

Agentic approach for neural theorem proving in Rocq

CANUM 2026

Guillaume Baudart, Emilio Jesús Gallego Arias, Marc Lelarge, Théo Stoskopf, et Jules Viennot

The outbreak of reasoning models



P. Erdős

ce Problem

num

+2 "

ionzero

cinmer,⁴
OpenAI
iversity ⁵Harvard University

ften presumed to vanish,
trations existing in Klein
l-form expression for the
ction of their momenta.
Weinberg's soft theorem.

³ The amplitude is divided into
ls are regions where sums of vari-
alf-collinear momenta are orthogo-
w. The (stripped) amplitudes are
egers in each chamber. These val-
rom the perturbative Berends–Giele
is equivalent to Feynman diagrams.
pecial kinematic region correspond-
ay into $n-1$ plus, we give a simple
this special region, the stripped
me values of $+1, -1, \text{ or } 0$.

³⁹) for the amplitude in this region
by GPT-5.2 Pro and then proved by
AI model. The solution was checked
rends–Giele recursion and was more-
vially obey the soft theorem, cyclic-
U(1) decoupling identities—none of
m direct inspection.
of these single-minus amplitudes in
mains to be understood. We note
ession is a dramatic simplification
n-diagram expression, it is entirely
mpler expression may be obtained
of analytic continuation, variables or
he single-minus decay channel. We
e more interesting insights to come
y and hope that this paper is a step
ore complete understanding of the
ltering amplitudes.

itudes also arise in self-dual Yang-
Mills theory (SDYM) [18], a restricted sector of Yang-
Mills, and potentially resolve a puzzle therein. In gen-
eral, the tree amplitudes of the Feynman expansion are
thought to be equivalent to the fully nonlinear classi-
cal theory. However, on the one hand the classical so-
lution space of SDYM is extremely nontrivial [19–21],
while the tree diagrams were previously supposed to yield
trivial two-point and three-point expressions. The latter
seem insufficient to reproduce the former. Potentially,
the single-minus tree amplitudes in SDYM found here
resolve this tension.

³ The half-collinear condition can be viewed as restricting the in-
going and outgoing momenta to a one-dimensional null circle on
the celestial torus at the boundary of Klein space.

EARLY EVIDENCE OF VIBE-PROVING WITH CONSUMER LLMs: A CASE STUDY ON SPECTRAL REGION CHARACTERIZATION WITH CHATGPT-5.2 (THINKING)

Brecht Verbeke^{1,2,3}
0000-0002-7506-3298

Brando Vagenende¹
0000-0002-0573-8093

Marie-Anne Guerry¹
0000-0001-5842-8905

Andres Algaba^{1,2}
0000-0002-0532-3066

Vincent Ginis^{1,2,3}
0000-0003-0063-9608

¹Data Analytics Lab, Vrije Universiteit Brussel, Pleinlaan 5, 1050 Brussel, Belgium
²imec-SMIT, Vrije Universiteit Brussel, Pleinlaan 9, 1050 Brussels, Belgium
³Department of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

February 24, 2026

ABSTRACT

Large Language Models (LLMs) are increasingly used as scientific copilots, but evidence on their
le in research-level mathematics remains limited, especially for workflows accessible to individual
researchers. We present early evidence for vibe-proving with a consumer subscription LLM through
an auditable case study that resolves Conjecture 20 of Ran and Teng (2024) on the exact nonreal
eutral region of a 4-cycle row-stochastic nonnegative matrix family. We analyze seven shareable
atGPT-5.2 (Thinking) threads and four versioned proof drafts, documenting an iterative pipeline
generate, referee, and repair. The model is most useful for high-level proof search, while human
perts remain essential for correctness-critical closure. The final theorem provides necessary and
ficient region conditions and explicit boundary attainment constructions. Beyond the mathematical
ult, we contribute a process-level characterization of where LLM assistance materially helps
d where verification bottlenecks persist, with implications for evaluation of AI-assisted research
rkflows and for designing human-in-the-loop theorem proving systems.

AI-assisted mathematics · Karpelevich region · large language models · stochastic matrices

duction

where programmers steer Large Language Models (LLMs) to generate software through high-level natural-
tent rather than line-by-line specification, has rapidly transformed code generation [1, 2, 3]. Underpinning
the sustained capability growth of frontier LLMs across successive systems and evaluations [4, 5, 6, 7],
odels grow more capable, they are increasingly deployed as scientific collaborators: generating candidate
as in controlled studies [8, 9], acting as agentic research assistants that plan, search, and iterate over literature
nents [10, 11, 12, 13, 14, 15, 16], and accelerating evidence synthesis and literature screening [17, 18, 19].
le, LLM tooling is associated with measurable shifts in scientific production and communication [20],
ins unclear how reliably these systems support auditable, research-level mathematical workflows under
cess conditions [21, 22].

pending author: brecht.verbeke@vub.be

1

arX

helicity particles and $n-2$ plus-helicity gluons, which
for generic (complexified) kinematics at tree level is the
maximally allowed number [4, 11–14]. This gives them a
privileged role in the theory, enabling their use as efficient
building blocks for the full Yang–Mills theory.

In general, $n-2$ is actually *not* the maximally al-
lowed number of plus gluons. In this paper, we show
that $n-1$ plus (or “single-minus”) amplitudes are in
fact allowed even at tree level² with restricted “half-

¹ The aforementioned agreement between theory and experiment
required over a half century of analytic and numerical work.

² Witten [5] notes that single-minus tree amplitudes are supported
at a point in twistor space; see also [6]. All-plus and single-minus
amplitudes are generically corrected at loop level [15, 16].

³ The half-collinear condition can be viewed as restricting the in-
going and outgoing momenta to a one-dimensional null circle on
the celestial torus at the boundary of Klein space.

p

?

THE MOTIVIC CLASS OF THE SPACE OF GENUS 0 MAPS TO THE FLAG VARIETY

JIM BRYAN, BALÁZS ELEK, FREDDIE MANNERS, GEORGE SALAFATINOS, AND RAVI VAKIL

ABSTRACT. Let \mathbb{F}_{n+1} be the variety of complete flags in \mathbb{A}^{n+1} and let $\Omega_{\beta}^2(\mathbb{F}_{n+1})$ be
the space of based maps $f: \mathbb{P}^1 \rightarrow \mathbb{F}_{n+1}$ in the class $f_*([\mathbb{P}^1]) = \beta$. We show that under
a mild positivity condition on β , the class of $\Omega_{\beta}^2(\mathbb{F}_{n+1})$ in $K_0(\text{Var})$, the Grothendieck
group of varieties, is given by

$$[\Omega_{\beta}^2(\mathbb{F}_{n+1})] = [\text{GL}_n \times \mathbb{A}^n].$$

The proof of this result was obtained in conjunction with Google Gemini and related
tools. We briefly discuss this research interaction, which may be of independent interest.
However, the treatment in this paper is entirely human-authored (aside from excerpts in an
appendix which are clearly marked as such).

1. INTRODUCTION

Let $\mathbb{F}_{n+1} = \text{GL}_{n+1}/B$ be the complete flag variety which parameterizes flags of
quotients:

$$\mathbb{k}^{n+1} \rightarrow E_n \rightarrow \dots \rightarrow E_1$$

where $\dim E_k = k$. There is a universal sequence of vector bundle quotients on \mathbb{F}_{n+1} :

$$\mathcal{O}_{\mathbb{F}_{n+1}}^{\oplus n+1} = \mathcal{E}_{n+1} \rightarrow \mathcal{E}_n \rightarrow \dots \rightarrow \mathcal{E}_1.$$

Let $\Omega_{\beta}^2(\mathbb{F}_{n+1})$ be the space of genus zero based maps in the class β :

$$\Omega_{\beta}^2(\mathbb{F}_{n+1}) = \{f: \mathbb{P}^1 \rightarrow \mathbb{F}_{n+1} : f_*([\mathbb{P}^1]) = \beta, f([1:0]) = [\text{Id}]\}$$

where $\beta \in A_1(\mathbb{F}_{n+1})$ and $[\text{Id}] \in \text{GL}_{n+1}/B$ is the base point in \mathbb{F}_{n+1} .

There is an isomorphism $A_1(\mathbb{F}_{n+1}) \cong \mathbb{Z}^n$ such that the n -tuple associated to β is given
by (d_1, \dots, d_n) where

$$d_i = \deg(f^* \mathcal{E}_i).$$

We say that $\beta = (d_1, \dots, d_n)$ is *strictly monotonic* if $0 < d_n < d_{n-1} < \dots < d_1$. The
main result of this paper is the following:

Theorem 1.1. *Suppose $\beta = (d_1, \dots, d_n)$ is strictly monotonic. Then we have the follow-
ing equality in $K_0(\text{Var}_{\mathbb{k}})$, the Grothendieck group of varieties over \mathbb{k} :*

$$[\Omega_{d_n, \dots, d_1}^2(\mathbb{F}_{n+1})] = [\text{GL}_n \times \mathbb{A}^{D-n^2}]$$

where $D = \sum_{k=1}^n 2d_k$.

In particular, Theorem 1.1 gives an equality of point counts over finite fields (see Corol-
lary A.1).

Our original interest in the variety $\Omega_{\beta}^2(\mathbb{F}_{n+1})$ is as an algebraic counterpart to the dou-
ble loop space in topology:

$$\Omega_{\beta, \text{top}}^2(\mathbb{F}_{n+1}) = \{f: S^2 \rightarrow \mathbb{F}_{n+1}, f_*([S^2]) = \beta, f([1:0]) = [\text{Id}]\}$$

with $x_0 \in \mathbb{R}^n$. Gradient descent is known to converge
 $\mathcal{O}(1/k)$ and

$$x_k \rightarrow x_{\infty} \in \arg \min$$

[27, 9, 14]. However, the $\mathcal{O}(1/k)$ rate is famously su-
Nesterov’s 1983 seminal paper [24] presented the Nes-
terov method

$$x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k), \\ y_{k+1} = x_{k+1} + \frac{t_k - 1}{t_{k+1}} (x_{k+1} - x_k)$$

with $x_0 = y_0 \in \mathbb{R}^n$, $t_0 = 1$, and $t_{k+1}^2 - t_{k+1} \leq t_k^2$. In
sense of $f(x_k) - \inf f \leq \mathcal{O}(1/k^2)$, an accelerated rate
or $t_k = (k+2)/2$ for $k = 0, 1, \dots$. However, whether

$$x_k \xrightarrow{?} x_{\infty} \in \arg \min$$

has remained an open problem [17, 4, 16, 10].

