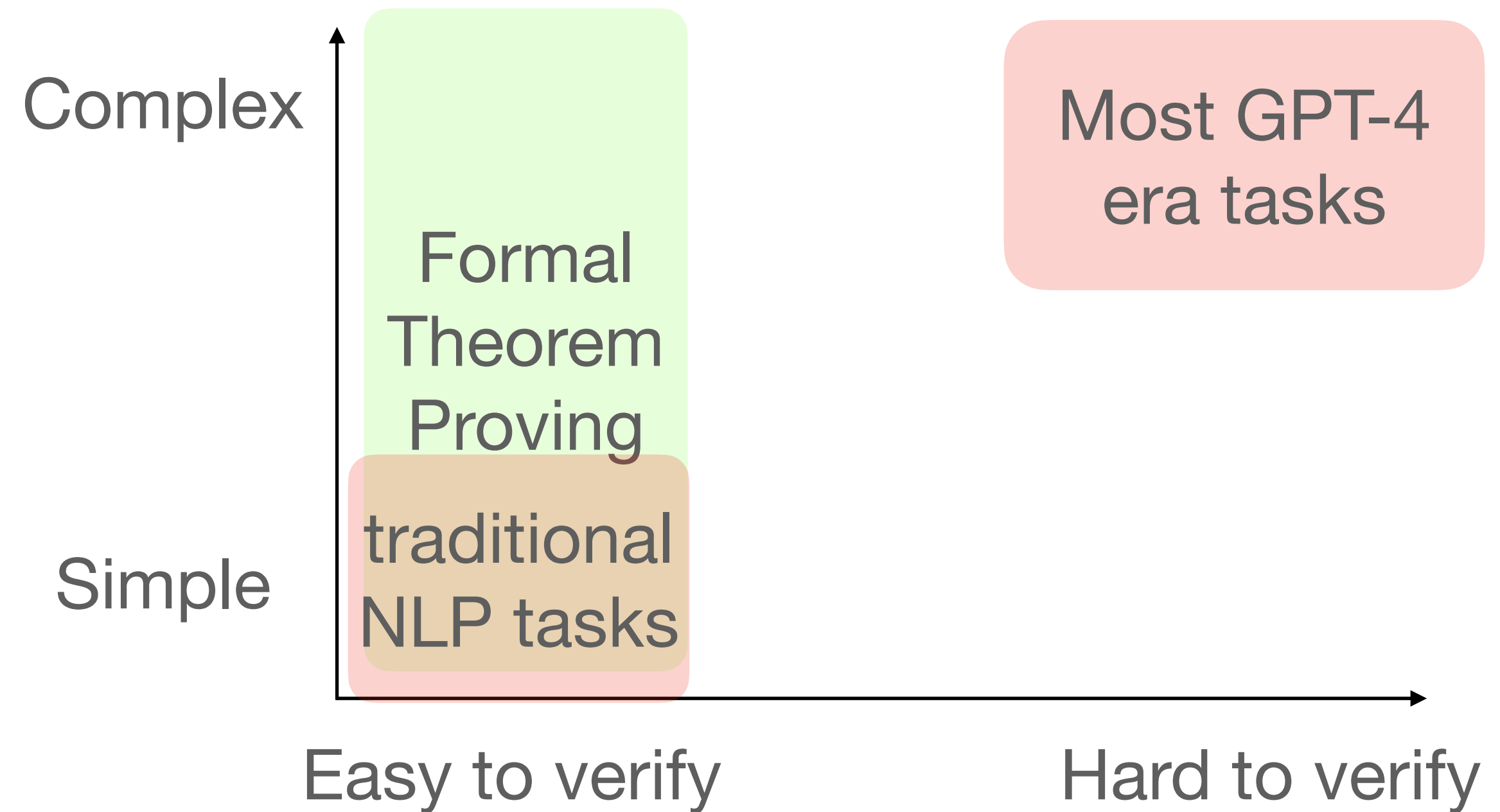
ML 1: reasoning

- Tests aspects of reasoning
- **Verifiable**, yet complex



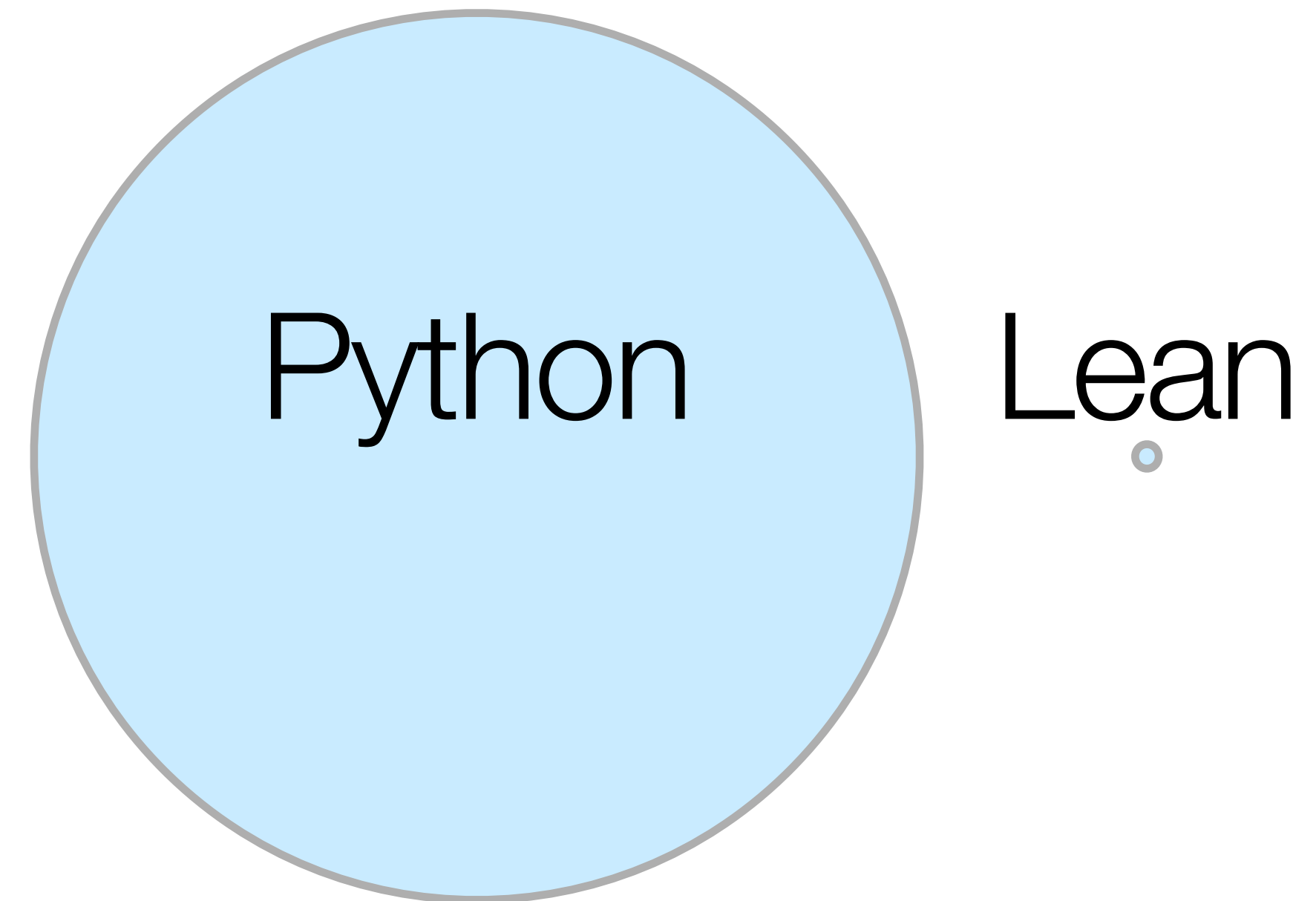
ML 1: reasoning

- Tests aspects of reasoning
- **Verifiable**, yet complex
- **Difficult** in GPT-4 era



ML 2: code generation

- Low-resource

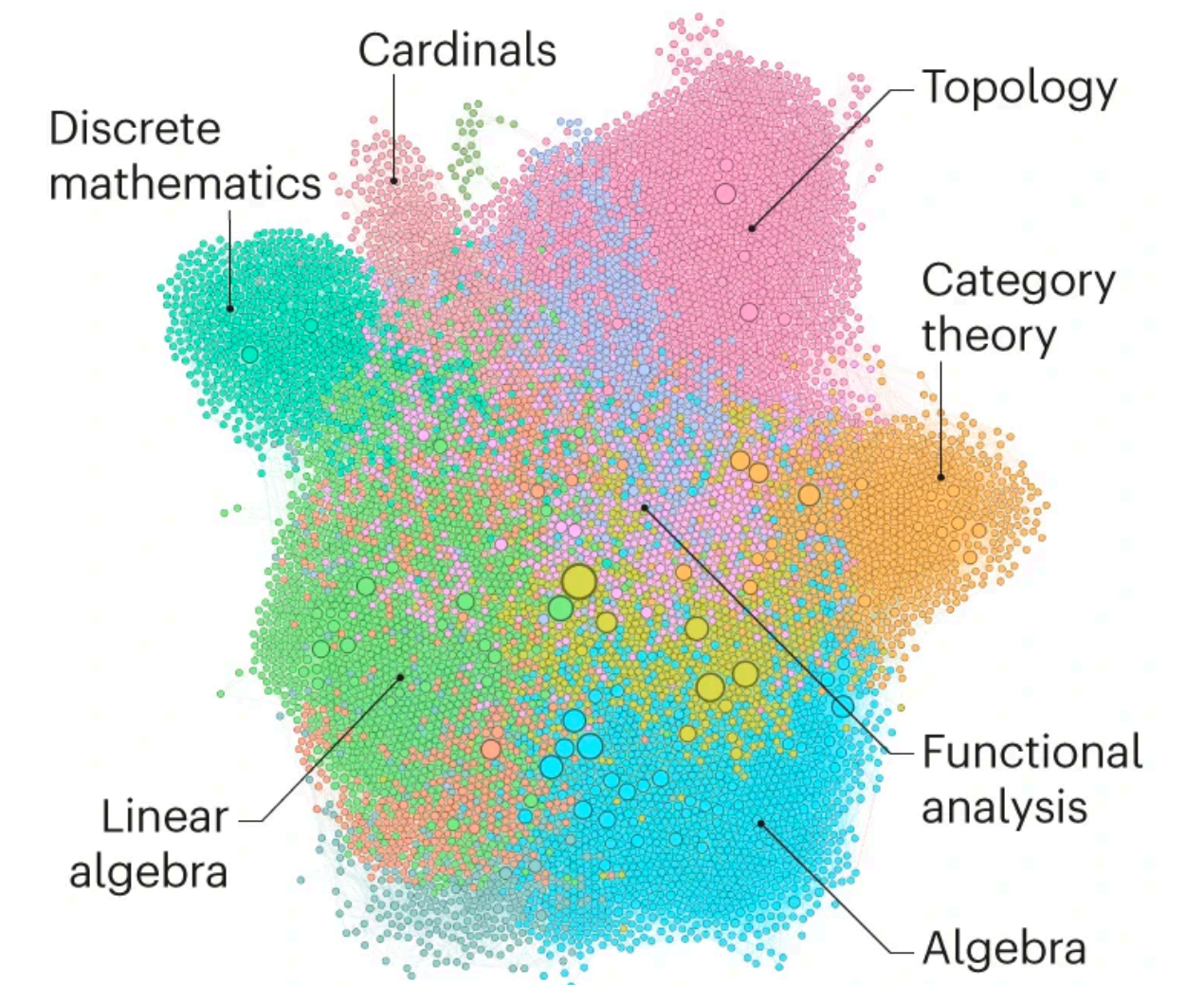


ML 2: code generation

- Low-resource
- Long-range, context-dependent

```
Code Blame 2755 Lines (2372 loc) · 149 KB
2724
2725 theorem norm_iteratedDeriv_clm_apply (f : E → F → L[k] G) (g : E → F) (N : ℕ) (n : ℕ)
2726 (hf : ContDiff k N f) (hg : ContDiff k N g) (x : E) (hx : ‖x‖ ≤ N) :
2727   iteratedDeriv k n (fun y : E => (f y) (g y)) x ≤ [1 in Finset.range (n + 1),
2728     (fun choose 1 => iteratedDeriv k 1 f x) * iteratedDeriv k (n - 1) g x] := by
2729     simp only [← iteratedDerivWithin_univ]
2730     exact norm_iteratedDerivWithin_clm_apply hf.contDiffOn hg.contDiffOn uniqueDiffOn_univ
2731     (Set.mem_univ x) hx
2732 #align norm_iterated_fderiv_clm_apply_norm_iteratedDeriv_clm_apply
2733
2734 theorem norm_iteratedDerivWithin_clm_apply_const (f : E → F → L[k] G) (g : E → F) (s : Set E) (x : E)
2735 (N : ℕ) (n : ℕ) (hf : ContDiff k N f) (hg : ContDiff k N g) (hx : ‖x‖ ≤ N) (hs : UniqueDiffOn s) (hx' : ‖x‖ ≤ N) :
2736   iteratedDerivWithin k n (fun y : E => (f y) (g y)) x ≤
2737     [C] * iteratedDerivWithin k n f s x := by
2738     let g : (F → L[k] G) → L[k] G := continuousLinearMap.apply f $ E
2739     have h := g.norm_contDiffWithin_mul [1] normMap (← iteratedDerivWithin k n f s x)
2740     rw [← g.iteratedDerivWithin_com_left h hx] at h
2741     refine' h.trans (mul_le_mul_of_nonneg_right (norm_nonneg _) _)
2742     rw [ContinuousLinearMap.apply_apply, mul_comm]
2743     exact f.le_op_norm c
2744 #align norm_iterated_fderiv_within_clm_apply_const norm_iteratedDerivWithin_clm_apply_const
2745
2746
2747 theorem norm_iteratedDeriv_clm_apply_const (f : E → F → L[k] G) (c : F) (x : E) (N : ℕ) (n : ℕ)
2748 (hf : ContDiff k N f) (hx : ‖x‖ ≤ N) :
2749   iteratedDeriv k n (fun y : E => (f y) c) x ≤ [c] * iteratedDeriv k n f x := by
2750     simp only [← iteratedDerivWithin_univ]
2751     exact norm_iteratedDerivWithin_clm_apply_const hf.contDiffOn uniqueDiffOn_univ
2752     (Set.mem_univ x) hx
2753 #align norm_iterated_fderiv_clm_apply_const norm_iteratedDeriv_clm_apply_const
2754
```

