

# AI and Isabelle: Experiences and Perspectives

Lawrence C Paulson, University of Cambridge

---

# “Autoformalisation is solved”

---

- ❖ Josef Urban, starting last November; by January had done a large portion of *Topology*, by James Munkres
- ❖ Urysohn’s lemma, Urysohn’s metrisation theorem, Tietze extension theorem
- ❖ 130K lines for Chad Brown’s higher-order set theory system Megalodon and its minimal library
- ❖ Afterwards, again in HOL Light and Isabelle/HOL
- ❖ “a long-running feedback loop” to ChatGPT or Claude

---

# The Becker–Mulligan Isabelle Toolset

---

- ❖ *I/Q*: an experimental Isabelle / jEdit plugin exposing proof editing capabilities as an MCP server, enabling MCP-capable agents to do interactive proofs, autonomously or collaboratively.
- ❖ *I/R* (Isabelle / REPL): interactive theory exploration outside jEdit, from the command line or programmatically via TCP and MCP.
- ❖ *Isabelle Assistant*: a set of AI-powered functions, e.g. to explain some Isabelle code, including a prompt box for general requests

*Free download* <https://github.com/aws-labs/AutoCorrode>

---

# Experiments using the toolset interactively

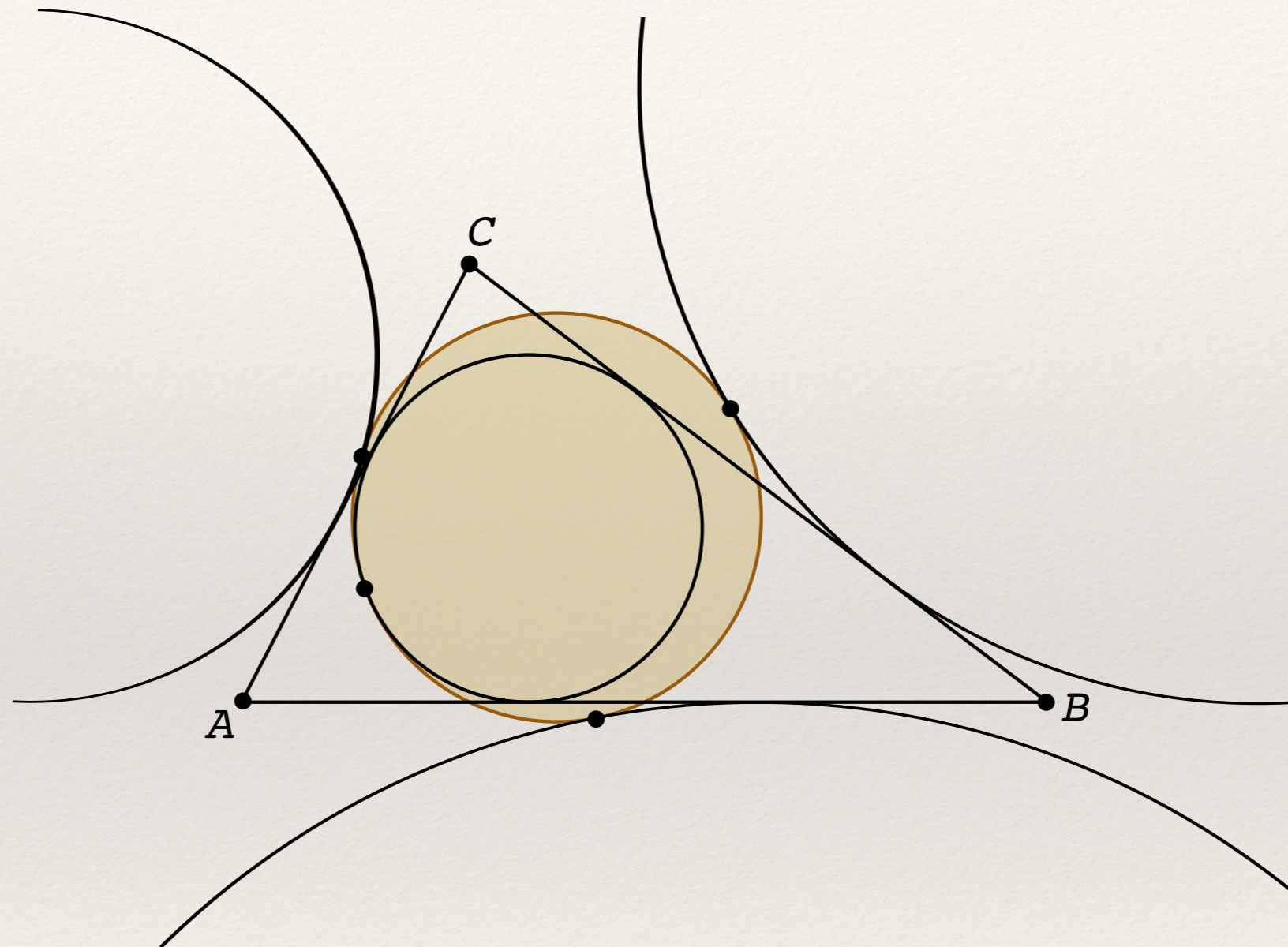
---

- ❖ Did Urysohn's metrisation theorem simply by pasting L<sup>A</sup>T<sub>E</sub>X a few lines at a time
