

Tao’s Equational Proof Challenge Accepted

Lydia Kondylidou¹, Jasmin Blanchette¹, and Marijn J.H. Heule²

¹ Ludwig-Maximilians-Universität München, Munich, Germany
{l.kondylidou,jasmin.blanchette}@lmu.de

² Carnegie Mellon University, Pittsburgh, United States
marijn@cmu.edu

Abstract. In the context of the Equational Theories Project, Terence Tao posed the challenge of finding alternatives to a complicated 62-step proof found by the Vampire superposition prover. We introduce a proof minimization tool called Krympa. Using a combination of brute force and heuristics, and exploiting both Vampire and the Twee equational prover, the tool reduces the 62-step proof to 20 steps, each corresponding to a rewrite. In an empirical evaluation, it also performs well on 1431 equational problems originating from the same project, reducing in particular a 151-step proof to only 10 steps.

Keywords: Theorem provers · Equational logic · Proof minimization.

1 Introduction

The Equational Theories Project [7], launched in September 2024 by Fields medalist Terence Tao, aims at exploring the relations between different equational theories of magmas. A *magma* is a basic algebraic structure consisting of a set equipped with a single binary operation \diamond closed on that set. The project’s first phase, concluded in April 2025, focused on equational laws for magmas that contain at most four applications of \diamond .

The project uses the Lean [23] proof assistant to express proofs and counter-examples but depends on automatic theorem provers and other external tools. The problems explored in the project’s first phase all fall within first-order logic’s unit equality fragment: They consist of a universally quantified equation as the sole axiom and a universally quantified equation as the proof goal, or conjecture.

For the problem 650 \implies 448, where 650 denotes the axiom $\forall x, y, z. x = x \diamond (y \diamond ((z \diamond x) \diamond y))$ and 448 denotes the conjecture $\forall x, y, z. x = x \diamond (y \diamond (z \diamond (x \diamond z)))$, the Vampire [5] superposition prover found a particularly complex proof, with 62 inference steps, excluding clausification and Skolemization. Given that the proof is unintelligible, Tao challenged the community to find “an alternate proof, by whatever means you wish—human, semi-automated, or automated” [29].

One idea could be to run a specialized equational prover, Twee [25], instead of Vampire, but this results in a very long, 137-step proof. Another approach would be to use Lean’s automation, such as the `aesop` [21], `canonical` [24], `duper` [9], and `grind` [1] tactics and the LeanHammer [31], to reconstruct and compress

consecutive superposition steps, in the style of Sledgehammer’s structured proof reconstruction [6, Sect. 6.3]. This would yield a shorter and more high-level proof, in which each step may combine multiple rewrites. Our approach is orthogonal. Our working hypothesis is that Vampire’s 62-step proof, which emerged as the byproduct of a saturation process, is likely suboptimal. By mixing and matching proofs generated by different automatic provers, as proposed by Sutcliffe et al. [28], we hope to achieve a shorter, simpler proof.

We introduce Krympa, a tool that minimizes equational proofs by decomposing them into independently provable components and reassembling them into more concise, intelligible proofs. Specifically, starting from a Vampire-generated proof, the tool transforms it into a direct proof (Sect. 3) and analyzes its inferences to break it down into intermediate results that serve as candidate lemmas. Each of these lemmas is then proved independently using Vampire and Twee (Sect. 4), the two leading systems in the unit equality division of CASC 2025 [27]. The resulting proofs are then combined into a single proof using heuristics that favor shorter derivations (Sect. 5).

Given the 62-step Vampire proof of $650 \implies 448$, our tool produces a 20-step proof, where 13 steps are generated by Twee (Sect. 6). In a larger empirical evaluation, we applied the tool to 1431 provable implications from the Equational Theories Project and obtained positive results (Sect. 7). In particular, the tool reduced a 151-step Vampire proof to 10 steps.

Our tool is implemented in Rust, OCaml, and Python. Its source code is available at <https://github.com/kondylidou/Krympa>. The files associated with Tao’s challenge and our empirical evaluation data are also available online [20].

2 Background

We briefly review the Vampire and Twee automatic provers and their associated proof formats.

2.1 Vampire and Superposition Proofs

Vampire is a saturation-based theorem prover for first-order logic with equality based on the superposition calculus [4]. It implements highly optimized search strategies and data structures, and integrates techniques such as literal selection, term orders, redundancy elimination, strategy scheduling, and portfolios.

Superposition works on implicitly universally quantified clauses. A preprocessor performs clausification and Skolemization. For example, the axiom $\forall x. f(x) = g(x)$ is transformed into $f(x) = g(x)$, where x is a free variable, and the conjecture $\forall x. f(x) = g(x)$ is negated and transformed into $f(\text{sk}) \neq g(\text{sk})$, where sk is a Skolem constant. The objective is to derive the contradictory clause \perp . For the unit equality fragment, the calculus’s two relevant inference rules are as follows:

$$\frac{t \neq u}{\perp} \text{equality resolution} \qquad \frac{t = t' \quad s[u] \bowtie s'}{\mu(s[t'] \bowtie s')} \text{superposition}$$

77 The equality resolution rule has one premise, $t \neq u$, one conclusion, \perp , and
 78 one side condition: that t and u are unifiable. The superposition rule has two
 79 premises and one conclusion. The \bowtie symbol denotes either $=$ or \neq throughout
 80 the rule. The $=$ and \neq operators are commutative; for example, the premise $t = t'$
 81 can match the equation $f(\mathbf{a}) = \mathbf{b}$ either as is or as $\mathbf{b} = f(\mathbf{a})$. The premises are
 82 assumed to have disjoint sets of variables, which can be achieved by renaming.
 83 Also in the rule, $s[\]$ is a term with a hole, the terms $s[u]$ and $s[t']$ are obtained
 84 by filling the hole in $s[\]$ with u and t' , and μ is a most general unifier of t
 85 and u . For example, the most general unifier of the terms $h(\mathbf{a}, y)$ and $h(x, \mathbf{b})$ is
 86 $\{x \mapsto \mathbf{a}, y \mapsto \mathbf{b}\}$; applying it on both terms yields $h(\mathbf{a}, \mathbf{b})$. Finally, the rule has
 87 further side conditions, not shown here, that restrict the search space.