62,000+ token
Lean file

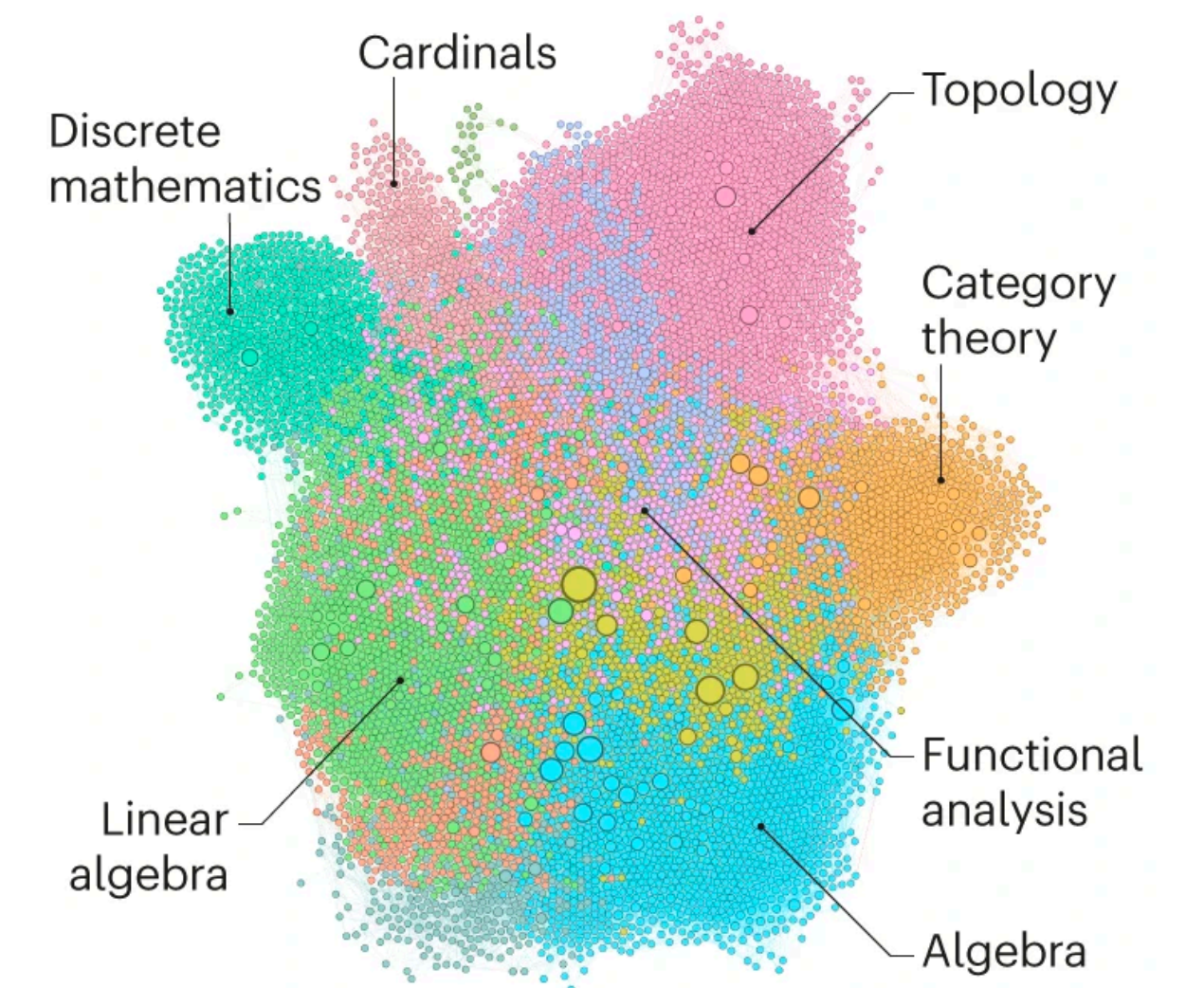


ML 2: code generation

- Low-resource
- Long-range, context-dependent
- **Verifiable**

```
Code Blame 2755 Lines (2372 loc) · 149 KB
2724
2725 theorem norm_iteratedDeriv_clm_apply (f : E → F → L[k] G) (g : E → F) (N : ℕ) (n : ℕ)
2726 (hf : ContDiff k N f) (hg : ContDiff k N g) (x : E) (hx : ‖x‖ ≤ N) :
2727   iteratedDeriv k n (fun y : E => (f y) (g y)) x ≤ [1 in Finset.range (n + 1),
2728     (fun i => iteratedDeriv k i f x) * iteratedDeriv k (n - i) g x] := by
2729     simp only [← iteratedDerivWithin_univ]
2730     exact norm_iteratedDerivWithin_clm_apply hf.contDiffOn hg.contDiffOn uniqueDiffOn_univ
2731     (Set.mem_univ x) hx
2732 #align norm_iterated_fderiv_clm_apply_norm_iteratedDeriv_clm_apply
2733
2734 theorem norm_iteratedDerivWithin_clm_apply_const (f : E → F → L[k] G) (g : E → F) (s : Set E) (x : E)
2735 (N : ℕ) (n : ℕ) (hf : ContDiff k N f) (hg : ContDiff k N g) (hx : ‖x‖ ≤ N) (hs : s ⊆ {x}) :
2736   iteratedDerivWithin k n (fun y : E => (f y) (g y)) x ≤
2737     [C] * iteratedDerivWithin k n f s x := by
2738     let g : (F → L[k] G) → L[k] G := continuousLinearMap.apply f s E
2739     have h := g.norm_contDiffWithin_mul [1 in Finset.range (n + 1), iteratedDerivWithin k n f s x]
2740     rw [← g.iteratedDerivWithin_const_left h hx] at h
2741     refine' h.trans (mul_le_mul_of_nonneg_right (norm_nonneg _) _)
2742     rw [ContinuousLinearMap.apply_apply, mul_comm]
2743     exact f.le_op_norm c
2744 #align norm_iterated_fderiv_within_clm_apply_const norm_iteratedDerivWithin_clm_apply_const
2745
2746
2747 theorem norm_iteratedDeriv_clm_apply_const (f : E → F → L[k] G) (c : F) (x : E) (N : ℕ) (n : ℕ)
2748 (hf : ContDiff k N f) (hx : ‖x‖ ≤ N) :
2749   iteratedDeriv k n (fun y : E => (f y) (c)) x ≤ [C] * iteratedDeriv k n f x := by
2750     simp only [← iteratedDerivWithin_univ]
2751     exact norm_iteratedDerivWithin_clm_apply_const hf.contDiffOn uniqueDiffOn_univ
2752     (Set.mem_univ x) hx
2753 #align norm_iterated_fderiv_clm_apply_const norm_iteratedDeriv_clm_apply_const
2754
```

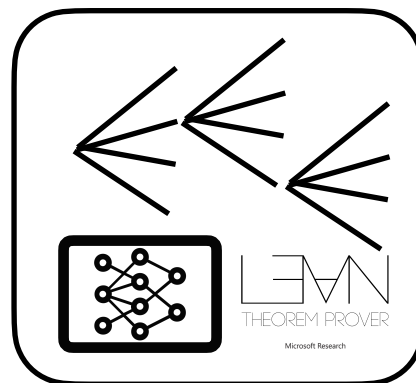
62,000+ token
Lean file



This tutorial

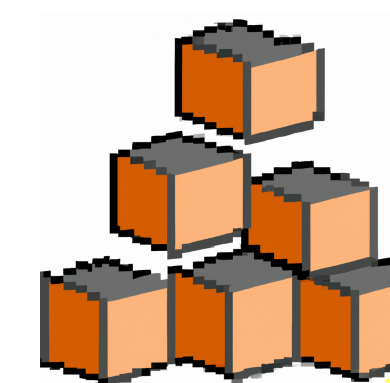
- Part I: next-step prediction

- Language model suggests next-proof-steps
- Tree search



- Part II: language cascades

- Compose language model functions
- Sketching, correction, tools



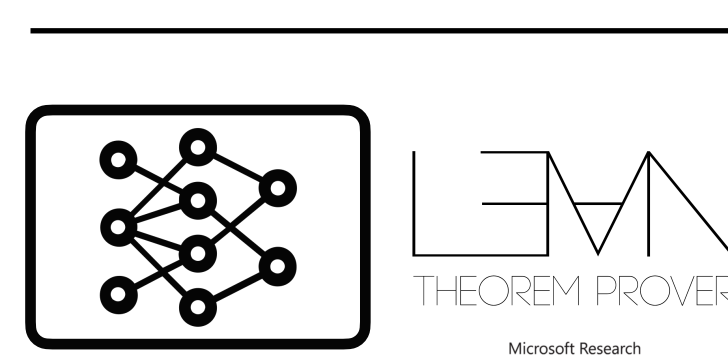
Tutorial code: <https://github.com/wellecks/ntptutorial>

Next-step prediction

- Build a “co-pilot” from scratch

```
example : ∀ (a: ℤ), a + 3 = 0 → a = -3 := by
  intro a ha
  llmstep
```

Proof state



▼ llmstep suggestions

Try this:

- `linarith`
- `rw [← sub_eq_zero] at ha`
- `apply eq_neg_of_add_eq_zero_left`
- `rw [← Int.negSucc_coe] at ha`

Next-step suggestions

Demo

```
42 example (f : ℕ → ℕ) : Monotone f → ∀ n, f n ≤ f (n + 1) := by
43 |
44
45
46
47
48
49
50
51
52
```

Lean Infoview ×

- ▼ Examples.lean:43:2
- ▼ Tactic state
- 1 goal**
- f** : ℕ → ℕ
- ┆ Monotone f → ∀ (n : ℕ), f n ≤ f (n + 1)

► All Messages (2)



Topic	Notebook
0. Intro	notebook
1. Data	notebook
2. Learning	notebook
3. Proof Search	notebook
4. Evaluation	notebook
5. <code>llmsuggest</code>	notebook

Tutorial code: <https://github.com/wellecks/ntptutorial>

0. Problem setup

- Proof: sequence of (state, step)
 - $(x_0, y_0), \dots, (x_t, y_t), \dots, (x_T, y_T)$



0. Problem setup

- Proof: sequence of (state, step)
- $(x_0, y_0), \dots, (x_t, y_t), \dots, (x_T, y_T)$
- x_t : proof state 
- y_t : proof step 



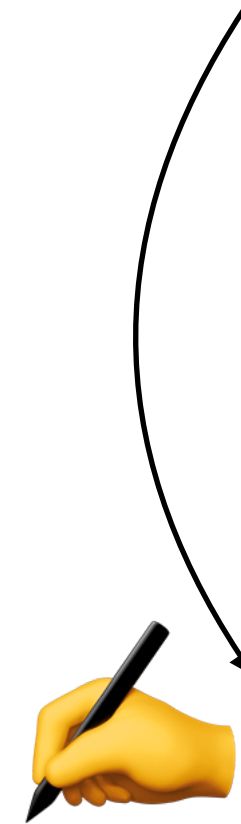
The screenshot shows a theorem prover interface with a dark background. At the top, a goal is displayed: `example : R ⊆ S → S ⊆ T → R ⊆ T := by`. Below this, a "Tactic state" panel is visible, containing the text: `1 goal`, `α : Type`, `R S T : Set α`, and `⊢ R ⊆ S → S ⊆ T → R ⊆ T`. To the right of the tactic state is the logo for "LEMN THEOREM PROVER" by Microsoft Research. A hand holding a pen is shown pointing to the `by` keyword in the goal, and another hand holding a pen is shown pointing to the `intro h1 h2` tactic in the input field below the goal.

0. Problem setup

- Proof: sequence of (state, step)
- $(x_0, y_0), \dots, (x_t, y_t), \dots, (x_T, y_T)$
- x_t : proof state 
- y_t : proof step 

```
example : R ⊆ S → S ⊆ T → R ⊆ T := by
```

```
▼Tactic state  
1 goal  
α : Type  
R S T : Set α  
⊢ R ⊆ S → S ⊆ T → R ⊆ T
```






```
example : R ⊆ S → S ⊆ T → R ⊆ T := by  
intro h1 h2
```

```
▼Tactic state  
1 goal  
α : Type  
R S T : Set α  
h1 : R ⊆ S  
h2 : S ⊆ T  
⊢ R ⊆ T
```

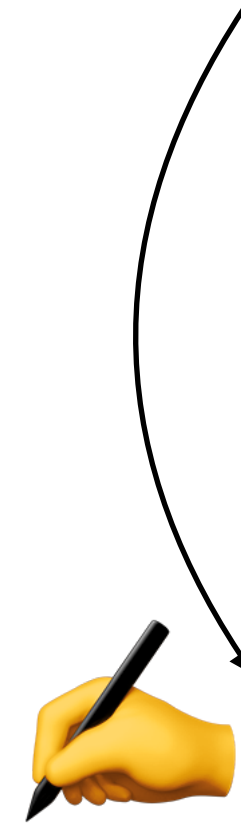


0. Problem setup

- Proof: sequence of (state, step)
- $(x_0, y_0), \dots, (x_t, y_t), \dots, (x_T, y_T)$
- x_t : proof state 
- y_t : proof step 
- x_T : “proof complete” 

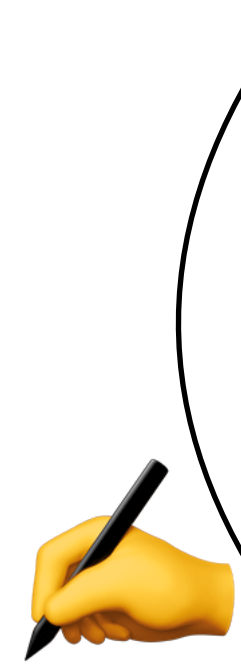
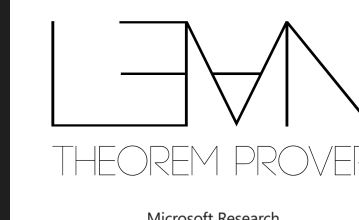
```
example : R ⊆ S → S ⊆ T → R ⊆ T := by
```

```
▼Tactic state  
1 goal  
α : Type  
R S T : Set α  
⊢ R ⊆ S → S ⊆ T → R ⊆ T
```




```
example : R ⊆ S → S ⊆ T → R ⊆ T := by  
intro h1 h2
```

```
▼Tactic state  
1 goal  
α : Type  
R S T : Set α  
h1 : R ⊆ S  
h2 : S ⊆ T  
⊢ R ⊆ T
```

