

Seventeen Provers under the Hammer

Martin Desharnais ✉ 

Graduate School of Computer Science, Saarland Informatics Campus, Saarbrücken, Germany

Max-Planck-Institut für Informatik, Saarland Informatics Campus, Saarbrücken, Germany

Petar Vukmirović ✉ 

Vrije Universiteit Amsterdam, Amsterdam, the Netherlands

Jasmin Blanchette ✉ 

Vrije Universiteit Amsterdam, Amsterdam, the Netherlands

Max-Planck-Institut für Informatik, Saarland Informatics Campus, Saarbrücken, Germany

Makarius Wenzel ✉ 

Augsburg, Germany

Abstract

One of the main success stories of automatic theorem provers has been their integration into proof assistants. Such integrations, or “hammers,” increase proof automation and hence user productivity. In this paper, we use Isabelle/HOL’s Sledgehammer tool to find out how useful modern provers are at proving formulas in higher-order logic. Our evaluation follows in the steps of Böhme and Nipkow’s Judgment Day study from 2010, but instead of three provers we use 17, including SMT solvers and higher-order provers. Our work offers an alternative yardstick for comparing modern provers, next to the benchmarks and competitions emerging from the TPTP World and SMT-LIB.

2012 ACM Subject Classification Computing methodologies → Theorem proving algorithms

Keywords and phrases Automatic theorem proving, interactive theorem proving, proof assistants

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

Funding This research has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 713999, Matryoshka). Blanchette has received funding from the Netherlands Organization for Scientific Research (NWO) under the Vidi program (project No. 016.Vidi.189.037, Lean Forward).

Acknowledgments We are grateful to the maintainers of StarExec for letting us use their service. Developers of automatic theorem provers, including Peter Baumgartner, Ahmed Bhayat, Pascal Fontaine, Konstantin Korovin, Giles Reger, Andrew Reynolds, Philipp Rümmer, Alexander Steen, and Martin Suda have helped us run their systems. Michael Färber, Thibault Gauthier, Konstantin Korovin, Lorenz Leutgeb, Jens Otten, Philipp Rümmer, Stephan Schulz, Hans-Jörg Schurr, Alexander Steen, Martin Suda, Mark Summerfield, Dmitriy Traytel, Josef Urban, and the anonymous reviewers suggested some textual improvements.

1 Introduction

Hammers [15] are tools that integrate automatic theorem provers in proof assistants, to automatically discharge proof obligations. CoqHammer [27], HOLyHammer [35], MizAR [58], and Sledgehammer (Section 2) are examples of hammers. They typically consist of three main components in addition to the automatic prover backends themselves:

- The *relevance filter* heuristically selects relevant lemmas (including definitions) based on the current goal, using either an iterative procedure or machine learning.
- The *problem translation* module encodes formulas from the proof assistant’s logic to the automatic provers’ usually less expressive logics. The encoding should ideally be sound and complete.



© Martin Desharnais, Petar Vukmirović, Jasmin Blanchette, and Makarius Wenzel; licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

44 ■ The *proof reconstruction* module translates proofs generated by the automatic provers to
45 proof texts expressed in the proof assistant’s language.

46 There is ample evidence that hammers can be effective. The landmark Judgment Day
47 study by Böhme and Nipkow [20] showed for the first time that a combination of three
48 automatic provers (E [49], SPASS [61], Vampire [38]) could discharge about half of the goals
49 that arise in typical Isabelle/HOL [43] developments. This came as a surprise to most Isabelle
50 users, who had not yet integrated Sledgehammer into their workflow. Since then, it has been
51 adopted by a large majority of users.

52 For over a decade, Judgment Day has been the most comprehensive evaluation of Sledge-
53 hammer, even though it considered only three provers and seven Isabelle theory files. Smaller
54 evaluations in later papers [7, 8, 11–13, 17, 45, 50, 53] cover more provers, but they mostly
55 predate the new generation of higher-order provers: cvc5 [1], Ehoh [60], Leo-III [51], Vam-
56 pire [9], and Zipperposition [59]. We do not have a clear picture of how well these provers
57 work in a hammer setting. We do not even know whether native support for higher-order
58 logic outperforms encodings.

