

# Proving Theorem Recursively

---

**Haiming Wang**<sup>1\*</sup> **Huajian Xin**<sup>1\*</sup> Zhengying Liu<sup>2†</sup> Wenda Li<sup>3</sup>

Yinya Huang<sup>4</sup> Jianqiao Lu<sup>5</sup> Zhicheng Yang<sup>6</sup> Jing Tang<sup>6,7</sup>

Jian Yin<sup>1†</sup> Zhenguo Li<sup>2</sup> Xiaodan Liang<sup>1,8†</sup>

<sup>1</sup>Sun Yat-sen University   <sup>2</sup>Huawei Noah's Ark Lab   <sup>3</sup>University of Edinburgh

<sup>4</sup>CityU   <sup>5</sup>HKU   <sup>6</sup>HKUST (Guangzhou)   <sup>7</sup>HKUST   <sup>8</sup>Pengcheng Laboratory

# Neural theorem proving

- Theorems/proofs represented formally as **computer codes** !

Theorem

Problem statement

$\sqrt{2}$  is irrational.

Proof

Assuming  $\sqrt{2} \in \mathbb{Q}$ , we have  $\sqrt{2}=a/b$ , and  $a, b$  are coprime. Then we have  $2 = a^2/b^2$  and  $2 \times b^2 = a^2$ . Thus, we know  $a$  is even,  $a = 2c$ . Substituting  $a$  into the previous equation, we have  $b^2 = (2 * c)^2$ . Thus, we know  $b$  is also even, and  $a, b$  are not coprime. This contradicts the original assumption. ■

lemma "sqrt 2  $\notin$  Q"

```

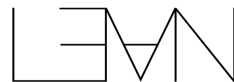
proof
  assume "sqrt 2  $\in$  Q"
  then obtain a b :: int where "sqrt 2 = a/b"
    "coprime a b" "b  $\neq$  0" sledgehammer
  then have c: "2 = a^2 / b^2"
    sledgehammer
  then have "b^2  $\neq$  0" sledgehammer
  then have *: "2*b^2 = a^2"
    sledgehammer
  then have "even a"
    sledgehammer
  then obtain c :: int where "a=2*c"
    sledgehammer
  with * have "b^2 = 2*c^2"
    sledgehammer
  then have "even b"
    sledgehammer
  with (coprime a b) (even a) (even b)
    show False sledgehammer
qed
    
```

Proof

Verification



Formal systems:

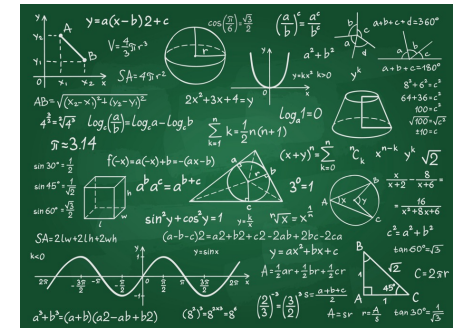


Coq



Mizar

Sources of Theorems



Mathematics



Software


# Previous approaches

---

LM + Search (GPT-f, Thor, DT-Solver):


- **Language model** suggests **action** given **current state**.
- **Formal system** executes action and updates state.
- **Search algorithm** finds correct action path.

lemma "sqrt 2  $\notin$   $\mathbb{Q}$ "

 **goals:** 1. sqrt 2  $\notin$   $\mathbb{Q}$


 **proof**


 **goals:** 1. sqrt 2  $\in$   $\mathbb{Q} \implies$  False

 **assume** "sqrt 2  $\in$   $\mathbb{Q}$ "

 **premise:** sqrt 2  $\in$   $\mathbb{Q}$

**goals:** 1. sqrt 2  $\in$   $\mathbb{Q} \implies$  False

 **then obtain** a b::int **where** "sqrt 2 = a/b"  
"coprime a b" "b  $\neq$  0" **sledgehammer**


 **premise:** sqrt 2 = real of int a / real of int b  
coprime a b

b  $\neq$  0

**goals:** 1. sqrt 2  $\in$   $\mathbb{Q} \implies$  False

 **then have** c: "2 = a<sup>2</sup> / b<sup>2</sup>"  
**sledgehammer**

 ...