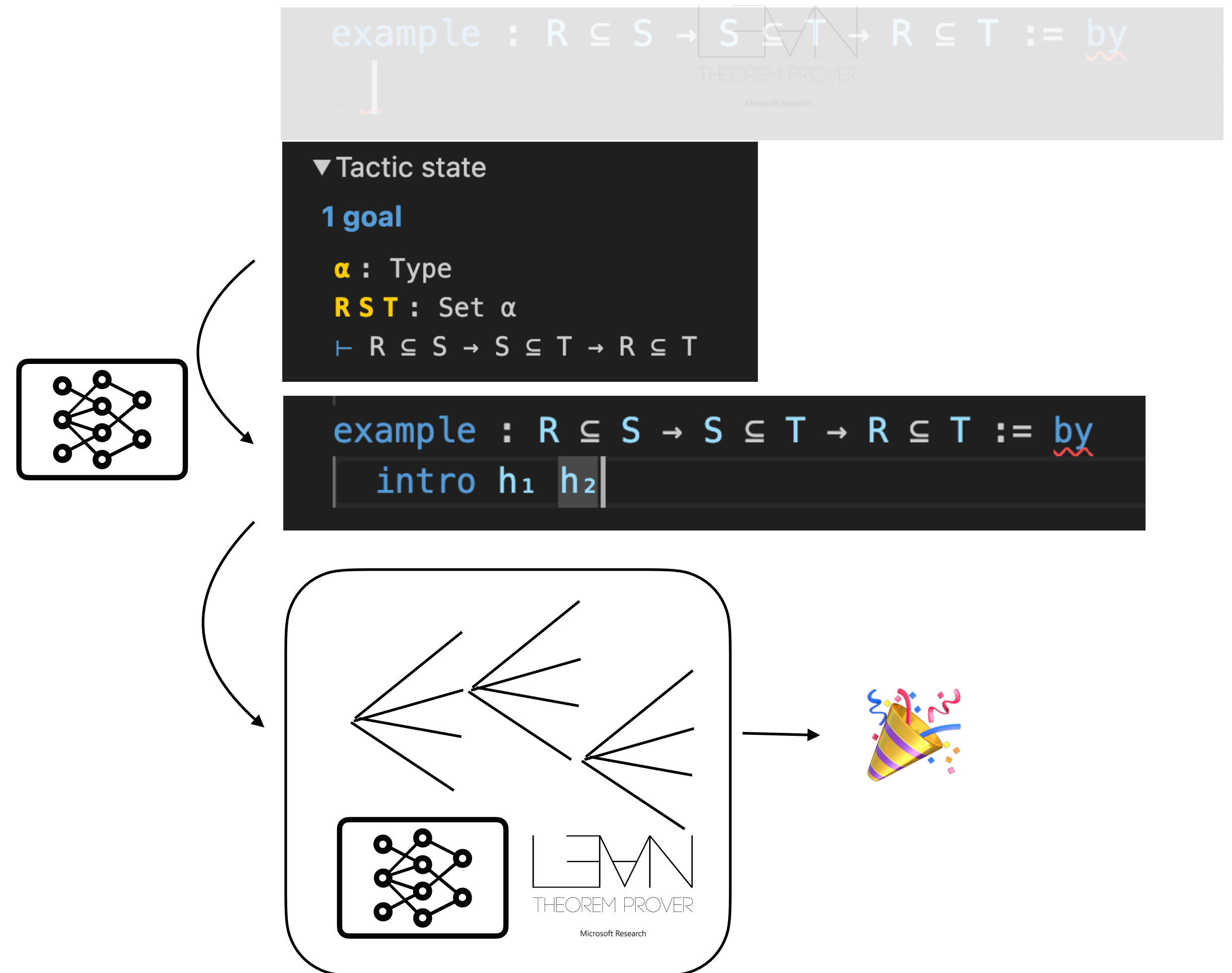


```
example : R ⊆ S → S ⊆ T → R ⊆ T := by  
intro h1 h2  
exact h1.trans h2
```

```
No goals 
```

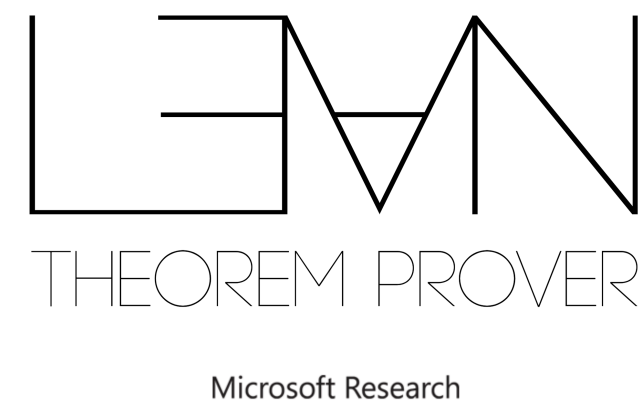
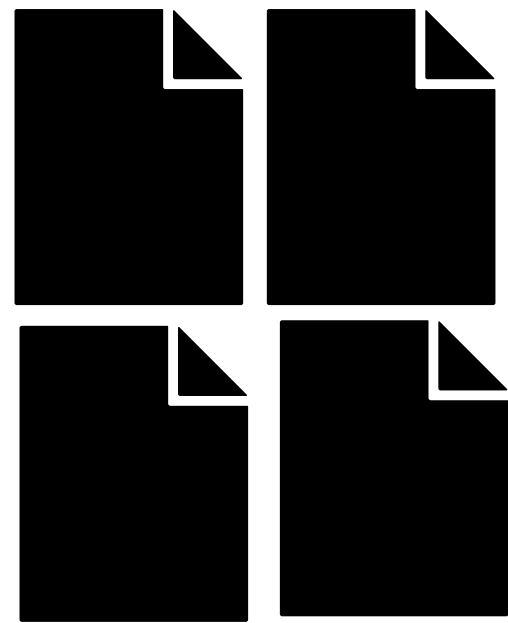

0. Problem setup

- Proof: sequence of (state, step)
 - $(x_0, y_0), \dots, (x_t, y_t), \dots, (x_T, y_T)$
- Language model:
 - $p_{\theta}(y_t | x_t)$
- Tree search to generate full proof



1. Data

- Extract (proofstate, next-step) pairs from human-written proofs

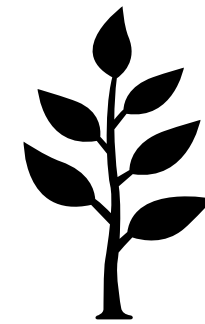
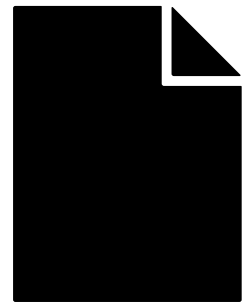


$$D = \{(x_t, y_t)\}$$

Data

Elaboration

Postprocess



$D = \{(x_t, y_t)\}$

```
example : R ⊆ S → S ⊆ T → R ⊆ T := by
  intro h1 h2
  exact h1.trans h2
```

```
{'args': [{'node': {'args': [...],
                    'info': 'none',
                    'kind': 'Lean.Parser.Command.declModifiers'}},
          {'node': {'args': [...],
                    'info': 'none',
                    'kind': 'Lean.Parser.Command.theorem'}}]},
'info': 'none',
'kind': 'Lean.Parser.Command.declaration'}
```

```
--- x1 ---
m n : ℕ
h : Nat.coprime m n
⊢ Nat.gcd m n = 1
--- y1 ---
rw [Nat.coprime] at h

--- x2 ---
m n : ℕ
h : Nat.gcd m n = 1
⊢ Nat.gcd m n = 1
--- y2 ---
exact h
```


Data

- Lean Dojo [Yang et al 2023]
 - repository $\rightarrow \{(x_t, y_t)\}$
 - Dependencies, multiple files, versioning

```
URL = "https://github.com/leanprover-community/mathlib4"  
COMMIT = "5a919533f110b7d76410134a237ee374f24eaaad"  
repo = LeanGitRepo(URL, COMMIT)  
traced_repo = trace(repo)
```

Data

- Mathlib:
 - 41,944 theorems + proofs

- => training data

train	169530
val	4053
test	3606

Data

- Mathlib:
 - 41,944 theorems + proofs
- => training data

train	169530
val	4053
test	3606

```
Number of non-empty training proofs: 41944
{'commit': '5a919533f110b7d76410134a237ee374f24eaaad',
 'end': [308, 76],
 'file_path': 'Mathlib/Analysis/BoxIntegral/Box/Basic.lean',
 'full_name': 'BoxIntegral.Box.withBotCoe_inj',
 'start': [307, 1],
 'traced_tactics': [{ 'state_after': 'no goals',
                      'state_before': 'ι : Type u_1\n'
                                      'I J : Box ι\n'
                                      'x y : ι → ℝ\n'
                                      'I J : WithBot (Box ι)\n'
                                      '⊢ ↑I = ↑J ↔ I = J',
                      'tactic': 'simp only [Subset.antisymm_iff, ← '
                                'le_antisymm_iff, withBotCoe_subset_iff]'},
                    ],
 'url': 'https://github.com/leanprover-community/mathlib4'}
```

2. Learning

- Standard supervised fine-tuning on $D = \{(x_t, y_t)\}$:

- $$\max_{\theta} \sum_{(x_t, y_t) \in D} -\log p_{\theta}(y_t | x_t)$$

2. Learning

Input:

[GOAL] ι : Type u_1

I† J† : Box ι

x_t x y : $\iota \rightarrow \mathbb{R}$

I J : WithBot (Box ι)

$\vdash \uparrow I = \uparrow J \leftrightarrow I = J$ [PROOFSTEP]

Output:

y_t simp only [Subset.antisymm_iff, ← le_antisymm_iff, withBotCoe_subset_iff]<|endoftext|>

2. Learning

<https://huggingface.co/wellecks/llmstep-mathlib4-pythia2.8b>

```
import transformers
```

```
MODEL = 'wellecks/llmstep-mathlib4-pythia2.8b'
```

```
model = transformers.GPTNeoXForCausalLM.from_pretrained(MODEL)
```

```
tokenizer = transformers.GPTNeoXTokenizerFast.from_pretrained(MODEL)
```

```
prompt = """[GOAL]m n :  $\mathbb{N}$ 
```

```
h : Nat.coprime m n
```

```
⊢ Nat.gcd m n = 1[PROOFSTEP]"""
```

```
input_ids = tokenizer.encode(prompt, return_tensors='pt')
```

```
out = model.generate(input_ids)
```

```
text = tokenizer.decode(out[0][input_ids.shape[1]:], skip_special_tokens=True)
```

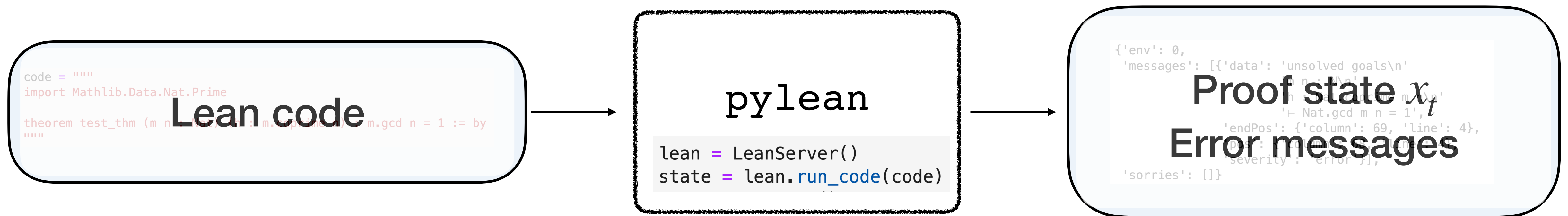
```
print(text)
```

```
rw [← h.gcd_eq_one]
```

3. Proof search

- Use next-step predictor $p_{\theta}(y_t | x_t)$ to generate a full proof y_1, \dots, y_T
- We need:
 - **Interaction** with Lean
 - **Algorithm** for search

Interaction | pylean



Interaction | pylean

```
[4]: # Generate a next step
prompt = f"[GOAL]{get_goal(state)}[PROOFSTEP]"

next_step = generate(prompt)
print(next_step)

rw [← h.gcd_eq_one]
```

Finally, we can give the generated next step to Lean and receive the next state.

```
[5]: code = """
import Mathlib.Data.Nat.Prime

theorem test_thm (m n : Nat) (h : m.coprime n) : m.gcd n = 1 := by

""" + next_step

lean = LeanServer()
state = lean.run_code(code)
lean.proc.close()

pprint(state)

{'env': 0, 'messages': [], 'sorries': []}
```

Interaction | Lean Dojo

- Lean Dojo
 - Traces entire Lean repository
 - Provides abstractions for interaction

Best-first search

type-checked candidates:

(-0.066) rintro rfl

(-0.307) rintro ⟨rfl, rfl⟩

(-0.035) intro h

(-0.230) rintro ⟨d, rfl⟩



Best-first search

type-checked candidates:

(-0.066) rintro rfl

(-0.307) rintro ⟨rfl, rfl⟩

(-0.035) intro h

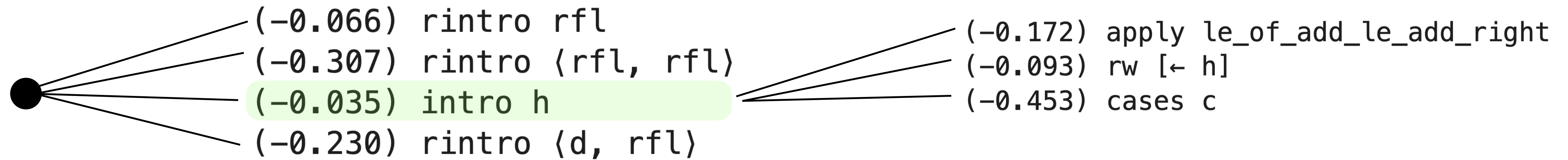
(-0.230) rintro ⟨d, rfl⟩

$$\frac{1}{Z} \sum_t \log p_{\theta}(y_t | x_t)$$

“value function”

Best-first search

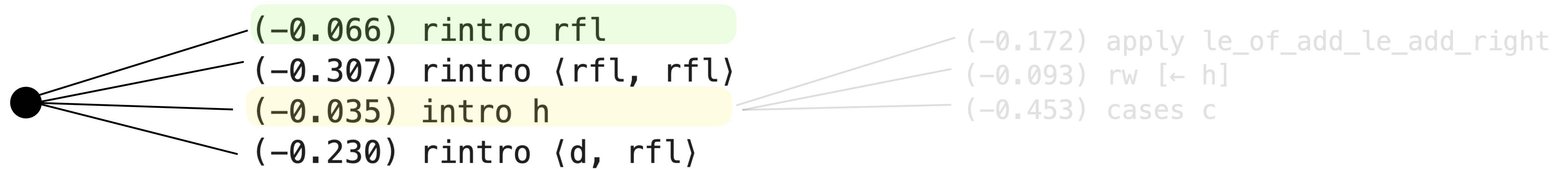
type-checked candidates:



$$\frac{1}{Z} \sum_t \log p_{\theta}(y_t | x_t)$$

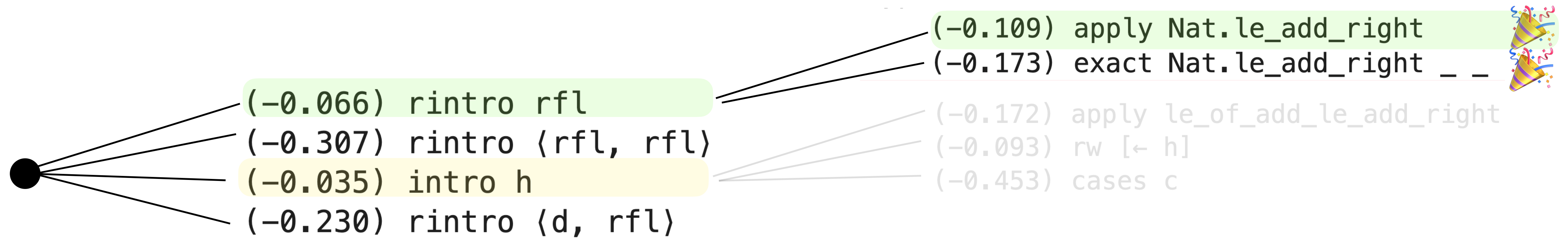
“value function”