88 **Example 1.** A subtle case of the superposition rule arises when both premises
 89 are the same clause. Consider the following rule instance, where the variable in
 90 the second premise has been renamed to avoid a clash:

$$\frac{f(f(x)) = g(x) \quad f(f(x')) \neq g(x')}{f(g(x)) \neq g(f(x))} \text{superposition}$$

91 This instance is obtained by taking $t := f(f(x))$, $t' := g(x)$, $\bowtie := \neq$, $s[\] := f([\])$,
 92 $u := f(x')$, $s' := g(x')$, and $\mu = \{x' \mapsto f(x)\}$. Applying the unifier μ to both
 93 premises yields the equation $f(f(x)) = g(x)$ and the disequation $f(f(f(x))) \neq$
 94 $g(f(x))$. The inference replaces the subterm $f(f(x))$ in the disequation with $g(x)$
 95 using the equation as a left-to-right rewrite rule, and derives the conclusion. ■

96 **Example 2.** Vampire implements *parallel superposition*, a variant of the super-
 97 position rule in which multiple subterms that match a term are replaced. The
 98 following inference illustrates this:

$$\frac{b = a \quad h(b, a, b) \neq h(a, b, a)}{h(a, a, a) \neq h(a, a, a)} \text{parallel superposition}$$

99
 100 Superposition proofs are expressed in a linear format. They are refutational
 101 and show how to derive \perp from the input axioms and the negated conjecture.

102 **Example 3.** The following is a superposition proof from clauses 1–3:

- | | | | |
|-----|---------------------------------|--|---|
| | 1. $a = b$ | axiom | |
| | 2. $f(x) = x$ | axiom | |
| | 3. $h(f(b), a) \neq h(a, f(b))$ | negated conjecture | |
| 103 | 4. $h(b, a) \neq h(a, b)$ | by parallel superposition from 2 and 3 | |
| | 5. $h(a, a) \neq h(a, a)$ | by parallel superposition from 1 and 4 | |
| 104 | 6. \perp | by equality resolution from 5 | ■ |

105 2.2 Twee and Structured Equational Chain Proofs

106 Twee is an automatic prover specialized for equational reasoning. It is based on
 107 the unfailing completion procedure [3], an extension of Knuth–Bendix completion
 108 [19]. As with superposition, quantifiers are eliminated by a preprocessor. In the
 109 DISCOUNT and Waldmeister tradition [8], Twee's proofs are structured as a

110 sequence of lemmas, where the lemmas and the conjecture are proved by chains
 111 of equalities. Twee introduces lemmas if they are needed more than once. Twee
 112 proofs are arguably more readable than Vampire proofs.

113 **Example 4.** The following is a Twee-style proof of goal 1 from axioms 1 and 2:

114 Axiom 1: $a = b$ Axiom 2: $f(x) = x$ Lemma 3: $f(b) = a$ Proof: $f(b)$ $= \{ \text{by axiom 1 right-to-left} \}$ $f(a)$ $= \{ \text{by axiom 2} \}$ a	Goal 1: $h(f(b), a) = h(a, f(b))$ Proof: $h(f(b), a)$ $= \{ \text{by lemma 3} \}$ $h(a, a)$ $= \{ \text{by lemma 3 right-to-left} \}$ $h(a, f(b))$
---	--

■

115 3 Conversion to Direct Proofs

116 Vampire generates proofs by refutation, whereas our mix-and-match approach
 117 requires direct proofs. To bridge this gap, we transform Vampire proofs into
 118 direct proofs. In the following sections, we will always use direct proofs.

119 In equational logic, to produce a direct proof, we first introduce existential
 120 quantifiers for Skolem constants and universal quantifiers for variables. For ex-
 121 ample, $h(x, sk) \neq x$ is transformed into $\exists z. \forall x. h(x, z) \neq x$. Then we apply
 122 the contrapositive to all inferences in which a premise and the conclusion are
 123 disequations to obtain positive equations. Thus, the inference

$$\frac{h(a, y) = b \quad h(x, sk) \neq x}{b \neq a} \text{superposition}$$

124 becomes

$$\frac{\forall y. h(a, y) = b \quad b = a}{\forall z. \exists x. h(x, z) = x}$$

125 Equality resolution inferences from a premise $t \neq t$ are omitted since their con-
 126 trapositives derive trivial equations.

127 **Example 5.** The following is a direct proof obtained from Example 3's proof
 128 by refutation.

1. $a = b$	axiom
2. $\forall x. f(x) = x$	axiom
129 3. $h(b, a) = h(a, b)$	from 1 and $h(a, a) = h(a, a)$
130 4. $h(f(b), a) = h(a, f(b))$	from 2 and 3

■

131 4 Proof Generation for Intermediate Lemmas

132 Our approach starts by translating the main theorem into a TPTP [13] input
 133 problem and running Vampire to produce a baseline proof. This proof is turned

134 into a direct proof, then decomposed into intermediate lemmas. For each lemma,
 135 we generate corresponding problems, with the objective of proving them using
 136 Vampire and Twee. Three problem variants are generated:

- 137 1. *Big-step problems* contain the axioms together with the lemma as the con-
 138 jecture, and nothing else. This allows us to investigate whether a radically
 139 new proof, with different intermediate steps, can be found.
- 140 2. *Small-step problems* contain the axioms together with the lemma as the con-
 141 jecture, and all lemmas derived prior to this lemma in the baseline proof as
 142 additional axioms. This allows us to investigate whether a somewhat similar
 143 variant of the original derivation can be found.
- 144 3. *Abstracted problems* are variants of big-step problems that contain the ax-
 145 ioms together with an abstracted version of the lemma as the conjecture.
 146 Specifically, selected subterms of the lemma—for example, expressions such
 147 as $x \diamond y$ that do not contain nested applications—are replaced by fresh vari-
 148 ables. This allows us to investigate whether a more general version of the
 149 lemma is provable, ideally with a shorter, more abstract proof.

150 Each problem is submitted to the two provers. If a proof is found for a small-step
 151 problem, we expand it to recursively include the shortest proofs of the lemmas
 152 used as axioms for the axioms referenced in the proof. Ties are broken arbitrarily.
 153 Note that abstracted problems might be unprovable.

154 Next, we compare the proofs of the three problem variants corresponding to
 155 the same lemma. If the abstracted problem has the shortest proof, the lemma
 156 it proves is replaced in all small-step problems where it appears as an axiom
 157 with the generalized lemma from the abstracted problem. Each updated small-
 158 step problem is then re-proved, and if the result has fewer steps, we replace the
 159 small-step problem's proof with it.

160 The length of a Vampire-generated proof is the number of steps of its direct
 161 proof, excluding preprocessing. For Twee, the length of a proof is the cumu-
 162 lative number of equalities in the equality chains. Thus, the Vampire proof in
 163 Example 5 has two steps, and the Twee proof in Example 4 has four steps.

164 5 Proof Construction for the Main Theorem

165 Based on the intermediate lemmas' proofs generated in the previous phase, our
 166 approach constructs a proof of the main theorem. The proof generally consists
 167 of three segments. The first segment starts with the axioms and ends with the
 168 derivation of a so-called *departure lemma*. The second segment derives a so-called
 169 *arrival lemma*. The third segment derives the conjecture. Different candidates
 170 are considered as the departure and arrival lemmas, yielding different proofs.
 171 The proof with the fewest steps is chosen.