- ❖ Proving the Ramsey number  $R(3,6) = 18$ : *much harder!*
  - ❖ Worse: proving  $R(3,6) > 17$  requires reading a diagram and analysing it for cliques. *Didn't try!*
- ❖ Elementary group theory from general knowledge
- ❖ Translating whole proofs from Lean and HOL Light

---

# Feuerbach's theorem

---



---

# Porting Feuerbach's theorem

---

- ❖ From a HOL Light proof that uses the WLOG tactics
- ❖ The port was largely automatic, but needed hints including HOL Light proof states
- ❖ Proof blew up from 214 lines to nearly 1500 due to verbose and needless calculations
- ❖ I got it down to 900 lines, partly automatically



---

# The good, the bad, the ugly

---

- ❖ Claude knows Isabelle conventions:
  - ❖ natural-sounding names; reasonable abbreviations
  - ❖ mapping  $\text{real}^1$  to  $\text{real}$  and  $\text{real}^N$  to `'a::euclidean_space` and much more
  - ❖ creating *idiomatic*, **readable** structured proofs
- ❖ But the proofs are needlessly fine-grained
- ❖ ... 1000s of lines often with duplication/symmetries

---

# The Fair Game theorem

---

- ❖ Aka *optional stopping theorem*, it's about martingales
- ❖ Ported from Lean, and it needed two Mathlib entries
- ❖ The whole job took just hours & was fully automatic
- ❖ I have no idea whether it is correct — clearly, correctness is crucial!