Best-first search



$$\frac{1}{Z} \sum_t \log p_{\theta}(y_t | x_t)$$

Best-first search



$$\frac{1}{Z} \sum_t \log p_{\theta}(y_t | x_t)$$

Best-first search

```
proofsearch.best_first_search(  
  model, tokenizer, header, theorem_statement,  
  max_iters=32,  
  num_samples=4,  
  temperatures=[0.0],  
  verbose=True  
)
```



```
--- current:  
theorem thm1 (a b c : Nat) : a + b = c → a ≤ c := by  
100%|██████████| 4/4 [00:03<00:00, 1.10it/s]  
--- type-checked candidates:  
(-0.066) rintro rfl  
(-0.307) rintro (rfl, rfl)  
(-0.035) intro h  
(-0.230) rintro (rfl, rfl)  
--- current:  
theorem thm1 (a b c : Nat) : a + b = c → a ≤ c := by  
intro h  
100%|██████████| 4/4 [00:03<00:00, 1.11it/s]  
--- type-checked candidates:  
(-0.173) exact Nat.le_add_right _ _  
(-0.099) apply Nat.le_add_right  
(-0.455) case h  
--- current:  
theorem thm1 (a b c : Nat) : a + b = c → a ≤ c := by  
rintro rfl  
100%|██████████| 4/4 [00:03<00:00, 1.10it/s]  
--- type-checked candidates:  
(-0.109) apply Nat.le_add_right  
(-0.173) exact Nat.le_add_right _ _
```

**Search
Trajectories**



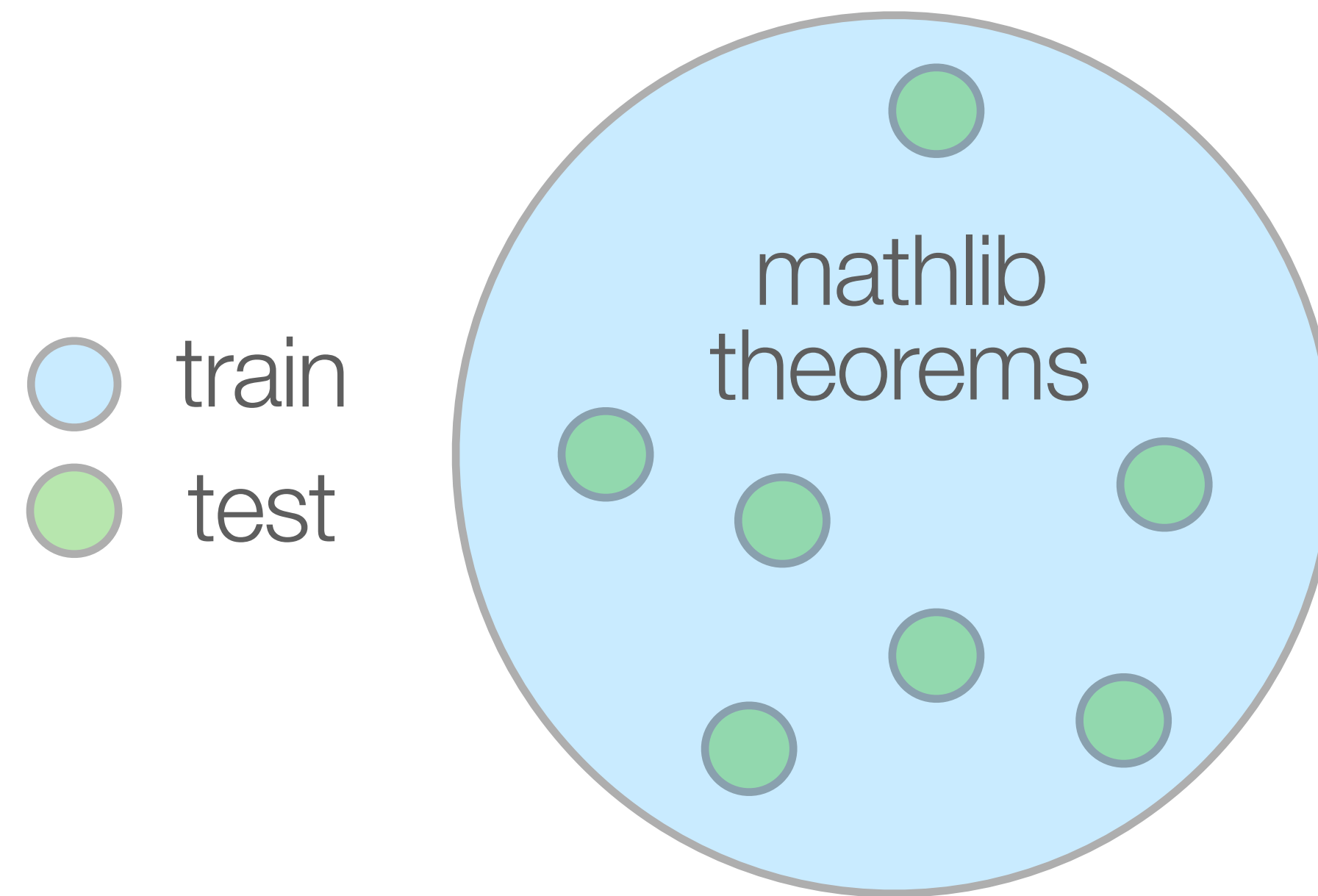
**Generated
Proof**

SUCCESS!

```
theorem thm1 (a b c : Nat) : a + b = c → a ≤ c := by  
  rintro rfl  
  apply Nat.le_add_right
```


4. Evaluation | in-domain

- Held-out theorems from training distribution



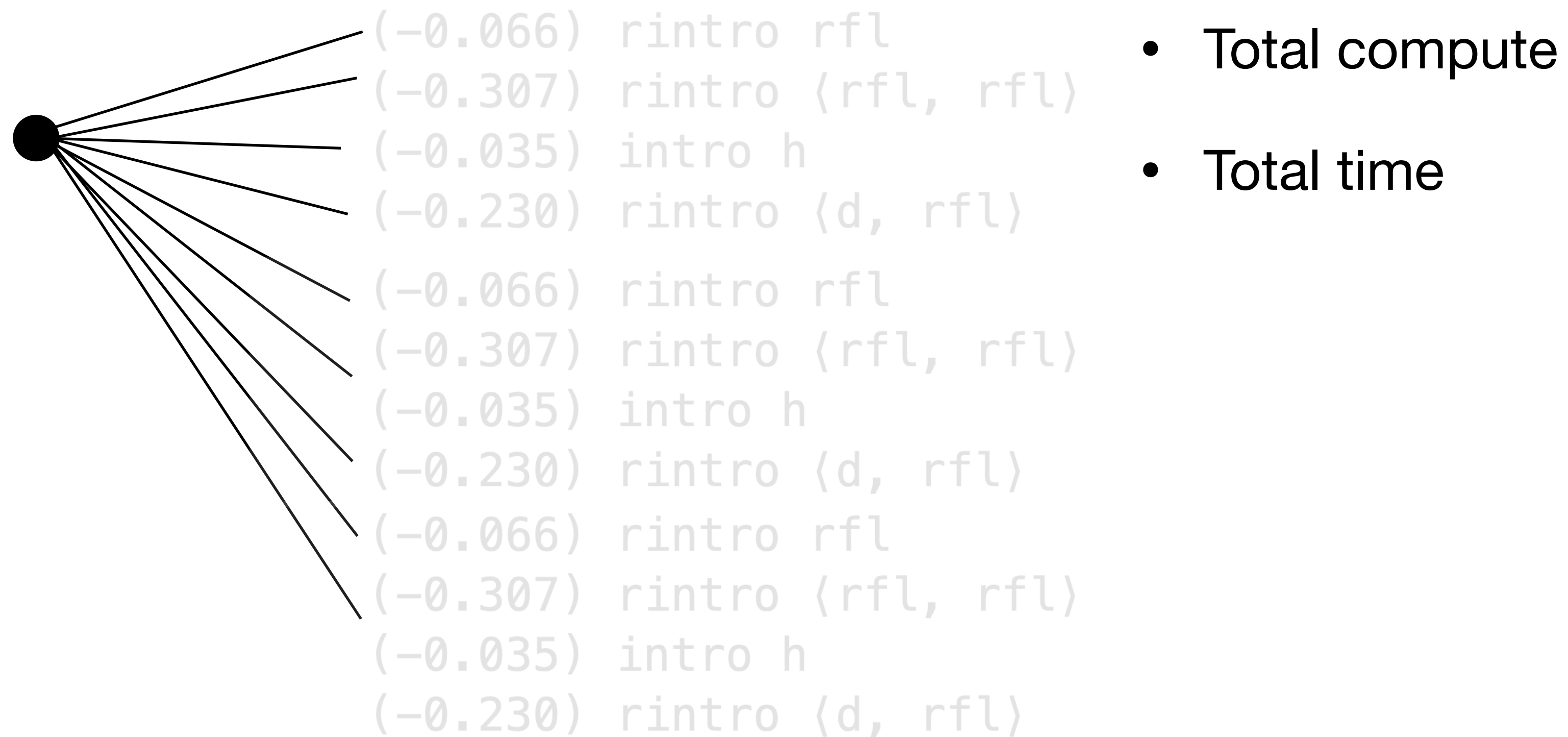
Performance depends on search

- pass rate = $f(p_\theta, \text{search}, \text{budget})$

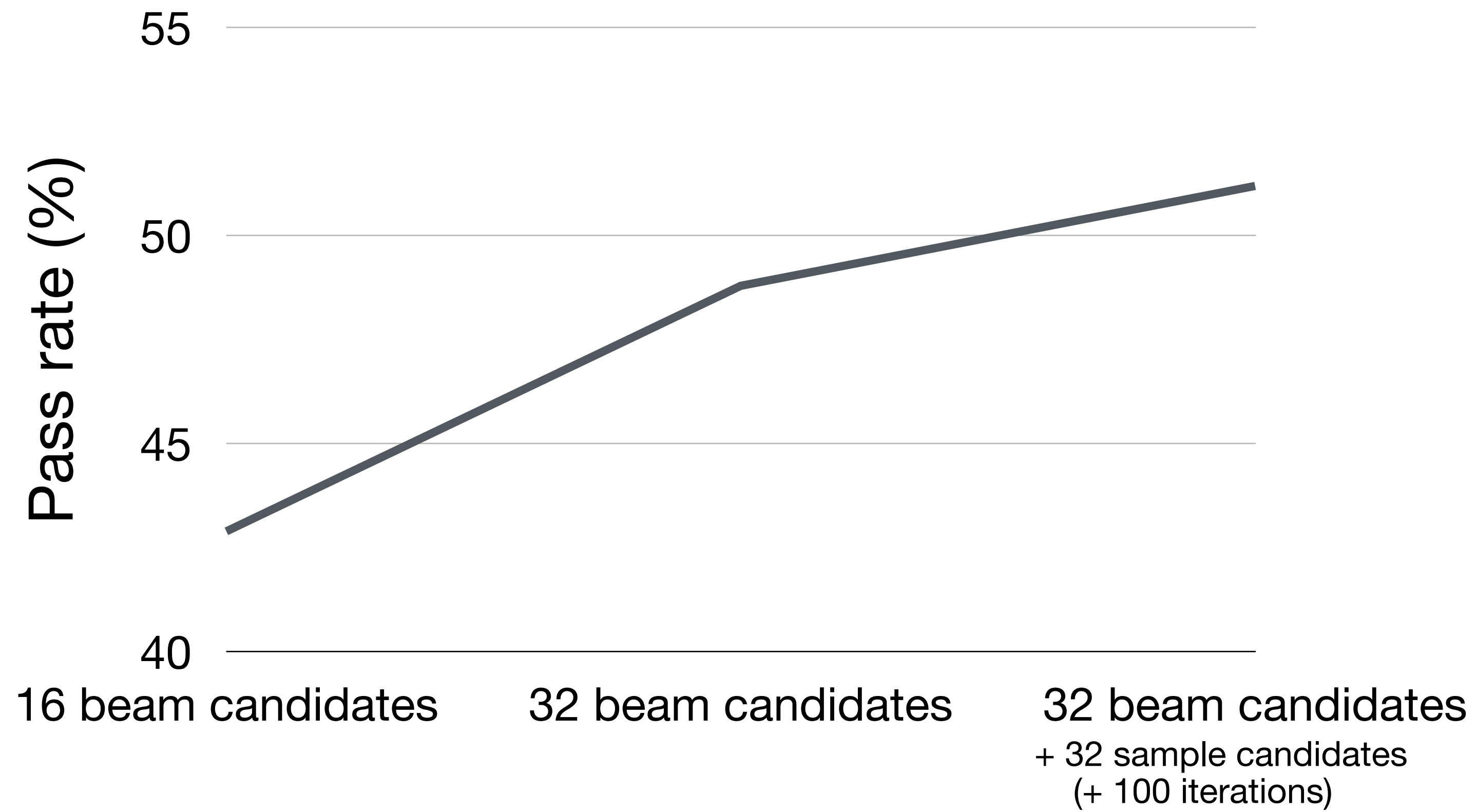


Performance depends on search

- pass rate = $f(p_\theta, \text{search}, \text{budget})$



- mathlib4 LeanDojo validation set

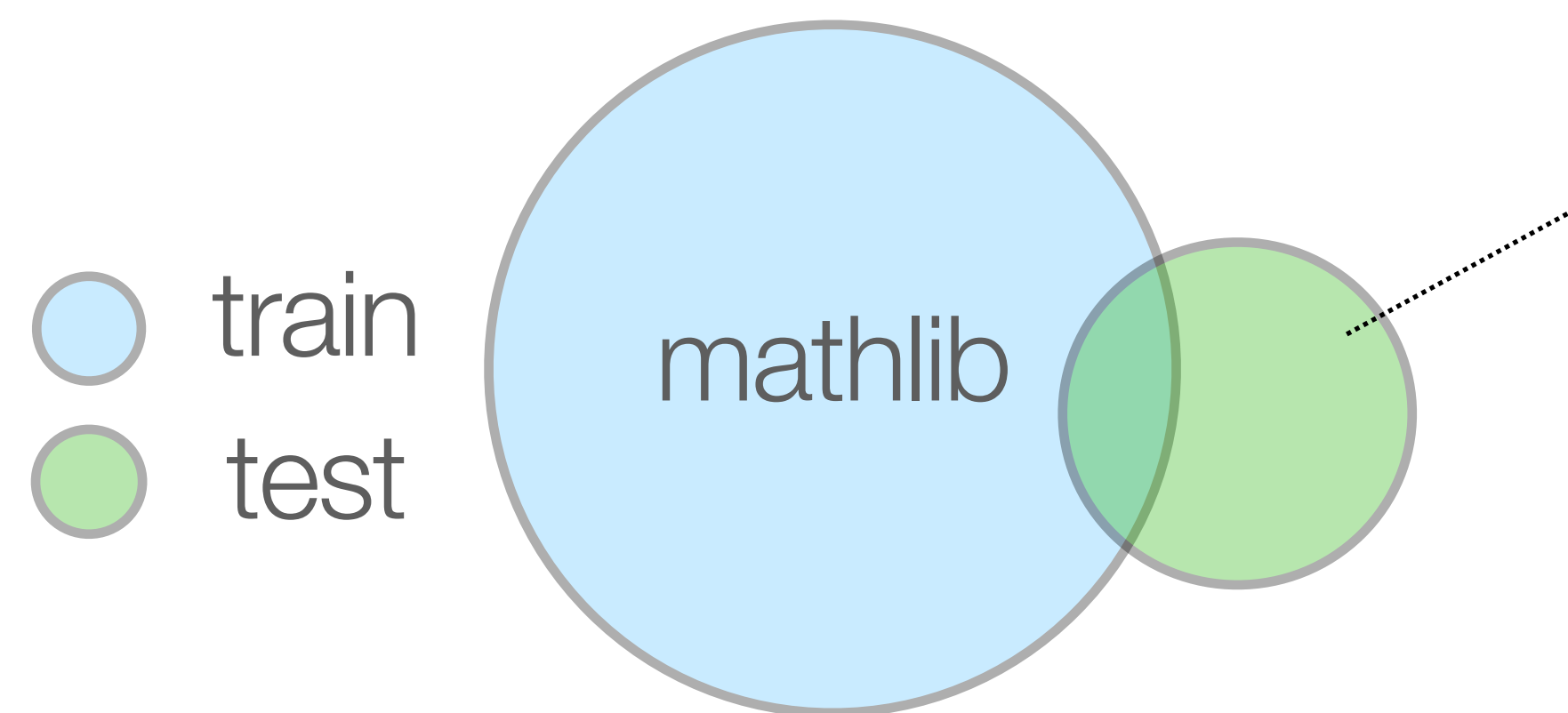


- recent results (Mathlib 3):