172 Specifically, we first identify up to six intermediate lemmas that arise close to
 173 the end of the baseline proof, including the conjecture, and consider them as can-
 174 didate arrival lemmas. For each of these, we consider its transitive dependencies

175 as candidate departure lemmas. Then, for each candidate departure lemma, we
 176 construct a problem with the axioms and the departure lemma’s dependencies
 177 as the axioms and the departure lemma itself as the conjecture. We run both
 178 provers and, if at least one succeeds, we use the shorter result as the proof of the
 179 first segment, unless an even shorter proof was generated in the previous phase.

180 Next, for each pair of candidate departure and arrival lemmas, we generate
 181 a new problem with the original axioms, the departure lemma, and its depen-
 182 dencies as axioms and the arrival lemma as the conjecture. We run both provers
 183 and, if at least one succeeds, we use the shorter result as the proof of the second
 184 segment, unless an even shorter proof was generated earlier. Finally, we generate
 185 a new problem with the original axioms, the departure lemma, its dependencies,
 186 and the arrival lemma as axioms and the original conjecture as the conjecture.
 187 We run both provers and, if at least one succeeds, we use the shorter result as the
 188 proof of the third segment, unless an even shorter proof was generated earlier.

189 Without the segmentation, proof minimization could be intractable due to
 190 combinatorial explosion. We chose to work with three segments, and six candi-
 191 date arrival lemmas, as a trade-off between performance and flexibility.

192 **Example 6.** Before we review the three-segment proof construction approach
 193 in detail, let us look at an example. The following sketch represents a baseline
 194 seven-step Vampire direct proof of a theorem $A \implies C$:

195	A	axiom
196	L_1	from A and A
197	L_2	from A and L_1
198	L_3	from L_1 and L_2
199	L_4	from L_2 and L_3
200	L_5	from L_3 and L_4
201	L_6	from A and L_5
202	C	from L_5 and L_6

203 Here, A denotes the axiom, and L_1, \dots, L_6 are the lemmas used to derive the
 204 conjecture C .

205 In the first phase, for each lemma L_1, \dots, L_6 , we construct big-step, small-
 206 step, and abstracted problems and try to prove them using Vampire and Twee,
 207 retaining the shortest proof for each lemma. Suppose the following: The shortest
 208 proof of L_1 has one step and is obtained from its big-step problem using Vampire;
 209 for L_2 and L_3 , the shortest proofs are obtained from their small-step problems
 210 using Twee; for L_4 , the shortest proof is obtained from its abstracted problem
 211 using Twee; and for L_5 and L_6 , the shortest proofs are obtained from their
 212 small-step problems using Vampire.

213 In the next phase, the last five lemmas, L_2, \dots, L_6 , and the conjecture C are
 214 considered as candidate arrival lemmas. We focus on L_6 . The proof below, found
 215 by Vampire for L_6 ’s small-step problem, is the shortest proof for L_6 :

216	A	axiom
217	L_1	from A and A
218	L_2	from A and L_1

219 L_3 from L_1 and L_2
 220 L_4 from L_2 and L_3
 221 L_5 from L_3 and L_4
 222 L_6 from A and L_5

223 This proof happens to be identical to the first six steps of the baseline proof,
 224 but in general it could be different.

225 Next, lemmas L_1 to L_5 are considered as candidate departure lemmas. We
 226 focus on L_3 . The proof of conjecture C is constructed by concatenating three
 227 segments. For the first segment, we create a new problem with A , L_1 , and L_2
 228 as axioms, since they are dependencies of the departure lemma L_3 in the above
 229 proof of L_6 , and L_3 as the conjecture. We run both provers on this problem
 230 and obtain a two-step Vampire proof of L_3 from A , L_1 , and a new lemma L'_2 .
 231 Since L_1 is treated as an axiom, we must include its proof to obtain a complete
 232 proof of L_3 . In the first phase, we found a one-step Vampire proof of L_1 from
 233 the axiom A , so we use it. In summary, the proofs of L_1 and L_3 form the first
 234 segment, which consists of one step for L_1 and two steps for L_3 .

235 For the second segment, we create a new problem with A , L_1 , L'_2 , and L_3 as
 236 axioms and the arrival lemma L_6 as the conjecture. We run both provers on this
 237 problem and obtain a two-step Twee proof of L_6 from L_1 and L_3 . Together with
 238 the first segment, this yields a five-step proof of L_6 . Since this proof is shorter
 239 than the six-step proof of L_6 presented above, it is used as the second segment.

240 For the third segment, we create a new problem with A , L_1 , L'_2 , the departure
 241 lemma L_3 , and the arrival lemma L_6 as axioms and C as the conjecture. We run
 242 both provers on this problem and obtain a two-step Twee proof of C from L'_2
 243 and L_3 . Since this proof does not use the arrival lemma L_6 , the second segment
 244 is excluded from the result. Concatenating the first and third segments yields a
 245 new five-step proof of C , which is two steps shorter than the baseline proof:

246 A axiom
 247 L_1 from A
 248 L'_2 from A and L_1
 249 L_3 from L_1 and L'_2
 250 C by a two-step equality chain using L'_2 and L_3

251 Finally, other combinations of candidate departure and arrival lemmas are
 252 considered, and the shortest proof is retained. ■

253 5.1 Construction of the Dependency Graph

254 We identify lemmas occurring close to the end of the derivation as candidate
 255 arrival lemmas. Different candidates typically depend on substantially different
 256 subsets of earlier lemmas. Each candidate therefore induces its own dependency
 257 chain, and different choices can lead to substantially different proof lengths.
 258 We consider six candidate arrival lemmas extracted from the baseline proof,
 259 including the conjecture itself, since our approach may produce a shorter proof
 260 of the conjecture by reproving it directly from a minimized dependency set.

261 For every candidate, we build a dependency graph that captures the lemmas
 262 required to derive it. Dependencies are determined from the shortest Vampire
 263 or Twee proof obtained for each lemma. Given that we generate three problem
 264 variants and run two provers, up to six proofs per lemma are considered. A
 265 lemma ℓ is considered to directly depend on a lemma ℓ' if the shortest proof
 266 of ℓ uses ℓ' as an axiom. Thus, for big-step and abstracted problems, only the
 267 original axioms can be dependencies. For small-step problems, each intermediate
 268 step in a Vampire proof and each lemma in a Twee proof is considered a lemma.

