

Isabelle/RL progress report: reaping the first fruits of the project



Jonathan Julián Huerta y Munive
huertjon@cvut.cz
Czech Technical University in Prague
April, 2025



**CZECH INSTITUTE
OF INFORMATICS
ROBOTICS AND
CYBERNETICS
CTU IN PRAGUE**



**Funded by
the European Union**

General Takeaway

Me: Come here as a user of interactive theorem provers (ITPs)

DeepIsaHOL Objective: create tools that serve the Machine Learning (ML) and the ITP community for advancing the state of the art in ML and ITP

Let's collaborate!

Example



Chengsong Tan

Hi Jonathan,

Thanks very much for the reply!

The problem I am faced with is automatic repair of proofs. Say we have



Jonathan Julian Huerta y Munive

writing my own project (<https://github.com/yonoteam/DeepIsaHOL/>) but I have not documented how to set it up. I can read a theory, and find the places where an apply-style-proof is made (all with ML code). Both my



*Accepted at
FomaliSE2025:*

The Burden of Proof: Automated Tooling for Rapid Iteration on Large Mechanised Proofs

Chengsong Tan^{*†}, Alastair F. Donaldson^{*}, Jonathan Julián Huerta y Munive[‡], and John Wickerson^{*}

^{*}Imperial College London, UK. Email: {c.tan, alastair.donaldson, j.wickerson}@imperial.ac.uk

[†]Kaihong, China. Email: tanchengsong@kaihong.com

[‡]Czech Technical University, Czech Republic. Email: huertjon@cvut.cz

Summary of this talk

- The plan is to showcase contributions with lessons learned and views to the future:
 1. Data extraction algorithm
 2. Training and evaluation loops for flan-T5 and Gemma
 3. Read-eval-print loops (REPL) for interacting with Isabelle
 4. Depth-first search (DFS) evaluation algorithm
 5. Small scale experiment comparing training set ups
 6. Proof fixing tool for large developments
 7. Tooling for calling an external next-step predictor from Isabelle

Approximate model of Isabelle

An Isabelle apply-style proof

User action: a command and its arguments

lemma

```
fixes P :: "('a::finite)  $\Rightarrow$  nat  $\Rightarrow$  bool"
assumes "∀i. ∃N::nat. ∀n ≥ N. P i n"
shows "∃N. ∀i. ∀n ≥ N. P i n"
apply (rule_tac x="Max {Inf {N. ∀n≥N. P i n} | i::'a. i ∈ UNIV}" in exI)
apply clarify
subgoal for i n
  using wellorder_Inf_lemma[of "λn. P i n"] assms
  Max_ge[of "{Inf {N. ∀n≥N. P i n} | i::'a. i ∈ UNIV}", OF finite_image]
  apply -
  apply (erule_tac x=i in allE)
  apply (subgoal_tac "∀n≥Inf {N. ∀n≥N. P i n}. P i n")
```

```
proof (prove)
goal (1 subgoal):
1. ∃N. ∀i n. N ≤ n  $\longrightarrow$  P i n
```

User states: the output the user sees in-between actions

```
proof (prove)
goal (1 subgoal):
1. ∀i n. Max {Inf {N. ∀n≥N. P i n} | i. i ∈ UNIV} ≤ n  $\longrightarrow$  P i n
```

```
proof (prove)
goal (1 subgoal):
1.  $\wedge$ i n. Max {Inf {N. ∀n≥N. P i n} | i. i ∈ UNIV} ≤ n  $\implies$  P i n
```

```
proof (prove)
goal (1 subgoal):
1. Max {Inf {N. ∀n≥N. P i n} | i. i ∈ UNIV} ≤ n  $\implies$  P i n
```

An Isabelle Isar-style Proof

```

lemma
  fixes P :: "('a::finite)  $\Rightarrow$  nat  $\Rightarrow$  bool"
  assumes "  $\forall i. \exists N::nat. \forall n \geq N. P i n$ "
  shows "  $\exists N. \forall i. \forall n \geq N. P i n$ "
proof-
  let "?bound i" = "Inf {N .  $\forall n \geq N. P i n$ }"
  let ?N = "Max {?bound i | i. i  $\in$  UNIV}"
  {fix n::nat and i::'a
    from assms
    obtain M where "  $\forall n \geq M. P i n$ "
      by blast
    hence obs: "  $\forall m \geq ?bound i. P i m$ "
      using wellorder_Inf_lemma[of "  $\lambda n. P i n$ "]
      by blast
    assume "  $n \geq ?N$ "
    also have "?N  $\geq$  ?bound i"
      using finite_image
      by(fastforce intro: Max_ge)
    ultimately have "  $n \geq ?bound i$ "
      using order.trans
      by blast
    hence "  $P i n$ "
      using obs
      by blast
  }
  thus "  $\exists N. \forall i n. N \leq n \longrightarrow P i n$ "
    by blast
qed

```

```

proof (prove)
goal (1 subgoal):
  1.  $\exists N. \forall i n. N \leq n \longrightarrow P i n$ 

```

```

proof (prove)
using this:
   $\forall i. \exists N. \forall n \geq N. P i n$ 
goal (1 subgoal):
  1.  $(\bigwedge M. \forall n \geq M. P i n \implies thesis) \implies thesis$ 

```

```

proof (prove)
using this:
   $\forall n \geq M. P i n$ 
goal (1 subgoal):
  1.  $\forall m \geq \text{Inf } \{N. \forall n \geq N. P i n\}. P i m$ 

```

```

proof (prove)
using this:
   $\forall n \geq M. P i n$ 
   $\exists N. \forall n \geq N. P i n \implies P i (\text{Inf } \{N. \forall n \geq N. P i n\})$ 
   $\exists N. \forall n \geq N. P i n \implies \forall n \geq \text{Inf } \{N. \forall n \geq N. P i n\}. P i n$ 
goal (1 subgoal):
  1.  $\forall m \geq \text{Inf } \{N. \forall n \geq N. P i n\}. P i m$ 