Method	random	novel_premises
tidy	23.8	5.4
GPT-4	28.8	7.5
ReProver (ours)	51.4	26.2
w/o retrieval	47.5	22.9

Evaluation | out-of-domain

- **MiniF2F**: 480 competition problems



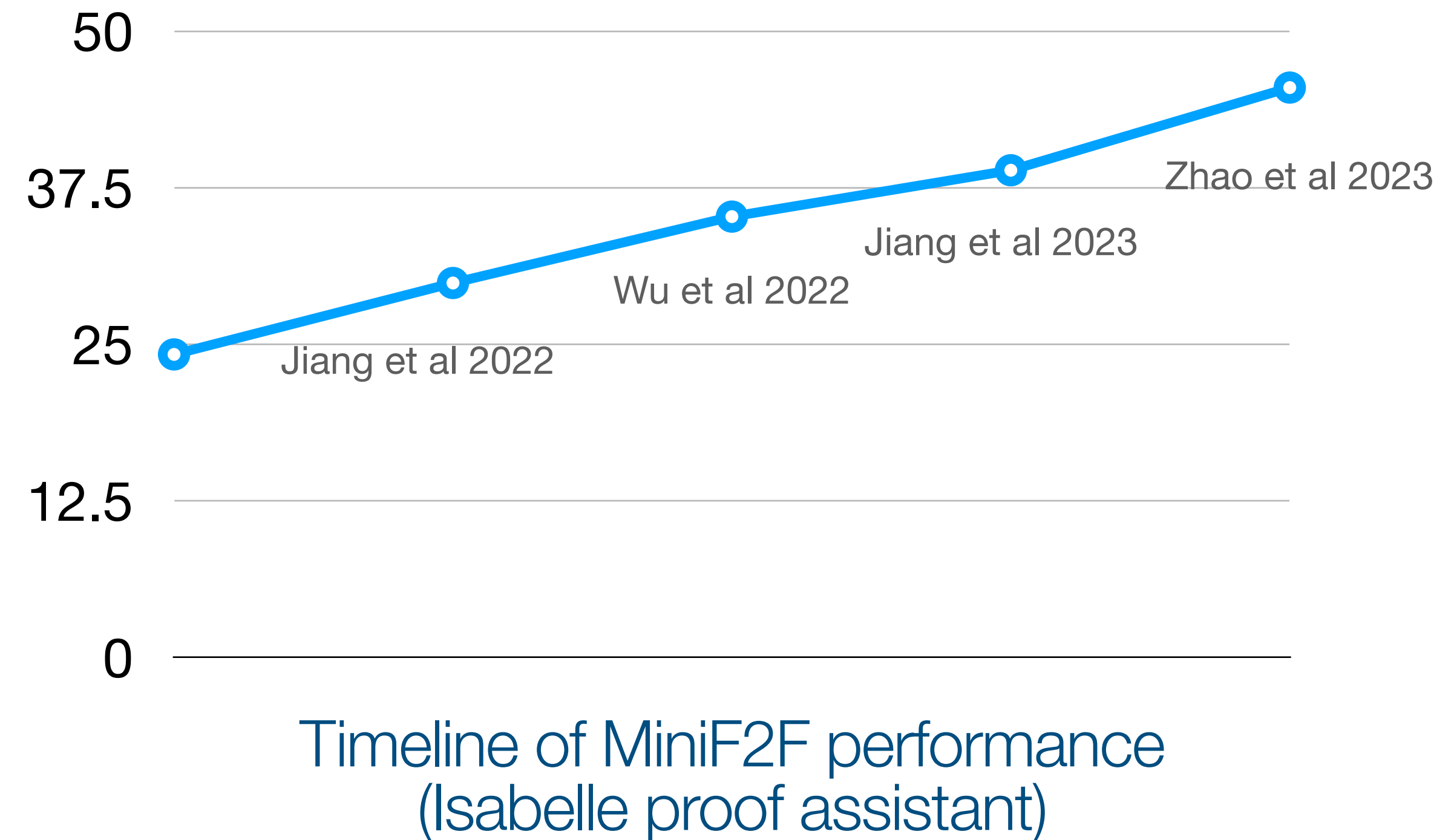
Problem 1959 IMO Problems/Problem 1

Prove that the fraction $\frac{21n + 4}{14n + 3}$ 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 :=
begin
```

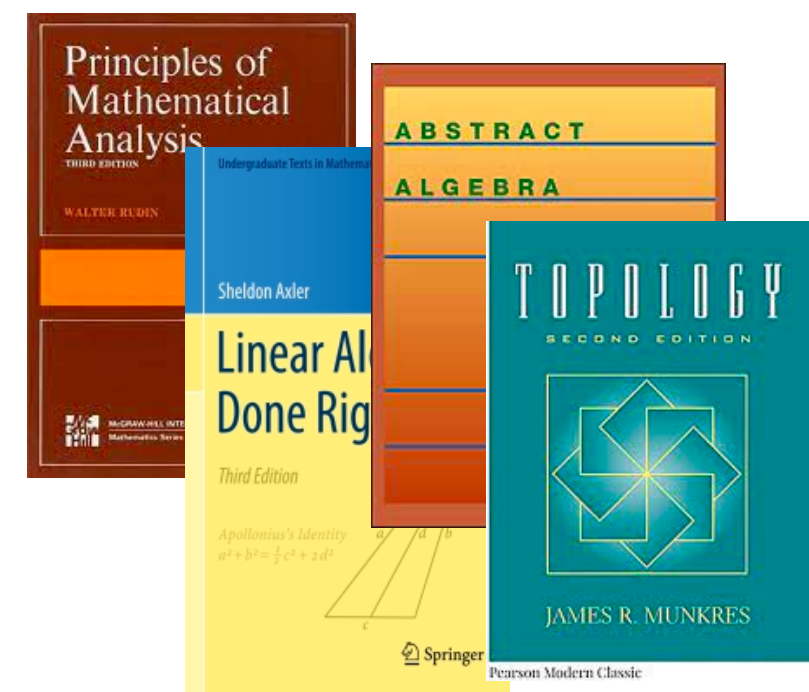
Evaluation | out-of-domain

- **MiniF2F**: 480 competition problems
 - Benchmark still challenging (e.g. olympiad problems)
 - **Lean** [Yang et al 2023]: **26.5%**



Evaluation | out-of-domain

- **ProofNet**: undergraduate textbooks
 - + informal statements/proofs



```
theorem exercise_4_5_14 {G : Type*}
  [group G] [fintype G]
  (hG : card G = 312) :
  ∃ (p : ℕ) (P : sylow p G), P.normal
```


5. interactive tool | llmstep

<https://github.com/wellecks/llmstep>

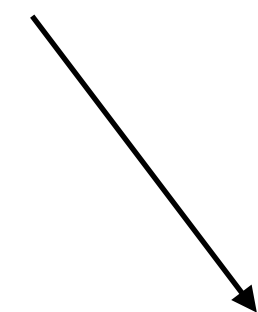
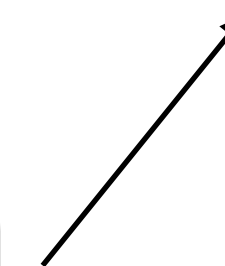
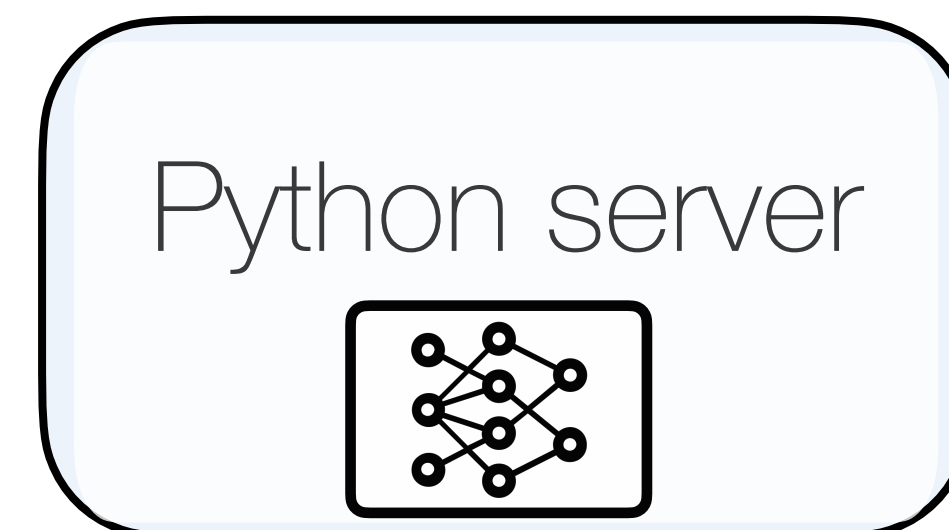
- [L]LM proof step suggestions
- Runs on your own device

```
example : ∀ (a: ℤ), a + 3 = 0 → a = -3 := by  
  intro a ha  
  llmstep
```

▼llmstep suggestions

Try this:

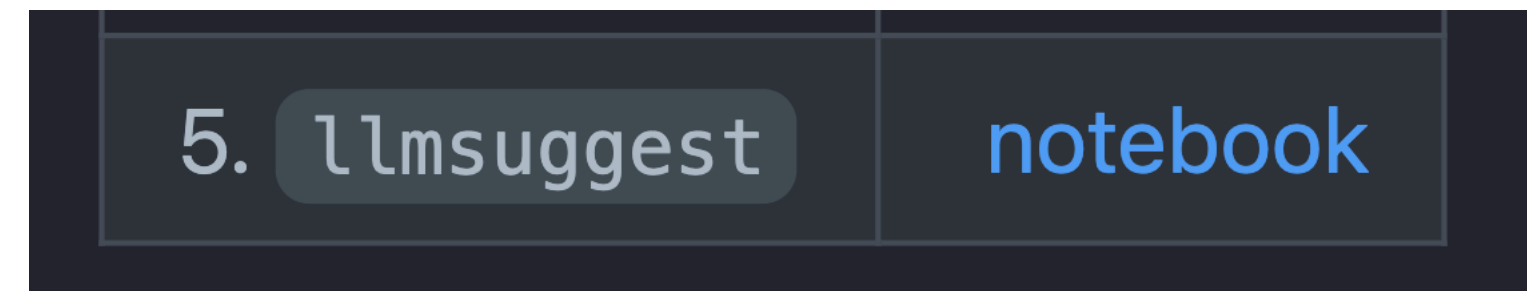
- 🚀 `linarith`
- `rw [← sub_eq_zero] at ha`
- `apply eq_neg_of_add_eq_zero_left`
- `rw [← Int.negSucc_coe] at ha`



5. interactive tool | llmstep

<https://github.com/wellecks/llmstep>

- [L]LM proof step suggestions
- Runs on your own device
- Simplified version:



<https://github.com/wellecks/ntptutorial>

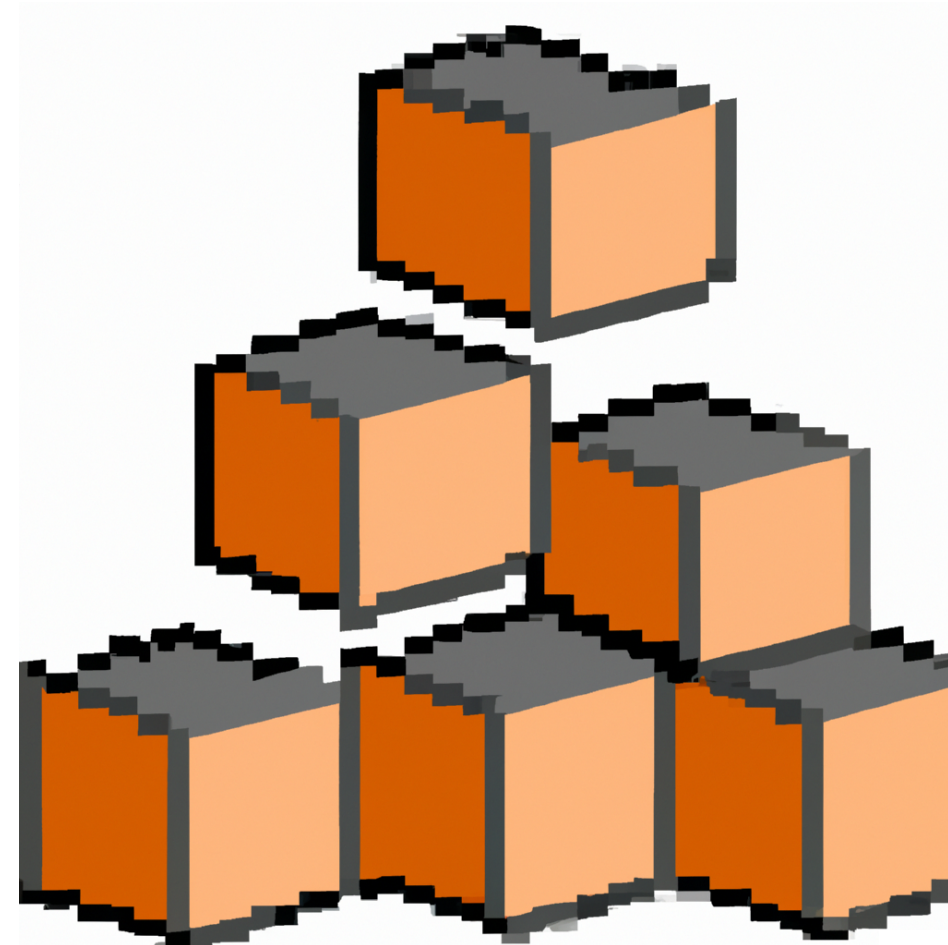
Summary

- Next-step suggestion $p_{\theta}(y_t | x_t)$
- Data, learning, search, “human-machine collaboration”

Part II: language cascades

Statement

If $\text{gcd}(n, 4) = 1$ and
 $\text{lcm}(n, 4) = 28$,
show that n is 7.



Verified formal proof