269 The dependency graph associated with a candidate arrival lemma is a di-
 270 rected acyclic graph (DAG) whose nodes correspond to lemmas and whose edges
 271 express derivability between them. Formally, let V be a finite set of lemmas, each
 272 represented by an equation and a set of dependencies on other lemmas. We con-
 273 struct a DAG (V, E) , where each vertex $\ell \in V$ corresponds to a lemma and each
 274 edge $(\ell, \ell') \in E$ indicates that lemma ℓ directly depends on lemma ℓ' . As an
 275 optimization, we merge lemmas that are identical up to the naming of variables,
 276 keeping the shortest proof.

277 5.2 Construction of the First Proof Segment

278 For each candidate arrival lemma, we investigate whether all lemmas included in
 279 its dependency graph are needed to derive it or whether a shorter proof can be
 280 obtained by choosing a departure lemma and recomputing parts of the derivation
 281 by combining proofs generated by the provers.

282 As candidate departure lemmas, we consider all lemmas in the DAG. Let ℓ
 283 be a candidate departure lemma. If ℓ depends only on the axioms, we take the
 284 shortest big-step, small-step, or abstracted proof previously found by Vampire or
 285 Twee. Otherwise, we build a problem that includes ℓ 's dependencies in the DAG
 286 as axioms and the departure lemma as the conjecture, and we run Vampire and
 287 Twee. If at least one of them succeeds, we choose the shorter proof as ℓ 's proof.
 288 This derivation, together with the shortest proofs of ℓ 's dependencies generated
 289 for the big-step, small-step, or abstracted problems, forms the first segment of
 290 the final proof. However, if we found an even shorter proof for the big-step, small-
 291 step, or abstracted problem, we use that proof instead. For small-step proofs, we
 292 must also include the proofs of the intermediate lemmas encoded as axioms.

293 5.3 Construction of the Remaining Proof Segments

294 To construct the second segment, we generate a problem with the departure
 295 lemma and its dependencies as axioms and the arrival lemma as the conjecture,
 296 and run both provers. If at least one of them succeeds, we choose the shorter
 297 proof as the proof of the arrival lemma. As above, we fall back on the proof of
 298 a big-step, small-step, or abstracted problem if it is even shorter.

299 Finally, to construct the third segment, we generate a problem with the
 300 departure lemma, the arrival lemma, and their dependencies as axioms and the
 301 original conjecture as the conjecture, and invoke both provers. If at least one of

302 them succeeds, we choose the shorter proof as the proof of the original conjecture.
 303 As above, we fall back on a previously derived proof if it is even shorter.

304 The final proof is obtained by concatenating the three segments. The proof
 305 might contain unreferenced lemmas; these are pruned.

306 5.4 Proof Output

307 Our tool generates the minimized proof in a native format, from which two Lean
 308 outputs are produced. The first Lean output is a step-by-step formalization using
 309 the `calc` tactic to reconstruct chains of equalities. It applies the `duper` tactic
 310 to fill in the subproofs. For example, a proof of $t_1 = t_2 = t_3 = t_4$ would be
 311 represented by

```
312 calc
313   t1 = t2 := by duper ...
314   _ = t3 := by duper ...
315   _ = t4 := by duper ...
```

316 where the ellipses stand for `duper`'s arguments. The second Lean output is a more
 317 compact Lean formalization in which each lemma is proved directly using Lean's
 318 automation without including the intermediate steps in chains of equalities.

319 6 Application to Tao's Challenge

320 We implemented our approach and tried the resulting tool, Krympa, on Tao's
 321 challenge theorem 650 \implies 448:

$$(\forall x, y, z. x = x \diamond (y \diamond ((z \diamond x) \diamond y))) \implies \forall x, y, z. x = x \diamond (y \diamond (z \diamond (x \diamond z))).$$

322 Our tool first ran Vampire to obtain a baseline 62-step superposition proof. Then
 323 it constructed 62 problems of each variant (big-step, small-step, and abstracted)
 324 and tried to prove them using Vampire and Twee. Among the six candidate
 325 arrival lemmas, the shortest proof was found by selecting

$$\forall x, y, z. x = x \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x). \quad (\text{lemma 9})$$

326 The coloring highlights repeating patterns. Next, our tool constructed the de-
 327 pendency graph for this lemma. The DAG contained 37 lemmas. It was based
 328 on big- and small-step proofs.

329 Among the 37 candidate departure lemmas, our tool found the shortest proof
 330 by selecting

$$\forall x, y, z, w. (x \diamond ((y \diamond x) \diamond x)) \diamond z = \\ ((x \diamond ((y \diamond x) \diamond x)) \diamond z) \diamond (w \diamond ((x \diamond ((y \diamond x) \diamond x)) \diamond w)). \quad (\text{lemma 7})$$

331 According to the DAG, the shortest proof of this lemma was found by running
 332 Vampire on the small-step problem consisting of the axiom and the following
 333 lemma dependencies:

$$\forall x, y, z, w. x \diamond ((y \diamond z) \diamond x) = (x \diamond ((y \diamond z) \diamond x)) \diamond (w \diamond (z \diamond w)) \quad (\text{lemma 1})$$

$$\forall x, y, z, w, v, u. x \diamond ((y \diamond ((z \diamond w) \diamond y)) \diamond x) = (x \diamond ((y \diamond ((z \diamond w) \diamond y)) \diamond x)) \diamond (v \diamond ((u \diamond (w \diamond u)) \diamond v)) \quad (\text{lemma 2})$$

$$\forall x, y, z, w, v. x \diamond (y \diamond x) = (x \diamond (y \diamond x)) \diamond (z \diamond ((w \diamond ((v \diamond y) \diamond w)) \diamond z)) \quad (\text{lemma 3})$$

$$\forall x, y, z, w, v. x \diamond (y \diamond x) = (x \diamond (y \diamond x)) \diamond ((z \diamond (y \diamond z)) \diamond (w \diamond ((v \diamond y) \diamond w))) \quad (\text{lemma 4})$$

$$\forall x, y, z, w. x \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x) = (x \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x)) \diamond (w \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond w)) \quad (\text{lemma 5})$$

$$\forall x, y, z, w. (x \diamond ((y \diamond x) \diamond x)) \diamond z = ((x \diamond ((y \diamond x) \diamond x)) \diamond z) \diamond ((w \diamond ((x \diamond ((y \diamond x) \diamond x)) \diamond w)) \diamond (z \diamond ((x \diamond ((y \diamond x) \diamond x)) \diamond z))). \quad (\text{lemma 6})$$

334 Following the inference steps of the baseline Vampire proof, our tool derived
 335 lemma 1 by applying a superposition inference with the axiom $x = x \diamond (y \diamond ((z \diamond$
 336 $x) \diamond y))$ as the first premise and a renamed copy $x' = x' \diamond (y' \diamond ((z' \diamond x') \diamond y'))$ as
 337 the second premise. The most general unifier of the first premise's right-hand side
 338 and the subterm $z' \diamond x'$ of the second premise is $\{x' \mapsto y \diamond ((z \diamond x) \diamond y), z' \mapsto x\}$.
 339 Applying the unifier to both premises yields the equations $x = x \diamond (y \diamond ((z \diamond x) \diamond$
 340 $y))$ and $y \diamond ((z \diamond x) \diamond y) = (y \diamond ((z \diamond x) \diamond y)) \diamond (y' \diamond ((x \diamond (y \diamond ((z \diamond x) \diamond y))) \diamond y'))$.
 341 The superposition inference replaced the subterm $x \diamond (y \diamond ((z \diamond x) \diamond y))$ in the
 342 second equation with x using the first equation as a right-to-left rewrite rule,
 343 and thus derived lemma 1, up to the naming of variables. Lemmas 2 to 7 were
 344 derived similarly following the steps of the baseline Vampire proof.