```

```

proof (state)
this:
   $\text{Max } \{\text{Inf } \{N. \forall n \geq N. P i n\} \mid i. i \in \text{UNIV}\} \leq ?n2 \implies P ?i2 ?n2$ 
goal (1 subgoal):
  1.  $\exists N. \forall i n. N \leq n \longrightarrow P i n$ 

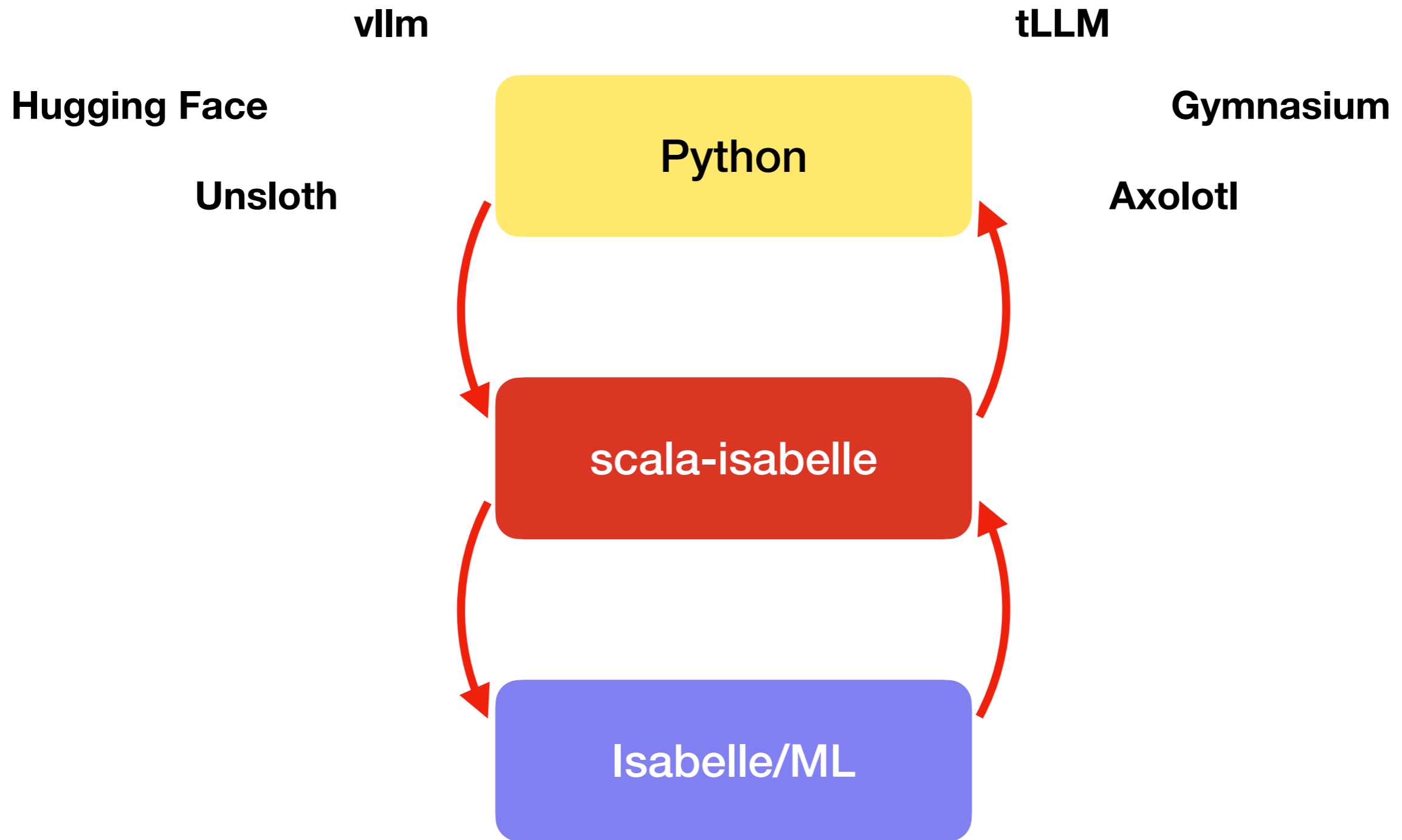
```

```

goal:
No subgoals!

