

1 Tao’s Equational Proof Challenge Accepted

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

3 ¹ Ludwig-Maximilians-Universität München, Munich, Germany

4 {l.kondylidou,jasmin.blanchette}@lmu.de

5 ² Carnegie Mellon University, Pittsburgh, United States

6 marijn@cmu.edu

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

16 **Keywords:** Theorem provers · Equational logic · Proof minimization.

17 1 Introduction

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

24 The project uses the Lean [22] proof assistant to express proofs and counter-
25 examples but depends on automatic theorem provers and other external tools.
26 The problems explored in the project’s first phase all fall within first-order logic’s
27 unit equality fragment: They consist of a \forall -quantified equation as the sole axiom
28 and a \forall -quantified equation as the proof goal, or conjecture.

29 For the problem 650 \implies 448, where 650 denotes the axiom $\forall x, y, z. x = x \diamond$
30 $(y \diamond ((z \diamond x) \diamond y))$ and 448 denotes the conjecture $\forall x, y, z. x = x \diamond (y \diamond (z \diamond (x \diamond z)))$,
31 the Vampire [5] superposition prover found a particularly complex proof, with 62
32 inference steps, excluding clausification and Skolemization. Given that the proof
33 is unintelligible, Tao challenged the community to find “an alternate proof, by
34 whatever means you wish—human, semi-automated, or automated” [28].

35 One idea could be to run a specialized equational prover, Twee [24], instead of
36 Vampire, but this results in a very long, 137-step proof. Another approach would
37 be to use Lean’s automation, such as the `aesop` [20], `canonical` [23], `duper` [9],
38 and `grind` [1] tactics and the LeanHammer [30], to reconstruct and compress

39 consecutive superposition steps, in the style of Sledgehammer’s structured proof
 40 reconstruction [6, Sect. 6.3]. This would yield a shorter and more high-level proof,
 41 in which each step may combine multiple rewrites. Our approach is orthogonal.
 42 Our working hypothesis is that Vampire’s 62-step proof, which emerged as the
 43 byproduct of a saturation process, is likely suboptimal. By mixing and match-
 44 ing proofs generated by different automatic provers, as proposed by Sutcliffe et
 45 al. [27], we hope to achieve a shorter, simpler proof.

46 We introduce Krympa, a tool that minimizes equational proofs by decompos-
 47 ing them into independently provable components and reassembling them into
 48 more concise, intelligible proofs. Specifically, starting from a Vampire-generated
 49 proof, the tool transforms it into a direct proof (Sect. 3) and analyzes its infer-
 50 ences to break it down into intermediate results that serve as candidate lemmas.
 51 Each of these lemmas is then proved independently using Vampire and Twee
 52 (Sect. 4), the two leading systems in the unit equality division of CASC 2025 [26].
 53 The resulting proofs are then combined into a single proof using heuristics that
 54 favor shorter derivations (Sect. 5).

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

60 Our tool is implemented in Rust, OCaml, and Python. Its source code is avail-
 61 able at <https://github.com/kondylidou/Krympa>. The files associated with
 62 Tao’s challenge and our empirical evaluation data are also available online [19].

63 2 Background

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

66 2.1 Vampire and Superposition Proofs

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

71 Superposition works on implicitly \forall -quantified clauses. A preprocessor per-
 72 forms clausification and Skolemization. For example, the axiom $\forall x. f(x) = g(x)$
 73 is transformed into $f(x) = g(x)$, where x is a free variable, and the conjecture
 74 $\forall x. f(x) = g(x)$ is negated and transformed into $f(\text{sk}) \neq g(\text{sk})$, where sk is a
 75 Skolem constant. The objective is to derive the contradictory clause \perp . For the
 76 unit equality fragment, the calculus’s two relevant inference rules are as follows:

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

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 $f(a) = b$ either as is or as $b = f(a)$. 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 $h(a, y)$ and $h(x, b)$ is
 86 $\{x \mapsto a, y \mapsto b\}$; applying it on both terms yields $h(a, b)$. 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:

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

100 Superposition proofs are represented 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 linear 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 |
| 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 |
| 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 comple-
 108 tion [18]. In the DISCOUNT and Waldmeister tradition [8], Twee's proofs are
 109 structured as a sequence of lemmas, where each lemma and the conjecture are

110 proved by a chain of equalities. Twee introduces lemmas if they are needed more
 111 than once. Twee proofs are arguably more readable than Vampire proofs. As
 112 with superposition, quantifiers are eliminated by a preprocessor.