345 Next, from the axiom and lemma 7, our tool proved the arrival lemma
 346 (lemma 9) using Twee. For this proof, Twee introduced the intermediate step

$$\forall x, y, z, w. (y \diamond ((z \diamond y) \diamond y)) \diamond w = ((y \diamond ((z \diamond y) \diamond y)) \diamond w) \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x). \quad (\text{lemma 8})$$

347 Finally, assuming all the lemmas derived so far, our tool proved the conjecture
 348 from lemmas 5 and 9 using Twee. The resulting proof has 20 steps, including
 349 three Twee-generated chains of equalities.

350 Below we present the final proof adapted from our tool's detailed Lean out-
 351 put. Instead of relying on proof automation, we use the `nth_rw` tactic, which
 352 performs a single rewrite step, where the numeric argument indicates which
 353 matching occurrence should be rewritten. In one case, two numbers are sup-
 354 plied, corresponding to a parallel rewrite.

```
355 theorem Equation650_implies_Equation448 (G : Type _) [Magma G]
356   (op_law : ∀ x y z : G, x = x ◊ (y ◊ ((z ◊ x) ◊ y))) :
357   ∀ x y z : G, x = x ◊ (y ◊ (z ◊ (x ◊ z))) :=
```

```

358 have lemma1 (x y z w : G) :
359   x ◊ ((y ◊ z) ◊ x) = (x ◊ ((y ◊ z) ◊ x)) ◊ (w ◊ (z ◊ w)) := by
360   nth_rw 3 [op_law z x y]
361   exact op_law (x ◊ ((y ◊ z) ◊ x)) w z
362
363 have lemma2 (x y z w v u : G) :
364   x ◊ ((y ◊ ((z ◊ w) ◊ y)) ◊ x) =
365   (x ◊ ((y ◊ ((z ◊ w) ◊ y)) ◊ x)) ◊ (v ◊ ((u ◊ (w ◊ u)) ◊ v)) := by
366   nth_rw 1 2 [lemma1 y z w u]
367   exact lemma1 x (y ◊ ((z ◊ w) ◊ y)) (u ◊ (w ◊ u)) v
368
369 have lemma3 (x y z w v : G) :
370   x ◊ (y ◊ x) = (x ◊ (y ◊ x)) ◊ (z ◊ ((w ◊ ((v ◊ y) ◊ w)) ◊ z)) := by
371   nth_rw 1 [lemma1 w v y x]
372   exact op_law (x ◊ (y ◊ x)) z (w ◊ ((v ◊ y) ◊ w))
373
374 have lemma4 (x y z w v : G) :
375   x ◊ (y ◊ x) = (x ◊ (y ◊ x)) ◊ ((z ◊ (y ◊ z)) ◊ (w ◊ ((v ◊ y) ◊ w))) := by
376   nth_rw 1 [lemma1 w v y z]
377   exact lemma3 x y (z ◊ (y ◊ z)) w v
378
379 have lemma5 (x y z w : G) :
380   x ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x) =
381   (x ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x)) ◊ (w ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ w)) := by
382   nth_rw 1 [lemma2 w y z y x ((z ◊ y) ◊ y)]
383   exact lemma4 x (y ◊ ((z ◊ y) ◊ y)) w x ((z ◊ y) ◊ y)
384
385 have lemma6 (x y z w : G) :
386   (x ◊ ((y ◊ x) ◊ x)) ◊ z =
387   ((x ◊ ((y ◊ x) ◊ x)) ◊ z) ◊ ((w ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ w)) ◊
388   (z ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ z))) := by
389   nth_rw 1 [lemma5 z x y w]
390   exact op_law ((x ◊ ((y ◊ x) ◊ x)) ◊ z) (w ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ w)) z
391
392 have lemma7 (x y z w : G) :
393   (x ◊ ((y ◊ x) ◊ x)) ◊ z =
394   ((x ◊ ((y ◊ x) ◊ x)) ◊ z) ◊ (w ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ w)) := by
395   nth_rw 1 [lemma5 w x y z]
396   exact lemma6 x y z w
397
398 have lemma8 (x y z w : G) :
399   ((x ◊ ((y ◊ x) ◊ x)) ◊ z) ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ w) =
400   (x ◊ ((y ◊ x) ◊ x)) ◊ z := by
401   let T := x ◊ ((y ◊ x) ◊ x)
402   calc
403     (T ◊ z) ◊ (T ◊ w) =
404     ((T ◊ z) ◊ ((T ◊ w) ◊ ((T ◊ (T ◊ w)) ◊ ((w ◊ (T ◊ w)) ◊ (T ◊ (T ◊ w)))))) := by
405     nth_rw 1 [←op_law]
406     - = ((T ◊ z) ◊ ((T ◊ w) ◊ ((T ◊ (T ◊ w)) ◊ ((w ◊ (T ◊ w)) ◊
407       (T ◊ ((T ◊ w) ◊ (w ◊ (T ◊ w))))))) := by