```

The communication process

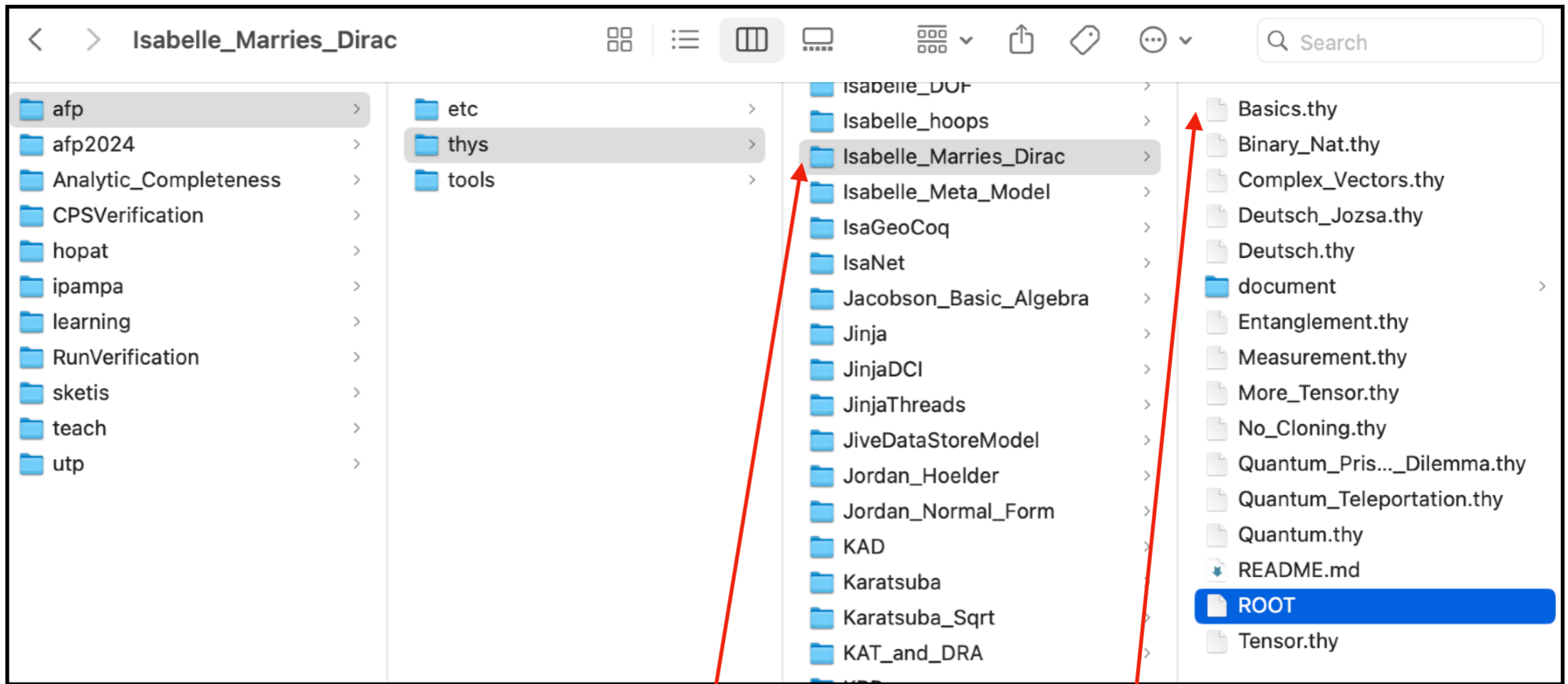


* plans to remove scala-isabelle in the future

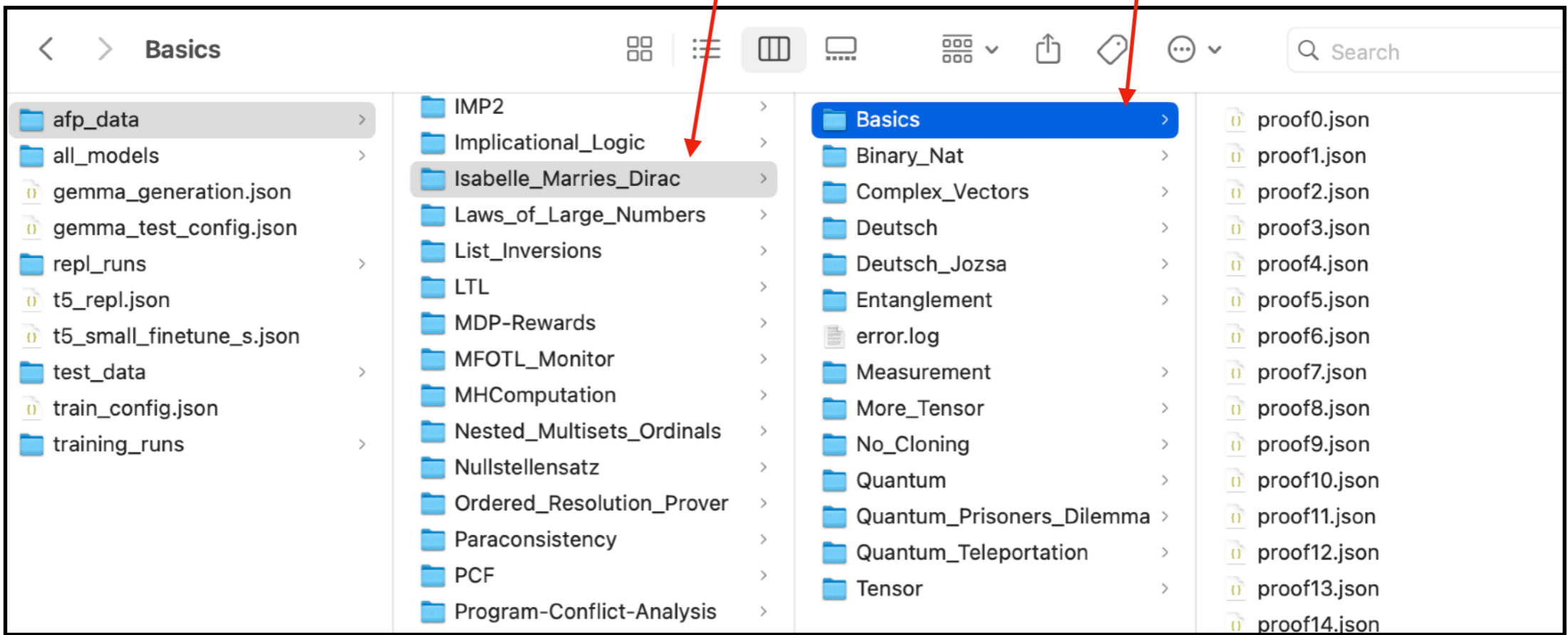
Data Extraction Algorithm

Data extraction process

AFP



Generated Data



Content of the generated data

Contents:

- user action
- user state
- Isabelle state/term
- variables
- constants
- *type constants*
- *type variables*
- apply-keywords
- isar-keywords
- proof methods
- lemma dependencies

Easily extensible to include:

- classes
- locales
- definitions
- conjecturing
- Sledgehammer premises

Programs > deepIsaHOL > all_data > afp_data > Isabelle_Marries_Dirac > Basics > {} proof0.

```
1 {
2   "proof": {
3     "steps": [
4       {
5         "step": {
6           "action": "lemma set_2_atLeast0 [simp]: \"{0..
7             <2>::nat} = {0,1}\"",
8           "user_state": "",
9           "term": "{0..<2} = {0, 1} \\<Longrightarrow> ({0..
10            <2} = {0, 1})\"",
11          "hyps": [
12            ],
13          "proven": [
14            ],
15          "variables": [
```

{} proof0.json ×

> deepIsaHOL > all_data > afp_data > Isabelle_Marries_Dirac > Basics > {} proof0.json > ...

```
2   "proof": {
241     "deps": [
389       {"thm": {"name": "Code_Numeral.
arity_unique_euclidean_semiring_natural", "term":
"OFCLASS(natural, unique_euclidean_semiring_class)
"}},
390       {"thm": {"name": "Code_Numeral.
arity_discrete_linordered_semimod_natural", "term":
"OFCLASS(natural, discrete_linordered_semimod_class)
"}},
391       {"thm": {"name": "Quickcheck_Narrowing.
arity_type_narrowing_type_IITN_narrowing_type",
"term": "OFCLASS(narrowing_type.
narrowing_type_IITN_narrowing_type, type_class)"}]}
392     ]
393   }
394 }
```

How to use it?

Simple interface for data extraction either as a single command:

```
sbt "run /path/to/a/read/directory/ /path/to/a/writing/directory/"
```

or via the Scala or Python REPL:

```
import isabelle_rl._
import de.unruh.isabelle.control.Isabelle

val logic = "LTL" // Isabelle session name

val writer = new Writer(
  "/path/to/a/read/directory/", // it can have a ROOT, or a .thy file
  "/path/to/a/writing/directory/",
  logic)

writer.write_data("your/file.thy")

writer.write_all()

val minion = writer.get_minion()

implicit val isabelle: Isabelle = minion.isabelle

isabelle.destroy()
```

Read-eval-print-loops (REPL)

How to use it?

Start a Py4j server running a scala-isabelle process with
`sbt "runMain isabelle_rl.Py4j_Gateway_Main" `

```
>>> import sys
>>> sys.path.append('/path/to/this/project/src/main/python')
>>> from repl import REPL

>>> repl = REPL("HOL")
REPL and minion initialized.

>>> repl.apply("lemma \"\\<forall>x. P x \\<Longrightarrow> P c\"")
'proof (prove) goal (1 subgoal):  1. \\<forall>x. P x \\<Longrightarrow> P c'

>>> repl.apply("apply simp")
'proof (prove) goal:  No subgoals!'

>>> repl.apply("done")
''

>>> repl.shutdown_gateway() # do not forget to close the Isabelle process
```

Small-scale experiment

Experimental setup

- Use small language models to predict the **next user action** in Isabelle proofs
- Text-to-text transfer transformer (T5) language models (encoder-decoders) requiring the SentencePiece tokenizer:
 - **small** (77 million parameters)
 - **base** (248 million parameters)
- To compare:
 - training from **scratch** vs **fine-tuning** them (T5 models are pre-trained with question-answering datasets on math *word* problems, coding problems for popular programming languages, and various English-language novels)
 - benefits of extra data: proof up to a given and user state (**s**) vs augmenting to that the premises to prove the theorem and the keywords available at that point (**spk**)