 ...

# Motivation

## Challenges:

- Step-by-step methods fail to find long proofs.
- Search space grows exponentially, leading to getting lost.
- High need for value functions to guide the search.

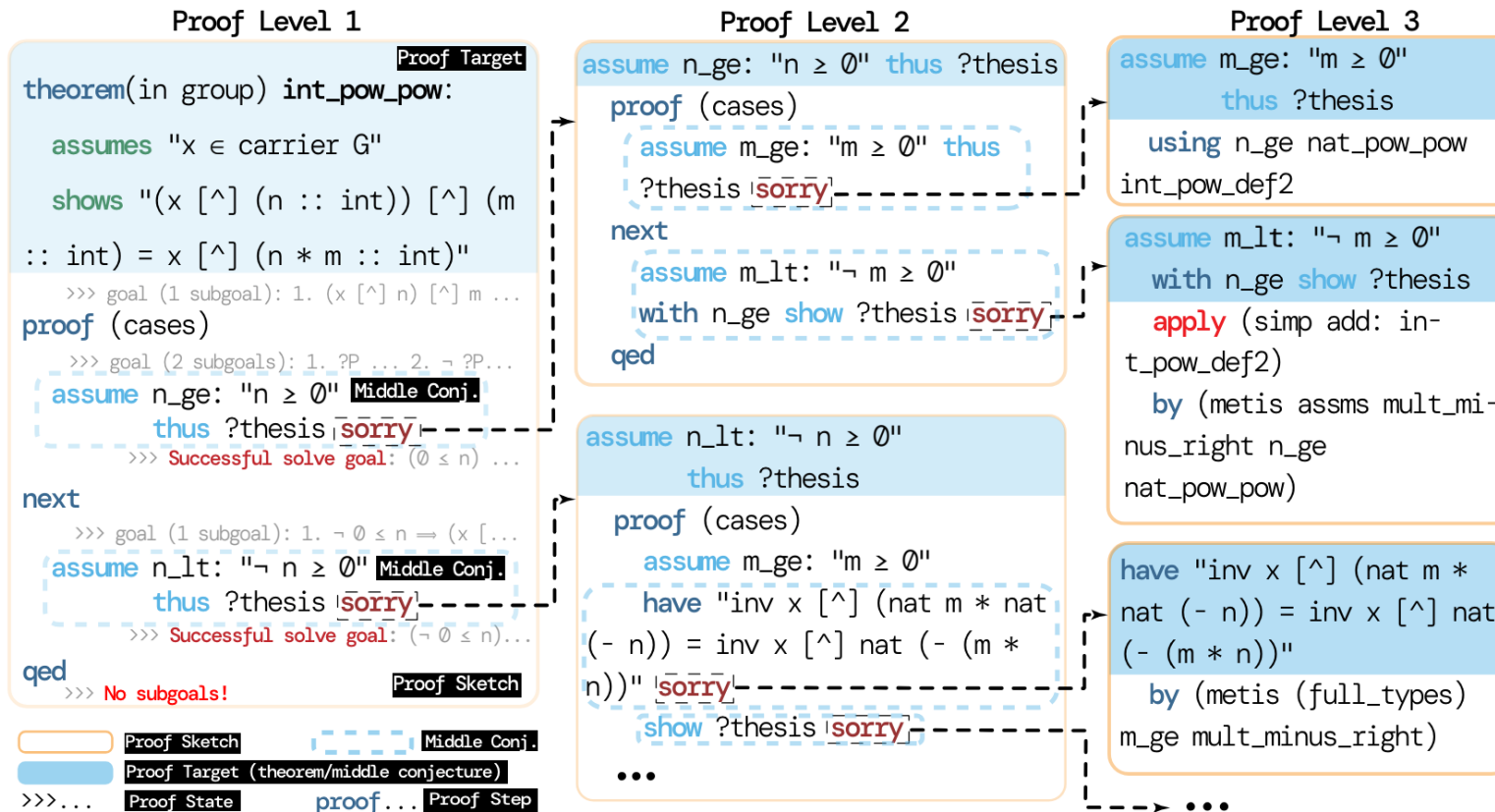
## Solution:

- Think like humans:

Plan → Verify → Plan → Verify → ...

```
theorem(in group) int_pow_pow:
  assumes "x ∈ carrier G"
  shows "(x [^] (n :: int)) [^] (m ::
int) = x [^] (n * m :: int)"
  >>> goal (1 subgoal): 1. (x [^] n) [^] m...
proof (cases)
  >>> goal (2 subgoals): 1. ?P ⇒ (x... 2...
  assume n_ge: "n ≥ 0" thus ?thesis
  >>> using this: 0 ≤ n goal (1 subgoal):...
  proof (cases)
    >>> goal (2 subgoals): 1. 0 ≤ n ⇒ ...
    assume m_ge: "m ≥ 0" thus ?thesis
    >>> using this: 0 ≤ m goal (1 subgoal)...
    using n_ge nat_pow_pow in
t_pow_def2
    >>> Successful solve goal (m ≥ 0) ...
  ...
qed
>>> No subgoals! Complete proof
```

# Prove theorem recursively



(b) Recursive Proof

## First recursive proving framework!

- Search proof sketch (plan for the proof at each stage)
- Verify the proof sketch by formal system!
- Proceed to deeper sketches after verified to be correct

# Experiments: main results

## Thor (Cambridge, NeurIPS 2022)

- LM is train on single step state action pair and finds proof with best first search algorithm

## Thor + expert iteration (Google + Cambridge, NeurIPS 2022)

- Extend Thor with extensive proof data generated by Codex LLM.

## Thor + Magnushammer (Cambridge, ICLR 2023)

- Extend Thor with neural enhanced sledgehammer.

## Our approaches (proofGPT-1.3B\*):

### GPT-f Baseline

- ablation setting which use step-by-step approach to prove theorem.

### POETRY

- Our recursive proving method

Table 1: **Comparing with baseline.** The table displays the pass@1 success rates of the baselines and POETRY, The highest success rates for each set are highlighted in bold.

Success rate	miniF2F-valid	miniF2F-test	PISA	single-level	multi-level
Thor w/o sledgehammer	25.0%	24.2%	39.0%	-	-
GPT-f Baseline	39.3%	37.3%	48.9%	<b>65.5%</b>	11.1%
– with sampling decoding	30.3%	31.5%	43.2%	57.8%	9.8%
POETRY	<b>42.2%</b>	<b>42.2%</b>	<b>49.6%</b>	65.4%	<b>13.6%</b>

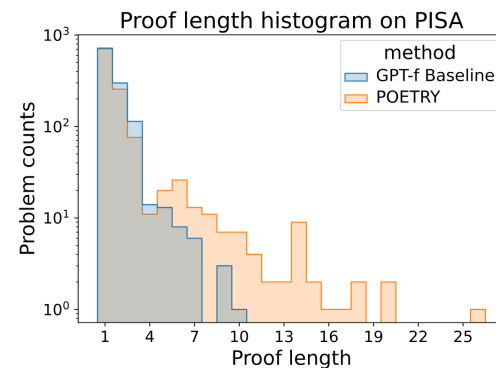
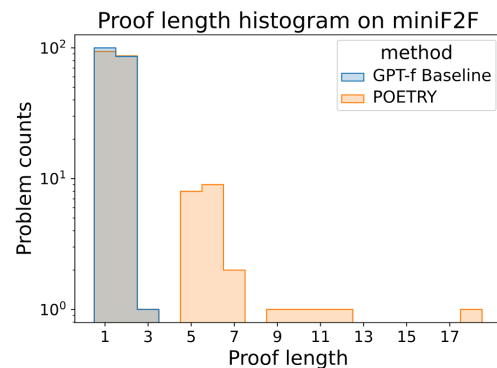
Table 2: **Comparing with state-of-the-art search-based methods on the miniF2F dataset.** The table displays the pass@1 success rates of previous works and POETRY, The highest success rates for each set are highlighted in bold.