In this work, we resolve this longstanding open
showing point convergence for NAG in the sense of

$$x_k \rightarrow x_{\infty}, \quad y_k \rightarrow x_{\infty}, \quad x_{\infty}$$

The discovery of the proof was heavily assisted by
language model, and we describe the process through-

^{*}Submitted to the editors DATE.

¹Department of Mathematics, University of California, Los Angeles (uijeongjang@math.ucla.edu,
eryu@math.ucla.edu).

Google DeepMind

2026-5-22

Advancing Mathematics Research with AI-Driven Formal Proof Search

George Tsoukalas¹, Anton Kovsharov¹, Sergey Shirobokov¹, Anja Surina¹, Moritz Firsching¹, Gergely
Bérczi², Francisco J. R. Ruiz¹, Arun Suggala¹, Adam Zsolt Wagner¹, Eric Wieser¹, Lei Yu¹, Aja Huang¹, Miklós
Z. Horváth¹, Andrew Ferraiuolo¹, Henryk Michalewski¹, Codrut Grosu³, Thomas Hubert¹, Matej Balog¹,
Pushmeet Kohli¹ and Swarat Chaudhuri¹

¹Google DeepMind¹, ²Aarhus University, ³Google

Large language models (LLMs) increasingly excel at mathematical reasoning, but their unreliability
limits their utility in mathematics research. A mitigation is using LLMs to generate formal proofs
in languages like Lean. We perform the first large-scale evaluation of this method’s ability to solve
open problems. Our most capable agent autonomously resolved 9 of 353 open Erdős problems at the
per-problem cost of a few hundred dollars, proved 44/492 OEIS conjectures, and is being deployed in
combinatorics, optimization, graph theory, algebraic geometry, and quantum optics research. A basic
agent alternating LLM-based generation with Lean-based verification replicated the Erdős successes but
proved costlier on the hardest problems. These findings demonstrate the power of AI-aided formal proof
search and shed light on the agent designs that enable it.

1. Introduction

Large language models (LLMs) have recently shown remarkable promise in solving complex
mathematics problems [21, 65], but *unreliability* remains a primary barrier to their integra-
tion into mathematics research. Because LLM-generated natural language proofs can contain
subtle logical errors, or “hallucinations,” they require expensive expert review. Mistakes in
unreviewed intermediate steps can cascade through a proof, limiting the complexity of tasks
that can be delegated to AI.

Recent efforts [29, 1] mitigate these issues by using AI to generate proofs in formal
languages like Lean [43], in which a compiler automatically verifies every logical step. So
far, successes of this paradigm have been concentrated in competition mathematics and
the human-aided formalization of natural language arguments [28]. In this paper, we
demonstrate its broader potential through a large-scale evaluation on open research-level
problems.

To this end, we developed a framework, AlphaProofNexus, for LLM-aided proof generation
and used it to build a basic agent in which a set of subagents independently searches for
proofs with feedback from the Lean compiler. We also developed a “full-featured” agent
in which subagents are coordinated using an evolutionary algorithm [46] and can use
AlphaProof [29], a system for olympiad-level Lean theorem-proving based on reinforcement
learning, as a focused proof tool.

Our full-featured agent autonomously solved 9 Erdős problems out of 353 attempted,
including two questions that had been open for 56 years [54, 7, 17], at the inference cost

¹Equal contributions. The first three authors are in random order. Correspondence to pushmeet@google.com
and swarat@google.com.

arXiv:2511.02864v3 [cs.NE] 22 Dec 2025

arXiv:2601.07222v1 [math.AG] 12 Jan 2026

arXiv:2510.23513v2 [ma

MATHEMATICAL EXPLORATIO

BOGDAN GEORGIEV, JAVIER GÓMEZ-SERRANO

ABSTRACT. AlphaEvoive, introduced in [234], is a gener-
al-purpose solver for all types to
attack problems in which a key objective is to construct
quantitative properties, such as obeying a certain inequ-
ity paper, we report on our experiments testing the per-
formance, primarily in the areas of analysis, combinat-
ics provided by AlphaEvoive were not merely numerical
in mathematics, by other tools such as Deep Think, an
able to match or exceed previous results in all cases, as
achieve could likely also have been matched by more trad-
itional tools that can autonomously explore mathematical
[291]. AlphaEvoive (see [224]) represents a step in this
when combined with evolutionary computation and rig-
orous constructions that either match or improve upon the best-known
large scales.

1. INTRO

The landscape of mathematical discovery has been fun-
damental tools that can autonomously explore mathematical
[291]. AlphaEvoive (see [224]) represents a step in this
when combined with evolutionary computation and rig-
orous constructions that either match or improve upon the best-known
large scales.

AlphaEvoive is not a general-purpose solver for all types
to attack problems in which a key objective is to construct
quantitative properties, such as obeying a certain inequ-
ity paper, we report on our experiments testing the per-
formance, primarily in the areas of analysis, combinat-
ics provided by AlphaEvoive were not merely numerical
in mathematics, by other tools such as Deep Think, an
able to match or exceed previous results in all cases, as
achieve could likely also have been matched by more trad-
itional tools that can autonomously explore mathematical
[291]. AlphaEvoive (see [224]) represents a step in this
when combined with evolutionary computation and rig-
orous constructions that either match or improve upon the best-known
large scales.

We have also seen that in many cases, besides the scaling, in order to get AlphaEvoive to output co-
results to the literature and in contrast to traditional ways of doing mathematics, very little overhead


The authors are listed in alphabetical order.

1

https://

Are humans still useful?

YES

- To **guide** Large Language Models (LLMs)
 - by hinting
 - by discussing
 - by building environments
- To **review** the generated proofs  Can it be automated?



A trustworthy, industrial-strength interactive theorem prover and dependently-typed programming language for mechanised reasoning in mathematics, computer science and more.

[Install](#)

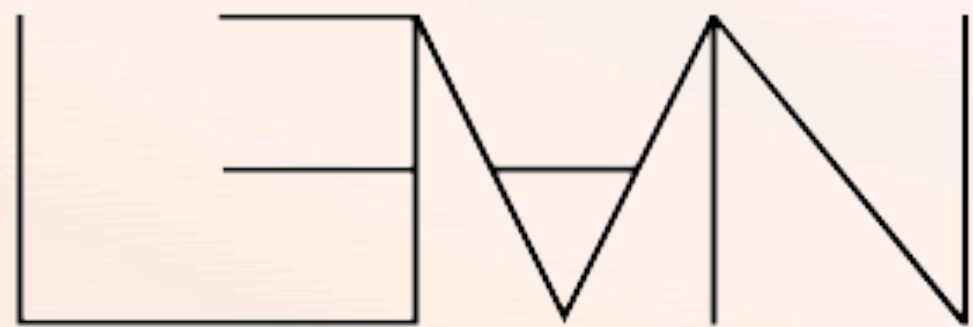
[About The Rocq Prover](#)

[Why Rocq?](#)

Latest Rocq Prover release: 8.20.1

Latest Rocq Prover release candidate: 9.0+rc1

Latest Rocq Platform release: 2024.10.1



The Rocq Prover was formerly known as the Coq Proof Assistant (see more on **the name evolution**).

Formalizing mathematics



Write definitions:
Peano natural numbers

```
ex1.v
1  Inductive nat : Set :=
2    | 0 : nat
3    | S : nat → nat.
```

▼ex1.v:3:20
No goals at this point!
►Messages (0)

Specialized models

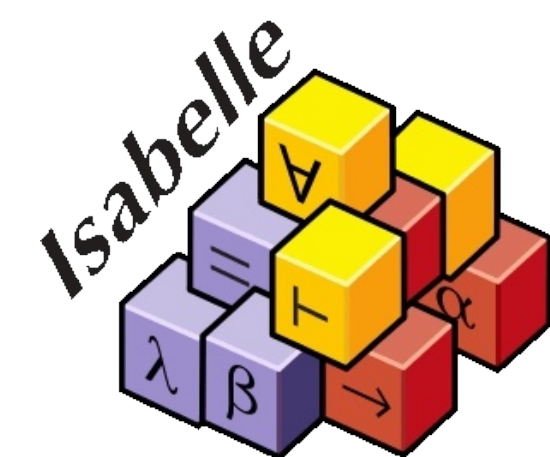
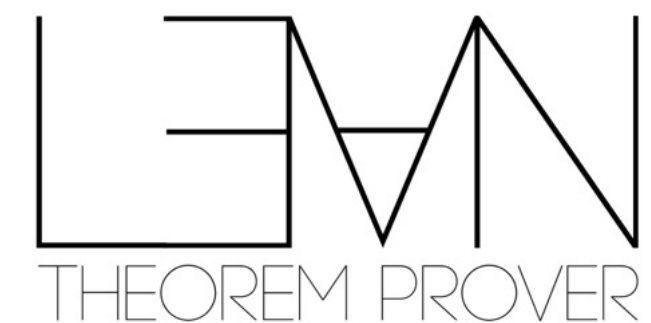
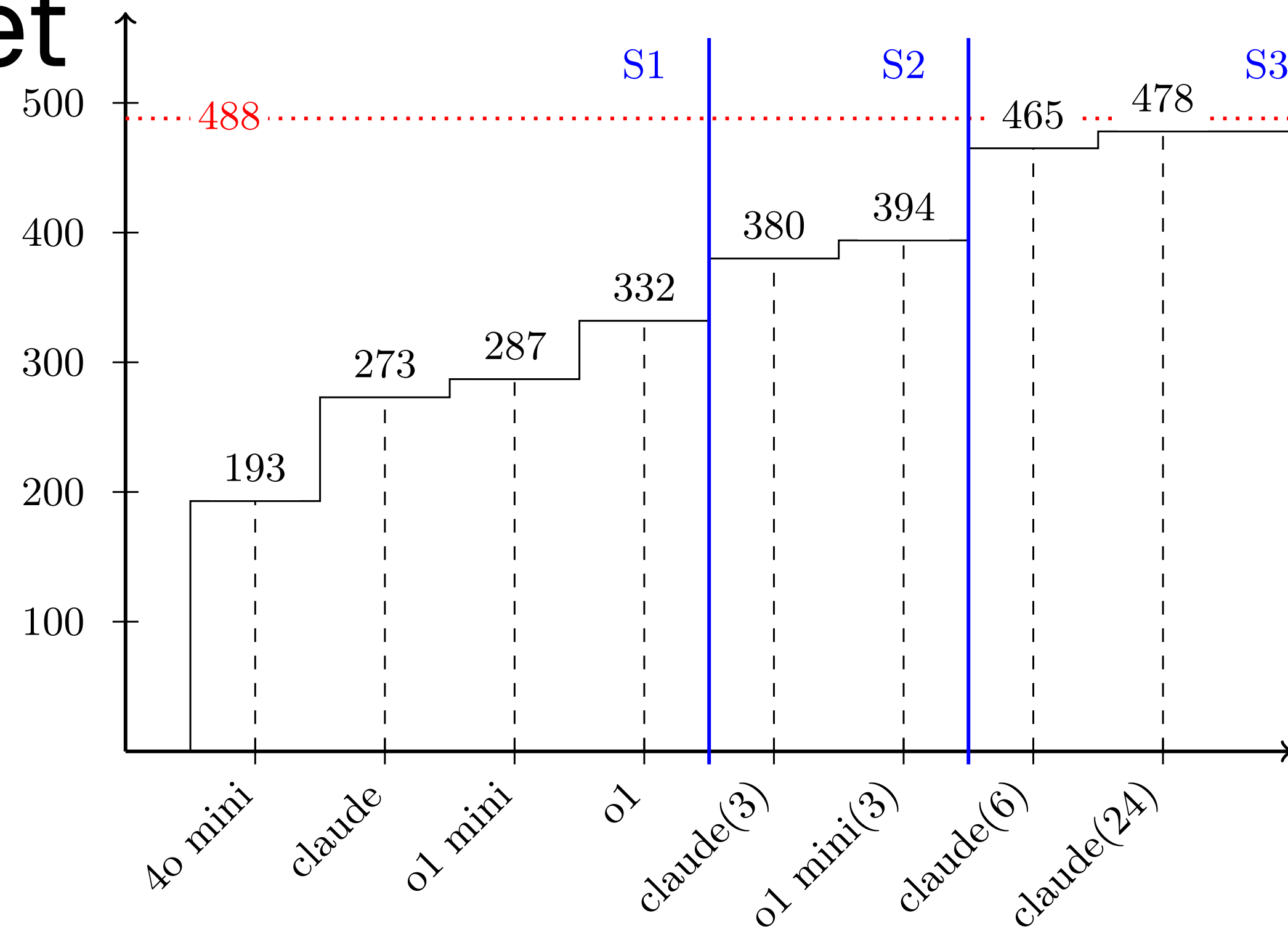
Can general models write in Rocq?