59 In this paper, we evaluate 17 modern provers (Section 3), both first- and higher-order.
60 They support a wide range of formats from the TPTP [54] and SMT-LIB [2] families. We first
61 generate problems using a tool called Mirabelle (Section 4), which invokes Sledgehammer on
62 goals originating from 50 Isabelle proof developments. We use the MePo [42] relevance filter,
63 which iteratively selects lemmas in a deterministic fashion based on the goal. To translate
64 the problem, we use Sledgehammer’s TPTP and SMT-LIB modules, targeting a wide range
65 of provers, including SMT (satisfiability modulo theories) solvers. Once the problems are
66 generated, we run the provers on them (Section 5). The evaluation yields a large collection
67 of new benchmarks that can be used to tune automatic provers—much larger than Judgment
68 Day, which currently consists of 1240 goals. Our collection of problems, called Seventeen,
69 and the raw evaluation data are archived online.¹

70 The evaluation is first and foremost an assessment of the provers and their features and of
71 various encoding schemes. Our setup does not attempt to reconstruct proofs in Isabelle. To
72 guard against incorrect proofs, we use only encodings known to be sound [12] and provers that
73 have shown themselves to be trustworthy over the years. The question of how to reconstruct
74 higher-order proofs in Isabelle is still partly open.

75 In addition, our evaluation provides raw data that can be used to fine-tune how Sledge-
76 hammer uses automatic provers. The evaluation also gives insights into which translation
77 methods are useful in hammers, guiding the development of hammers.

78 Although it may be tempting to view this paper as a prover competition, there are
79 important differences. In a competition, the developers have full control over how their
80 provers are invoked, and there is a clear protocol for interacting with the organizers. We
81 tried to pass meaningful command-line options and consulted prover developers, but there
82 are no guarantees that we succeeded in finding the best options for all provers. Of course,
83 the success rates we report are only lower bounds on what can be achieved.

84 **2 Sledgehammer**

85 When we invoke Sledgehammer on a goal (i.e., a proof obligation, which might be provable
86 or not), Sledgehammer’s relevance filters and problem translation modules come into play.

¹ <https://doi.org/10.5281/zenodo.5940084>



■ **Table 1** Output formats supported by Sledgehammer

Format	Description
FOF	Untyped first-order logic
TF0	Many-sorted first-order logic
TX0	Many-sorted first-order logic with Booleans, if-then-else, and let
SMT2	Many-sorted first-order logic with Booleans, if-then-else, let, and linear arithmetic
TF1	Rank-1 polymorphic first-order logic
TX1	Rank-1 polymorphic first-order logic with Booleans, if-then-else, and let
TH0 ⁻	Many-sorted higher-order logic
TH0 ⁺	Many-sorted higher-order logic with Hilbert choice, if-then-else, and let
SMT3	Many-sorted higher-order logic with if-then-else, let, and linear arithmetic
TH1 ⁻	Rank-1 polymorphic higher-order logic
TH1 ⁺	Rank-1 polymorphic higher-order logic with Hilbert choice, if-then-else, and let

87 We briefly describe them below. The literature covers them in more detail.

88 2.1 The Relevance Filter

89 Given a goal and a cutoff number n , the relevance filter selects a subset of n lemmas among
 90 all the lemmas that are currently loaded (typically numbering in the thousands) and ranks
 91 them from 1 to n in order of decreasing expected relevance for proving the goal. Isabelle
 92 includes three relevance filters.

- 93 ■ MePo [42], named after *Meng* and *Paulson*, works iteratively starting from the goal. It is
 94 superficially similar to the SInE [31] algorithm implemented in several automatic provers
 95 (E, Vampire, Zipperposition) but was developed independently.
- 96 ■ MaSh [13], the *Machine learner for Sledgehammer*, implements two fast machine learning
 97 algorithms: naive Bayes and k -nearest neighbors. It learns which lemmas are useful to
 98 reason about which symbols and ranks the lemmas accordingly.
- 99 ■ MeSh [13] is a combination of *MePo* and *MaSh*.

100 While MaSh and MeSh give higher success rates than MePo [13], they are tricky to use
 101 properly in our complicated evaluation setup, with 50 different Isabelle developments. They
 102 also introduce nondeterminism. Since we are more interested in the relative performance of
 103 the provers than in the absolute success rates, we decided to use only MePo for this paper.
 104 The MaSh paper [13] has a detailed evaluation of the three filters.

105 MePo operates on a set of *known symbols*, initialized to consist of the symbols that occur
 106 in the goal. Roughly speaking, MePo works as follows:

- 107 1. Rank the (thousands of) available lemmas. The more symbols a lemma shares with the
 108 goal, and the fewer other symbols it contains, the higher its weight.
- 109 2. Select a number of lemmas based on their weights, removing them from the set of available
 110 lemmas. If no lemmas have suitably high weights, stop.
- 111 3. Enlarge the set of known symbols to include all the symbols occurring in the newly
 112 selected lemmas.
- 113 4. If n lemmas have been selected, stop; otherwise, go to step 1.