```

```

408     nth_rw 1 [←lemma7]
409   _ = ((T ◊ z) ◊ ((T ◊ w) ◊ ((T ◊ (T ◊ w)) ◊ ((w ◊ (T ◊ w)) ◊
410     ((T ◊ ((T ◊ w) ◊ (w ◊ (T ◊ w)))) ◊ ((T ◊ w) ◊ (w ◊ (T ◊ w))) ◊
411     (T ◊ ((T ◊ w) ◊ (w ◊ (T ◊ w))))))))) := by
412     nth_rw 2 [←lemma7]
413   _ = ((T ◊ z) ◊ ((T ◊ w) ◊ ((T ◊ (T ◊ w)) ◊ (w ◊ (T ◊ w)))) := by
414     nth_rw 1 [←op_law]
415   _ = ((T ◊ z) ◊ ((T ◊ w) ◊ (T ◊ (T ◊ w)))) := by
416     nth_rw 1 [←lemma7]
417   _ = ((x ◊ ((y ◊ x) ◊ x)) ◊ z) := by
418     nth_rw 1 [←lemma7]
419
420 have lemma9 (x y z : G) :
421   (x ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x)) = x := by
422   calc
423     (x ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x)) =
424     (x ◊ (((y ◊ ((z ◊ y) ◊ y)) ◊ x) ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x))) := by
425     nth_rw 1 [lemma8]
426   _ = (x ◊ (((y ◊ ((z ◊ y) ◊ y)) ◊ x) ◊ ((y ◊ ((z ◊ y) ◊ y)) ◊ x) ◊
427     ((y ◊ ((z ◊ y) ◊ y)) ◊ x))) := by
428     nth_rw 2 [lemma8]
429   _ = x := by
430     nth_rw 1 [←op_law]
431
432 show _ by
433   intros x y z
434   calc
435     x = x ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ x) := by
436     nth_rw 1 [lemma9]
437   _ = (x ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊ x) ◊ ((y ◊ (z ◊ (x ◊ z))) ◊
438     ((x ◊ ((y ◊ x) ◊ x)) ◊ (y ◊ (z ◊ (x ◊ z))))) := by
439     nth_rw 1 [←lemma5]
440   _ = x ◊ ((y ◊ (z ◊ (x ◊ z))) ◊ ((x ◊ ((y ◊ x) ◊ x)) ◊
441     (y ◊ (z ◊ (x ◊ z)))) := by
442     nth_rw 1 [lemma9]
443   _ = x ◊ (y ◊ (z ◊ (x ◊ z))) := by
444     nth_rw 1 [lemma9]

```

7 Experiments on Other Equational Proofs

446 To assess the general potential of our approach, we evaluated our tool on a set of
447 equational theorems obtained from the Equational Theories Project repository
448 [7]. We selected all problems in the 13 Lean files `Proofs1` to `Proofs13` that have
449 a proof and translated them to TPTP problem files, yielding 1431 benchmarks.

450 For each file, we invoked our tool's TPTP problem generator, which parses
451 the Lean theorems and produces corresponding TPTP problem files. For each
452 problem, our tool was given 2700 seconds to produce a minimized proof using
453 Vampire to find the baseline proof and Vampire and Twee to find subproofs; on

454 failure, the baseline Vampire proof was output. A time limit of 10 seconds was
 455 used for each prover invocation. The experiments were conducted on a server
 456 equipped with a dual-socket AMD EPYC 9965 system (384 cores, 768 threads)
 457 running at 2.25–3.70 GHz with 3 TiB of DDR5 ECC RAM, and running Debian
 458 GNU/Linux 13 (kernel 6.17.13+deb13-amd64).

459 Overall, proofs for the 13 Lean files have an average length of 6.6 steps before
 460 minimization and 4.5 steps after minimization using the combination of small-
 461 step and abstracted problems and both provers. This corresponds to a 31.5%
 462 decrease, showing that even short proofs can often be made shorter.

463 Since longer proofs present more opportunities for minimization, we now
 464 focus on the 117 benchmarks whose baseline proofs have at least 15 steps. Table 1
 465 compares proof lengths before and after minimization. The “Avg. before” column
 466 shows the average number of steps in the baseline proofs. The remaining columns
 467 report the average proof length after minimization under four configurations,
 468 which differ in which problem variants are used: “BA” denotes the combination
 469 of the big-step and abstracted variants; “SA” denotes the combination of the
 470 small-step and abstracted variants; “BS” denotes the combination of the big- and
 471 small-step variants; and “BSA” denotes the combination of all three variants.

472 The results show an often substantial reduction in proof length. SA generally
 473 yielded the shortest proofs. Across all problems for the 13 Lean files, the average
 474 reduction with SA is 56.7%. BS and BSA also produced substantial reductions,
 475 whereas BA generally yielded the least improvements.

476 It might seem counterintuitive that SA, which does not consider big-step
 477 problems, outperforms BSA. However, the nonmonotonicity is to be expected.

Table 1. Comparison of proof lengths before and after minimization for problems with baseline proofs of at least 15 steps

File	Num. problems	Avg. before	Avg. after			
			BA	SA	BS	BSA
Proofs1	9	17.3	16.0	13.1	13.3	13.3
Proofs2	8	16.9	14.4	11.5	11.5	11.5
Proofs3	7	19.3	15.6	10.9	10.7	10.9
Proofs4	11	19.1	14.3	10.9	11.4	11.4
Proofs5	9	20.1	17.6	12.8	11.9	11.9
Proofs6	11	25.6	18.9	12.5	12.6	12.6
Proofs7	11	37.2	19.7	11.8	11.9	11.9
Proofs8	7	24.4	15.6	12.3	13.1	13.1
Proofs9	14	39.8	29.0	13.1	14.1	14.1
Proofs10	2	21.5	16.0	8.0	11.0	11.0
Proofs11	5	25.4	22.4	13.0	14.0	14.0
Proofs12	13	24.6	16.5	8.0	8.5	8.5
Proofs13	10	35.3	27.7	9.1	10.1	10.1
Total	117	26.3	19.2	11.4	11.8	11.8

478 Provers are nondeterministic, especially when invoked with a time limit. More
 479 importantly, our approach makes different heuristic choices when constructing
 480 the three proof segments depending on which problem variants are used. As a
 481 result, SA might find a short proof that escapes BSA.

482 The reduction in proof length is especially noticeable in individual cases. The
 483 problem $2666 \implies 3460$ has a Vampire proof with 51 inference steps, which our
 484 tool reduces to only 12 single rewrite steps, and $2923 \implies 2628$ is reduced from
 485 180 steps to only 34. The problem $3569 \implies 3957$ is reduced from 92 to 23 steps
 486 and, even more dramatically, $3957 \implies 3971$ is reduced from 141 steps to only 23.
 487 Furthermore, $2860 \implies 2660$ is reduced from 44 to 14 steps, and $723 \implies 872$ goes
 488 from 57 to 13 steps. Finally, $947 \implies 3897$ underwent the largest reduction, from
 489 151 to 10 steps. Overall, these results demonstrate that our approach produces
 490 shorter proofs across a diverse set of equational theorems.