namespace MeasureTheory

```
variable {Ω : Type*} {m0 : MeasurableSpace Ω} {μ : Measure Ω} {G : Filtration ℕ m0} {f : ℕ → Ω → ℝ}
  {τ π : Ω → ℕ}
```

```
set_option backward.isDefEq.respectTransparency false in
```

```
/-- Given a submartingale `f` and bounded stopping times `τ` and `π` such that `τ ≤ π`, the
expectation of `stoppedValue f τ` is less than or equal to the expectation of `stoppedValue f π`.
This is the forward direction of the optional stopping theorem. -/
theorem Submartingale_expected_stoppedValue_mono {E : Type*} [NormedAddCommGroup E]
  [NormedSpace ℝ E] [CompleteSpace E] [PartialOrder E] [IsOrderedAddMonoid E]
  (hadd : StronglyAdapted G f)
  [IsOrderedModule ℝ E] [ClosedCofTopology E] [SigmaFiniteFiltration μ G] {f : ℕ → Ω → E}
  (hf : Submartingale f G μ) (hτ : IsStoppingTime G τ) (hπ : IsStoppingTime G π) (hle : τ ≤ π) :
  {N : ℕ} (hbdd : ∀ ω, π ω ≤ N) → μ[stoppedValue f τ] ≤ μ[stoppedValue f π] := by
```

```
rw [← sub_nonneg, ← integral_sub', stoppedValue_sub_eq_sum' hle hbdd]
· simp only [Finset.sum_apply]
have : ∀ i, MeasurableSet[G i] {ω : Ω | τ ω ≤ i ∧ i < π ω} := by
  specialize hf (s.pieces (fun _ => i) fun _ => j) _ (isStoppingTime_pieces_const hij hs)
  intro i
  refine (ht i).inter ?_
  convert (hπ i).compl using 1
  ext x
  simp; rfl
  split ifs with hw
  · refine Finset.sum_nonneg fun i _ => ?_
  rw [integral_indicator (G.le _ (this _)), integral_sub', sub_nonneg]
  · exact hf.setIntegral_le (Nat.le_succ i) (this _)
  · exact (hf.integrable _).integrableOn (G.le _ hs) (hint _).integrableOn (hint _).integrableOn, ←
  · exact (hf.integrable _).integrableOn (G.le _ hs) (hint j), add_le_add_iff_right] at hf
intro i _
```

```
exact IntegrableIndicator (integrable_sub (hf.integrable _) (hf.integrable _)) (hf.integrable _)
```

```
· exact hf.integrable_stoppedValue ht hbdd
· exact hf.integrable_stoppedValue ht fun ω => le_trans (hle ω) (hbdd ω)
  (hadd : StronglyAdapted G f) (hint : ∀ i, Integrable (f i) μ) :
```

```
Submartingale f G μ ↔ ∀ τ π : Ω → ℕ, IsStoppingTime G τ → IsStoppingTime G π →
  τ ≤ π → (∃ N : ℕ, ∀ x, π x ≤ N) → μ[stoppedValue f τ] ≤ μ[stoppedValue f π] :=
  ⟨fun hf _ _ ht hπ hle (hN) => hf.expected_stoppedValue_mono ht hπ hle hN,
  submartingale_of_expected_stoppedValue_mono hadd hint⟩
```

---

# The isoperimetric theorem\*

---

- ❖ “among all closed curves in a plane with a fixed perimeter, the circle encloses the largest area”
- ❖ needs mountains of libraries, e.g. *bounded variation*
- ❖ co-operative proof including **helping the LLM:**
  - ❖ noting when it's gone wrong & cancelling
  - ❖ giving hints or doing part of the job manually
  - ❖ giving it *refactoring* suggestions

Claude seems fine at porting Lean proofs

It's also good with L<sup>A</sup>T<sub>E</sub>X; less good with PDF

For HOL Light it needs lots of hints: proofs have less structure and are **much** longer

What do we do with our heaps of  
machine-generated proofs?

---

# An analogy with compilers?

---

## The FORTRAN Automatic Coding System

J. W. BACKUS†, R. J. BEEBER†, S. BEST‡, R. GOLDBERG†, L. M. HAIBT†,  
H. L. HERRICK†, R. A. NELSON†, D. SAYRE†, P. B. SHERIDAN†,  
H. STERN†, I. ZILLER†, R. A. HUGHES§, AND R. NUTT||

### INTRODUCTION

THE FORTRAN project was begun in the summer of 1954. Its purpose was to reduce by a large factor the task of preparing scientific problems for IBM's next large computer, the 704. If it were possible for the 704 to code problems for itself and produce as good programs as human coders (but without the errors), it was clear that large benefits could be achieved.

system is now complete. It has two components: the FORTRAN language, in which programs are written, and the translator or executive routine for the 704 which effects the translation of FORTRAN language programs into 704 programs. Descriptions of the FORTRAN language and the translator form the principal sections of this paper.