■ **Table 2** Type encodings supported by Sledgehammer’s TPTP module

Encoding	Description
\mathbf{g}	Polymorphism-preserving encoding based on type guards (predicates)
$\mathbf{g}?$	Lightweight variant of \mathbf{g}
$\mathbf{g}??$	Featherweight variant of \mathbf{g}
$\mathbf{g}@$	Alternative polymorphism-preserving encoding based on type guards
\mathbf{t}	Polymorphism-preserving encoding based on type tags (functions)
$\mathbf{t}?$	Lightweight variant of \mathbf{t}
$\mathbf{t}??$	Featherweight variant of \mathbf{t}
$\mathbf{t}@$	Alternative polymorphism-preserving encoding based on type guards
$\tilde{\mathbf{g}}, \tilde{\mathbf{g}}?, \tilde{\mathbf{g}}??, \tilde{\mathbf{t}}, \tilde{\mathbf{t}}?, \tilde{\mathbf{t}}??$	Monomorphizing variants of $\mathbf{g}, \mathbf{g}?, \mathbf{g}??, \mathbf{t}, \mathbf{t}?, \mathbf{t}??$, respectively

114 A further refinement is that symbols are annotated with their types, to increase precision.
 115 For example, if the symbol `nil` of type *nat list* is known, lemmas containing `nil` of polymorphic
 116 type α *list* will be considered relevant, unlike lemmas containing `nil` of type *bool list*.

117 2.2 The TPTP Translation

118 The TPTP translation module [41] generates problems in TPTP formats FOF, TF0, TX0,
 119 TF1, TX1, TH0, and TH1. Moreover, we draw an unofficial distinction between TH0⁻ and
 120 TH0⁺ variants of TH0 and similarly for TH1. Table 1 presents a brief overview of all these
 121 syntaxes, as well as the SMT-LIB syntaxes described below. Note that all higher-order
 122 formats support interpreted Booleans. Moreover, although some of the syntaxes provide
 123 interpreted arithmetic types and symbols, the TPTP translation module does not exploit
 124 this and maps Isabelle’s arithmetic operators to uninterpreted symbols.

125 When translating Isabelle’s logic to the automatic provers’ possibly weaker logics, four
 126 types of constructs may need to be encoded: types, partial application, λ -abstractions, and
 127 Booleans. The encodings are naturally not used in formats that support the respective
 128 features; for example, TX0, TX1, TH0, and TH1 support interpreted Booleans, so there is
 129 no need to encode them.

- 130 ■ For the *types* and the type classes, some translation schemes encode Isabelle’s entire
 131 rank-1 polymorphic type system using terms and clauses [12, 41]. An alternative that
 132 works particularly well in conjunction with the many-sorted formats TF0, TX0, and
 133 TH0 is to iteratively monomorphize the problem [12, Section 5.6]. Monomorphization is
 134 generally incomplete [18, Section 2]. Sledgehammer’s sound type encodings are listed in
 135 Table 2 and described by Blanchette et al. [12].
- 136 ■ The *Boolean* constants `False`, `True`, the logical connectives, and the quantifiers are either
 137 mapped to their TPTP equivalents or translated to uninterpreted “proxy” symbols.
 138 Axioms are provided for reasoning about the proxies (e.g., `False` \neq `True`).
- 139 ■ *Partial application* is encoded using a distinguished binary symbol `app`. For example,
 140 `rev` and `rev xs` are mapped to `rev` and `app(rev, xs)`, respectively. As an optimization, if all
 141 occurrences of `rev` take at least one argument, it can be passed directly (e.g., `rev(xs)`).
- 142 ■ *λ -abstractions* are encoded using SKBCI combinators [57] or λ -lifting [33]. The generated
 143 problem may include axioms that define the introduced symbols, or may leave them as
 144 *opaque* symbols.



145 ► **Example 1.** Consider the Isabelle symbol filter of type $(\alpha \Rightarrow \text{bool}) \Rightarrow \alpha \text{ list} \Rightarrow \alpha \text{ list}$ such
 146 that filter $p \ xs$ is defined as the sublist of xs that consists of the elements x for which $p \ x$.
 147 The following lemma is part of Isabelle’s list library:

148
$$\text{filter } P (\text{filter } Q \ xs) = \text{filter } (\lambda x. Q \ x \wedge P \ x) \ xs$$

149 If the target is TF0, and both monomorphization and SKBCI are used, the encoding is

150
$$\text{filter}_{nat}(P, \text{filter}_{nat}(Q, \ xs)) = \text{filter}_{nat}(S_{nat, bool, bool}(\mathbf{B}_{nat, bool, bool \Rightarrow bool}(\text{conj}, Q), P), \ xs)$$

151 where the subscripts identify monomorphic instances. In addition, axioms are included
 152 to characterize the combinators $S_{nat, bool, bool}$ and $\mathbf{B}_{nat, bool, bool \Rightarrow bool}$ and the proxy `conj`. By
 153 contrast, if the target is TH1, the lemma is encoded as itself.