491 8 Related Work

492 At least two other researchers took on Tao’s challenge. Kinyon [18] found a 24-
 493 step proof (excluding preprocessing) of $650 \implies \forall x, y. x = x \diamond y$ using Prover9
 494 [22], from which $650 \implies 448$ follows by instantiation. Later, Le Floch [11] devel-
 495 oped a pen-and-paper proof and translated it to Lean. The Lean proof relies on
 496 only 14 rewrite steps but includes additional reasoning as proof terms, and two of
 497 the rewrite steps are parallel, so the overall length is similar to ours. The proof is
 498 “loosely based” on the output of multiple Prover9 runs “with intermediate results
 499 thrown in as assumptions or as goals”—in essence, a manual approximation of
 500 our approach. Also in the context of the Equational Theories Project, Janota [16]
 501 evaluated Vampire on the project’s problems and showed that combining super-
 502 position with finite model finding can solve almost all problems. One reason for
 503 this success might be the work on term ordering diagrams by Hajdu et al. [15].

504 We are aware of little work on automated proof minimization. Stachniak
 505 [26] designed an algorithm for constructing resolution proofs in propositional
 506 logics known as strongly finite logics. Amjad [2] and Cotton [10] introduced
 507 techniques for minimizing propositional resolution proofs. Vyskočil et al. [30]
 508 proposed to compress proofs by inventing new definitions using a heuristic based
 509 on substitution trees. Gu et al. [14] developed ProofOptimizer, which uses large
 510 language models to simplify Lean proofs.

511 Some SAT (satisfiability) and SMT (satisfiability modulo theories) solvers
 512 can minimize the number of axioms needed for a proof, but the result can be a
 513 longer proof. SAT solvers commonly interleave search, which can be expressed
 514 as resolution steps, with formula-rewriting techniques that go beyond resolution.
 515 This interleaving, known as inprocessing [17], is highly effective and often yields
 516 both faster solving times and shorter proofs than either approach in isolation.

517 The idea of automatically mixing and matching proofs is not new. Sutcliffe
 518 et al. [28] introduced a method for combining automatically generated proofs to
 519 generate new ones. Their proofs are represented as DAGs, enabling the identifi-
 520 cation and replacement of subproofs across different proofs. Proof combination

521 is guided by heuristics that measure structural similarity, and a greedy search
 522 strategy is used to explore alternative combinations that yield proofs differing
 523 from the originals. In contrast to our approach, the main objective is to increase
 524 proof diversity rather than minimize proof length.

525 9 Conclusion

526 Historically, more research has gone into finding proofs automatically than into
 527 improving and presenting them. We introduced an approach for minimizing equa-
 528 tional proofs by mixing and matching the output of separate runs of Vampire
 529 and Twee, and implemented it in a new tool, Krympa. We used the tool to min-
 530 imize the proof of problem 650 \implies 448 from the Equational Theories Project
 531 from 62 to 20 steps, thereby providing a fully automatic solution to a challenge
 532 posed by Tao. We also obtained remarkable reductions on other problems origi-
 533 nating from the project. The shorter proofs are arguably easier to understand
 534 by humans and sometimes more general. Our work shows that proof automation
 535 and readability can go hand in hand.

536 Our approach could be extended in several ways. First, it could be generalized
 537 to support full first- or higher-order logic. Second, alternative lemma abstraction
 538 strategies could be explored. Third, proofs with more than three segments could
 539 be synthesized. Fourth, we might want to consider not only the number of steps
 540 but also term size when measuring proofs, as suggested by Le Floch [12]. Fifth,
 541 we could try to translate Vampire’s superposition steps to Twee’s structured
 542 equality chain format.

543 Some possible extensions specifically concern the implementation. First, we
 544 could explore whether nondefault Vampire strategies can produce shorter proofs,
 545 following a private suggestion by Martin Suda. Second, since proof generation
 546 relies heavily on external provers, performance could benefit from better schedul-
 547 ing of prover invocations, using adaptive time limits. Third, as the number of
 548 possible lemma combinations grows rapidly, exploiting parallelism at multiple
 549 levels—such as lemma re-proving, dependency graph construction, and proof
 550 construction—would be a natural extension of the current architecture. Finally,
 551 we could integrate other provers than Vampire and Twee.

552 **Acknowledgments.** We thank Bruno Le Floch, Andrew Reynolds, Stephan
 553 Schulz, and Uwe Waldmann for fruitful discussions. We thank Laura Kovács,
 554 Luca Maio, Martin Suda, and Mark Summerfield for helpful textual suggestions.

555 Blanchette’s research was cofunded by the European Union (ERC, Nekoka,
 556 101083038). Views and opinions expressed are however those of the authors only
 557 and do not necessarily reflect those of the European Union or the European
 558 Research Council. Neither the European Union nor the granting authority can
 559 be held responsible for them.

560 Heule’s research is supported by the NSF under grant DMS-2434625 and
 561 funding from AFRL and DARPA under Agreement FA8750-24-9-1000.