The experience of the FORTRAN group in using the system has confirmed the original expectations con-

---

# Is the output of AI like object code?

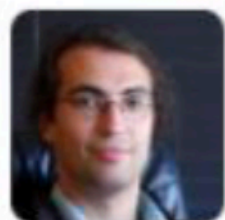
---

- ❖ In the 1950s, compilers “wrote code” for us
- ❖ Nobody looked at this code or edited it, but only ran it
- ❖ Should we use AI the same way?
  - ❖ Simply trust those proofs without looking at them?
  - ❖ *Shouldn't formal proofs offer more?*

# In Math, Rigor Is Vital. But Are Digitized Proofs Taking It Too Far?



*The quest to make mathematics rigorous has a long and spotty history — one mathematicians can learn from as they push to formalize everything in the computer program Lean.*



**Patrick Massot**  MOVED



I think the situation is pretty clear: AI companies, and especially MathInc, will indeed thoroughly bomb this area to turn it in a giant radioactive wasteland that will never be able to sustain life again, so we will never get the benefits expected from formalization (improved understanding and accessibility). I strongly advise young people to contribute to less shiny projects that are less likely to be destroyed.



5



17



5



3



---

# Proofs are not compiled code!

---

- ❖ Programming languages are precise; mathematical texts contain gaps and errors.
- ❖ Assembler is unreadable; formal proofs don't have to be.
- ❖ What if the LLM exploits some soundness bug?
- ❖ We *demand* more from a proof than “computer says OK”

*We must inspect auto-generated proofs and make them readable.*

```

lemma continuous_on_path_length_subpath_right:
  assumes "rectifiable_path g" "s ∈ {0..1}"
  shows "continuous_on {0..1} (λt. path_length (subpath s t g))"
proof -
  have g_bv: "has_bounded_variation_on g {0..1}"
    using assms(1) unfolding rectifiable_path_def by auto
  have g_cont: "continuous_on {0..1} g"
    using assms(1) unfolding rectifiable_path_def path_def by auto
  define V where "V t = vector_variation {0..t} g" for t :: real
  have V_cont: "continuous_on {0..1} V"
    unfolding V_def continuous_on_eq_continuous_within
  proof
    fix c :: real assume c01: "c ∈ {0..1}"
    have "continuous (at c within {0..1}) (λx. vector_variation {0..x} g) ↔
      continuous (at c within {0..1}) g"
      using vector_variation_continuous[OF g_bv c01] .
    moreover have "continuous (at c within {0..1}) g"
      using g_cont c01 unfolding continuous_on_eq_continuous_within by auto
    ultimately show "continuous (at c within {0..1}) (λt. vector_variation {0..t} g)"
      by simp
  qed
  have V_mono: "V x ≤ V y" if "x ∈ {0..1}" "y ∈ {0..1}" "x ≤ y" for x y
  proof -
    have bv_0y: "has_bounded_variation_on g {0..y}"
      using has_bounded_variation_on_subset[OF g_bv] that by auto
    have "V y = V x + vector_variation {x..y} g"
      using vector_variation_combine[OF bv_0y] that unfolding V_def by auto
    moreover have "vector_variation {x..y} g ≥ 0"

```

---

# Which is why you should use Isabelle

---

- ❖ Isabelle has the best of AUTOMATH (a weak logical framework), the HOL system, Mizar, Haskell (type classes)
- ❖ *nested structured proofs*, which can be legible
- ❖ Best-in-class automation
- ❖ Sledgehammer can repair or simplify AI slop.
- ❖ *No dependent types*

---

# Don't we need dependent types?

---

- ❖ Who needs 4GB proof objects?
- ❖ Universe polymorphism complicates proofs
- ❖ They don't even work:  $T(n + 1)$  and  $T(1 + n)$  are different types!
- ❖ So Mathlib *discourages* dependent types, as do Rocq's SSRreflect and Mathematical Components