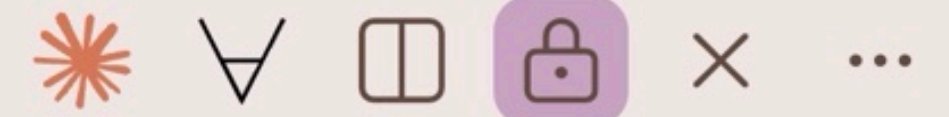
MiniF2F dataset

Category	Subcategory
IMO	
AIME	
AMC	
Other	Algebra
	Number Theory
	Induction
Total	

mathd_numbertheory_227

One morning each member of Angela's family drank a certain amount of coffee and milk. The amounts of coffee and milk varied. Angela drank a quarter of the total amount. How many people are in the family?



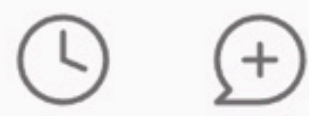


```

1  Require Import Reals Lia Lra Psatz.
2
3  Theorem mathd_numbertheory_227:
4    forall (x y n : nat),
5      x < 0 →
6      y < 0 →
7      n < 0 →
8      (INR x / 4 + INR y / 6
9       = (INR x + INR y) / INR n)%R →
10     n = 5.
11 Proof.
12   intros x y n Hx Hy Hn Heq.
13   assert (Hxpos : (0 < INR x)%R) by (apply lt_0_INR; lia).
14   assert (Hypos : (0 < INR y)%R) by (apply lt_0_INR; lia).
15   assert (Hnpos : (0 < INR n)%R) by (apply lt_0_INR; lia).
16   assert (HnR : (INR n < 0)%R) by lra.
17   assert (key : (3 * INR x * (INR n - 4) = 2 * INR y * (6 - INR n))%R).
18   {
19     assert (H : (12 * INR n * (INR x / 4 + INR y / 6) =
20              12 * INR n * ((INR x + INR y) / INR n))%R)
21       by (f_equal; exact Heq).
22     assert (H1 : (12 * INR n * (INR x / 4 + INR y / 6) =
23                3 * INR x * INR n + 2 * INR y * INR n)%R)
24       by (field; lra).
25     assert (H2 : (12 * INR n * ((INR x + INR y) / INR n) =
26                12 * INR x + 12 * INR y)%R)
27       by (field; lra).
28     lra.
29   }
30   destruct n as [ [] [] [] [] [] n ]].
31   - contradiction.
32   - exfalso; rewrite S_INR, INR_0 in key; lra.
33   - exfalso; rewrite !S_INR, INR_0 in key; lra.
34   - exfalso; rewrite !S_INR, INR_0 in key; lra.
35   - exfalso; rewrite !S_INR, INR_0 in key; lra.
36   - reflexivity.
37   - exfalso; rewrite !S_INR in key.
38     pose proof (pos_INR n) as Hnn.
39     nra.
40 Qed.

```

Prove Angela's family theorem in Rocq



`</> ex2.v`

Can you propose a proof of the theorem in `@ex2.v` please? Here is the informal statement of the problem: "One morning each member of Angela's family drank an

● `nlinarith` isn't in scope — `Psatz` provides `nra` instead.

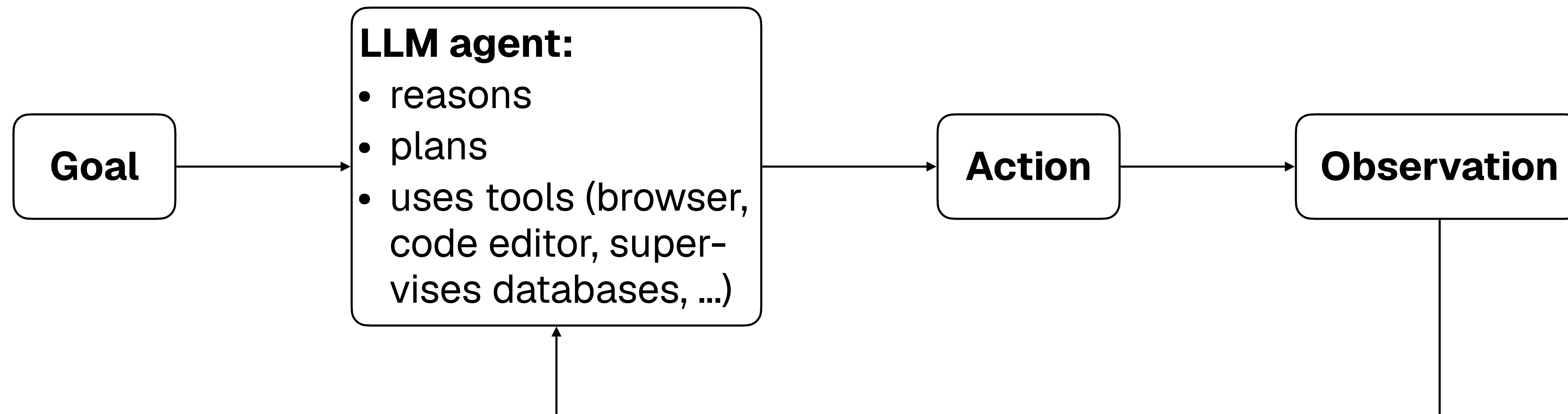
● **Edit** `ex2.v`
Modified

You've used 91% of your session limit · resets in 3h · [View usage](#) ×

⌘ Esc to focus or unfocus Claude

Agents

What are agents?



Classical LLM:

- one question = one answer
- no internal state
- no effects outside of the conversation

LLM agent:

- one goal leads to a sequence of actions
- maintains context overtime
- creates files, writes code, searches on the web, ...

Agents on Putnam 2025

```

/-- The sum of the first n squares formula -/
theorem sum_of_squares (n : ℕ) :
  (∑ i in Finset.range (n + 1), (i : ℚ)^2) =
    n * (n + 1) * (2 * n + 1) / 6 :=
  
```



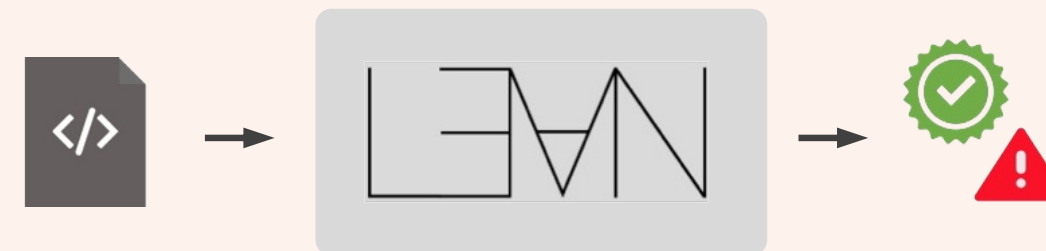
Claude Code

```

/-- The sum of the first n squares formula -/
theorem sum_of_squares (n : ℕ) :
  (∑ i in Finset.range (n + 1), (i : ℚ)^2) =
    n * (n + 1) * (2 * n + 1) / 6 := by
  induction n with n ih
  · simp
  · rw [Finset.sum_range_succ, ih]
    field_simp
    ring
  
```

Numina-Lean-MCP

Lean-LSP-MCP



LeanDex

Simplify rational expressions with polynomials

1. `Mathlib.Tactic.FieldSimp.reduceExpr`
 In: `Mathlib/Tactic/FieldSimp.lean`

DOCSTRING:

Given `x` in a commutative group-with-zero, construct a new expression in the standard form `*** / ***` (all denominators at the end) which is equal to `x`.

DESCRIPTION:

This procedure takes an expression `x` in a commutative group-with-zero structure and rewrites it into a standard fractional form where all denominators appear at the end, producing both the simplified expression and a proof of its equality to `x`. It operates by first analyzing the expression's structure to ensure it contains only standard arithmetic operations, then utilizes a normalization algorithm that transforms the expression into a sign-adjusted multiplicative factorization, simplifies this factorization by removing zero-exponent terms with nonzero bases, and finally evaluates it to produce the result.

Informal Prover

Prove the sum of squares formula for the first n natural numbers: $0^2 + 1^2 + 2^2 + \dots + n^2 = n(n+1)(2n+1) / 6$

Detailed Proof:

1. Base Case: $n = 0$ Left-Hand Side (LHS) = $0^2 = 0$. Right-Hand Side (RHS) = $[0(0+1)(2*0+1)] / 6 = (0 * 1 * 1) / 6 = 0$. LHS = RHS. The base case holds.
2. Inductive Hypothesis (IH) Assume that...