References

- 562
- 563 1. The Lean Language Reference (2025), [https://lean-lang.org/doc/reference/](https://lean-lang.org/doc/reference/latest/)
564 [latest/](https://lean-lang.org/doc/reference/latest/)
 - 565 2. Amjad, H.: Compressing propositional refutations. In: Merz, S., Nipkow, T. (eds.)
566 AVoCS 2006. Electronic Notes in Theoretical Computer Science, vol. 185, pp. 3–15.
567 Elsevier (2006)
 - 568 3. Bachmair, L., Dershowitz, N., Plaisted, D.A.: Completion without failure. In: Ait-
569 Kaci, H., Nivat, M. (eds.) Rewriting Techniques, pp. 1–30. Academic Press (1989)
 - 570 4. Bachmair, L., Ganzinger, H.: Strict basic superposition. In: Kirchner, C., Kirchner,
571 H. (eds.) CADE 1998. LNCS, vol. 1421, pp. 160–174. Springer (1998)
 - 572 5. Bártek, F., Bhayat, A., Coutelier, R., Hajdú, M., Hetzenberger, M., Hozzová, P.,
573 Kovács, L., Rath, J., Rawson, M., Reger, G., Suda, M., Schoisswohl, J., Voronkov,
574 A.: The Vampire diary. In: Piskac, R., Rakamaric, Z. (eds.) CAV 2025. LNCS, vol.
575 15933, pp. 57–71. Springer (2025)
 - 576 6. Blanchette, J.C., Böhme, S., Fleury, M., Smolka, S.J., Steckermeier, A.: Semi-
577 intelligible Isar proofs from machine-generated proofs. *J. Automated Reas.* **56**,
578 155–200 (2016)
 - 579 7. Bolan, M., Breitner, J., Brox, J., Carlini, N., Carneiro, M., van Doorn, F., Dvorak,
580 M., Goens, A., Hill, A., Husum, H., Mejia, H.I., Kocsis, Z.A., Floch, B.L., Bar-on,
581 A.L., Luccioli, L., McNeil, D., Meiburg, A., Monticone, P., Nielsen, P., Osazuwa,
582 E.O., Paolini, G., Petracci, M., Reinke, B., Renshaw, D., Rossel, M., Roux, C.,
583 Scanvic, J., Srinivas, S., Tadipatri, A.R., Tao, T., Tsyrlkevich, V., Vaquerizo-Villar,
584 F., Weber, D., Zheng, F.: The Equational Theories Project (2025)
 - 585 8. Buch, A., Hillenbrand, T.: WALDMEISTER: Development of a high performance
586 completion-based theorem prover (1996)
 - 587 9. Clune, J., Qian, Y., Bentkamp, A., Avigad, J.: Duper: A proof-producing superposi-
588 tion theorem prover for dependent type theory. In: Bertot, Y., Kutsia, T., Norrish,
589 M. (eds.) ITP 2024. LIPIcs, vol. 309, pp. 1–20. Leibniz-Zentrum für Informatik
590 (2024)
 - 591 10. Cotton, S.: Two techniques for minimizing resolution proofs. In: Strichman, O.,
592 Szeider, S. (eds.) Theory and Applications of Satisfiability Testing - SAT 2010,
593 13th International Conference, SAT 2010, Edinburgh, UK, July 11-14, 2010. Pro-
594 ceedings. Lecture Notes in Computer Science, vol. 6175, pp. 306–312. Springer
595 (2010)
 - 596 11. Floch, B.L.: Zulip post in “Machine Learning for Theorem Proving:
597 A (Semi)-Autoformalization Challenge (650 → 448)”. Zulip (2025),
598 available at [https://leanprover.zulipchat.com/#narrow/channel/](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/518970125)
599 [219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/518970125)
600 [29-autoformalization.20challenge.3A.20650.3D.3E448/near/518970125](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/518970125)
 - 601 12. Floch, B.L.: Zulip post in “Machine Learning for Theorem Proving:
602 A (Semi)-Autoformalization Challenge (650 → 448)”. Zulip (2025),
603 available at [https://leanprover.zulipchat.com/#narrow/channel/](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/568313777)
604 [219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/568313777)
605 [29-autoformalization.20challenge.3A.20650.3D.3E448/near/568313777](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/568313777)
 - 606 13. Geoff Sutcliffe: Stepping stones in the TPTP World. In: Benzmüller, C., Heule,
607 M., Schmidt, R. (eds.) IJCAR 2024. pp. 30–50. LNCS (2024)
 - 608 14. Gu, A., Piotrowski, B., Gloeckle, F., Yang, K., Markosyan, A.H.: ProofOptimizer:
609 Training language models to simplify proofs without human demonstrations. *CoRR*
610 [abs/2510.15700](https://arxiv.org/abs/2510.15700) (2025)

- 611 15. Hajdu, M., Coutelier, R., Kovács, L., Voronkov, A.: Term ordering diagrams. In:
612 Barrett, C.W., Waldmann, U. (eds.) CADE-30. LNCS, Springer, to appear
- 613 16. Janota, M.: Experimental results for Vampire on the Equational Theories Project
614 (2025)
- 615 17. Järvisalo, M., Heule, M., Biere, A.: Inprocessing rules. In: Gramlich, B., Miller, D.,
616 Sattler, U. (eds.) IJCAR 2012. LNCS, vol. 7364, pp. 355–370. Springer (2012)
- 617 18. Kinyon, M.: Zulip post in “Machine Learning for Theorem Proving:
618 A (Semi)-Autoformalization Challenge (650 → 448)”. Zulip (2025),
619 available at [https://leanprover.zulipchat.com/#narrow/channel/
620 219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.
621 29-autoformalization.20challenge.3A.20650.3D.3E448/near/518961204](https://leanprover.zulipchat.com/#narrow/channel/219941-Machine-Learning-for-Theorem-Proving/topic/A.20.28semi.29-autoformalization.20challenge.3A.20650.3D.3E448/near/518961204)
- 622 19. Knuth, D.E., Bendix, P.B.: Simple word problems in universal algebras. In: Leech,
623 J. (ed.) Computational Problems in Abstract Algebra, pp. 263–297. Pergamon
624 Press (1970)
- 625 20. Kondylidou, L., Blanchette, J., Heule, M.: Tao's equational proof challenge accepted. Zenodo (2026), <https://doi.org/10.5281/zenodo.18624123>
- 626 21. Limperg, J., From, A.H.: Aesop: White-box best-first proof search for Lean. In:
627 CPP 2023. pp. 253–266. Association for Computing Machinery (2023)
- 628 22. McCune, W.: Prover9 and Mace4 (2005–2010)
- 629 23. de Moura, L., Ullrich, S.: The Lean 4 theorem prover and programming language.
630 In: Bertot, Y., Kutsia, T., Norrish, M. (eds.) CADE 2021. LNCS, vol. 12699, pp.
631 625–635. Springer (2021)
- 632 24. Norman, C., Avigad, J.: Canonical for automated theorem proving in Lean. In:
633 ITP 2025. LIPIcs, vol. 352, pp. 1–20. Leibniz-Zentrum für Informatik (2025)
- 634 25. Smallbone, N.: Twee: An equational theorem prover. In: Platzer, A., Sutcliffe, G.
635 (eds.) CADE 2021. LNCS, vol. 12699, pp. 602–613. Springer (2021)
- 636 26. Stachniak, Z.: Minimization of resolution proof systems. *Fundam. Informaticae*
637 **14**(1), 129–146 (1991)
- 638 27. Sutcliffe, G.: The 12th IJCAR Automated Theorem Proving System
639 Competition—CASC-J12. *AI Communications* **38**, 3–20 (2025)
- 640 28. Sutcliffe, G., Chang, C., McGuinness, D., Lebo, T., Ding, L., da Silva, P.P.: Com-
641 bining proofs to form different proofs. In: Fontaine, P., Stump, A. (eds.) PxTP
642 2011. pp. 60–73. LNCS (2011)
- 643 29. Tao, T.: Machine learning for theorem proving: A (semi)-autoformalization
644 challenge (650 → 448). [https://leanprover-community.github.io/archive/
645 stream/219941-Machine-Learning-for-Theorem-Proving/](https://leanprover-community.github.io/archive/stream/219941-Machine-Learning-for-Theorem-Proving/), Lean Zulip thread,
646 created May 16, 2025
- 647 30. Vyskočil, J., Stanovský, D., Urban, J.: Automated proof compression by invention
648 of new definitions. In: Clarke, E.M., Voronkov, A. (eds.) LPAR-16. LNCS, vol. 6355,
649 pp. 447–462. Springer (2010)
- 650 31. Zhu, T., Clune, J., Avigad, J., Jiang, A.Q., Welleck, S.: Premise selection for a
651 Lean hammer. arXiv preprint arXiv:2506.07477 (2025)
- 652