Success rate	environment	miniF2F-valid	miniF2F-test
<i>Baselines</i>			
PACT [Han et al., 2022]	Lean	23.9%	24.6%
Leandojo [Yang et al., 2023]	Lean	-	26.5%
FMSCL [Polu et al., 2022]	Lean	33.6%	29.6%
COPRA [Thakur et al., 2024]	Lean	-	30.7%
Thor [Jiang et al., 2022a]	Isabelle	28.3%	29.9%
Thor + expert iteration [Wu et al., 2022]	Isabelle	37.3%	35.2%
Thor + Magnushammer [Mikuła et al., 2023]	Isabelle	36.9%	37.3%
<i>Ours</i>			
POETRY	Isabelle	<b>42.2%</b>	<b>42.2%</b>

# Experiments: analysis

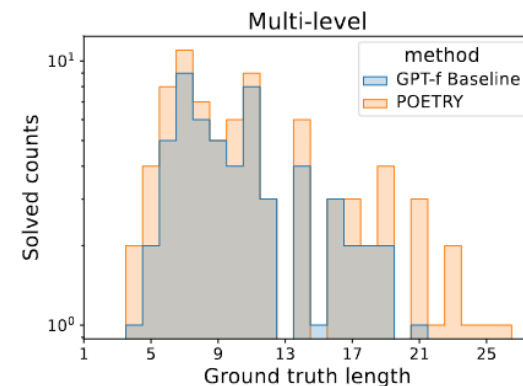
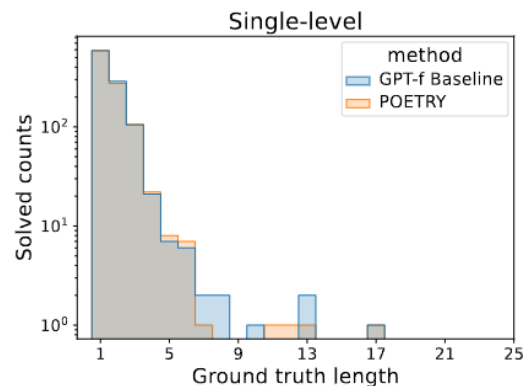
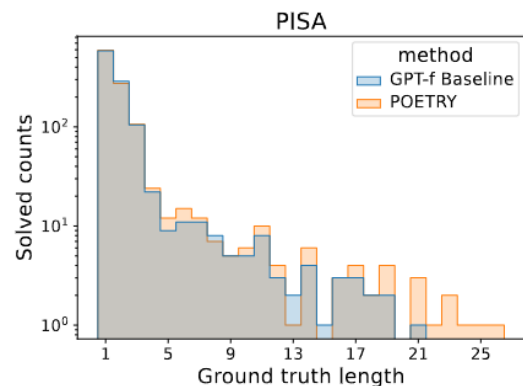
## POETRY capable of finding long proof

- GPT-f Baseline found maximum length of 3 (10) steps in miniF2F (PISA)
- POETRY found maximum length of 18 (26) steps in miniF2F (PISA)



## POETRY capable of finding harder proof

- POETRY and GPT-f Baseline have similar performance in single-level problem
- POETRY excels at solving problems requires structural reasoning.



# Experiments: case study

```
lemma(in UP_cring) n_mult_closed:
  assumes "f ∈ carrier P"
  shows "n_mult f ∈ carrier P"

proof(rule UP_car_memI[of "deg R f"])
  show "∧n. deg R f < n ⇒ n_mult f n = 0"
  unfolding n_mult_def
  using assms
  unfolding P_def
  by (simp add: UP_car_memE(2)) Proof level 2
  show "∧n. n_mult f n ∈ carrier R"
  using assms
  unfolding n_mult_def
  by (simp add: assms cfs_closed) Proof level 2
qed Proof level 1
```

(a)

```
lemma(in UP_cring) n_mult_closed:
  assumes "f ∈ carrier P"
  shows "n_mult f ∈ carrier P"
proof(rule UP_car_memI[of "deg R f"])
  fix n
  assume A: "deg R f < n"
  show "n_mult f n = 0"
    unfolding n_mult_def
    proof -
      X Timeout after 600 seconds
    end
  X Never explored
end
```

(b)

## Case comparison between POETRY and GPT-f Baseline.

- Recursive proof found by POETRY in 71.2 seconds, the proof contains two proof levels.
- Failure-proof paths found by the GPT-f Baseline. GPT-f Baseline failed to find proof due to timeout after 600 seconds. We select two different failure proof paths found by GPT-f Baseline.



**Thanks**

---