Discussion Partner

The `ring` tactic is failing. What's missing?

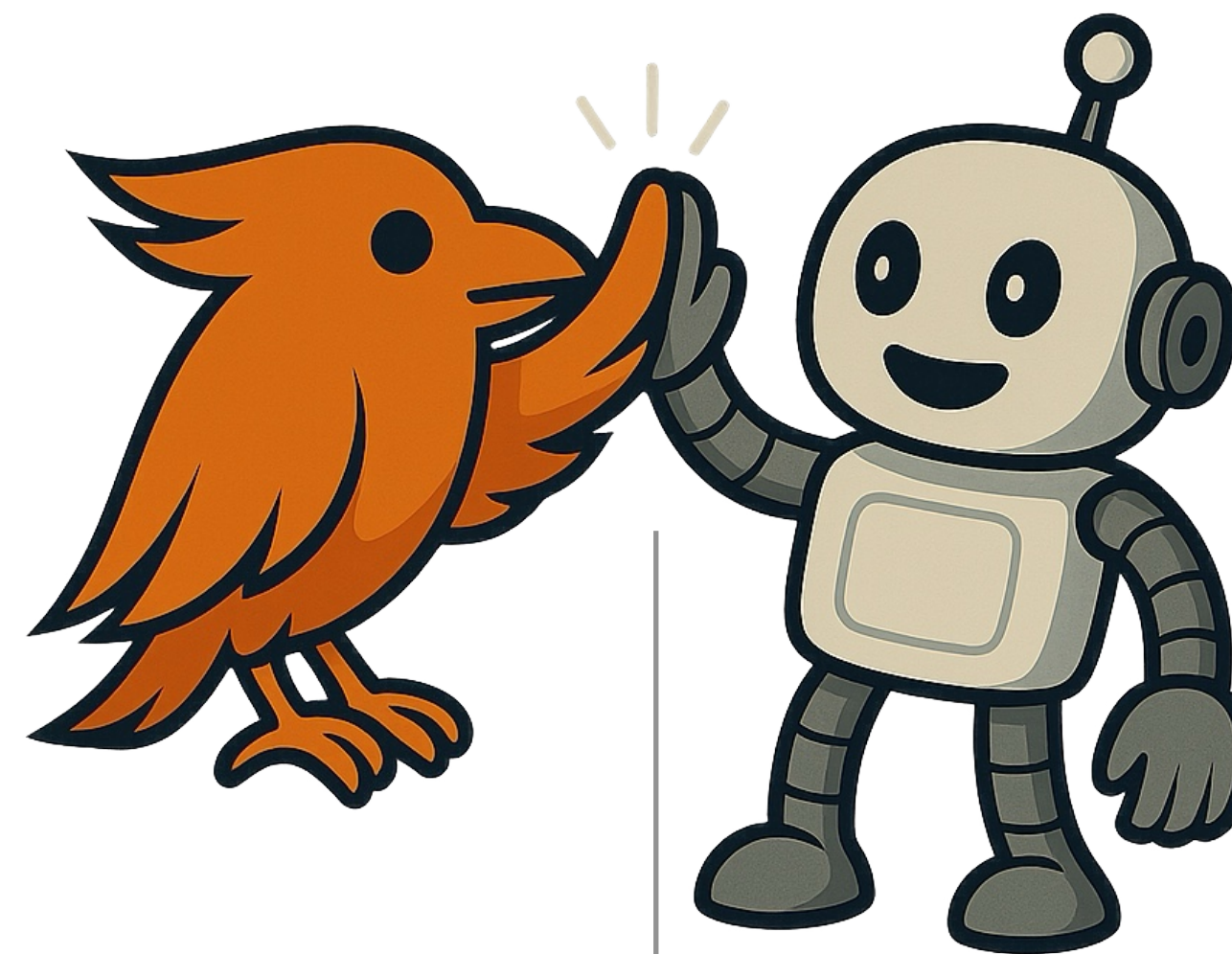


Apply `field_simp` first. You need to clear the denominators before `ring` can solve the polynomial.

What do agents need?

coq-isp

Flèche »→



petanque / pytanque

rocq-mcp

NLIR

LLM4Docq

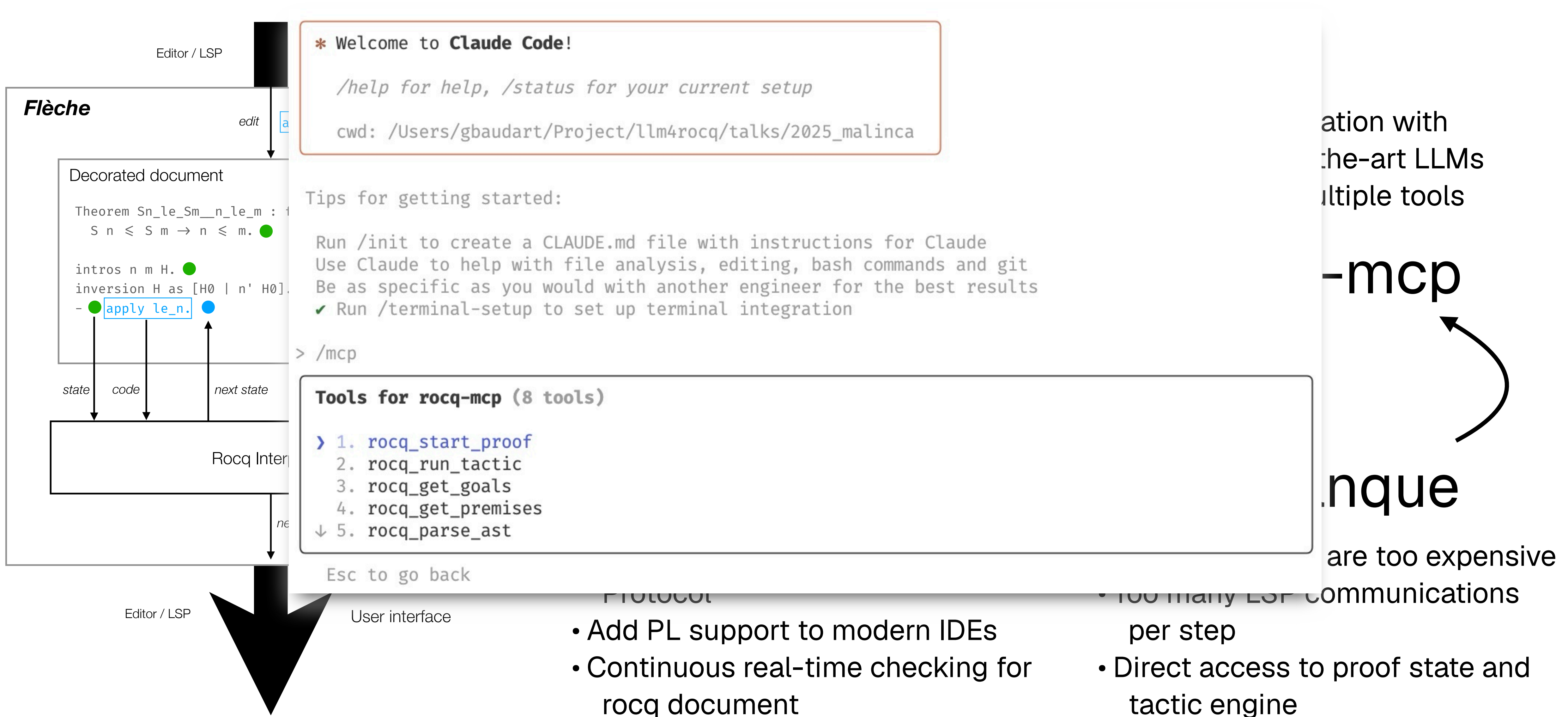
Crrrocq

Babel

Tacq



Interacting with Rocq



Tacq:

importance of context for LLMs

Overview

Current goal

```
R : fieldType
h : {poly R}
p : {poly %/ h}
hQM : hQ \is monic
hCp : coprimep hQ p
p_neq0 : p ≠ 0
F : (egcdp hQ p).1 * hQ
    + (egcdp hQ p).2 * p %= 1
```

```
((egcdp hQ p).2 * p) %% hQ %= 1 %% hQ →
(lead_coef (((egcdp hQ p).2 * p) %% hQ))^-1
 *: (((egcdp hQ p).2 * p) %% hQ)
 = 1:P
```

Extracting notations in Rocq

```
Record group : Type := { ...
```

```
Notation "a + b" := (group_operation _ a b).
```

```
Lemma morph_id :
```

```
forall G1 G2 : group,
```

```
forall morph : elements G1 → elements G2,
```

```
(forall a b, morph (a + b) = morph a + morph b) →
```

```
morph (identity_element G1) = identity_element G2.
```

```
Proof.
```

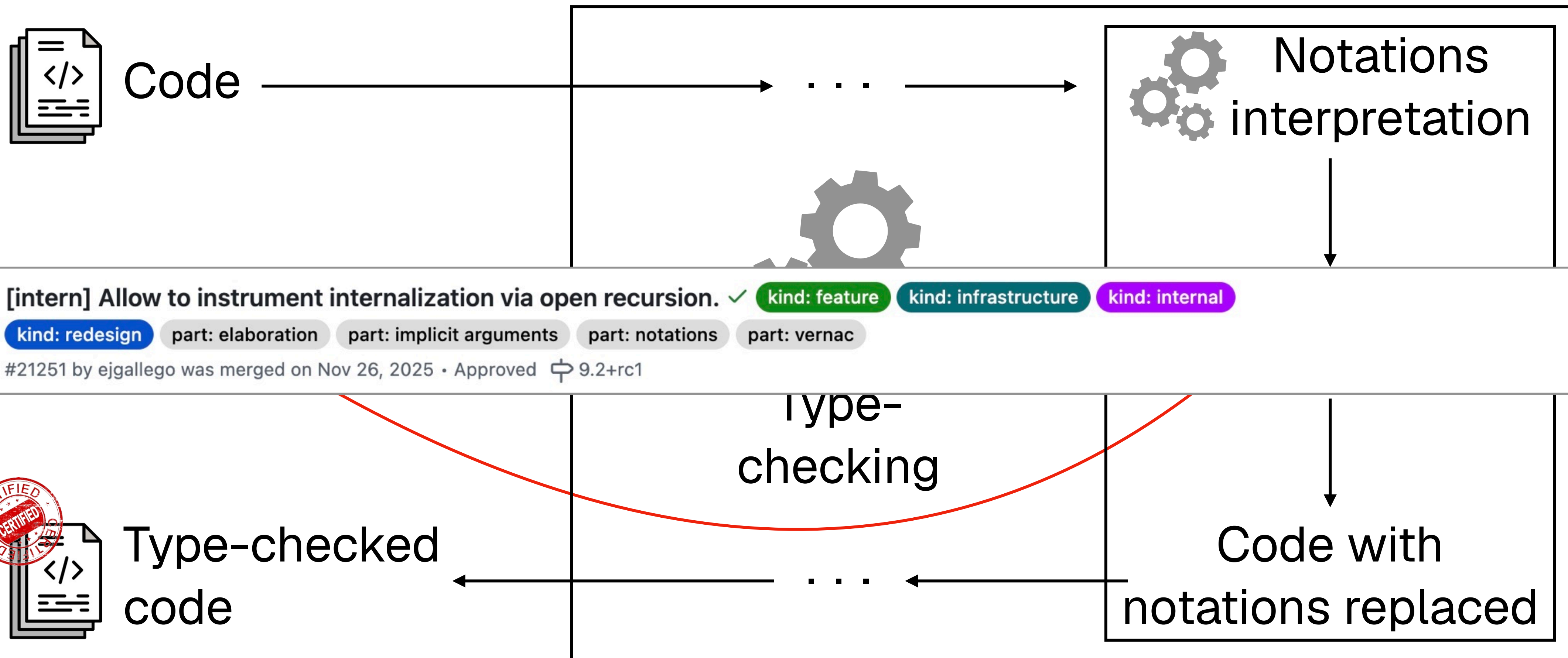


\neq

How to differentiate them?