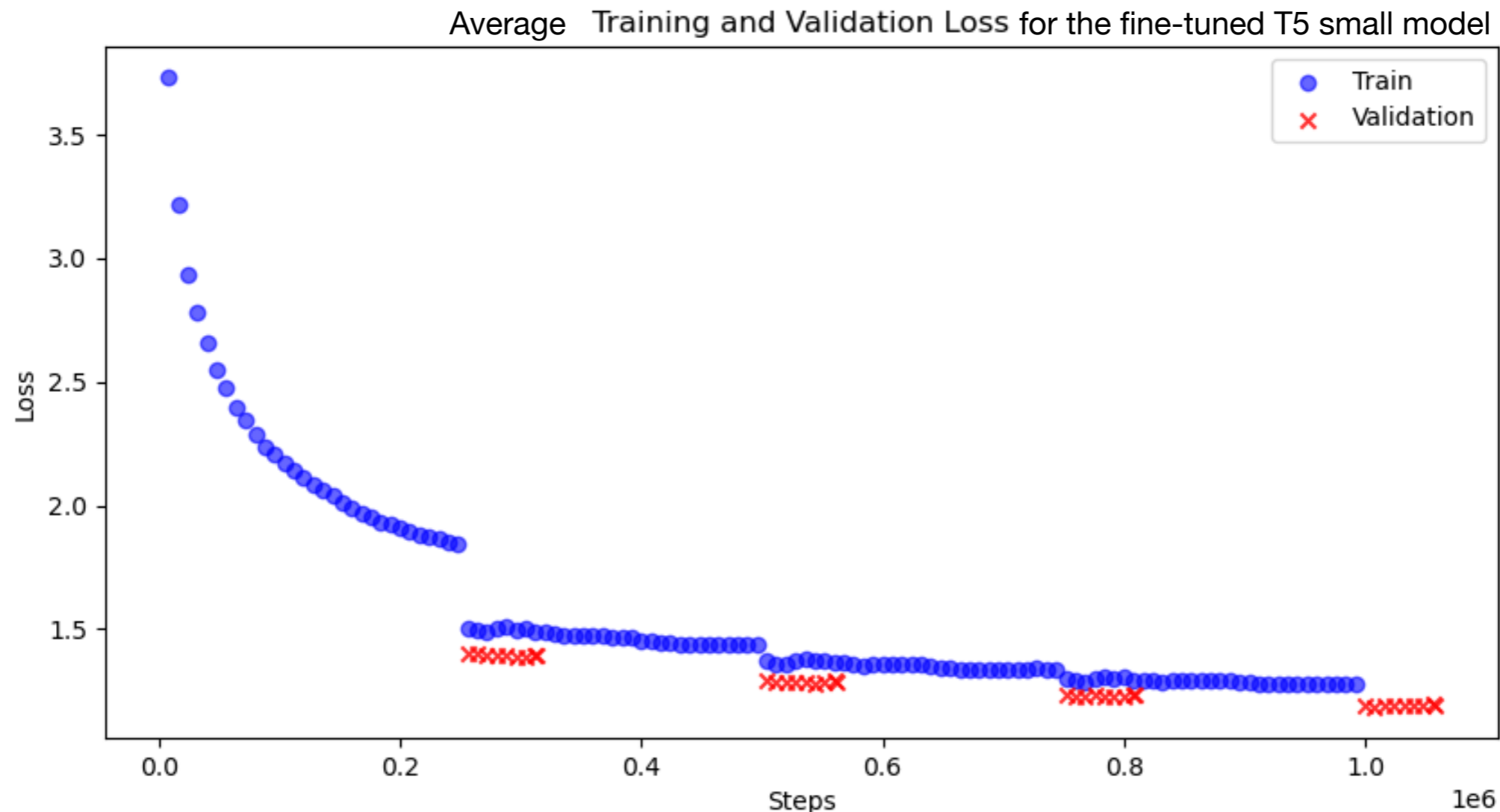
Experimental setup

AFP official numbers
~296,000 lemmas
926 AFP entries

2025 data
198,108 proofs (~66.93%)
922 AFP entries (99.57%)
1,555,653 training proof steps

2024 data
187,210 proofs (~63.24%)
861 AFP entries (92.98%)
1,854,901 training proof steps

64% of each .thy file is reserved for *training*, 16% for *validation*, and 20% for *testing*



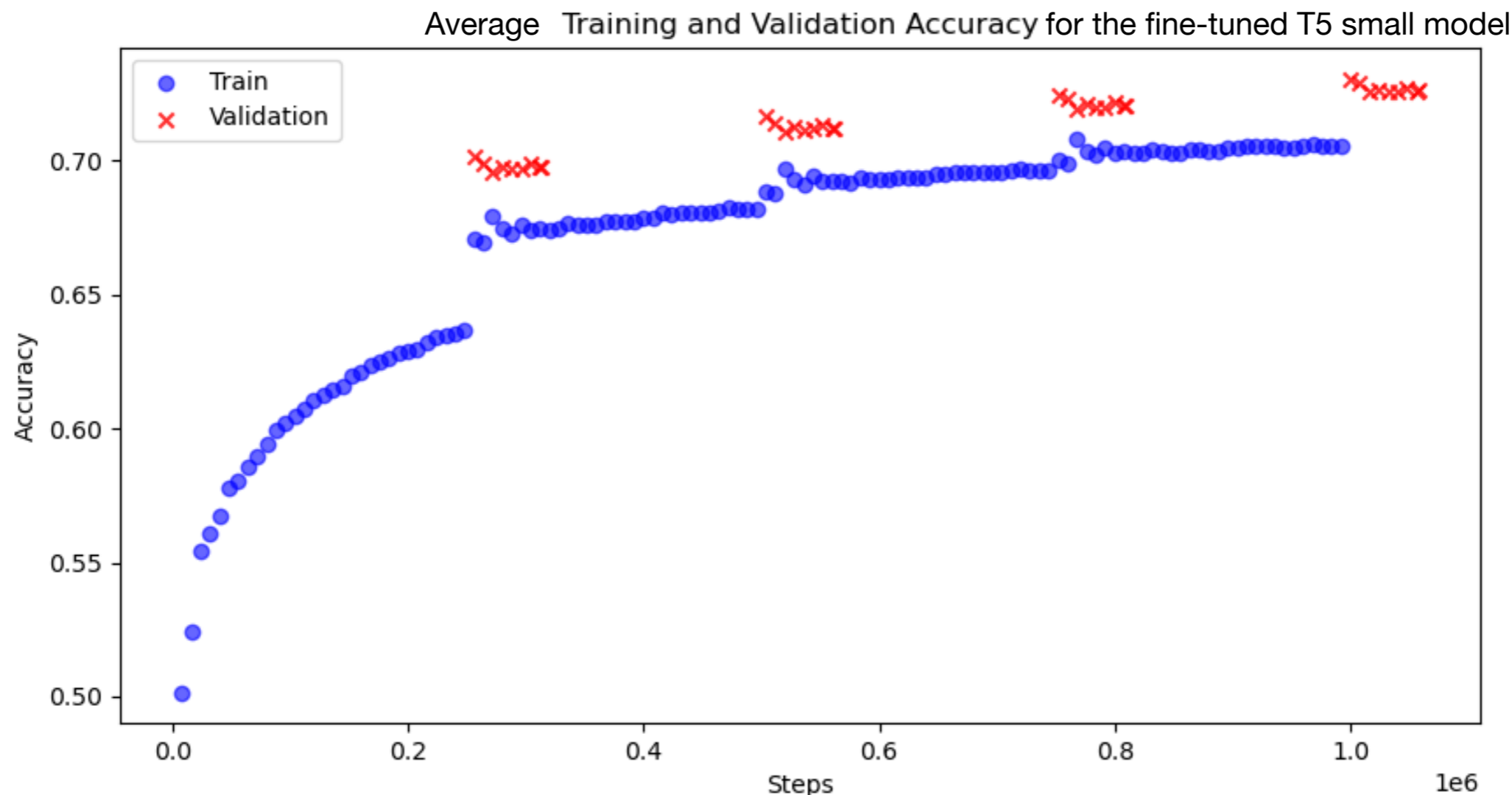
Experimental setup

AFP official numbers
~296,000 lemmas
926 AFP entries

2025 data
198,108 proofs (~66.93%)
922 AFP entries (99.57%)
1,555,653 training proof steps

2024 data
187,210 proofs (~63.24%)
861 AFP entries (92.98%)
1,854,901 training proof steps

64% of each .thy file is reserved for *training*, 16% for *validation*, and 80% for *testing*



Loss, accuracy, and other metrics (as LLMs)

Model	Train Loss	Train Acc.	Valid Loss	Valid Acc.
t5_small_s (4ep)	3.02	0.445	3.13	0.436
t5_small_finetune_s (4ep)	1.18	0.720	1.14	0.733
t5_small_spk (3ep)	2.96	0.417	3.29	0.375
t5_small_spk_trim (4ep)	2.71	0.438	3.21	0.371
t5_base_s (4ep)	2.84	0.474	3.06	0.441