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

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

115 3 Proof Redirection

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 To redirect a proof by refutation in equational logic, we first introduce \exists
 120 quantifiers for Skolem constants and \forall quantifiers for variables. For example,
 121 $h(x, sk) \neq x$ is transformed into $\exists z. \forall x. h(x, z) \neq x$. Then we apply the contra-
 122 positive to all inferences in which a premise and the conclusion are disequations
 123 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.

$$\begin{array}{ll}
 \text{129 1. } a = b & \text{axiom} \\
 2. \forall x. f(x) = x & \text{axiom} \\
 3. h(b, a) = h(a, b) & \text{from 1 and } h(a, a) = h(a, a) \\
 4. h(f(b), a) = h(a, f(b)) & \text{from 2 and 3} \blacksquare
 \end{array}$$

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

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 redirected
 161 proof, excluding preprocessing. For Twee, the length of a proof is the
 162 cumulative number of equalities in the equality chains. Thus, the Vampire proof
 163 in 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 initial 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 separation into segments, proof minimization could be intractable
 190 due to combinatorial explosion. We chose to work with three segments as a trade-
 191 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 an initial
 194 seven-step Vampire-generated redirected 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 initial proof, but
 224 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 :

| | | |
|-----|--------|---|
| 246 | A | axiom |
| | L_1 | from A |
| | L'_2 | from A and L_1 |
| | L_3 | from L_1 and L'_2 |
| | C | by a two-step equality chain using L'_2 and L_3 |

247 Finally, other combinations of candidate departure and arrival lemmas are
 248 also considered, and the shortest proof is retained. ■

249 5.1 Construction of the Dependency Graph

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

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

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

273 5.2 Construction of the First Proof Segment

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

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

289 5.3 Construction of the Remaining Proof Segments

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

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

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

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

302 **5.4 Proof Output**

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

```
308 calc
309    $t_1 = t_2 := \text{by duper } \dots$ 
310    $\_ = t_3 := \text{by duper } \dots$ 
311    $\_ = t_4 := \text{by duper } \dots$ 
```

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

315 **6 Application to Tao's Challenge**

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

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

318 Our tool first ran Vampire to obtain an initial 62-step superposition proof. Then
 319 it constructed 62 problems of each variant (big-step, small-step, and abstracted)
 320 and tried to prove them using Vampire and Twee. Among the six candidate
 321 arrival lemmas, the shortest proof was found by selecting

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

322 The coloring highlights repeating patterns. Next, our tool constructed the
 323 dependency graph for this lemma. The DAG contained 37 lemmas. It was based
 324 on big- and small-step proofs.

325 Among the 37 candidate departure lemmas, our tool found the shortest proof
 326 by selecting

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

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

$$\begin{aligned}
& \forall x, y, z, w. x \diamond ((y \diamond z) \diamond x) = \\
& \quad (x \diamond ((y \diamond z) \diamond x)) \diamond (w \diamond (z \diamond w)) & \text{(lemma 1)} \\
& \forall x, y, z, w, v, u. x \diamond ((y \diamond ((z \diamond w) \diamond y)) \diamond x) = \\
& \quad (x \diamond ((y \diamond ((z \diamond w) \diamond y)) \diamond x)) \diamond (v \diamond ((u \diamond (w \diamond u)) \diamond v)) & \text{(lemma 2)} \\
& \forall x, y, z, w, v. x \diamond (y \diamond x) = \\
& \quad (x \diamond (y \diamond x)) \diamond (z \diamond ((w \diamond ((v \diamond y) \diamond w)) \diamond z)) & \text{(lemma 3)} \\
& \forall x, y, z, w, v. x \diamond (y \diamond x) = \\
& \quad (x \diamond (y \diamond x)) \diamond ((z \diamond (y \diamond z)) \diamond (w \diamond ((v \diamond y) \diamond w))) & \text{(lemma 4)} \\
& \forall x, y, z, w. x \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x) = \\
& \quad (x \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond x)) \diamond \\
& \quad (w \diamond ((y \diamond ((z \diamond y) \diamond y)) \diamond w)) & \text{(lemma 5)} \\
& \forall x, y, z, w. (x \diamond ((y \diamond x) \diamond x)) \diamond z = \\
& \quad ((x \diamond ((y \diamond x) \diamond x)) \diamond z) \diamond ((w \diamond ((x \diamond ((y \diamond x) \diamond x)) \diamond w)) \diamond \\
& \quad (z \diamond ((x \diamond ((y \diamond x) \diamond x)) \diamond z)). & \text{(lemma 6)}
\end{aligned}$$