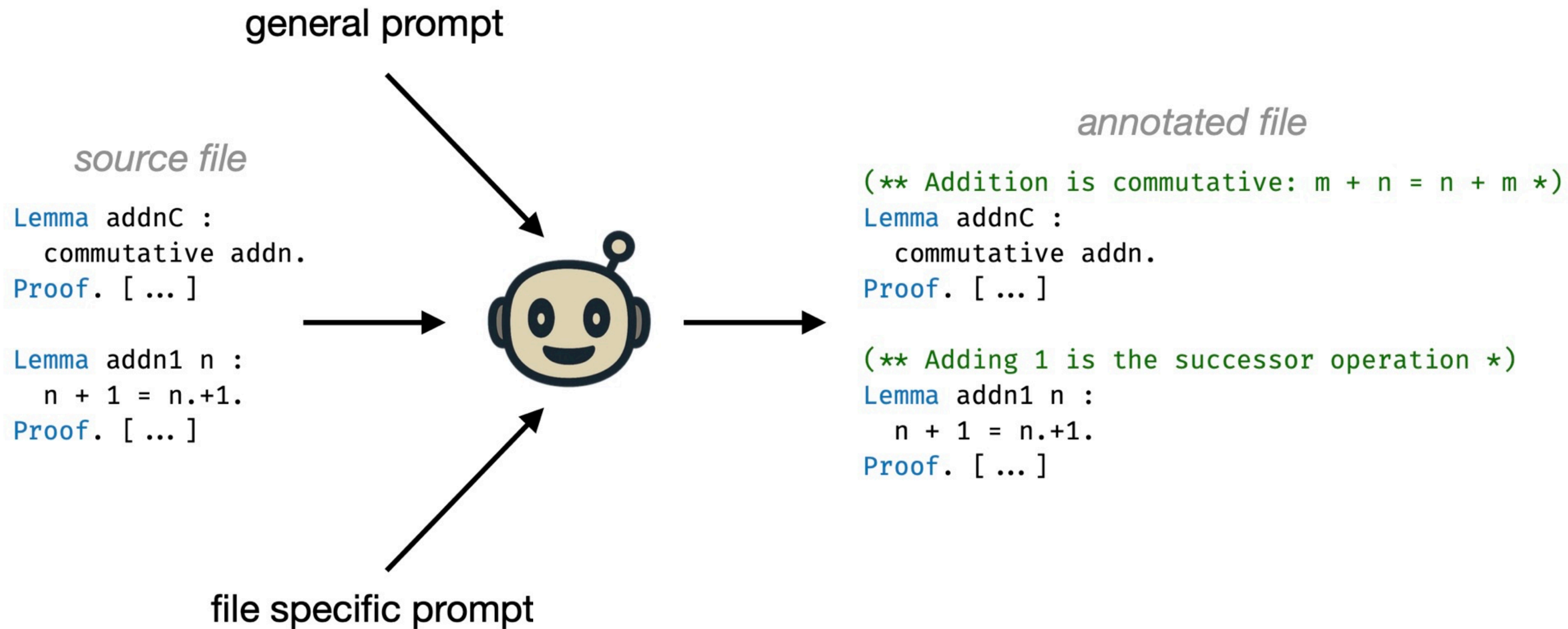
Extracting notations in Rocq

Modified Rocq interpreter



LLM4Docq: generating docstrings

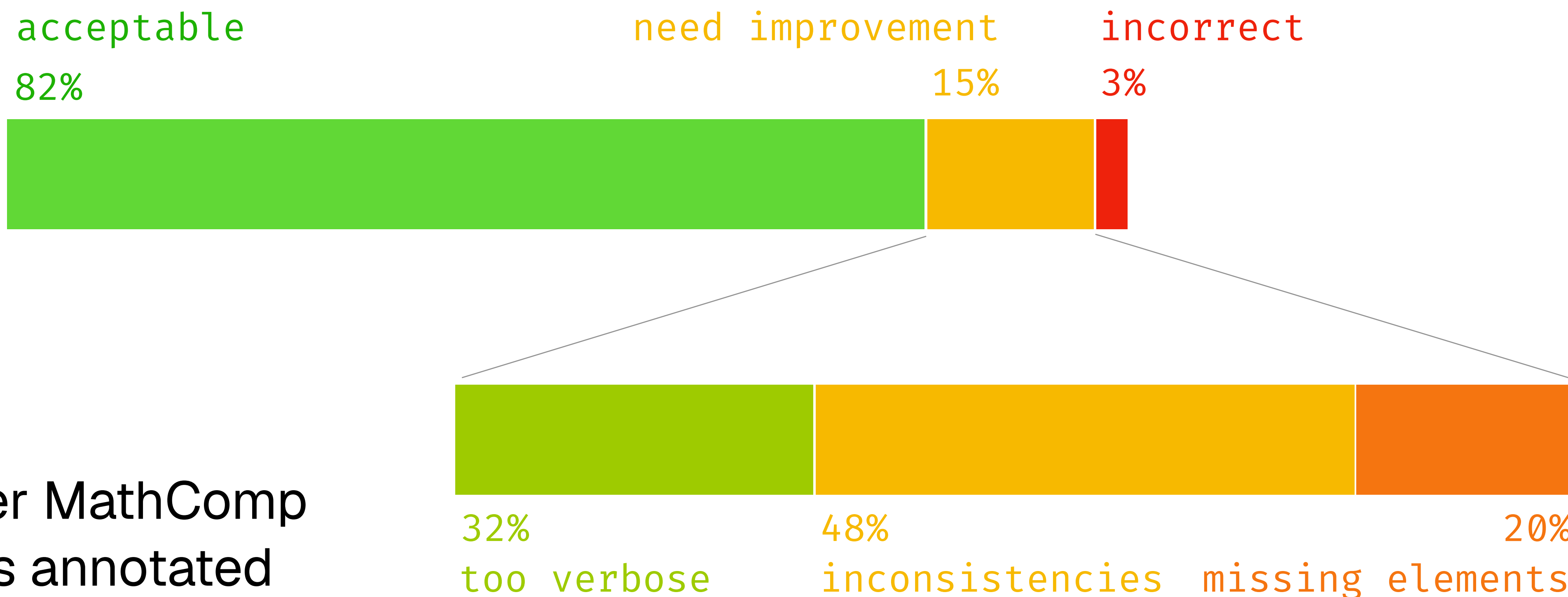
MathComp: library of ~15,000 mathematical theorems in Rocq
→ use LLMs to annotate the library



LLM4Docq: generating docstrings

MathComp: library of ~15,000 mathematical theorems in Rocq

→ use LLMs to annotate the library



Results

- 2 iterations over MathComp
- 20,000 objects annotated
- 600 annotations reviewed

Experiments

- **Configurations:** G — $G + D + L$ — $G + D \& N$ — $G + D \& N + L$

Example:

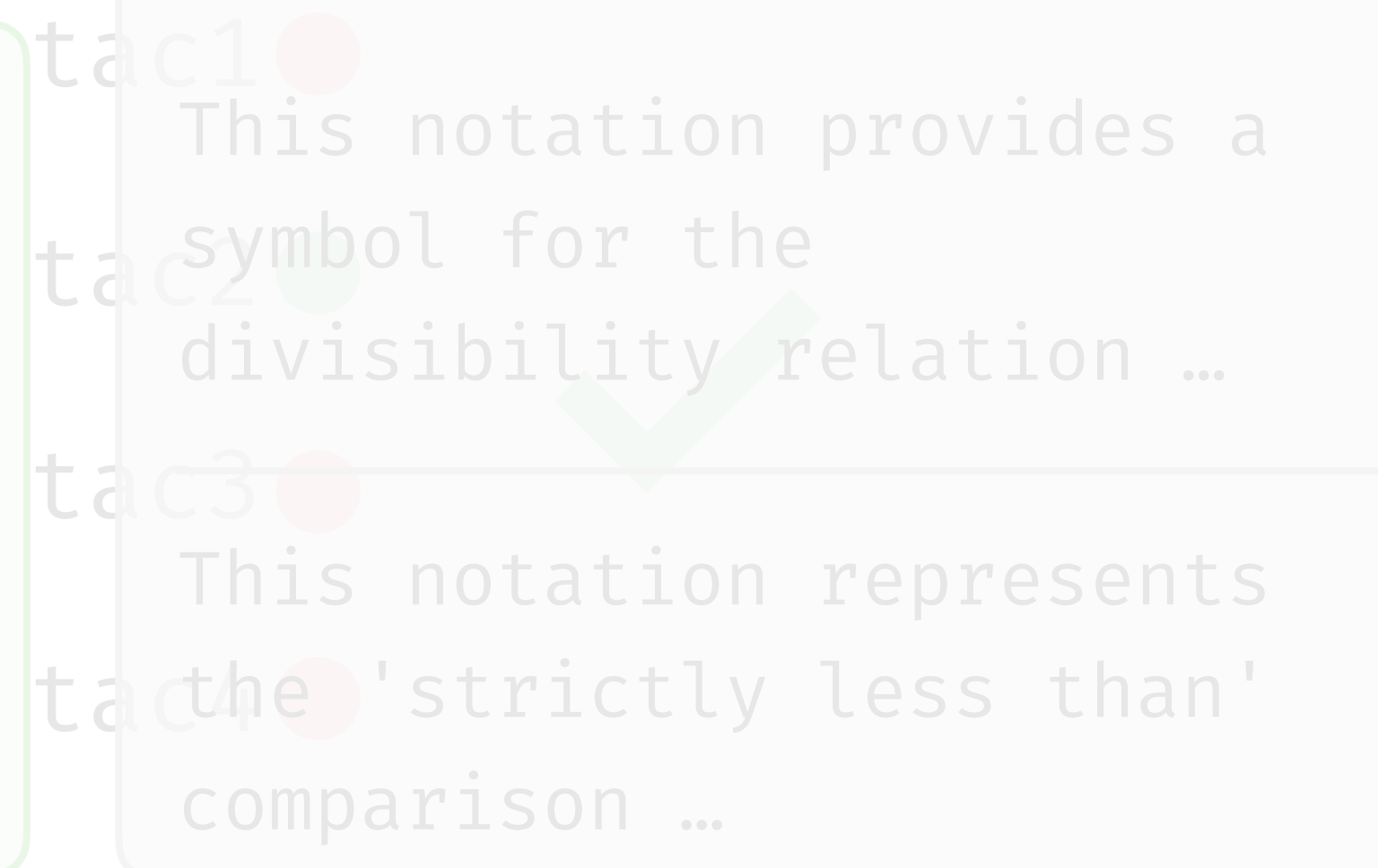
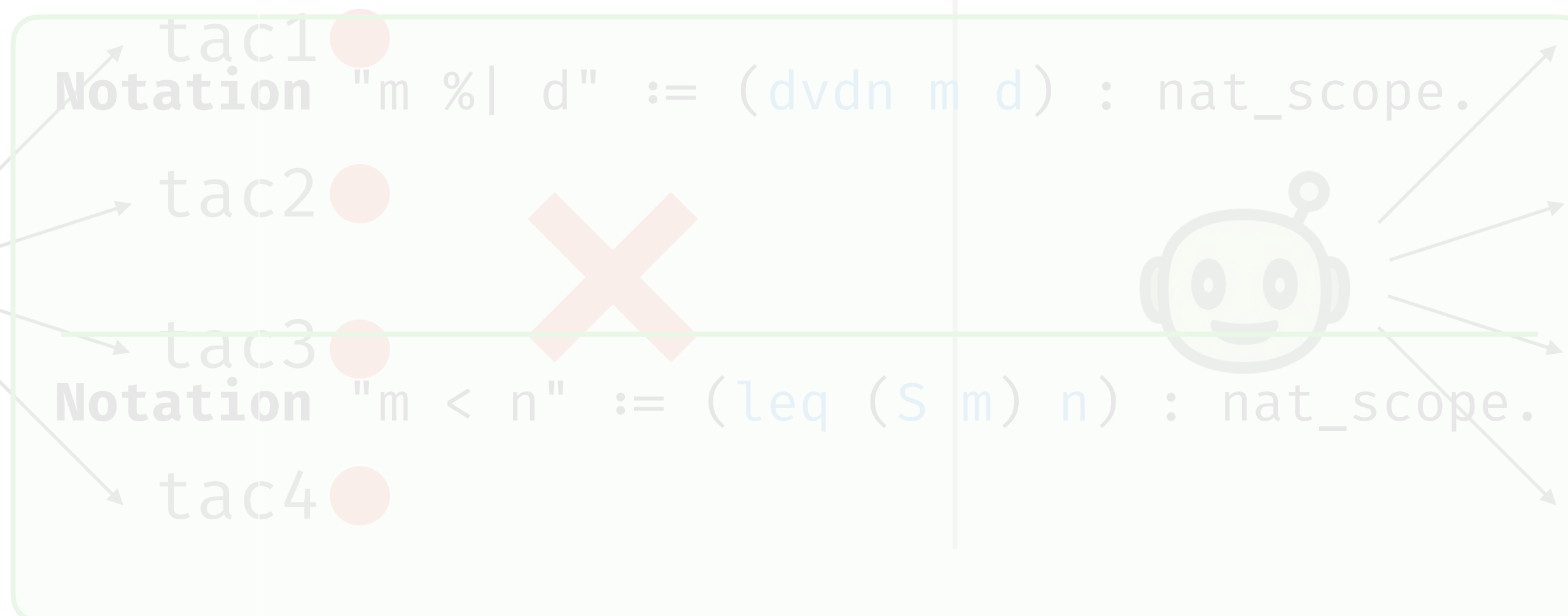
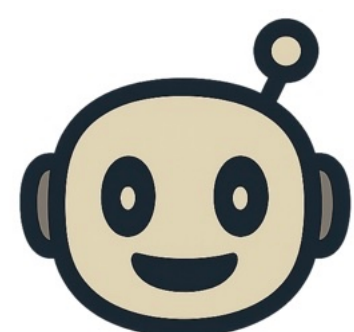
```
n : nat
fact div :
Pass@k:  $k \in \{1, 4, 8, 32\}$ 
  n % | factorial n
probability of one success among k generated samples
Example: pass@4
```

```
Inductive nat : Set :=
  | 0 : nat
  | S : nat → nat.

Fixpoint factorial n :=
  if n is n'.+1 then n * factorial n'
  else 1.
```

This is an inductive type representing the set of natural numbers, ...

This definition computes the factorial of a natural number, ...



Results

Is the generated tactic valid?

1. Bigger models benefit the most from additional context

Model	Configuration	pass@1	pass@4	pass@8	pass@32
GPT 4o	G	32.7	52.4	61.2	—
	G + D + L	35.2	55.2	64.3	—
	G + D & N	41.2	60.7	68.0	—
	G + D & N + L	42.1	61.7	70.4	—
CLAUDE SONNET 4	G	47.8	59.4	64.5	—
	G + D + L	43.5	57.7	62.8	—
	G + D & N	45.8	57.7	62.6	—
	G + D & N + L	50.1	62.9	67.8	—
QWEN3 32B	G	25.5	47.6	58.2	74.4
	G + D + L	27.8	50.3	60.7	77.2
	G + D & N	29.8	53.2	63.0	77.4
	G + D & N + L	29.7	53.0	63.1	77.8
QWEN3 14B	G	26.2	43.4	51.8	67.1
	G + D + L	23.8	40.9	47.5	65.1
	G + D & N	23.5	41.2	49.8	65.7
	G + D & N + L	24.1	41.8	50.4	65.7
QWEN3 4B	G	11.1	22.6	30.3	48.5
	G + D + L	10.9	22.7	30.2	47.6
	G + D & N	9.4	19.9	27.1	42.9
	G + D & N + L	11.0	22.5	29.8	46.8

**Back to Putnam 2025:
do models know their needs?**



Tristan Stérin

Computer Scientist



Contents

- [Summary](#)
- [Setup](#)
- [Tasks](#)



Tristan Stérin EDITED

With this same setup we achieved 81.1% of solved problems (198/244) on the test set of [miniF2F-rocq](#) in about 6h of active time (25h total but stalled a lot of the time because of waiting for session limits to reset).

20 FEB

19:03

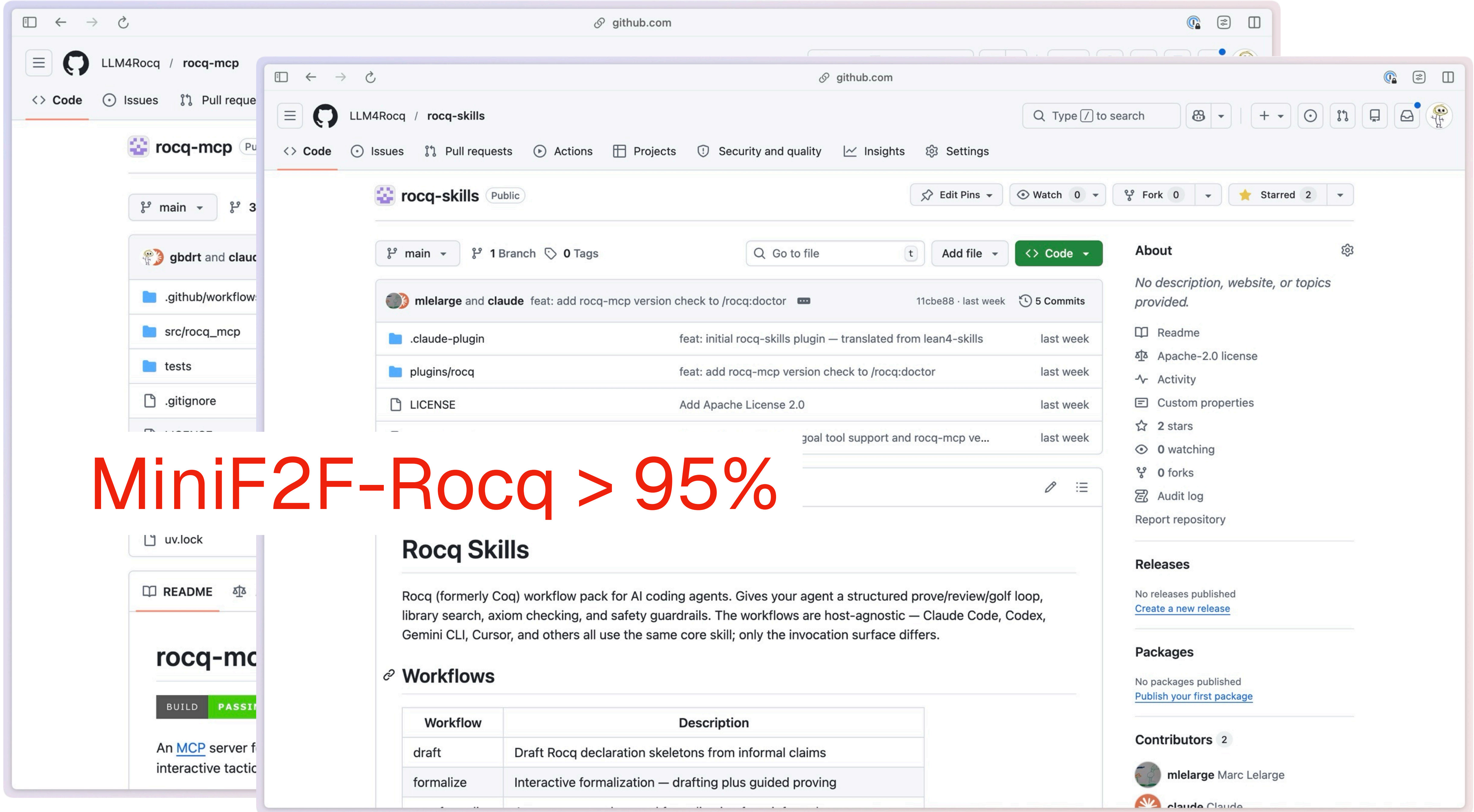
[\(4\)_proof](#)

Ventures Rese

Opus 4.6 is great at formal proofs

February 17th 2026

What started as an excuse to try out Anthropic's new [agent teams](#) feature led me down a rabbit hole exploring how good Opus 4.6 is at formal proofs in Rocq and Lean4.