Model	Test Exact Acc.	Test 1st-Tok Acc.	Test All Wrong
t5_small_s (4ep)	0.137	0.316	0.303
t5_small_finetune_s (4ep)	0.150	0.302	0.231
t5_small_spk (2ep)	0.042	0.254	0.553
t5_small_spk_trim (4ep)	0.022	0.194	0.371
t5_base_s (4ep)	0.170	0.307	0.299

Model performance metrics after 4 epochs over the training data set
(except for the spk which was only trained over 3 epochs)

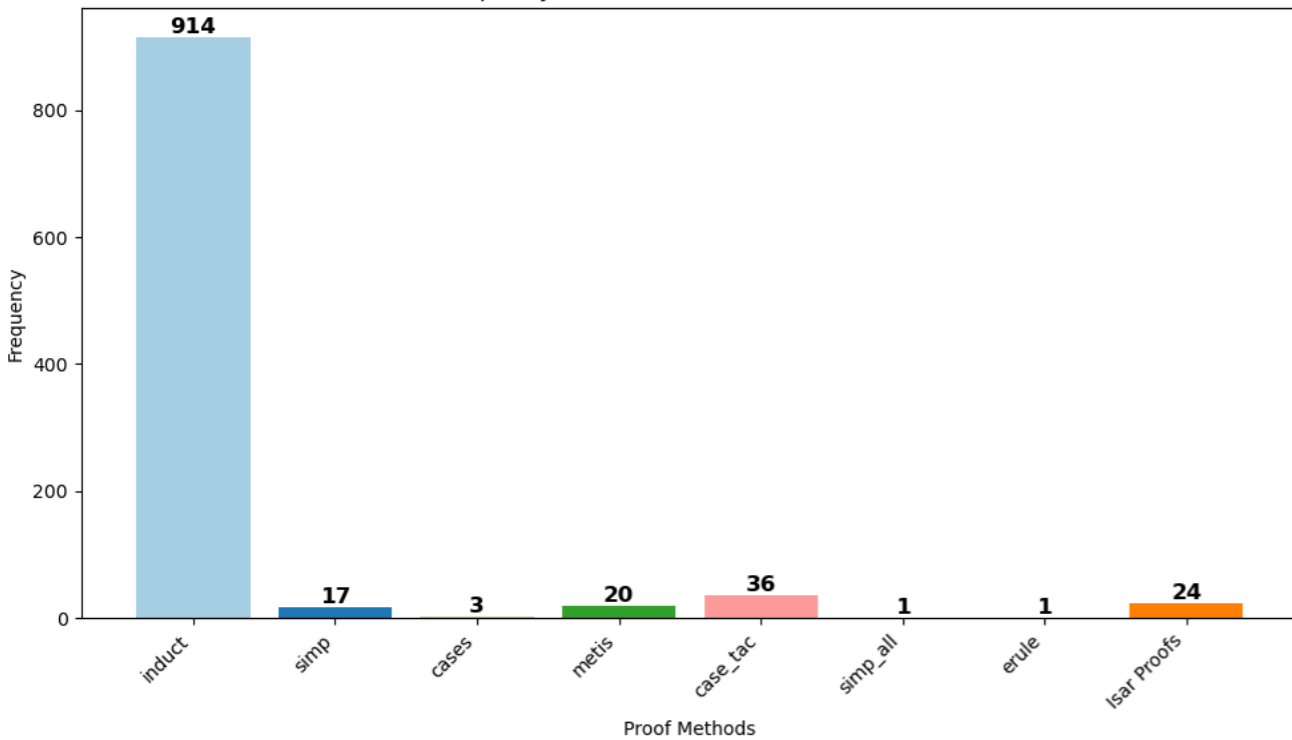
3. Training of T5 LLMs for simple experiments

Let the models suggest the next user-action (5 opportunities) in a depth-first search (DFS) fashion for at most 5 consecutive user actions

Model	Predictions	Progress	by commands	Correct by	Proofs	Finished proofs	Avg. time
small_s	7,426,213	61.58 %	187,683	42.02 %	43,538	7.13 %	48.92 s
small_fs	3,508,777	37.79 %	162,986	8.40 %	42,408	11.28 %	40.09 s
small_spk	2,939,773	48.82 %	55,195	14.76 %	43,218	8.81 %	29.84 s
base_s	5,955,871	59.45 %	101,268	25.36 %	43,001	9.51 %	50.33 s
base_fspk	2,315,016	52.57 %	118,959	36.70 %	38,532	16.71 %	53.28 s