```
have c1: "1*28 = n*4"  
using assms  
  by (smt (z3) prod_gcd_lcm_nat)  
then have c2: "n = 1*28/4"  
  by auto  
then show ?thesis  
  by auto
```

Part II: language cascades

Topic	Notebook
1. Language model cascades	notebook
2. Draft, Sketch, Prove	notebook

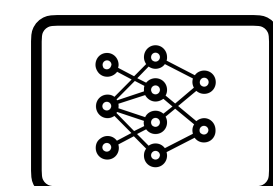
Tutorial code: <https://github.com/wellecks/ntptutorial>

Language model cascade

Prompted language model implements a function

$$y \sim p_{\theta}(y | x; P)$$

```
python_code = gpt4(  
    prompt,  
    """There are 800+72 apples in a barrel.  
    How many apples in 2233 barrels?""")  
)
```



$(800+72)*2233$

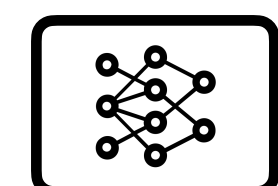
Language model cascade

Chain together functions to form a “cascade”

$$y \sim p_{\theta}(y | x; P)$$

$$z = f(y)$$

```
python_code = gpt4(  
    prompt,  
    """There are 800+72 apples in a barrel.  
    How many apples in 2233 barrels?""")  
)
```



$(800+72)*2233$



1947176

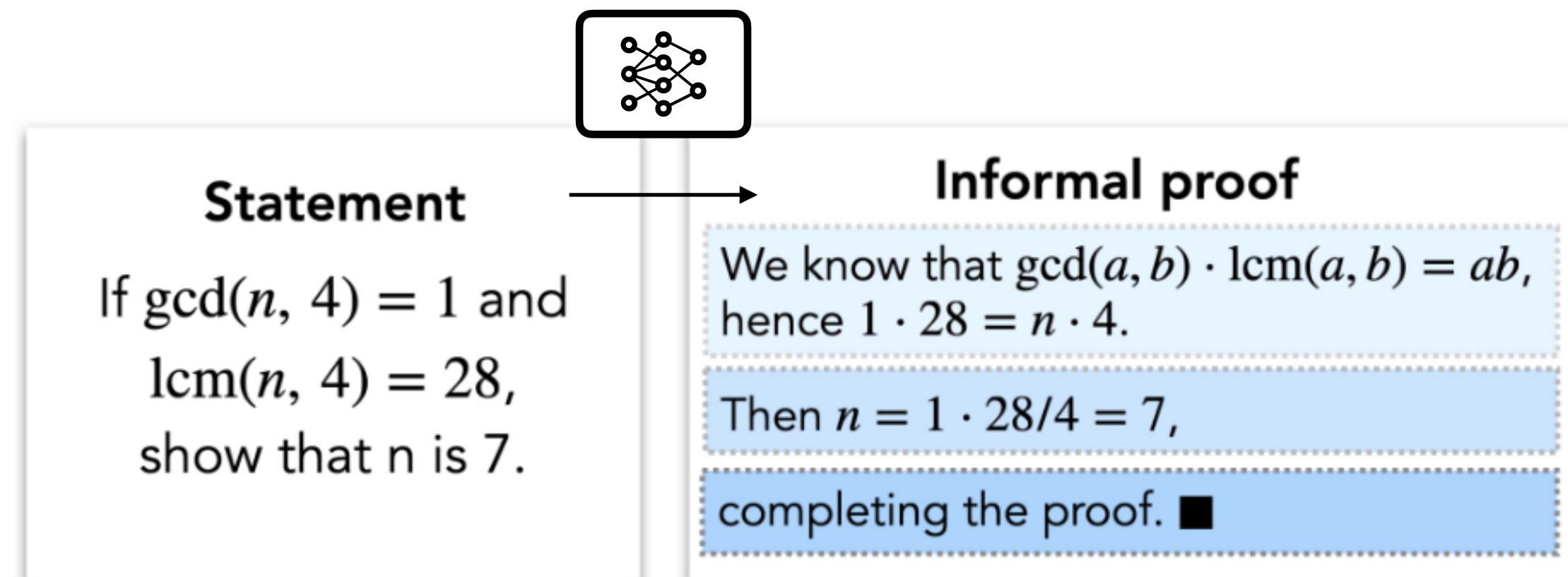
Draft, Sketch, Prove

Given informal theorem x_I
formal theorem x_F

Draft, Sketch, Prove

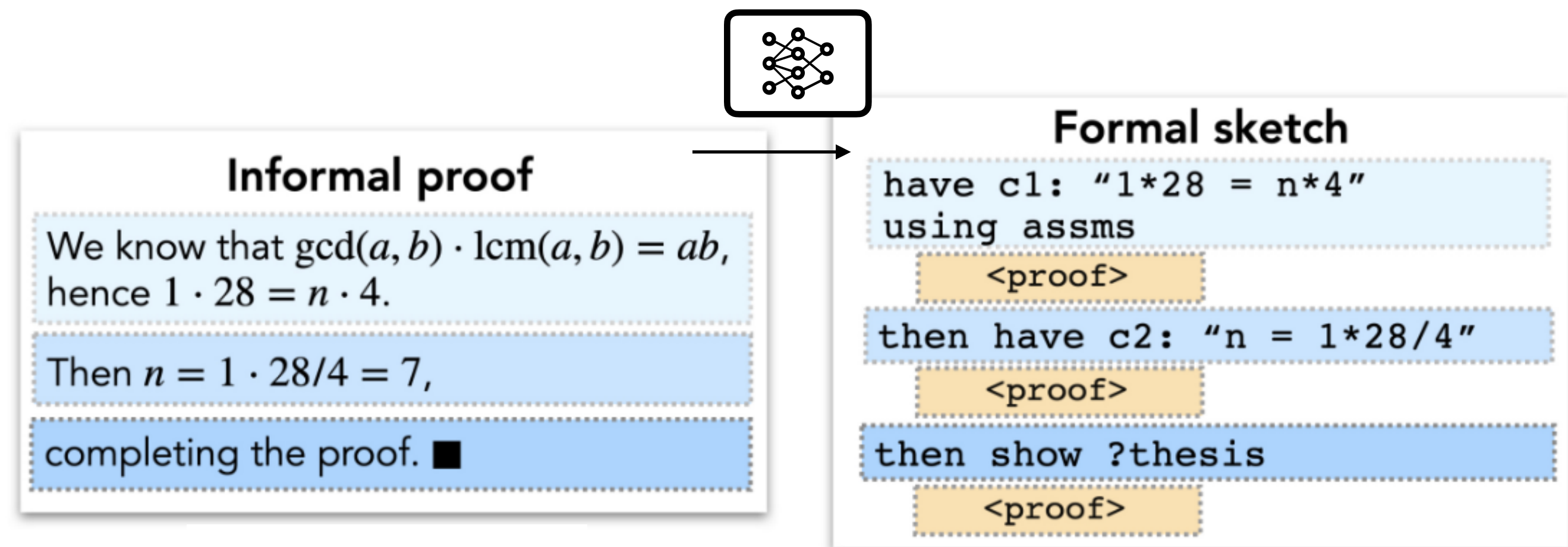
Given informal theorem x_I
formal theorem x_F

1. Draft $y_I \sim p(\cdot | x_I)$



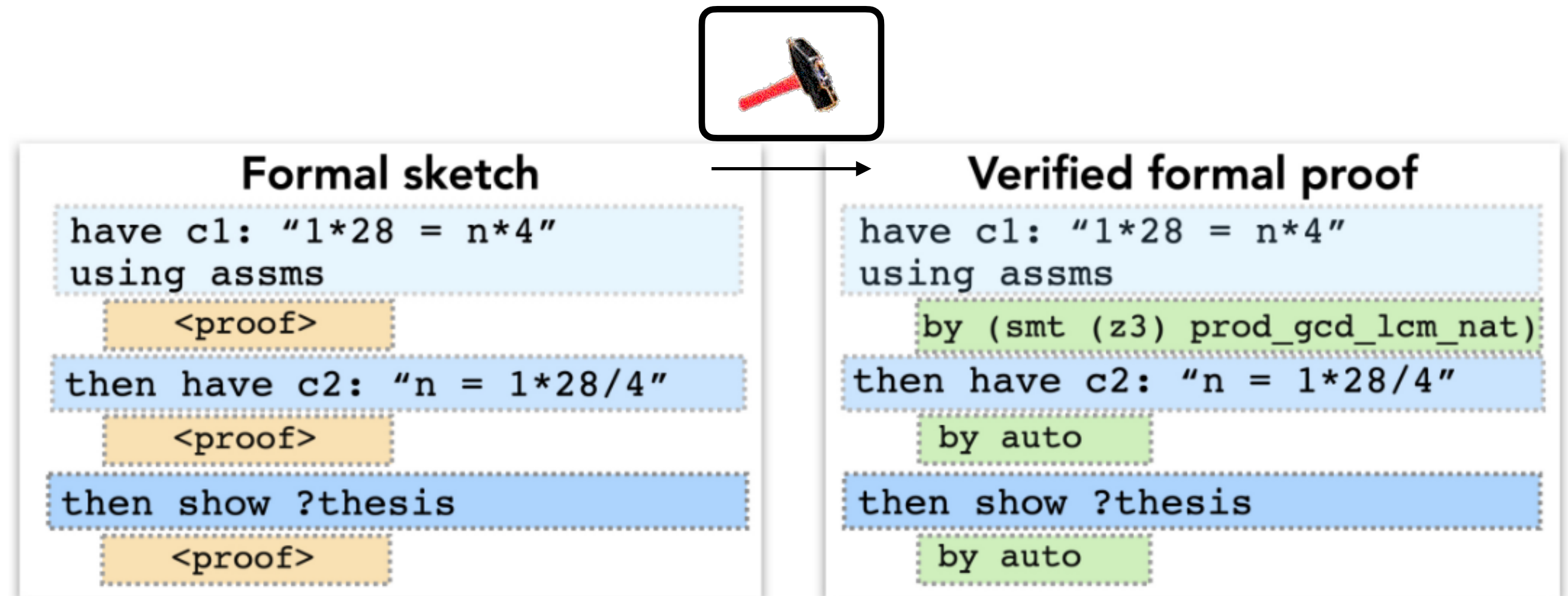
Given informal theorem x_I
formal theorem x_F

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



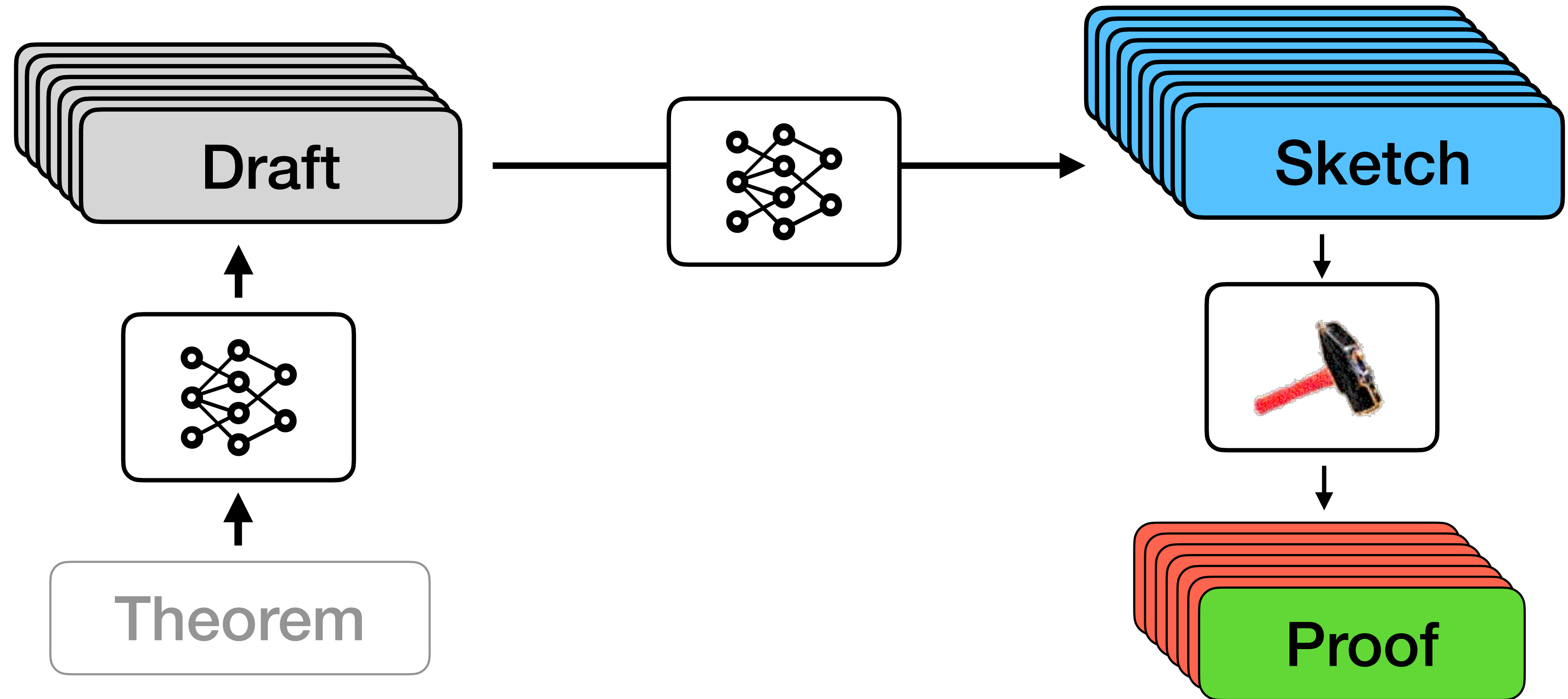
Given informal theorem x_I
formal theorem x_F

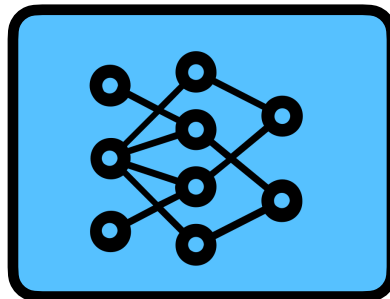
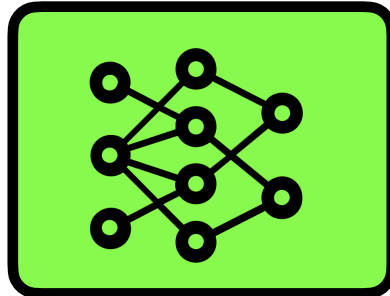
1. Draft $y_I \sim p(\cdot | x_I)$
2. Sketch $z_F \sim p(\cdot | x_F, x_I, y_I)$
3. Prove $y_F = f(x_F, z_F)$



“Classical” prover
Sledgehammer

Proof search





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

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_p1:

fixes n :: nat

