

Automated Sketching and Repairing of Large Mechanised Proofs



Repository



Video demo

C. Tan^{1,2}, A. F. Donaldson¹, J. Wickerson¹, J. J. Huerta y Munive³

¹Imperial College London, UK ²Kaihong, China ³Aalborg University, Denmark

{alastair.donaldson, j.wickerson}@imperial.ac.uk tanchengsong@kaihong.com jjhymuni@cs.aau.dk

Motivation & Background

- Formal verification with Interactive Theorem Proving (ITP) provides high assurance of software quality, but it is time-consuming which is a barrier to industrial adoption.
- We aim to address the automation challenges that hinder verification workflows by developing tools that streamline interactions among the human expert, the proof assistant, and automation tools.
- Case study:** A formal model of CXL cache in Isabelle/HOL (74 theory files, ≈ 310 k LoC), proving **Single-Writer-Multiple-Reader (SWMR)** cache coherence via an inductive invariant P .
- P is an conjunction comprised of $n = 796$ conjuncts that must be preserved by $m = 68$ rules, yielding a **proof obligation matrix** with over 54,000 proof obligations.

$$\begin{pmatrix} \ddots & & & & \ddots \\ & \left(\begin{array}{c} P \sigma \wedge \\ \text{guard}_{R_i} \sigma \wedge \\ \sigma \rightarrow_{R_i} \sigma' \\ \implies \\ \phi_j \sigma' \end{array} \right) & & & \\ & & (i,j) & & \\ & & & & \ddots \\ \ddots & & & & \ddots \end{pmatrix}_{m \times n}$$

super_sketch

Generates Isar proof skeletons without generative AI, inspired by Haftmann's `sketch` tool.

How it works:

- Retrieves all proof goals after user-suggested methods
- Concurrently invokes Sledgehammer and user-supplied heuristics
- Applies goal splitting and reduction for harder subgoals before retrying

Results on the 68×796 matrix:

- One full rule lemma (>700 conjuncts) regenerated in **<30 min**
- Full matrix: ≈ 1 machine-day of compute
- <100 failures** out of $>50,000$ goals

super_fix

Automated proof debugging and repair, emulating a proof engineer's debugging workflow.

How it works:

- 1st pass: Parses the `.thy` file and process the commands sequentially. Detects and fixes **timeouts**, **misaligned goals**, **failed methods** with `sorry` as a placeholder and prioritising upstream errors
- 2nd pass: Invokes Sledgehammer on each `sorry` to find alternative proofs

Repair results (after `super_sketch`):

- Pass 1: 18 files with errors \rightarrow 9 files (1 error each)
- Pass 2: \rightarrow **proof complete (0 errors)**

try_sketch example

```
(* example proof produced by try_sketch *)
lemma "( $\sum i \leq (n::nat). 2*i-1) = n^2$ "
try_sketch SORRYS[simp] (induct n)
proof(induct n)
  show goal0:"( $\sum i::nat \leq 0. 2*i-1) = 0^2$ "
  by simp
next
  show goal1:"( $\sum i \leq \text{Suc } n. 2*i-1) = (\text{Suc } n)^2$ "
  if "( $\sum i \leq n. 2*i-1) = n^2$ "
  for n :: "nat"
  using that
  sorry
qed
```

llm_try_proof

```
55 lemma (in category) homE:
56   assumes "A ∈ Ob" and "B ∈ Ob" and "f ∈ Hom A B"
57   shows "Dom f = A" and "Cod f = B"
58
59
60 llm_try_proof 59
61 proof-
62   show "Dom f = A" using assms by (simp add: hom_def)
63   show "Cod f = B" using assms by (simp add: hom_def)
64 qed
65
66 lemma (in category) id_arrow [intro]:
67   assumes "A ∈ Ob"
68   shows "id A ∈ Hom A A"
```

Ongoing Extensions

meta_sketch: recursive DSL-driven generalisation of `super_sketch`.

- Users specify a hierarchy of goal splitters and a prioritised solver portfolio
- Three-stage leaf engine: *racing*, *rescuing*, *last-resort splitting*
- Isabelle/ML Future orchestration replaces ad-hoc Python scripting

LLM integration in Isabelle:

- Python client-server bridge to Hugging Face models, Ollama, ChatGPT, and Gemini
- `llm_recommend`: suggests the next proof step
- `llm_try_proof`: DFS-based incremental proof completion with configurable depth/breadth

Fixer functions:

- Treats tool actions as functions $f: S \rightarrow \Sigma^*$ over ITP states
- `try_sketch` reimplements `super_sketch` inside this framework
- Provides a modular basis for combining LLM actions with classical automation

Future work

Evaluation. Stress-tested on 54,128 obligations (68×796): avg. ≈ 0.90 s/goal. Excessive parallelism causes *timeout cascades* (16%–40% failure rate). Next: adaptive timeouts & concurrency control.

ITP Automation Survey identifying which proof engineering tasks experts most want to automate informing case studies and tooling directions.

Acknowledgements. EPSRC grant EP/R006865/1; Horizon MCSA 2022 Postdoctoral Fellowship (project 101102608); National S&T Major Project of China (Grant No. 2024ZD0803002).

super_fix example

```
1 theory Super_Fix_Demo
2   imports Complex_Main
3
4 begin
5
6
7 lemma real_sqrt_le_iff: "x ≥ 0 ⇒ y ≥ 0
8 ⇒ sqrt x ≤ y ↔ x ≤ y ^ 2"
9   apply blast
10  done
11
12 lemma sqrt_le_itself: "1 ≤ x ⇒ sqrt x ≤ x"
13   sorry
14
15 lemma sqrt_real_nat_le: "sqrt (real n) ≤ real n"
16   apply simp
17  done
18
19 lemma semiring_factor_left:
20   fixes b:: "a::semiring"
21   shows "a * b + a * c = a * (b + c)"
22   apply simp
```