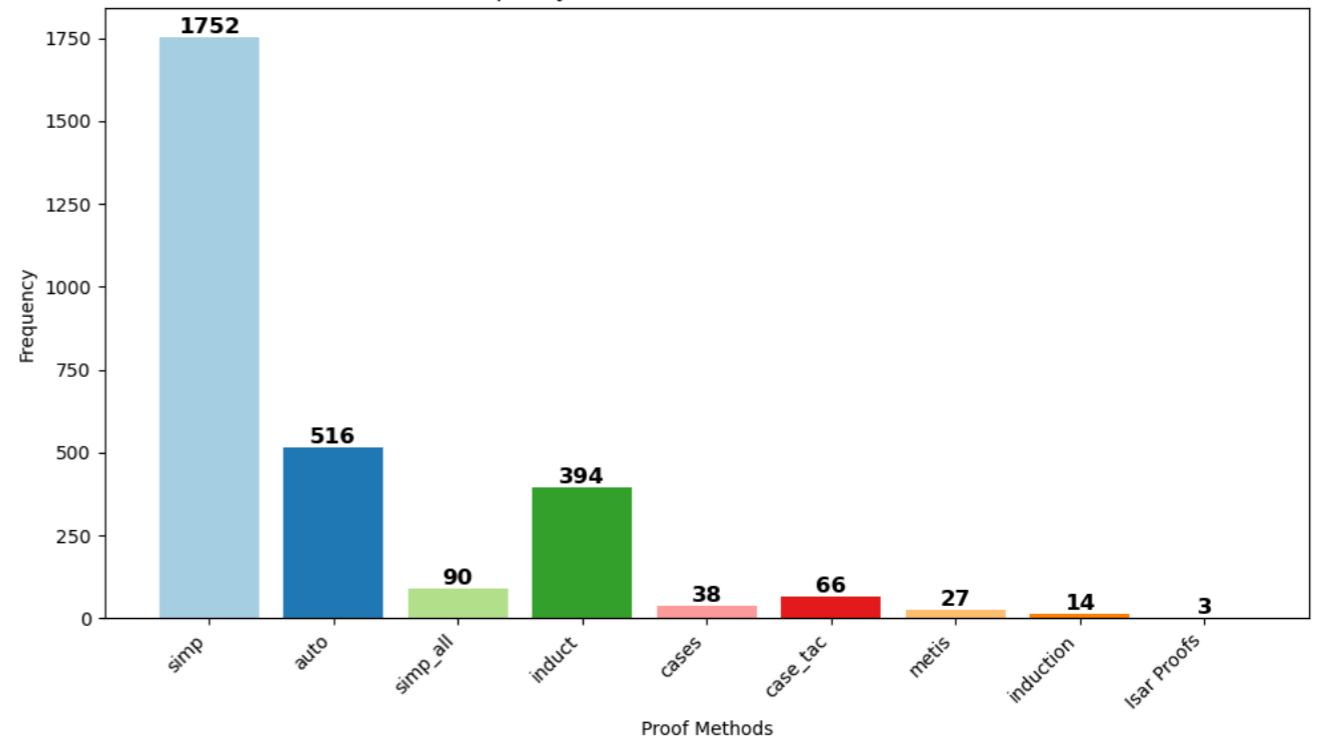
Extras on the small-scale experiment 3

Frequency of Proof Methods in Isabelle Proofs



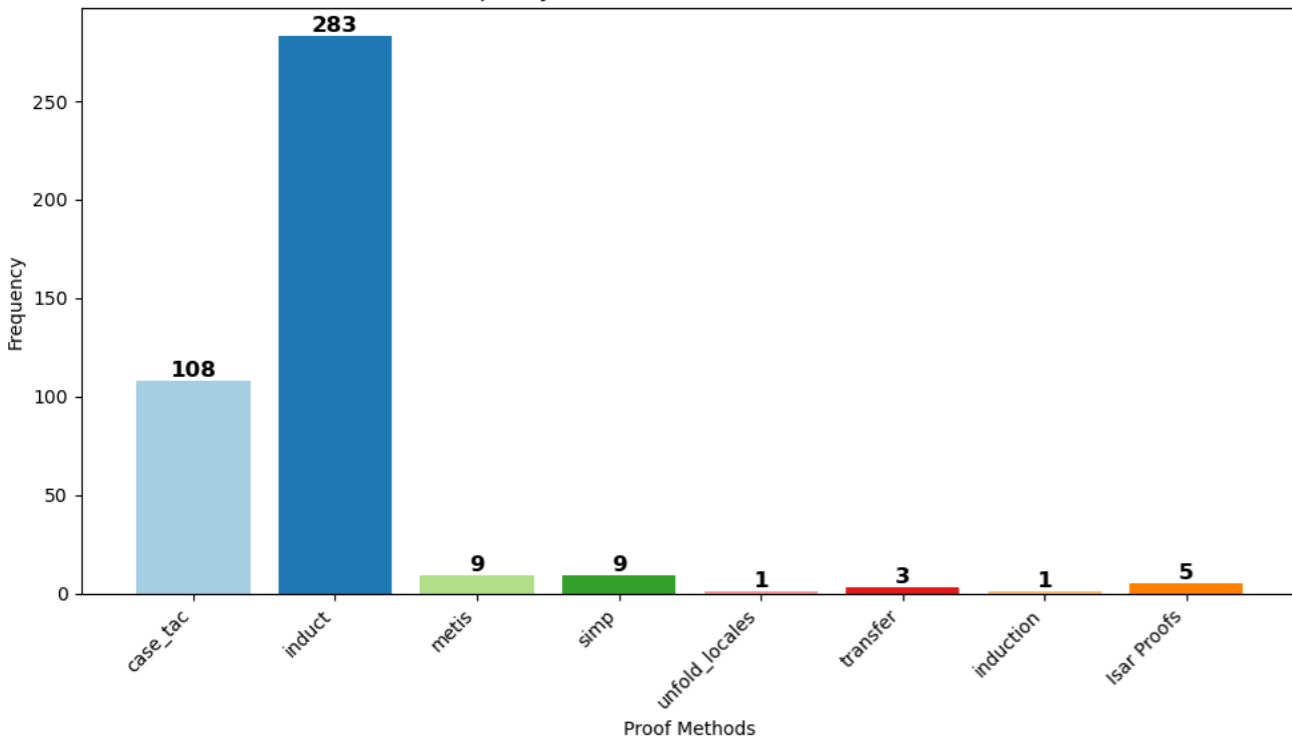
t5_small (s)

Frequency of Proof Methods in Isabelle Proofs



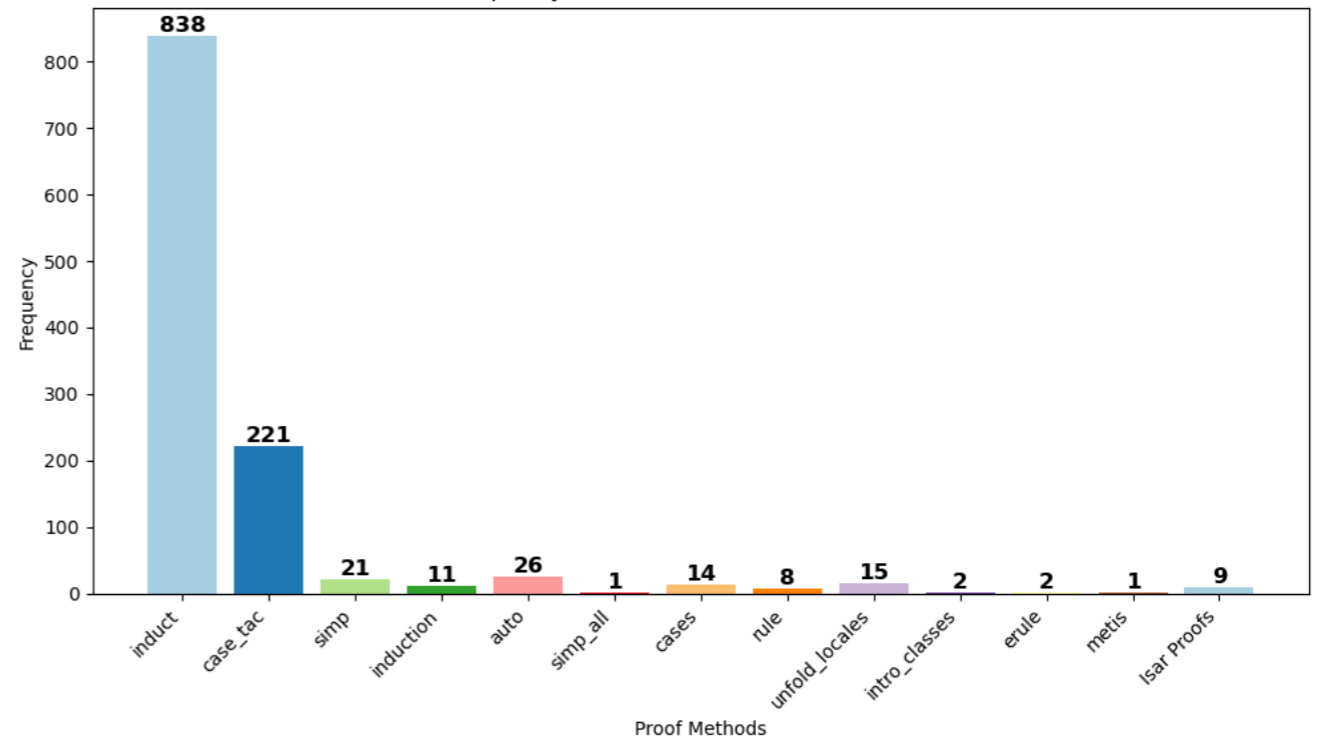
t5_small (fs)

Frequency of Proof Methods in Isabelle Proofs



t5_small (spk)