154 2.3 The SMT-LIB Translation

155 The SMT-LIB translation module [19, Chapter 2] was first developed as part of the *smt*
 156 proof method [21] and later added to Sledgehammer [11]. The SMT-LIB 2 standard roughly
 157 amounts to TPTP TX0 with arithmetic and other theories. The forthcoming SMT-LIB 3 is
 158 expected to support higher-order logic and polymorphism and thus be a superset of TPTP
 159 TH1. Both *cvc5* and a fork of *veriT* support a preliminary version of the standard.

160 The SMT-LIB translation module monomorphizes polymorphic types. This is done even
 161 for SMT-LIB 3 output. In addition, for SMT-LIB 2, partial application is encoded using `app`,
 162 and λ -lifting is used. Isabelle’s basic arithmetic operators on integers and reals are mapped
 163 to the corresponding SMT operators. We use the names SMT2 and SMT3 to refer to the
 164 fragments of SMT-LIB 2 and 3 used by Sledgehammer, respectively.

165 3 The Automatic Provers

166 For our evaluation, we selected 17 provers that support at least one of Sledgehammer’s
 167 formats. The formats supported by each prover are listed in Table 3. All of the provers
 168 participated in the CASC [56] or the SMT-COMP [4] competition at some point. Except when
 169 mentioned otherwise, the provers were not specifically tuned for Seventeen-style problems.

- 170 ■ *agsyHOL* [39] is a higher-order prover based on an intuitionistic sequent calculus and a
 171 narrowing engine. It was developed to showcase the technology behind the Agsy [40]
 172 proof tool for Agda. We use version 1.0.
- 173 ■ *Beagle* [5] is a first-order prover based on hierarchic superposition, a calculus that
 174 generalizes standard superposition (which in turn generalizes resolution) with theory
 175 reasoning—in Beagle’s case, linear arithmetic reasoning. The main developer, Peter
 176 Baumgartner, points out that the prover is stronger on problems that require arithmetic
 177 reasoning and that it usually performs worse than state-of-the-art superposition provers
 178 on problems without theories. We use version 0.9.50.
- 179 ■ *cocATP²* is a Coq-inspired higher-order automatic theorem prover based on the calculus
 180 of constructions without inductive constructions. We use version 0.2.0.
- 181 ■ *cvc5* is a first-order SMT solver based on CDCL(T) with some support for higher-order
 182 logic [1]. It is CVC4’s [3] successor. The developers provided us with a portfolio of
 183 configurations that are tuned for low timeouts. We use version 0.0.4.

² <http://www.tptp.org/CASC/J7/SystemDescriptions.html#cocATP---0.2.0>



■ **Table 3** Input formats supported by the provers

Prover	FOF	TF0	TX0	SMT2	TF1	TX1	TH0 ⁻	TH0 ⁺	SMT3	TH1 ⁻	TH1 ⁺
agsyHOL							✓				
Beagle	✓	✓		✓							
cocATP							✓				
cvc5	✓	✓		✓			✓		✓		
E	✓	✓	✓				✓				
ENIGMA	✓										
iProver	✓	✓	✓	✓							
leanCoP	✓										
Leo-III	✓	✓	✓		✓	✓	✓	✓		✓	✓
Princess	✓	✓		✓							
Satallax							✓				
SPASS	✓										
Vampire	✓	✓	✓	✓	✓	✓	✓			✓	
veriT				✓							
Z3	✓	✓		✓							
Zenon	✓										
Zipperposition	✓	✓			✓		✓			✓	

184 ■ E [49] is a first-order prover based on the superposition calculus. An extension called
 185 Ehoh supports some higher-order constructs [60]. We use E 2.6 with Ehoh. This version
 186 can parse TH0 but does not yet implement higher-order unification. Note that Vukmirović
 187 is a developer of E.

188 ■ *ENIGMA* [34] is a variant of E that uses machine learning for determining the order in
 189 which clauses are processed by E, thereby steering the search space traversal. ENIGMA’s
 190 models were trained on the TPTP library, which could give suboptimal performance on
 191 Sledgehammer benchmarks. We use version 0.5.1.