---

# An explanation from set theory

---

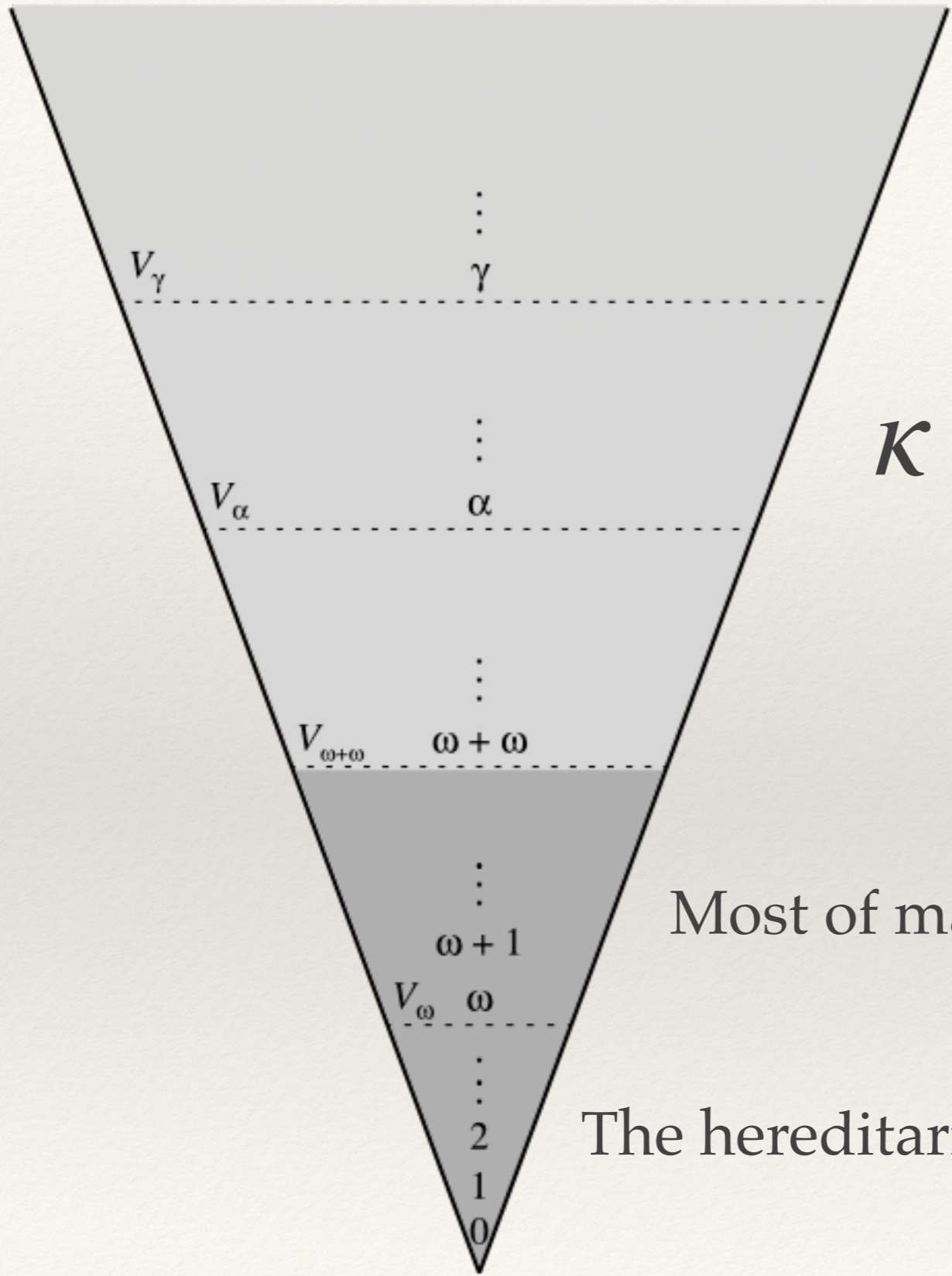
- ❖ Multi-universes  $\approx$  Tarski-Grothendieck set theory
- ❖ Most maths, including structures such as schemes, live in a tiny bit of set theory  $\approx$  *simple type theory*
- ❖ In Isabelle/HOL, we can add ZF (and more) if needed:
  - ❖ to use the ordinal / cardinal hierarchy
  - ❖ for foundational constructions e.g. *surreal numbers*

*Here be dragons*

$$\mathcal{K} = \mathcal{H}_{\mathcal{K}}$$

Most of maths, likely including FLT

The hereditarily finite sets



*Inaccessible cardinals are called  
“inaccessible” for a reason.*

*Equipping yourself with a proper class of  
them is like pushing a wheelbarrow filled  
with trillion dollar coins.*