Frequency of Proof Methods in Isabelle Proofs

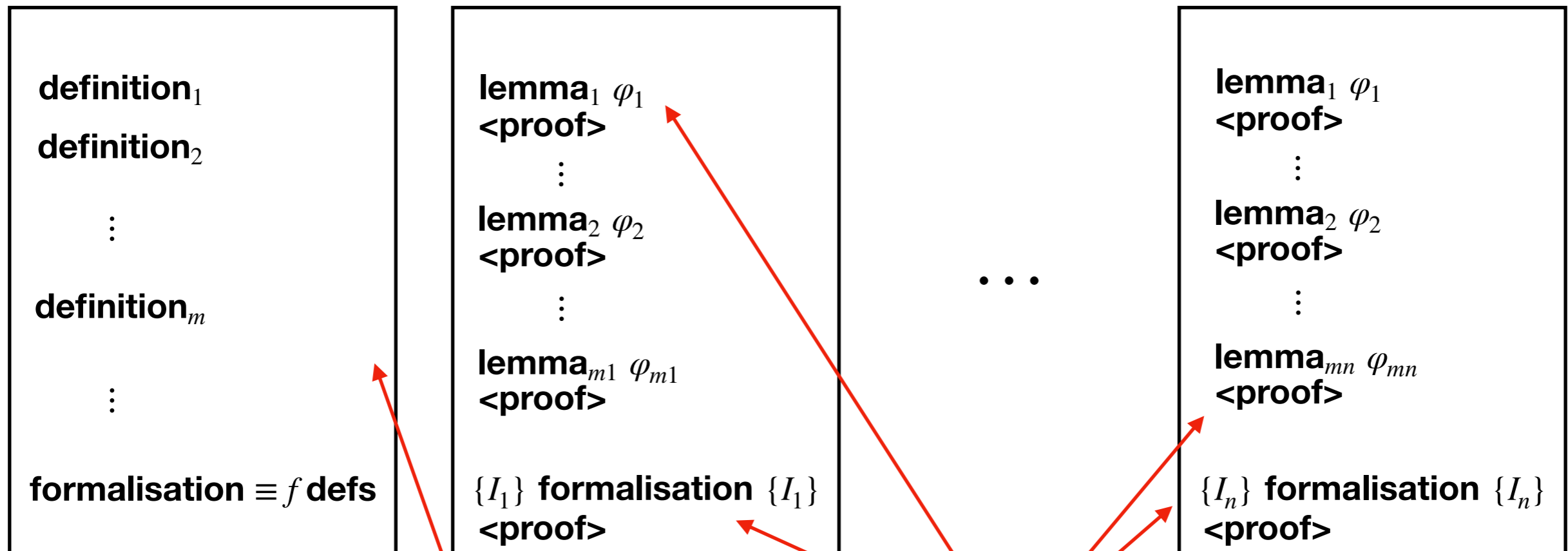


t5_base (s)

Proof-fixing tool

General problem

- Prove the partial correctness of a protocol: $\{ \mathbf{inv} \}$ **formalisation** $\{ \mathbf{inv} \}$
- The invariant is really a huge conjunction: $\mathbf{inv} \triangleq I_1 \wedge I_2 \wedge I_3 \wedge \dots \wedge I_n$
- Each (minor) edition in the formalisation, breaks downstream proof progress:

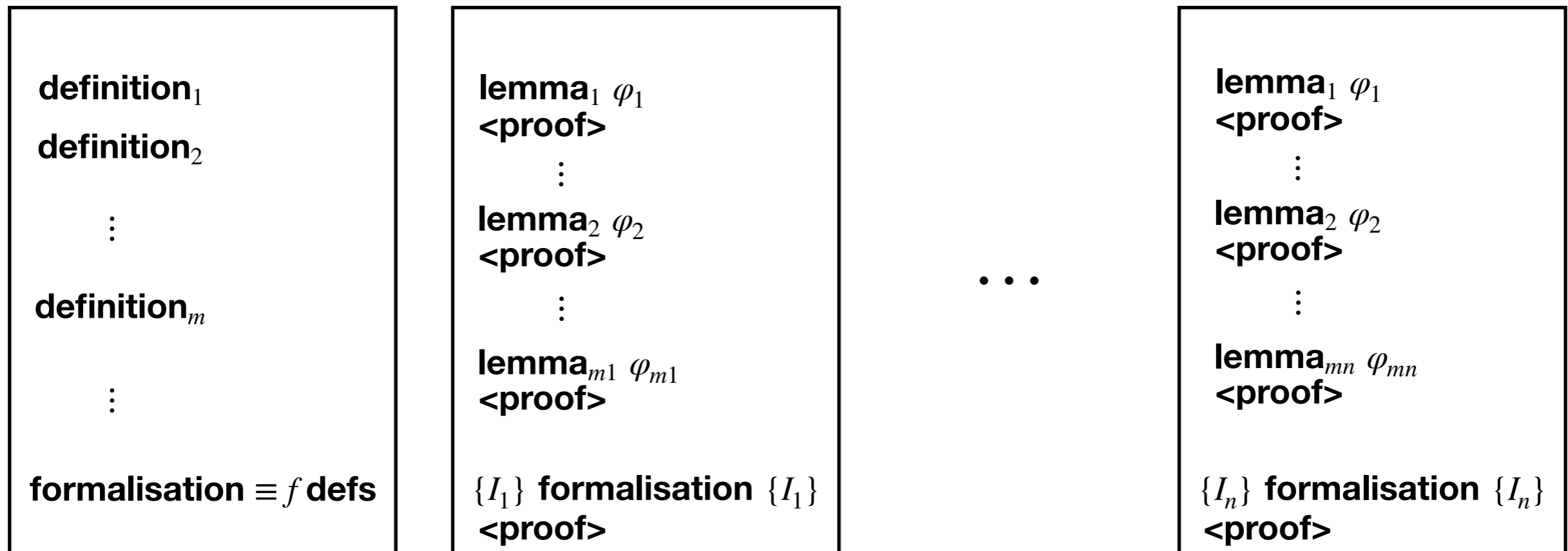


(minor) change here
(e.g. adding an invariant)

produces (trivial) errors there
(solved by calling Sledgehammer)

General solution

super_fix: traverses your .thy file, it fixes misaligned proof obligations and replaces failed (or infinitely looping) proof methods with sorry's. Then, it traverses the .thy file again and calls Sledgehammer where each sorry appears. If Sledgehammer finds a proof, it replaces the sorry with that proof. Otherwise, it leaves the sorry to the user.



this helped complete the large scale proof of a cache-coherence protocol!

Conclusion

Next steps

- Call an external tool via an Isabelle tactic
- Continue the small scale experiment
 - Compare against Sledgehammer
 - Variation of the model architecture (from encoder-decoder to autoencoder)
 - Usage of the latest optimised models (Gemma or DeepSeek-R1)
 - Train on small-scale premise selection
 - Train on “conjecturing” predicting the next state
 - Variations on the proof-exploration tree (best-first search, MCTS)
- Complete the infrastructure for reinforcement learning experiments
- Focus on cleverly engineered premise selection to improve the models’ performance

Previous ML approaches on ITPs

Prover\Strategy	Premise selection	Reinforcement Learning	Next step prediction for proof completion	Autoformalisation
HOL4	TacticToe (kNN)	Tactic Zero	TacticToe (kNN)	
Coq (now Rocq)			CoqPilot, Tactician (kNN, GNN, and LLM)	
Isabelle	Magnushammer		Thor, DeepIsaHOL	Draft-Sketch-Prove, LEGO-Prover
Lean			LeanDojo's ReProver	Deep Seek Prover
HOL Light		<i>HOL-list (RL)</i>	<i>HOL-list (GNN)</i>	

Final invitation

DeepIlsaHOL is an open project intended to serve the machine learning community and the ITP community.