192 ■ *iProver* [37] is a first-order prover based on instantiation. Recently, a superposition
 193 module was added [29]. The developers prepared a version of the prover trained using
 194 HOL-ML [32] on Sledgehammer benchmarks—version post-3.5 (git revision 04c55471083).
 195 Compared with 3.5, this version adds an extended parser that supports TF0, TX0, and
 196 SMT2 benchmarks.

197 ■ *leanCoP* [44] is a first-order prover based on the clausal connection calculus, a goal-
 198 oriented refinement of the clausal tableau calculus. Implemented in a few lines of Prolog,
 199 it combines extreme minimalism with decent performance. We use version 2.2.

200 ■ *Leo-III* [51] is a higher-order prover based on the higher-order paramodulation calculus.
 201 Paramodulation is a variant of superposition. The provers E and CVC4 are invoked as
 202 “end-game” backends at intervals on a first-order encoding of the current clause set.
 203 The main developer, Alexander Steen, provided command-line options suited for our
 204 experiments and fixed a few issues we discovered. We use version 1.6.6.

205 ■ *Princess* [48] is a first-order SMT solver based on a tableau calculus and with support
 206 for several arithmetic and nonarithmetic theories. We use version 2021-11-15 with
 207 command-line options provided by the developer.



- 208 ■ *Satallax* [24] is a higher-order prover based on a tableau calculus guided by a SAT solver.
 209 Ehoh is invoked at regular intervals as an “end-game” backend on an applicative first-order
 210 (or λ -free higher-order) encoding of the current clause set. We use version 3.5.
- 211 ■ *SPASS* [61] is a first-order prover based on superposition. We use version 3.9.
- 212 ■ *Vampire* [38] is a higher-order prover based on superposition and instantiation. Its
 213 superposition calculus uses *SKBCI* combinators to represent λ -abstractions [9]. We use
 214 version 4.6.1.sl,³ a customized version provided by the developers, with command-line
 215 options also provided by them. This version uses *Z3* as a backend.
- 216 ■ *veriT* [23] is a first-order SMT solver based on the *CDCL(T)* calculus. A prototype
 217 supports the SMT3 format [1], but we use the standard version 2021.06-rmx. We invoke
 218 veriT with command-line options that were derived by the developers using a custom
 219 tool [50, Section 5.1].
- 220 ■ *Z3* [28] is a first-order SMT solver based on the *CDCL(T)* calculus. We use version 4.8.12.
 221 The solver was not optimized for these benchmarks.
- 222 ■ *Zenon* [22] is a first-order prover based on tableaux. We use version 0.7.1.
- 223 ■ *Zipperposition* [59] is a higher-order prover based on λ -superposition, a higher-order
 224 variant of the superposition calculus. The E prover is invoked as an “end-game” backend
 225 at regular intervals on an applicative first-order (or λ -free higher-order) encoding of the
 226 current clause set. We run it in a custom mode designed for low timeouts. We use version
 227 2.1. Note that Vukmirović and Blanchette are developers of Zipperposition.

228 4 Mirabelle

229 Mirabelle is a testing and evaluation tool included with Isabelle/HOL since version 2010. It
 230 was initially developed Sascha Böhme, Blanchette, and other Isabelle developers as a Perl
 231 script and some Standard ML code. It was used for many evaluations of Sledgehammer and
 232 automatic provers [6–8, 11–13, 16, 17, 20, 45, 47, 50, 53, 60], starting with Judgment Day [20]. It
 233 was also used to generate TPTP, CASC, and SMT-LIB benchmarks.

234 For Isabelle version 2021-1, Mirabelle was reimplemented by Wenzel and Desharnais. The
 235 new tool is implemented as part of Isabelle/Scala and is properly integrated with modern
 236 Isabelle concepts such as sessions and parallel proof checking.

237 The new Mirabelle is a command-line tool that takes four arguments as input:

- 238 ■ *An Isabelle session*: A session is a collection of Isabelle theory files forming a project. For
 239 example, each entry of the *Archive of Formal Proofs* (AFP) [14] constitutes a session.
- 240 ■ *A theory filter*: The filter selects a subset of the theory files from the session for further
 241 processing. For example, a theory filter may keep all theory files or only a specified one.
- 242 ■ *A goal filter*: The filter selects a subset of the goals within the selected theory files for
 243 further processing. For example, a goal filter may keep all goals or only the first n goals
 244 from the selected theories.
- 245 ■ *A list of actions*: Each action is applied to each selected goal and produces a final report.