shows "gcd (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 $\gcd(7n+1, 1)=1$, we have $\gcd(21n+4, 14n+3)=1$. *)

then have "gcd (7*n + 1) 1 = 1"

using c1

by auto [ATP]

then have "gcd (21*n + 4) (14*n + 3) = 1"

using c1

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

numeral_One numeral_eq_Suc numerals(1) semiring_norm(3)) [ATP]

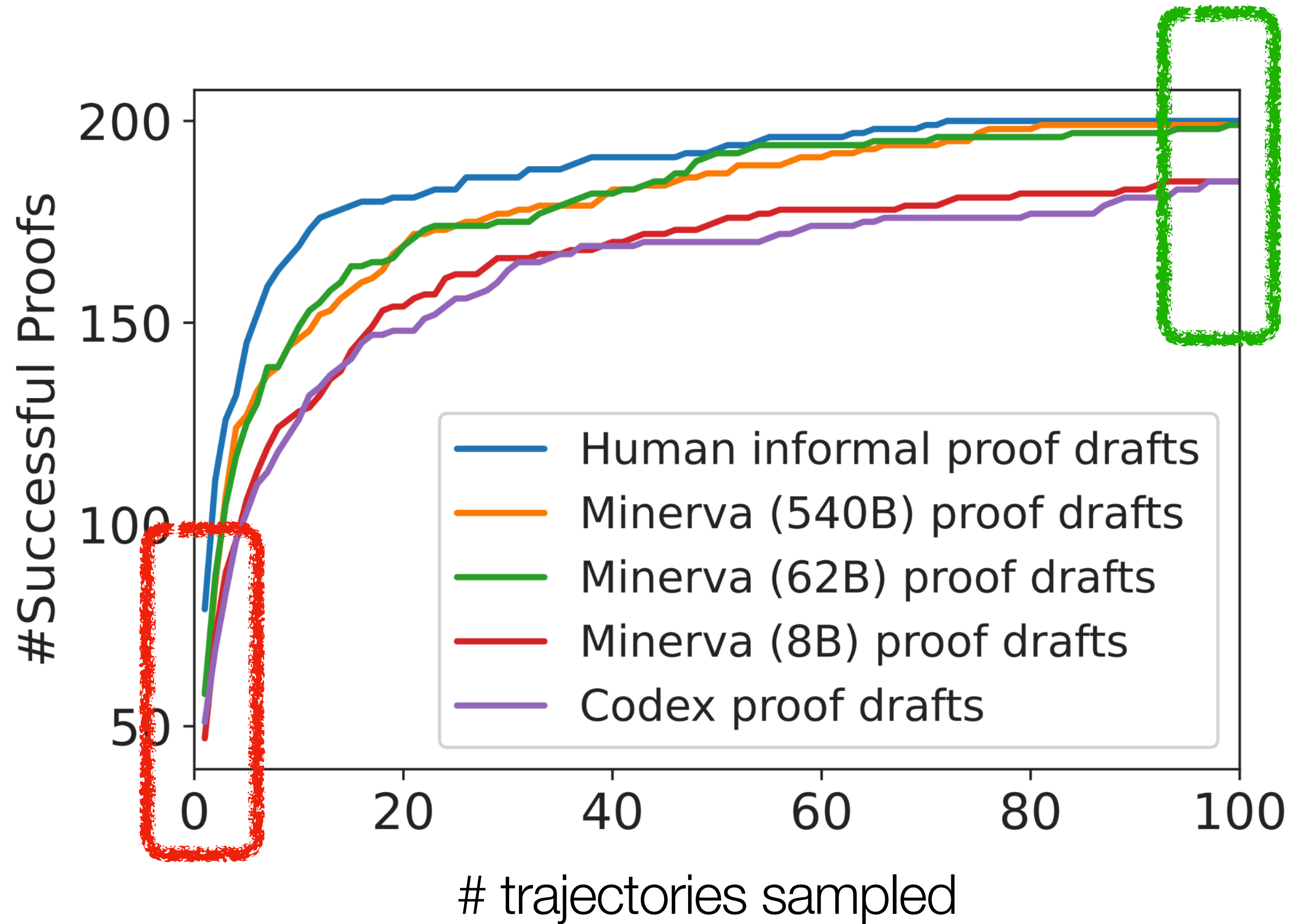
then show ?thesis

using c1

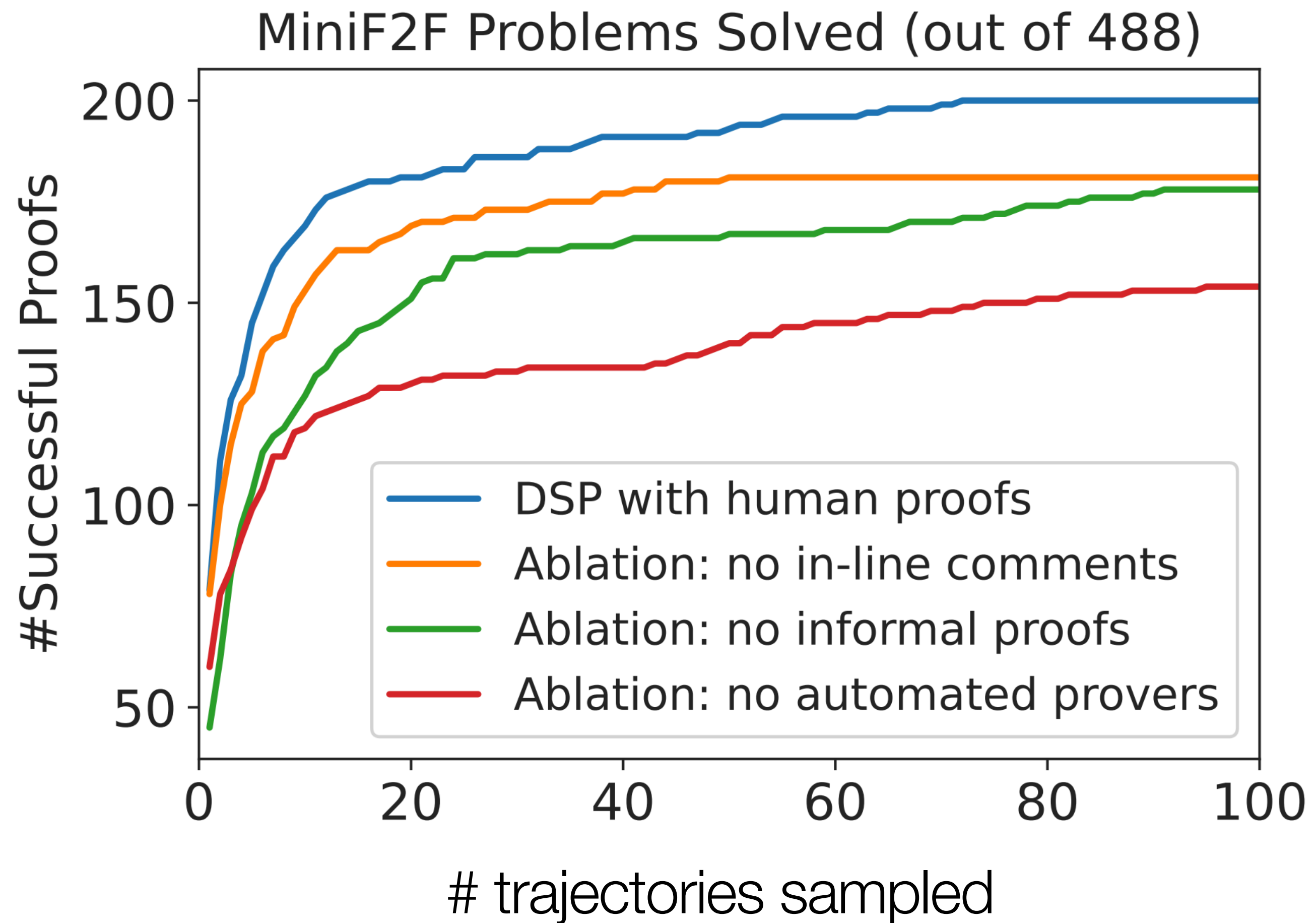
by blast [ATP]

qed

Scaling proof search

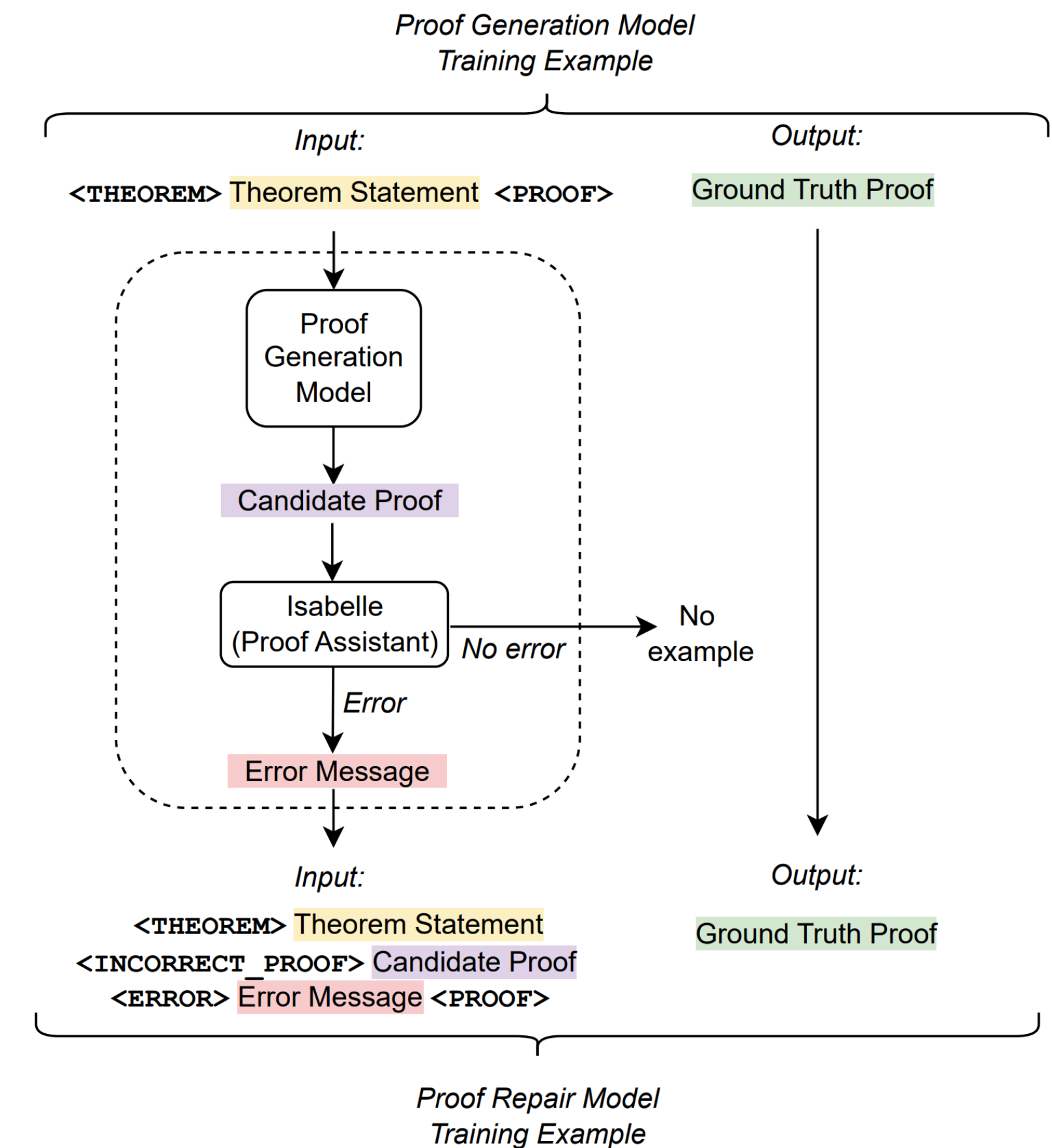


Ablations



Baldur

- Generate $y^{(1)} \sim p_{\theta}(\cdot | x)$
- Repair $y^{(2)} \sim p_{\theta}(\cdot | y^{(1)}, x, \text{errors})$
- Train repair module on generator's outputs



Baldur

Model	16 samples	64 samples
Baldur 8b generate	34.8%	40.7%
Baldur 8b generate + repair	36.3%*	—
Baldur 8b w/ context	40.9%	47.5%
Baldur 62b w/ context	42.2%	47.9%
Baldur 8b w/ context \cup Thor	—	65.7%

Table 4: Proof rate of different models.

***The repair approach uses half the number of samples, and then one repair attempt for each sample.**

Summary

- Flexible language models
- Chain together in a cascade

(Some) open challenges

- Data
- Context
- Efficiency

Data scarcity



**General
domain**
> 6T tokens

Lean
•
~200M tokens

Data scarcity

- Option 1: transfer

Expert
domain

~50B tokens

Lean

~200M tokens

**General
domain**

> 6T tokens

- Ongoing project: [EleutherAI/math-1m](#)

Data scarcity

- Option 2: synthesize

General domain
> 6T tokens

Expert domain

~50B tokens

Transfer

Synthetic

??

Lean

~200M tokens

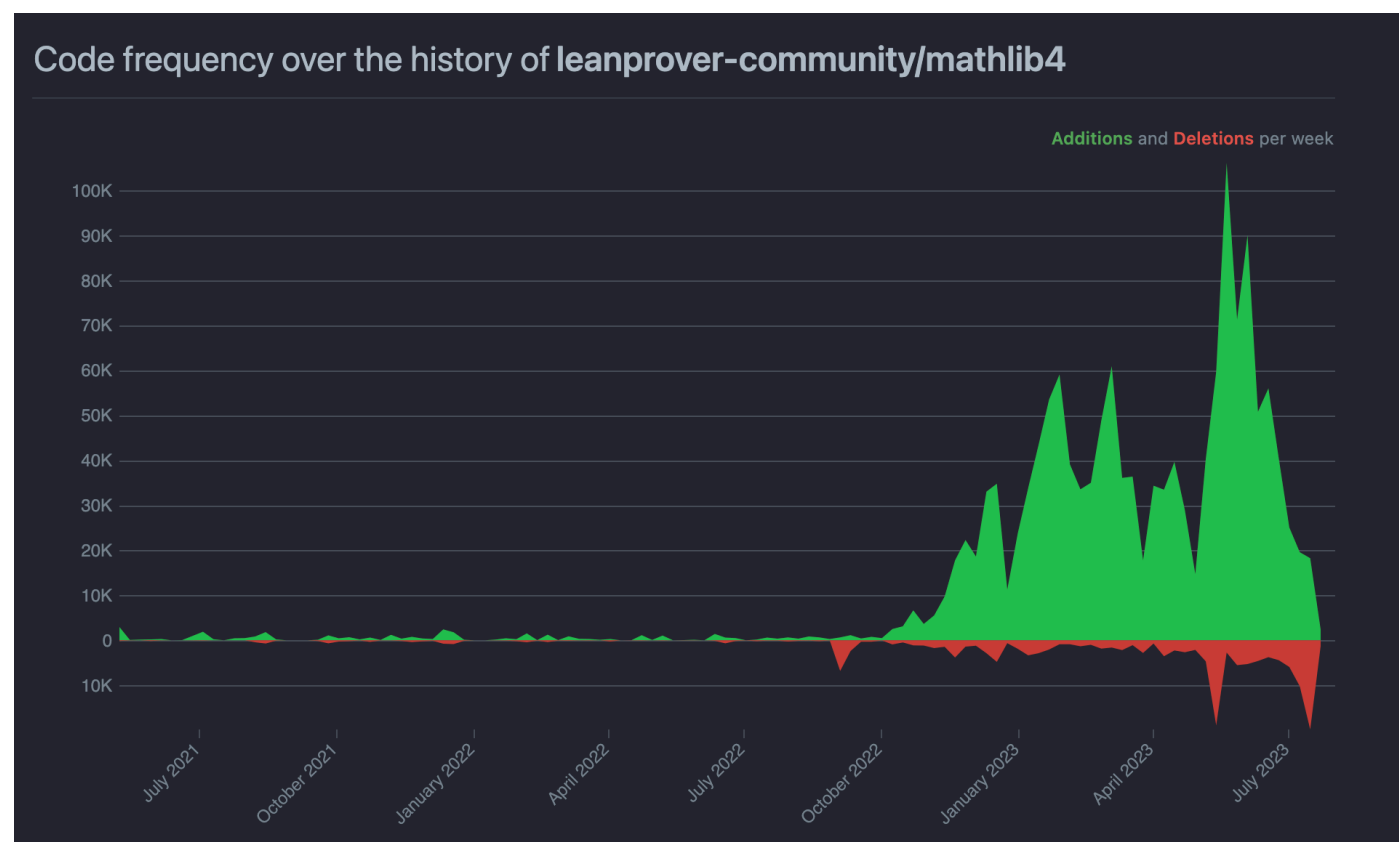
- E.g. see Expert Iteration [Polu et al ICLR 2023], Autoformalization [Wu et al Neurips 2022]

Context

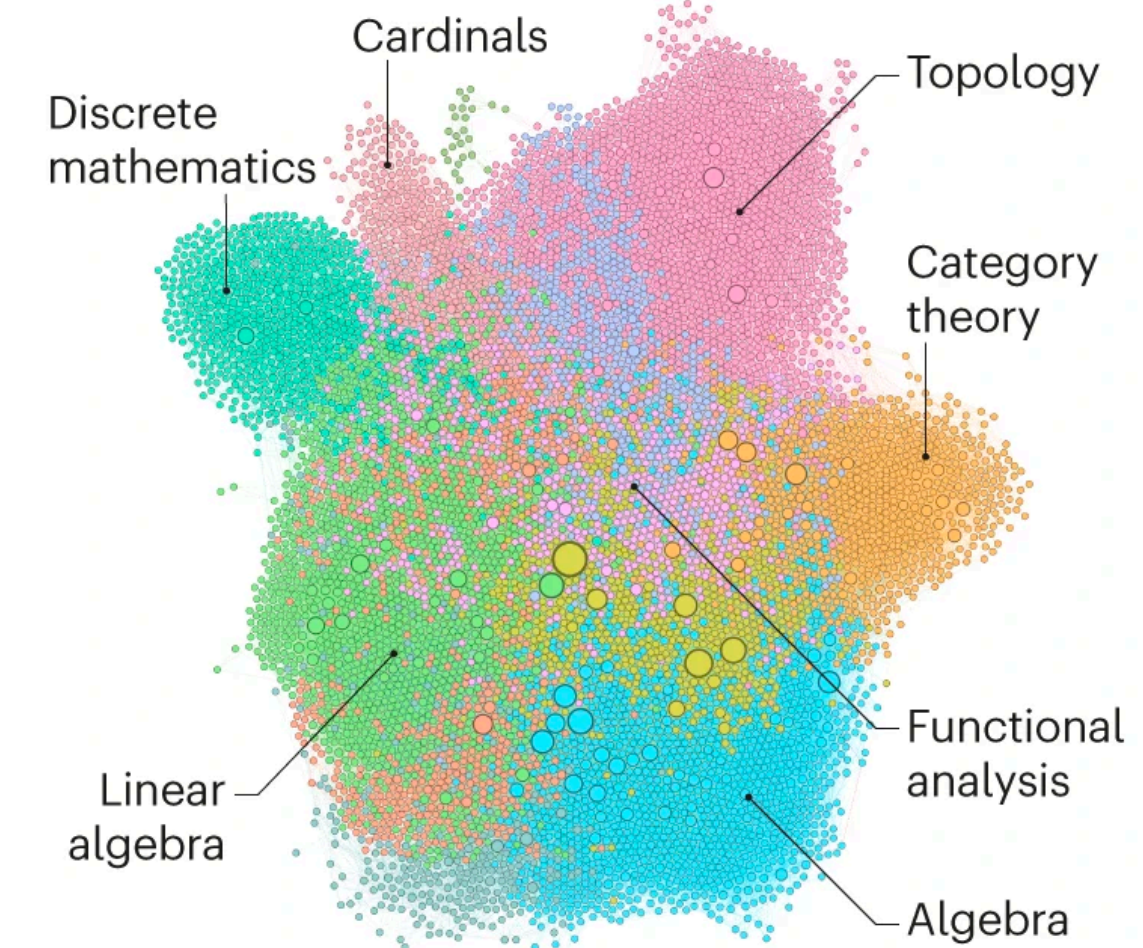
$$p_{\theta}(y_t | x_t)$$

Changing code

Intermediate definitions/lemmas

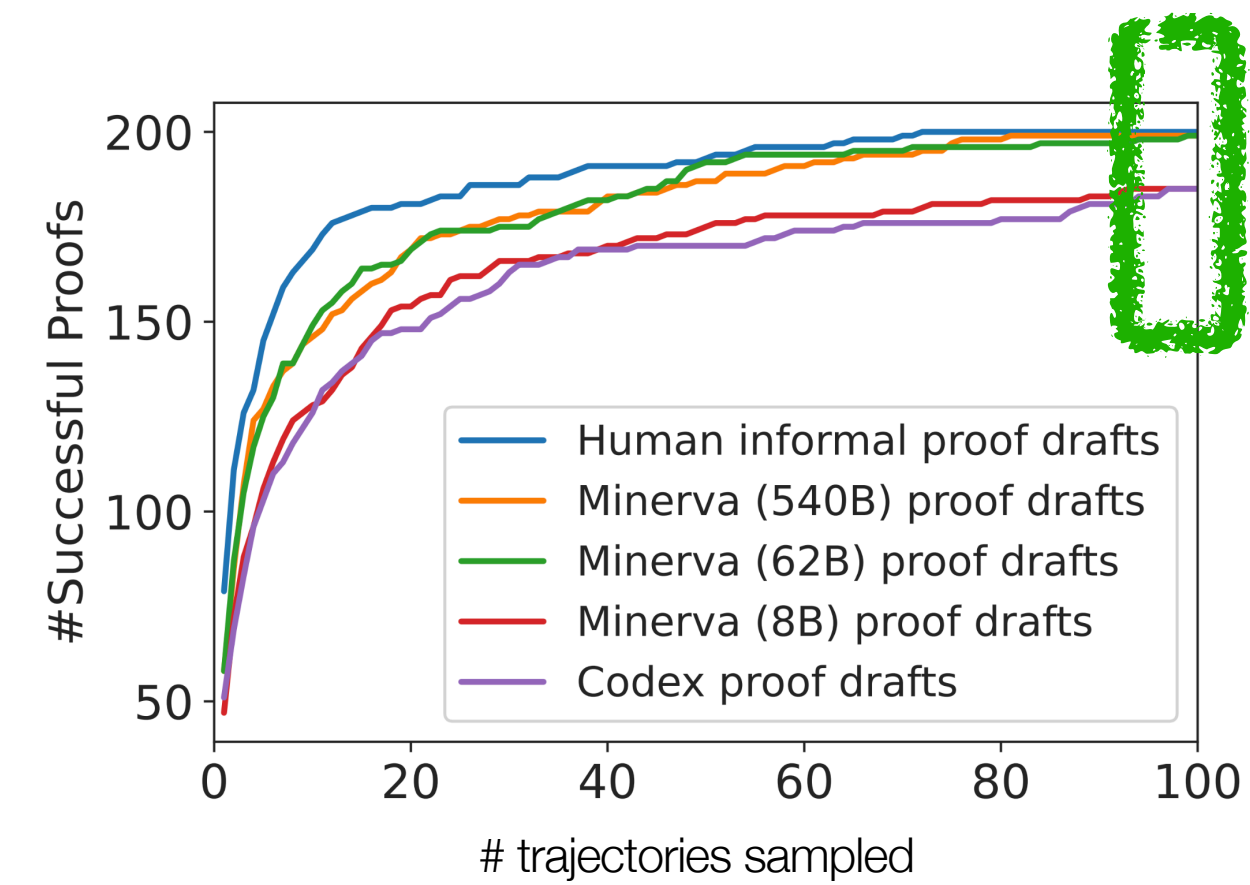
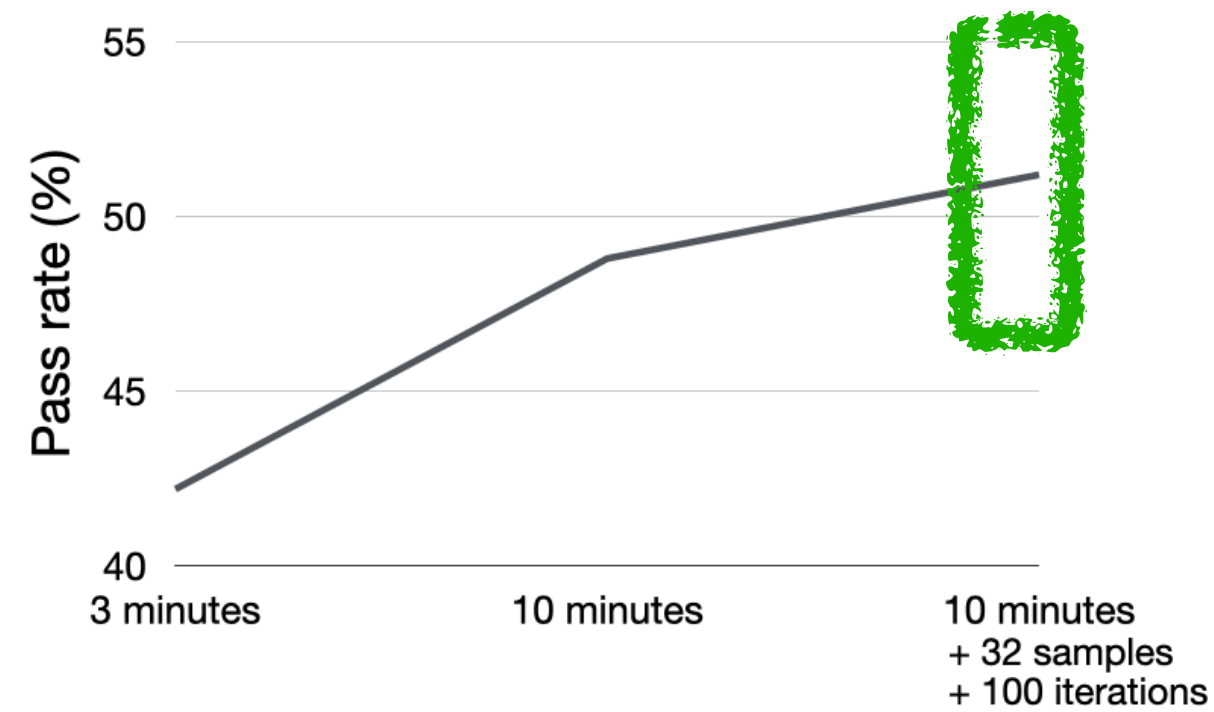