# A plea for weaker frameworks

N.G. de Bruijn

Eindhoven University of Technology

## Abstract

It is to be expected that logical frameworks will become more and more important in the near future, since they can set the stage for an integrated treatment of verification systems for large areas of the mathematical sci-

---

# AI as *interactive* assistant

---

- ❖ Translating from other formalisms (or informal text)
- ❖ Proving theorems from its *general knowledge*
  - ❖ *identifying* a theorem from its formula
  - ❖ ... suggesting simplifications to a given formula
  - ❖ ... noting that a claim *might be false* (as stated)
- ❖ You can **discuss** and **plan** the proof

```

Lemma absolutely_setcontinuous_on_imp_has_bounded_setvariation_on:
  fixes f :: "real set  $\Rightarrow$  'a::euclidean_space"
  assumes "operative (+) 0 f"
    "absolutely_setcontinuous_on f s"
    "bounded s"
  shows "has_bounded_setvariation_on f s"
proof -
  from assms(2)[unfolded absolutely_setcontinuous_on_def, rule_format, OF zero_less_one]
  obtain r where r_pos: <r > 0>
    and r_bound: < $\forall d$  t. d division_of t  $\wedge$  t  $\subseteq$  s  $\wedge$  ( $\sum k \in d$ . content k) < r  $\longrightarrow$  ( $\sum k \in d$ . norm (
  by auto
  from <bounded s> obtain a :: real where s_sub: <s  $\subseteq$  cbox (-a) a>
    using bounded_subset_cbox_symmetric by blast
  define  $\delta$  where < $\delta = \min$  r 1 / 3>
  have  $\delta$ _pos: < $\delta$  > 0> unfolding  $\delta$ _def using r_pos by auto
  obtain p where p_div: <p tagged_division_of {-a..a}> and p_fine: <( $\lambda x$ . ball x  $\delta$ ) fine p>
    using fine_division_exists_real[OF gauge_ball[OF  $\delta$ _pos]] by blast
  show ?thesis sorry
qed

```

```

proof (state)
  this:
  ▪ p tagged_division_of {- a..a}
  ▪ ( $\lambda x$ . ball x  $\delta$ ) fine p

goal (1 subgoal):
  1. has_bounded_setvariation_on f s

```

- p\_div: p tagged\_division\_of {-a..a}
- p\_fine: ( $\lambda x$ . ball x  $\delta$ ) fine p

Let me read the final proof state:

The proof now obtains p as a tagged division of {-a..a} that is fine with respect to the gauge  $\lambda x$ . ball x  $\delta$  where  $\delta = \min$  r 1 / 3.  
The key steps:

1. **Define**  $\delta = \min$  r 1 / 3 — the radius of the gauge balls, matching the HOL Light `min r 1 / 3 % vec 1` (which in 1D is just `min r 1 / 3`)
2. **Prove**  $\delta > 0$  — since  $r > 0$  and  $1 > 0$ , their min is positive, and dividing by 3 preserves positivity
3. **Obtain** p via `fine_division_exists_real[OF gauge_ball[OF  $\delta$ _pos]]` — this chains:
  - $\delta$ \_pos :  $0 < \delta$  into `gauge_ball` to get gauge  $(\lambda x$ . ball x  $\delta$ )
  - Then into `fine_division_exists_real` to get a tagged division p of {-a..a} with  $(\lambda x$ . ball x  $\delta$ ) fine p