³ <https://github.com/vprover/vampire/releases/tag/v4.6.1.sl>



246 Mirabelle runs the specified session using Isabelle, collecting all selected goals. At the
247 end, it applies the actions to the collected goals, using parallelism. Each action application
248 produces some output, which can contain more details than the final reports. A Mirabelle
249 *goal* consists of the Isabelle goals before and after a proof step, whether it is a one-line proof
250 (**by**) or not, the theory file, and the goal’s position in the file. A Mirabelle *action* consists of
251 two functions: *run* performs the action, and *finalize* produces the final report.

252 Mirabelle includes a number of predefined actions. The one we use is Sledgehammer.
253 Like the **sledgehammer** command, it supports many options. These can be specified on
254 the command line. The action can either run the automatic provers or only generate the
255 TPTP or SMT-LIB files without running the provers. It is this latter mode that we use for
256 our evaluation. Other Mirabelle actions include *arith*, *metis* [46], and Quickcheck [26]. Users
257 can register custom actions.

258 ► **Example 2.** The following command launches Mirabelle on the `Compiler.thy` file from
259 the VeriComp session and runs E for 30 seconds on the first 10 goals:

```
260 isabelle mirabelle -d '$AFP' -O output -T Compiler -m 10 \  
261 -A "sledgehammer[provers=e,timeout=30]" VeriComp
```

262 The Sledgehammer invocation on the first goal is identified by the label `0.sledgehammer`
263 `goal.by VeriComp.Compiler 61:1935`. The associated action output is

```
264 succeeded (26348+464) [e]:  
265 Try this: using L1.prog_behaves_def by auto (16 ms)  
266 none (sledgehammer): succeeded (0)
```

267 The second line gives an Isabelle proof that reconstructs the automatic prover’s proof in
268 Isabelle and the time needed for reconstruction.

269 5 Evaluation

270 For our evaluation, we employed Mirabelle to generate problems from 5000 goals originating
271 from 50 randomly selected entries of the AFP (100 goals per entry). Within an entry, the
272 goals were selected at regular intervals, alleviating the issue that consecutive goals tend to
273 be similar. We used the repository revisions `b87fcf474e7f` of Isabelle (15 January 2022) and
274 `e2ae9549a7b0` of the AFP (19 December 2021). The experiments were run on the StarExec
275 Iowa cluster [52] with a CPU timeout of 30 s per problem.

276 5.1 Base Configuration

277 Table 4 presents the number of proved problems (out of 5000) for the 17 provers and the
278 different supported input formats. In this and later tables, the maximum of each row is
279 shown in bold, and the maximum of each column is shown in italics.

280 The problems were generated using a *base configuration* of Sledgehammer, which sets the
281 following options to provide a reasonable baseline for the experiments:

- 282 ■ The *fact_filter* option, which specifies the relevance filter, is set to MePo.
- 283 ■ The *max_facts* option, which controls the target number of Isabelle facts (definitions,
284 lemmas, etc.) to include in generated problems, is set to 512—a large number chosen to
285 stress the provers.



■ **Table 4** Proved problems for each prover and each input format, using the base configuration

	FOF	TF0	TX0	SMT2	TF1	TX1	TH0 ⁻	TH0 ⁺	SMT3	TH1 ⁻	TH1 ⁺
agsyHOL							814				
Beagle	901	1064		853							
cocATP							584				
cvc5	<i>2619</i>	2779		<i>2631</i>			2522		<i>2557</i>		
E	2394	2456	2534				2557				
ENIGMA	2301										
iProver	2291	2418	2211	2144							
leanCoP	1487										
Leo-III	1938	2136	2084		503	459	1632	<i>1539</i>		608	<i>447</i>
Princess	1205	1353		1260							
Satallax							1695				
SPASS	1550										
Vampire	2495	2608	<i>2587</i>	2204	<i>2471</i>	<i>2212</i>	2179			<i>1814</i>	
veriT				2463							
Z3	2226	2249		2252							
Zenon	812										
Zipperposition	2551	2607			1702		2718			1674	

- 286 ■ The *induction_rules* option, which controls the translation of induction rules, is set so
287 that these rules are excluded from generated problems.
- 288 ■ The *uncurried_aliases* option, which controls the translation of curried function applica-
289 tions in the TPTP module, is set to “false.”
- 290 ■ The *lam_trans* option, which controls the translation scheme for λ -abstractions in the
291 TPTP module, is set to use λ -lifting with defining equations for first-order formats and
292 to keep the λ 's as is for higher-order formats.
- 293 ■ The *type_enc* option, which controls the translation scheme for types in the TPTP module,
294 is set as follows: For polymorphic formats, keep the types as is. For monomorphic formats,
295 monomorphize the types. For FOF, use Sledgehammer’s efficient $\tilde{g}??$ encoding [12].
- 296 ■ The *max_mono_iters* and *max_new_mono_instances* options, which limit the number
297 of axiom instances generated by monomorphization, are set to 3 and 100, respectively.