MiniF2F-Rocq > 95%

Workflow	Description
draft	Draft Rocq declaration skeletons from informal claims
formalize	Interactive formalization — drafting plus guided proving

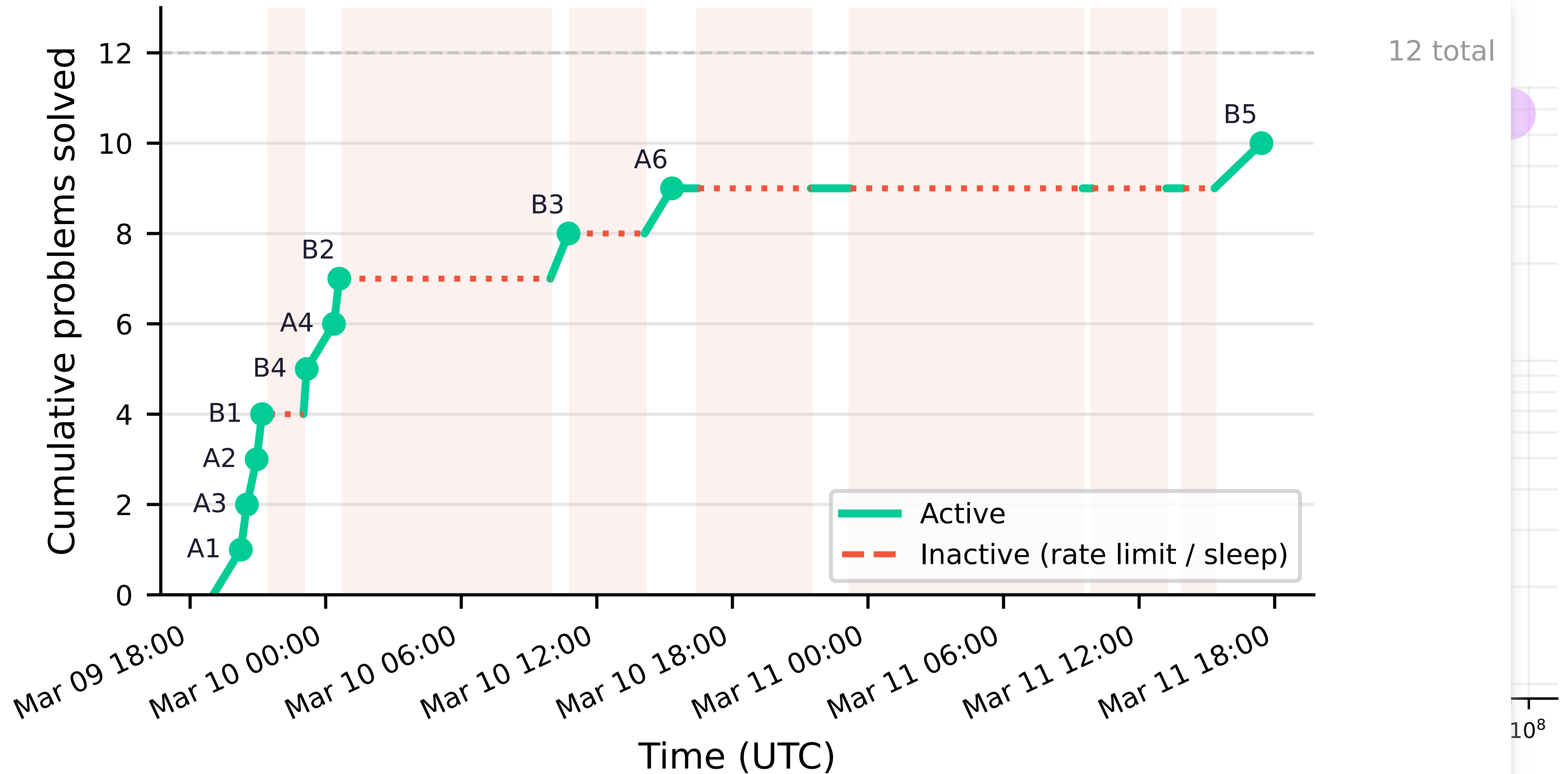
Results

Launch a

Formal statements

Problem
putnam_2025_a
putnam_2025_a
putnam_2025_a
putnam_2025_a
putnam_2025_a
putnam_2025_a
putnam_2025_b
putnam_2025_b
putnam_2025_b
putnam_2025_b
putnam_2025_b
putnam_2025_b
putnam_2025_b
putnam_2025_b

10/12 verified



Targeting the remaining problems

- **B6**: solved when given the Lean proof
- **A5**: lead to the creation of a whole library:
mathcomp-eulerian
 - Claude Opus 4.6 / 4.7 + rocq-mcp
 - exact correspondence with sections 1.4 and 1.6 of Stanley's *Enumerative Combinatorics*
 - 33 files \approx 17,600 lines
 - following the style of *mathcomp*, the biggest mathematical library in Rocq

```
000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
```

Library Before Proof:
Making LLM-Generated Rocq Usable by Mathematicians

Anonymous Authors¹

Abstract

LLM-assisted formalization now produces thousands of lines of Lean or Rocq for research-level mathematics, but the resulting code remains hard to read for a working mathematician who is not a Rocq expert. We argue that the bottleneck is no longer closing goals, but the library built *around* the proof terms. We propose a three-artifact approach: mathcomp-style Rocq backed by an explicit coding-style skill, a Lean-style blueprint cross-linking informal math to formal lemmas, and a verification PDF pairing each definition and theorem statement with its informal version. We demonstrate it on two case studies: a textbook chapter of Stanley's *Enumerative Combinatorics* and a Rocq formalization of Quasi-Borel Spaces. Both libraries contain zero `Admitted` and zero custom axioms; both were produced in a single Docker environment by Claude Opus driving Rocq through an MCP server, to which we contributed. The contribution is not the libraries themselves but a proposal for the artifacts that should accompany LLM-generated Rocq for the mathematicians, plus the tooling to produce them.

1. Introduction

Putnam 2025 Problem A5 is, on paper, a one-page combinatorics exercise: the sign sequences for which the descent pattern admits the most permutations of $\{1, \dots, n\}$ are the two alternating patterns, and the maximum is the n -th Euler/zigzag number (Stanley, 2011; Mathematical Association of America, 2026). The Axiom Math team's autonomous Lean prover failed this problem during the 2025 competition and, on a post-hoc attempt, produced a formal proof in 518 minutes using 9.1 M tokens, totalling 2,054 lines, 52 theorems and 3,074 tactics (Axiom Math, 2025). None of these numbers reflect the *informal* difficulty of A5. They reflect that the basic enumerative combinatorics of descents and Euler numbers is not packaged as a library in either Mathlib (Lean) or mathcomp (Rocq), so a goal-closing system has to re-derive it from scratch every time.

The question we explore is therefore not *can an LLM close goals in Lean or Rocq*, which is by now well-established, but *can the resulting library be read, audited and reused by a working mathematician without proof-assistant expertise?* We treat this as an open design question and propose an answer: a Rocq-MCP server and a coding-style skill keep the LLM mathcomp-compatible from the first line, while an auto-generated blueprint and a verification PDF expose the resulting code in the mathematician's own language. The blueprint and PDF are first attempts at such an interface. We demonstrate the stack end-to-end on two libraries built with Claude Opus and ship the supporting tooling so others can extend or replace any piece of it:

- **Recipe (§3)**. Three artifacts behind “readable formalization”: mathcomp-style Rocq (backed by an explicit style guide), a blueprint, and a verification PDF.
- **Case studies (§4)**. Library E (substantial portions of Stanley §1.4 and §1.6, centered on descents, Eulerian numbers, and the alternating-set maximum; \approx 17,600 lines, 609 named results, full recipe *including* verification PDF) and Library Q (Quasi-Borel Spaces, 8,913 lines, 412 proofs, first two artifacts only, no verification PDF). Both: zero `Admitted`, zero custom axioms, same Docker environment, same model, same MCP tools.
- **MCP-tool features (§5)**. Five upstream features that the recipe exposes as needed.
- **Scaling, and what we want to build next (§6)**. The full three-artifact recipe is feasible at the Eulerian scale; at the QBS scale the blueprint carries more of the reading load than a single PDF can. The verification PDF is itself a first step, not a final design: what we want is a richer mathematician-facing interface for navigating the same information.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

1

Proof agent:

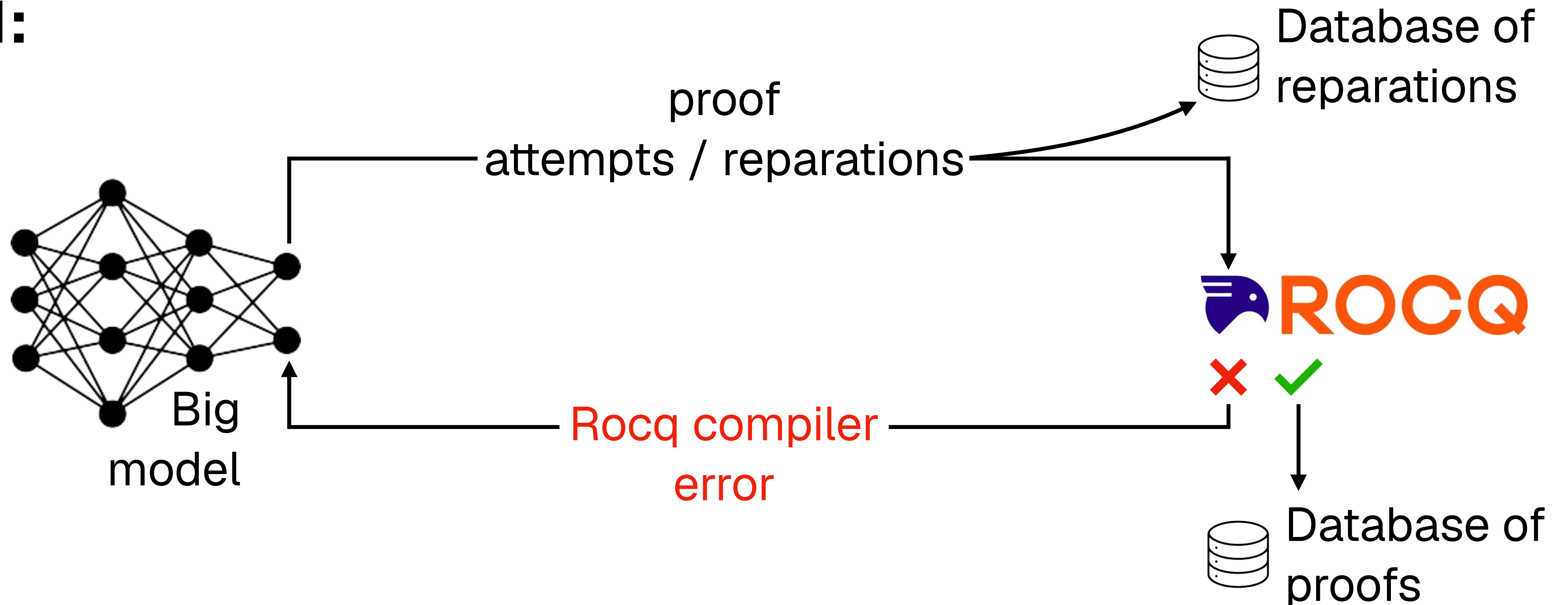
impact of good examples

Distillation by error recovery

Distillation: transferring knowledge from a bigger model to a smaller one

Goal: *limit the number of calls to big models*

Run 1:



Distillation by error recovery

Distillation: transferring knowledge from a bigger model to a smaller one

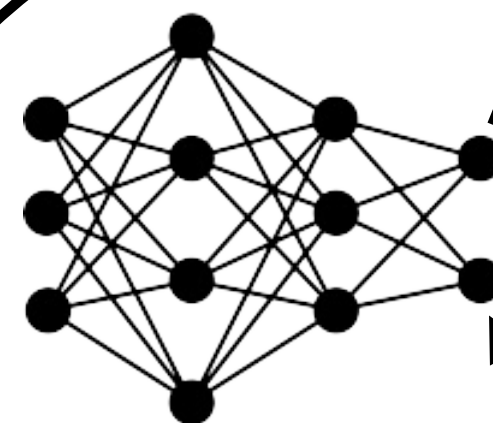
Goal: *limit the number of calls to big models*

Run 2:

 Database of reparations

proof

attempts / reparations



Small model



Rocq compiler error

 Database of proofs

Introduction to embedding models

Question: *can models represent meaning?*

The meaning is **implicit** in LLMs

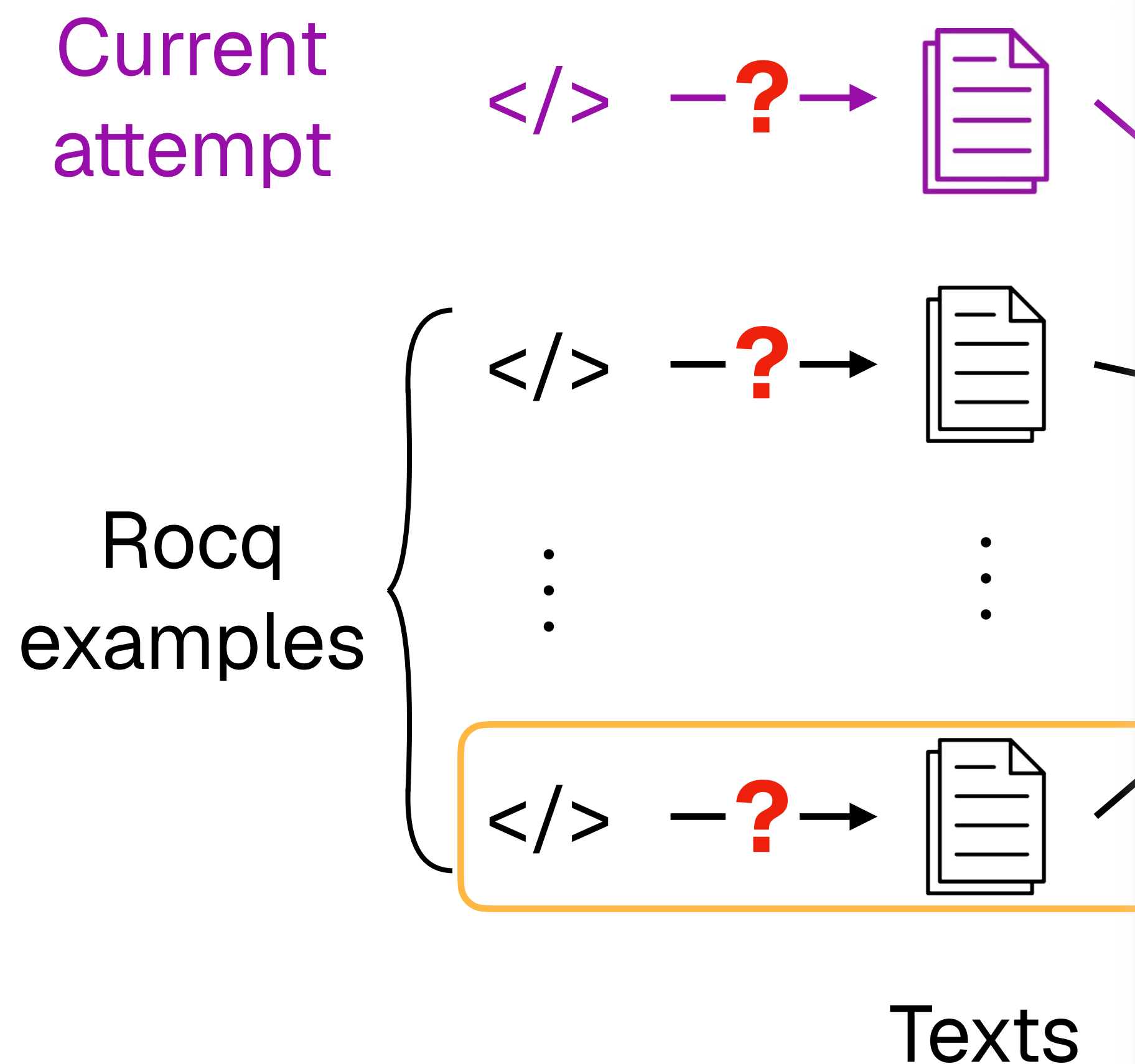
- The mouse ate the ...
 cheese ✓
 cat ✗ → choose the most probable
- Learning about the probability = learning about the meaning
- Inside LLMs: words = vectors → turn **language into geometry**

king - man + woman \approx queen

Embedding models make the meaning **explicit**

= they return vectors representing the meaning of a sentence

Retrieving relevant examples



Successful proofs:

- natural language description of the problem
- Rocq statement of the problem

Reparation examples:

- Rocq compiler error
- failing tactic
- proof up until the failing tactic

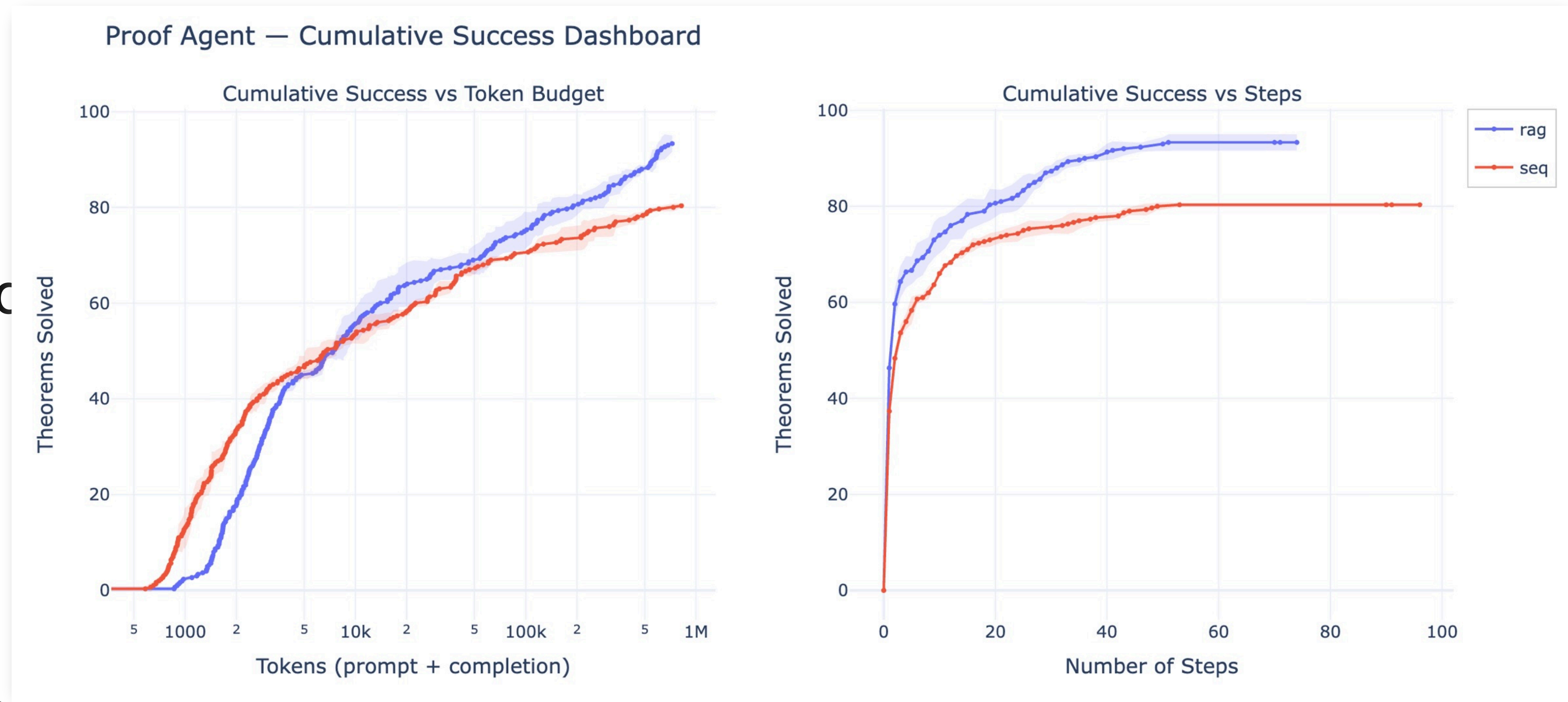
Results

Models:

- Big
- Small
- Embedd

Run 1 on

5 successful proofs examples +
5 reparation examples at each error



Thank you!

Contact: jules.viennot@gmail.com

GitHub: <https://github.com/LLM4Rocq>