330 Following the inference steps of the initial Vampire proof, our tool derived
331 lemma 1 by applying a superposition inference with the axiom $x = x \diamond (y \diamond ((z \diamond$
332 $x) \diamond y))$ as the first premise and a renamed copy $x' = x' \diamond (y' \diamond ((z' \diamond x') \diamond y'))$ as
333 the second premise. The most general unifier of the first premise's right-hand side
334 and the subterm $z' \diamond x'$ of the second premise is $\{x' \mapsto y \diamond ((z \diamond x) \diamond y), z' \mapsto x\}$.
335 Applying the unifier to both premises yields the equations $x = x \diamond (y \diamond ((z \diamond x) \diamond$
336 $y))$ and $y \diamond ((z \diamond x) \diamond y) = (y \diamond ((z \diamond x) \diamond y)) \diamond (y' \diamond ((x \diamond (y \diamond ((z \diamond x) \diamond y))) \diamond y'))$.
337 The superposition inference replaced the subterm $x \diamond (y \diamond ((z \diamond x) \diamond y))$ in the
338 second equation with x using the first equation as a right-to-left rewrite rule,
339 and thus derived lemma 1, up to the naming of variables. Lemmas 2 to 7 were
340 derived similarly following the steps of the initial Vampire proof.

341 Next, from the axiom and lemma 7, our tool proved the arrival lemma
342 (lemma 9) using Twee. For this proof, Twee introduced the auxiliary lemma

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

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

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

```

351 class Magma ( $\alpha$  : Type  $\_$ ) where
352   op :  $\alpha \rightarrow \alpha \rightarrow \alpha$ 
353
354   infix:65 " $\diamond$ " => Magma.op