The gauge  $\lambda x$ . ball x  $\delta$  is equivalent to the HOL Light gauge  $\lambda x$ . {x -  $\delta$  <..< x +  $\delta$ } via `ball_eq_greaterThanLessThan`.

```

have shrink_bounds: <Inf (?shrink k) = Inf k + (Sup k - Inf k) / 2^n>
                    <Sup (?shrink k) = Sup k - (Sup k - Inf k) / 2^n>
  if <k ∈ ?S> for k
proof -
  have ne: <Inf k + (Sup k - Inf k) / 2^n ≤ Sup k - (Sup k - Inf k) / 2^n>
    using k_props(5)[OF that] by (auto simp: box_real)
  show <Inf (?shrink k) = Inf k + (Sup k - Inf k) / 2^n>
    unfolding box_real using cInf_atLeastAtMost[OF ne]
  show <Sup (?shrink k) = Sup k - (Sup k - Inf k) / 2^n>
    unfolding box_real using cSup_atLeastAtMost[OF ne]
qed
— <Rewrite the d'-sum as a sum over ?S>
have d'_sum: <(∑K∈d'. norm (f (Sup K) - f (Inf K))) =
  (∑k∈?S. norm (f (Sup k - (Sup k - Inf k) / 2^n) - f (Inf k + (Sup k - Inf k) / 2^n))>
proof -
  have <(∑K∈d'. norm (f (Sup K) - f (Inf K))) =
    (∑k∈?S. norm (f (Sup (?shrink k)) - f (Inf (?shrink k))))>
    unfolding d'_eq using sum.reindex[OF inj_shrink] by (simp add: o_def)
  also have <... = (∑k∈?S. norm (f (Sup k - (Sup k - Inf k) / 2^n) -
    f (Inf k + (Sup k - Inf k) / 2^n)))>

```

```

proof (prove)
goal (1 subgoal):
1. absolutely_continuous_on s f

```

### Summary

The last sorry — the bound claim  $\forall n \exists F$  in sequentially.  $\sigma_n \leq \epsilon/2$  — is now fully proved. The proof connects the shrunken-sum sequence  $\sigma_n$  to the sum over the shrunken division  $d'$ , which is bounded by  $r_{int}$ . The key steps were:

- Injectivity of ?shrink on ?S** (lines 1761–1807): If  $?shrink\ k_1 = ?shrink\ k_2$  for distinct  $k_1, k_2 \in ?S$ , then  $?shrink\ k_1 \subseteq interior\ k_1$  and  $?shrink\ k_2 \subseteq interior\ k_2$ , but  $interior\ k_1 \cap interior\ k_2 = \{\}$  (from  $division\_of$ ), contradicting  $?shrink\ k_1 = ?shrink\ k_2$ .
- Inf/Sup of shrunken intervals** (lines 1808–1819): Using  $cInf\_atLeastAtMost$  and  $cSup\_atLeastAtMost$  with the non-emptiness condition from  $k\_props(5)$ .
- Rewriting the d'-sum** (lines 1820–1836): Via  $sum.reindex$  (using injectivity) and  $sum.cong$  (substituting  $shrink\_bounds$ ).
- Zero summands for degenerate intervals** (lines 1837–1857): For  $k \in d$  with content  $k = \emptyset$ , we have  $Inf\ k = Sup\ k$ , so the summand vanishes. Using  $sum.mono\_neutral\_right$ .
- Conclusion** (lines 1858–1861):  $\sigma_n = (\sum_{k \in d'} \|f(Sup\ k) - f(Inf\ k)\|) < \epsilon/2$ , hence  $\sigma_n \leq \epsilon/2$ .

---

# Other aspects of AI interaction

---

- ❖ Asking Claude why something doesn't work
- ❖ Asking it about the Isabelle libraries
- ❖ Asking it to read the HOL Light sources to discover exactly what some tactic does (RING\_CONV)
- ❖ ... then to read the Isabelle sources to find some equivalent (the algebra method)
- ❖ Helping it by Sledgehammer when it's "almost there"

---

# So where do things stand?

---

- ❖ We can formalise mathematics *interactively*, much faster
- ❖ We get proofs that are correct but messy, and should not trust them blind: *we must make them readable*
- ❖ (Sledgehammer makes this easy.)
- ❖ Source material? L<sup>A</sup>T<sub>E</sub>X, Lean, its own knowledge
- ❖ ... preferably not HOL Light or PDF

---

# And where are things going?

---

Given the piece of change, it's foolish to make detailed predictions

But things just got *much* easier

Let's take the time to create proofs we can be proud of.

Many thanks to Hanno and Dominic for these tools, also to Josef and many other research colleagues for discussions