```
Code Blame 2755 Lines (2372 Loc) · 149 KB  
2724  
2725 theorem norm_iteratedFderiv_clm_apply (f : E → F →L[k] G) (g : E → F) (N : ℕ) (n : ℕ)  
2726 (hf : ContDiff k N f) (hg : ContDiff k N g) (x : E) (hn : n ≤ N) :  
2727   ||iteratedFderiv k n (fun y : E => (f y) (g y)) x|| ≤ ∑ i in Finset.range (n + 1),  
2728     +(n.choose i) * ||iteratedFderiv k i f x|| * ||iteratedFderiv k (n - i) g x|| := by  
2729   simp only [- iteratedFderivWithin_univ]  
2730   exact norm_iteratedFderivWithin_clm_apply hf.contDiffOn hg.contDiffOn uniqueDiffOn_univ  
2731   (Set.mem_univ x) hn  
2732 #align norm_iterated_fderiv_clm_apply norm_iteratedFderiv_clm_apply  
2733  
2734 theorem norm_iteratedFderivWithin_clm_apply_const (f : E → F →L[k] G) (c : F) (s : Set E) (x : E)  
2735 (N : ℕ) (n : ℕ) (hf : ContDiffOn k N f s) (hs : UniqueDiffOn k s) (hx : x ∈ s) (hn : n ≤ N) :  
2736   ||iteratedFderivWithin k n (fun y : E => (f y) c) s x|| ≤  
2737     ||c|| * ||iteratedFderivWithin k n f s x|| := by  
2738   let g : (F →L[k] G) →L[k] G := ContinuousLinearMap.apply k G c  
2739   have h := g.norm_compContinuousMultilinearMap_le (iteratedFderivWithin k n f s x)  
2740   rw [- g.iteratedFderivWithin_comp_left hf hs hx hn] at h  
2741   refine' h.trans (mul_le_mul_of_nonneg_right _ (norm_nonneg _))  
2742   refine' g.op_norm_le_bound (norm_nonneg _) fun f => _  
2743   rw [ContinuousLinearMap.apply_apply, mul_comm]  
2744   exact f.le_op_norm c  
2745 #align norm_iterated_fderiv_within_clm_apply_const norm_iteratedFderivWithin_clm_apply_const  
2746  
2747 theorem norm_iteratedFderiv_clm_apply_const (f : E → F →L[k] G) (c : F) (x : E) (N : ℕ) (n : ℕ)  
2748 (hf : ContDiff k N f) (hn : n ≤ N) :  
2749   ||iteratedFderiv k n (fun y : E => (f y) c) x|| ≤ ||c|| * ||iteratedFderiv k n f x|| := by  
2750   simp only [- iteratedFderivWithin_univ]  
2751   exact norm_iteratedFderivWithin_clm_apply_const hf.contDiffOn uniqueDiffOn_univ  
2752   (Set.mem_univ x) hn  
2753 #align norm_iterated_fderiv_clm_apply_const norm_iteratedFderiv_clm_apply_const  
2754
```



Efficiency

- Real users:
 - run on own device (e.g. laptop)
- Large inference costs



Thank you

Tutorial code: <https://github.com/wellecks/ntptutorial>

Topic	Notebook
0. Intro	notebook
1. Data	notebook
2. Learning	notebook
3. Proof Search	notebook
4. Evaluation	notebook
5. <code>llmsuggest</code>	notebook

Topic	Notebook
1. Language model cascades	notebook
2. Draft, Sketch, Prove	notebook