```

```

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

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

For each file, we invoked our tool's TPTP problem generator, which parses 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 initial proof and Vampire and Twee to find subproofs; on
 454 failure, the initial 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 fo-
 464 cuses on problems whose initial proofs have at least 15 steps. Table 1 compares
 465 proof lengths before and after minimization. The “Avg. before” column shows the
 466 average number of inference steps in the initial Vampire proofs. The remaining
 467 columns report the average proof length after minimization under four configu-
 468 rations, which differ in which problem variants are used to generate candidate
 469 lemmas: “BA” denotes the combination of the big-step and abstracted variants;
 470 “SA” denotes the combination of the small-step and abstracted variants; “BS”
 471 denotes the combination of the big- and small-step variants; and “BSA” denotes
 472 the combination of all three variants.

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

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

| File | Num. problems | Avg. before | Avg. after | | | |
|----------|---------------|-------------|------------|-------------|-------------|-------------|
| | | | BA | SA | BS | BSA |
| Proofs1 | 135 | 17.3 | 16.0 | 13.1 | 13.3 | 13.3 |
| Proofs2 | 117 | 16.9 | 14.4 | 11.5 | 11.5 | 11.5 |
| Proofs3 | 108 | 19.3 | 15.6 | 10.9 | 10.7 | 10.9 |
| Proofs4 | 125 | 19.1 | 14.3 | 10.9 | 11.4 | 11.4 |
| Proofs5 | 116 | 20.1 | 17.6 | 12.8 | 11.9 | 11.9 |
| Proofs6 | 115 | 25.6 | 18.9 | 12.5 | 12.6 | 12.6 |
| Proofs7 | 117 | 37.2 | 19.7 | 11.8 | 11.9 | 11.9 |
| Proofs8 | 114 | 24.4 | 15.6 | 12.3 | 13.1 | 13.1 |
| Proofs9 | 112 | 39.8 | 29.0 | 13.1 | 14.1 | 14.1 |
| Proofs10 | 101 | 21.5 | 16.0 | 8.0 | 11.0 | 11.0 |
| Proofs11 | 110 | 25.4 | 22.4 | 13.0 | 14.0 | 14.0 |
| Proofs12 | 123 | 24.6 | 16.5 | 8.0 | 8.5 | 8.5 |
| Proofs13 | 38 | 35.3 | 27.7 | 9.1 | 10.1 | 10.1 |

477 It might seem counterintuitive that SA, which does not consider big-step
 478 problems, outperforms BSA. However, the nonmonotonicity is to be expected.
 479 Provers are nondeterministic, especially when invoked with a time limit. More
 480 importantly, our approach makes different heuristic choices when constructing
 481 the three proof segments depending on which problem variants are used. As a
 482 result, SA might find a short proof that escapes BSA.

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

492 8 Related Work

493 At least two other researchers took on Tao’s challenge. Kinyon [17] found a 24-
 494 step proof (excluding preprocessing) of $650 \Rightarrow \forall x, y. x = x \diamond y$ using Prover9
 495 [21], from which $650 \Rightarrow 448$ follows by instantiation. Later, Le Floch [11] devel-
 496 oped a pen-and-paper proof and translated it to Lean. The Lean proof relies on
 497 only 14 rewrite steps but includes additional reasoning as proof terms, and two
 498 of the rewrite steps are parallel, so the overall length is similar to ours. The proof
 499 idea is “loosely based” on the output of multiple Prover9 runs “with intermediate
 500 results thrown in as assumptions or as goals”—in essence, a manual approxima-
 501 tion of our approach. Also in the context of the Equational Theories Project,
 502 Janota [15] evaluated Vampire on the project’s problems and showed that com-
 503 bining superposition with finite model finding can solve almost all problems.

504 We are aware of little work on automated proof minimization. Stachniak
 505 [25] 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. [29]
 508 proposed to compress proofs by inventing new definitions using a heuristics 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 [16], 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. [27] 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 orig-
 533 inating 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, following a private suggestion by Martin Suda, explore whether nondefault
 545 Vampire strategies can produce shorter proofs. 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. Finally, 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.

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

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

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

561 References

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

610 15. Janota, M.: Experimental results for Vampire on the Equational Theories Project
 611 (2025)

612 16. Järvisalo, M., Heule, M., Biere, A.: Inprocessing rules. In: Gramlich, B., Miller, D.,
 613 Sattler, U. (eds.) IJCAR 2012. LNCS, vol. 7364, pp. 355–370. Springer (2012)

614 17. Kinyon, M.: Zulip post in “Machine Learning for Theorem Proving:
 615 A (Semi)-Autoformalization Challenge (650 → 448)”. Zulip (2025),
 616 available at <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>

617 18. Knuth, D.E., Bendix, P.B.: Simple word problems in universal algebras. In: Leech,
 618 J. (ed.) Computational Problems in Abstract Algebra, pp. 263–297. Pergamon
 619 Press (1970)

620 19. Kondylidou, L., Blanchette, J., Heule, M.: Tao's equational proof challenge ac-
 621 cepted. Zenodo (2026), <https://doi.org/10.5281/zenodo.18624123>

622 20. Limpert, J., From, A.H.: Aesop: White-box best-first proof search for Lean. In:
 623 CPP 2023. pp. 253–266. Association for Computing Machinery (2023)

624 21. McCune, W.: Prover9 and Mace4 (2005–2010)

625 22. de Moura, L., Ullrich, S.: The Lean 4 theorem prover and programming language.
 626 In: Bertot, Y., Kutsia, T., Norrish, M. (eds.) CADE 2021. LNCS, vol. 12699, pp.
 627 625–635. Springer (2021)

628 23. Norman, C., Avigad, J.: Canonical for automated theorem proving in Lean. In:
 629 ITP 2025. LIPICS, vol. 352, pp. 1–20. Leibniz-Zentrum für Informatik (2025)

630 24. Smallbone, N.: Twee: An equational theorem prover. In: Platzer, A., Sutcliffe, G.
 631 (eds.) CADE 2021. LNCS, vol. 12699, pp. 602–613. Springer (2021)

632 25. Stachniak, Z.: Minimization of resolution proof systems. Fundam. Informaticae
 633 **14**(1), 129–146 (1991)

634 26. Sutcliffe, G.: The 12th IJCAR Automated Theorem Proving System
 635 Competition—CASC-J12. AI Communications **38**, 3–20 (2025)

636 27. Sutcliffe, G., Chang, C., McGuinness, D., Lebo, T., Ding, L., da Silva, P.P.: Com-
 637 bining proofs to form different proofs. In: Fontaine, P., Stump, A. (eds.) PxTP
 638 2011. pp. 60–73. LNCS (2011)

639 28. Tao, T.: Machine learning for theorem proving: A (semi)-autoformalization
 640 challenge (650 → 448). <https://leanprover-community.github.io/archive/stream/219941-Machine-Learning-for-Theorem-Proving/>, Lean Zulip thread,
 641 created May 16, 2025

642 29. Vyskočil, J., Stanovský, D., Urban, J.: Automated proof compression by invention
 643 of new definitions. In: Clarke, E.M., Voronkov, A. (eds.) LPAR-16. LNCS, vol. 6355,
 644 pp. 447–462. Springer (2010)

645 30. Zhu, T., Clune, J., Avigad, J., Jiang, A.Q., Welleck, S.: Premise selection for a
 646 Lean hammer. arXiv preprint arXiv:2506.07477 (2025)