298 The other Sledgehammer options are left with their default values [10]. In the next subsections,
299 we will deviate from this baseline to study the effect of various options.

300 In Table 4, we see that the best prover in the base configuration is *cvc5*. Remarkably,
301 this success is achieved with the TF0 format and not with the arithmetic-capable SMT2.
302 The situation is similar for *iProver*, *Princess*, and *Vampire*, but not *Z3*.

303 Intuitively, we would also expect support for higher-order logic to be beneficial, but this
304 is not always the case. *E* and *Zipperposition* perform better with the higher-order TH0⁻
305 format than with the first-order TF0, but for *cvc5*, *Leo-III*, and *Vampire* the opposite is true.
306 Similar effects are visible with formats supporting if-then-else and let constructs (TX0, TX1,
307 TH0⁺, TH1⁺).

308 Is higher-order logic useful at all? It would appear so. The data underlying Table 4
309 reveals that among the 3321 goals that are collectively proved by all provers and formats,
310 146 goals are proved with higher-order formats but not with first-order formats.



■ **Table 5** Running time of each prover’s best input format

Prover	Format	Proved problems	Percentile						
			25	50	75	90	95	99	100
agsyHOL	TH0 ⁻	814	0.5	1.5	4.2	9.0	15.3	24.9	29.5
Beagle	TF0	1064	5.5	7.0	9.6	12.3	13.7	15.4	16.8
cocATP	TH0 ⁻	584	2.4	4.9	11.1	20.1	24.6	28.4	30.1
cvc5	TF0	2779	0.2	0.3	0.5	2.1	7.6	24.4	30.0
E	TH0 ⁻	2557	0.2	0.4	1.1	9.5	13.5	23.6	29.2
ENIGMA	FOF	2301	2.7	2.9	3.4	4.5	6.2	9.0	15.1
iProver	TF0	2418	1.3	1.3	4.3	7.8	12.8	25.9	29.9
leanCoP	FOF	1487	1.4	1.4	3.5	11.6	12.7	26.3	29.2
Leo-III	TF0	2136	5.5	6.5	7.6	8.7	9.2	10.3	12.0
Princess	TF0	1353	6.8	7.6	8.5	9.2	9.5	9.8	11.0
Satallax	TH0 ⁻	1695	2.6	5.8	12.7	20.0	21.1	24.4	28.6
SPASS	FOF	1550	1.7	3.7	10.0	18.1	22.9	28.8	29.7
Vampire	TF0	2608	0.2	0.3	5.3	11.0	12.8	23.0	29.4
veriT	SMT2	2463	0.1	0.1	0.2	1.1	4.6	19.6	29.2
Z3	SMT2	2252	0.2	0.2	0.3	0.6	2.3	17.4	29.8
Zenon	FOF	812	0.6	1.6	3.2	6.3	11.4	24.7	29.0
Zipperposition	TH0 ⁻	2718	0.9	2.1	4.4	6.9	10.2	23.1	30.0

311 5.2 Running Time

312 How quickly do the provers find the proofs? Table 5 answers this question for each prover,
 313 focusing on the most successful input format for the prover according to Table 4. The median,
 314 or 50th percentile, is shown, as well as other percentiles.⁴ For example, the number 12.3 in
 315 the Beagle row indicates that 90% of the 1064 problems proved by Beagle are proved within
 316 12.3 s. As in Table 4, the baseline configuration is used. Recall that the timeout is 30 s.

317 Table 5 shows that if a prover finds a proof, it will likely find it quickly. Some provers
 318 are dazzlingly fast; the median time it takes to find a proof is 0.1 s for veriT, 0.2 s for Z3,
 319 and 0.3 s for cvc5 and Vampire. Remarkably, Princess proves its 1353 problems within 11 s;
 320 the remaining 19 s are unnecessary.

321 5.3 Number of Facts

322 An important Sledgehammer option is `num_facts`, which governs the number of facts to be
 323 selected by the relevance filter. These facts are included as axioms in the generated problems.
 324 The base configuration specifies 512 facts, but Table 6 shows what happens when we change
 325 this option (and keep the other options as in the base configuration). Figure 1 depicts the
 326 same information graphically for six provers, using a logarithmic x -axis.