Everyone is invited to collaborate, use it, and provide feedback

Jonathan Julián Huerta y Munive
huertjon@cvut.cz
Czech Technical University in Prague



<https://github.com/yonoteam/DeepIlsaHOL>

Extras

Points of discussion

1. why did the project use the scala-isabelle library instead of the protocol between Isabelle/ML and Isabelle/Scala?
 - Understanding the protocol required expert knowledge while the scala-isabelle library was straightforward to understand.
 - Yet, there is an ongoing conversation with the Isabelle developers. The next Isabelle release includes an extension that will enable us to turn the data-generation algorithm (and other of the project's tools) into a more standard Isabelle tool.
2. why did the project not use Isabelle/dump?
 - Isabelle/dump does not provide all the required information with the relationships that the project needs (e.g. theorem to its proof-steps to their constants, variables, ...)
3. why did the project not compare with Sledgehammer Isabelle/mirabelle?
 - For evaluation against Sledgehammer, we are going to.
4. why did the project not use Lean/Coq/HOL4/Mizar/Metamath/...?
 - Isabelle chosen because of
 - its closeness to natural language via its scripting language Isar
 - its automation via Sledgehammer
 - its large repository of proofs
 - its growing number of users

Previous ML approaches on ITPs: TacticToe

k-nearest neighbours (kNN) and Monte-Carlo tree search (MCTS) to predict next proof-step
 Proves **66.4%** of 7164 top-level theorems from HOL4 standard library

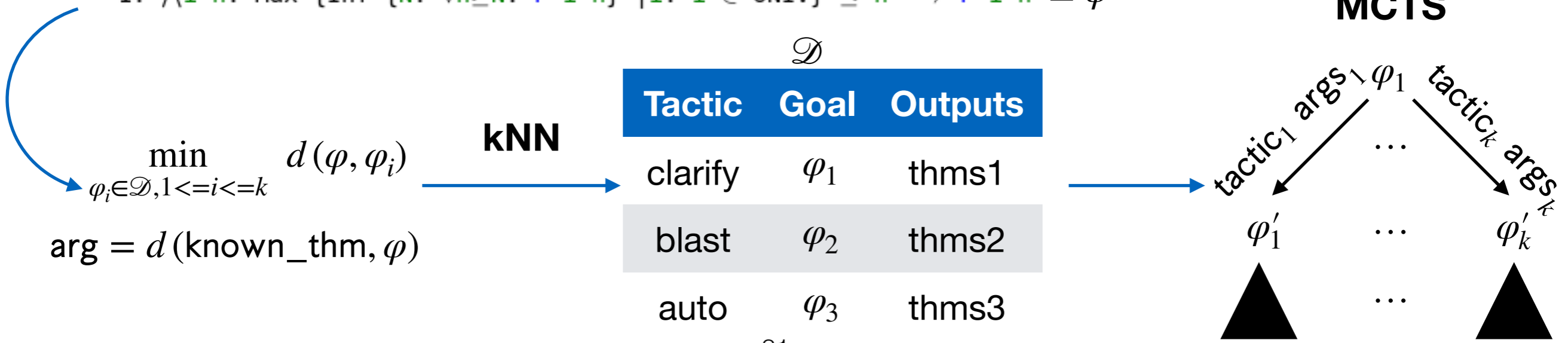
lemma

```

fixes P :: "('a::finite)  $\Rightarrow$  nat  $\Rightarrow$  bool"
assumes " $\forall i. \exists N::nat. \forall n \geq N. P i n$ "
shows " $\exists N. \forall i. \forall n \geq N. P i n$ "
apply (rule_tac x="Max {Inf {N.  $\forall n \geq N. P i n$ } | i::'a. i  $\in$  UNIV}" in exI)
apply clarify
subgoal for i n
  using wellorder_Inf_lemma[of " $\lambda n. P i n$ "] assms
  Max_ge[of "{Inf {N.  $\forall n \geq N. P i n$ } | i::'a. i  $\in$  UNIV}", OF finite_image]
  apply -
  apply (erule_tac x=i in allE)
  apply (subgoal_tac " $\forall n \geq \text{Inf } \{N. \forall n \geq N. P i n\}. P i n$ ")
  apply (metis (mono_tags, lifting) dual_order.trans iso_tuple_UNIV_I mem_Collect_eq)
  apply blast
  done
done
    
```

```

proof (prove)
goal (1 subgoal):
1.  $\bigwedge i n. \text{Max } \{ \text{Inf } \{ N. \forall n \geq N. P i n \} \mid i. i \in \text{UNIV} \} \leq n \implies P i n \equiv \varphi$ 
    
```



Previous ML approaches on ITPs: “Draft, Sketch, Prove” and Magnushammer

Isabelle

Formal proof
sketching

Informal to
LaTeX

L^AT_EX

```
lemma
  fixes P :: "('a::finite) ⇒ nat ⇒ bool"
  assumes "∀i. ∃N::nat. ∀n ≥ N. P i n"
  shows "∃N. ∀i. ∀n ≥ N. P i n"
proof-
  let ?bound i = "Inf {N . ∀n ≥ N. P i n}"
  let ?N = "Max {?bound i | i. i ∈ UNIV}"
  {fix n::nat and i::'a
    from assms
    obtain M where "∀n ≥ M. P i n"
      \<proof>
    hence obs: "∀ m ≥ ?bound i. P i m"
      \<proof>
    assume "n ≥ ?N"
    also have "?N ≥ ?bound i"
      \<proof>
    ultimately have "n ≥ ?bound i"
      \<proof>
    hence "P i n"
      \<proof>
  }
  thus "∃N. ∀i n. N ≤ n → P i n"
    \<proof>
qed
```

Call hammers to
complete the proofs

Since ATPs only serve for
premise selection, replace
them with neural networks (NNs)

Combined approach: Proof-rate of 71% on publicly
[available benchmark](#) of 183K Isabelle proofs from the AFP

Other previous ML approaches on ITPs

- **TacticZero for HOL4:** It is an RL approach using NNs. It achieves **~49% success** on a test set of 268 problems.
- **Tactician for Coq (now Rocq):** It is similar to TacticToe but extends the approach with a graph neural network (GNN) and compares it with the kNN.
 - It compares them with an LLM trained only on Coq data.
 - Intersection between kNN and GNN is thin.
 - GNN almost fully covers CoqHammer's results while
 - GNN proves some that CoqHammer does not prove.
- **HOL-list environment for HOL Light:** It is a light-weight tool for testing different approaches for tactic and premise selection with a Python interface.
 - It has tested GNNs and an RL algorithm.
 - The RL approach solves **56.3%** of a set of 3225 theorems of the HOL Light's mathematical library.
- **LeanDojo for Lean:** Open project that interfaces Python with the Lean ITP.
- **More LLMs in Lean:** Many recent articles on using LLMs for autoformalization, premise selection, tactic suggestion, and proof search.

Literature review lessons

Lessons learned

- Main approaches discussed:
 - ML (kNN or GNN or LLM) for predicting next tactic + premises
 - LLM for autoformalisation + hammer for premise selection
 - Medium LM as a replacement for hammer as premise selection
 - Reinforcement learning replacing the standard best-first search exploration
- Language models alone have not sufficed for full automation in theorem proving
- Premise selection is a very important aid for small and medium sized language models in these projects (takes them as far as 50-60% on already curated proofs)
- The approach of scaling the model to a LARGE language model for memorising proof structures and autoformalising also helps
- More work is needed to:
 - Combine complementary approaches to improve proof success rate
 - Compare these approaches in terms of resources needed (training time, testing time, computing power) and across provers
 - Create tools for users (only know of TacticToe, Tactician, LeanCopilot)

Isabelle tactic calling an external tool