327 The stronger provers tend to peak with more facts than the weaker ones. For example,
 328 the strongest prover, cvc5, peaks at 512 facts. Some of the stronger provers (E, Vampire,
 329 Zipperposition) implement the SInE [31] algorithm; others simply seem to scale well in the
 330 presence of extraneous axioms.

⁴ Percentiles are calculated using linear interpolation.

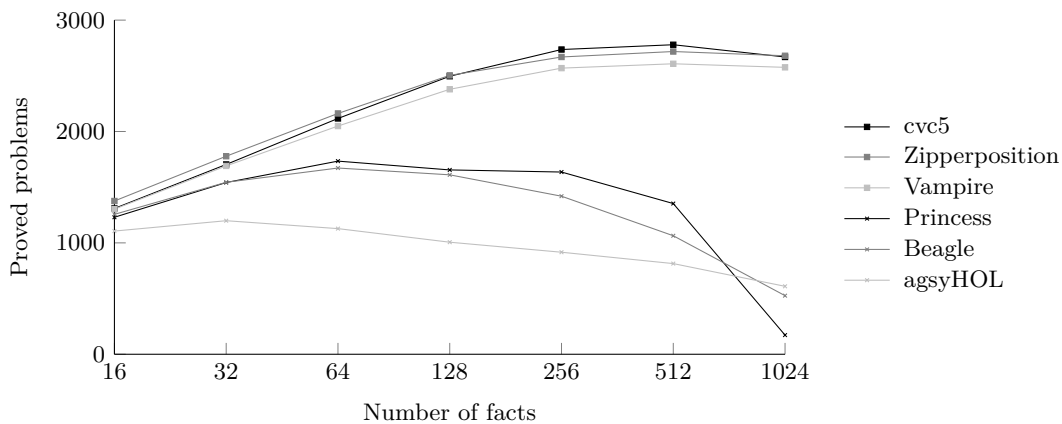
■ **Table 6** Proved problems for each prover’s best input format and different numbers of facts

	Format	16	32	64	128	256	512	1024
agsyHOL	TH0 ⁻	1106	1199	1128	1006	916	814	610
Beagle	TF0	1256	1544	1672	1611	1419	1064	526
cocATP	TH0 ⁻	745	830	837	793	706	584	440
cvc5	TF0	1309	1704	2117	2496	<i>2736</i>	2779	2669
E	TH0 ⁻	1346	1736	2116	2438	2590	2557	2368
ENIGMA	FOF	1297	1670	2004	2246	2347	2301	2137
iProver	TF0	1294	1667	2009	2300	2431	2418	2303
leanCoP	FOF	1150	1376	1510	1581	1562	1487	1356
Leo-III	TF0	1305	1680	2043	2333	2385	2136	1344
Princess	TF0	1230	1542	1734	1655	1636	1353	172
Satallax	TH0 ⁻	1233	1459	1596	1702	1725	1695	1632
SPASS	FOF	1273	1578	1778	1843	1787	1550	973
Vampire	TF0	1303	1691	2049	2379	2569	2608	2576
veriT	SMT2	1312	1683	2020	2307	2446	2463	2382
Z3	SMT2	1314	1668	1964	2207	2310	2252	2150
Zenon	FOF	955	1048	1067	1009	899	812	712
Zipperposition	TH0 ⁻	<i>1376</i>	<i>1778</i>	<i>2162</i>	<i>2504</i>	2699	2718	<i>2680</i>

331 5.4 Encoding of Lambdas

332 Not all provers support λ -abstractions. How practical are the encodings of λ ’s implemented in
 333 Sledgehammer, and is native support for λ ’s preferable when available? Table 7 presents a very
 334 fragmented picture: Some provers prefer λ -lifting (“lift”) to SKBCI combinators (“combs”),
 335 while others prefer SKBCI combinators or a combination of λ -lifting and combinators. Some
 336 prefer opaque definitions, while others work better when the defining equations are present.

337 Particularly striking is that only two higher-order provers, E and Zipperposition, work
 338 best with TH0⁻ and native λ ’s, whereas the other three prefer TF0 encodings. This is
 339 disappointing; we would have hoped that native higher-order support in a prover pays off.
 340 For Leo-III and Satallax, which rely on backends, there is a potential explanation: The



■ **Figure 1** Proved problems for the top provers’ best input formats and different numbers of facts

■ **Table 7** Proved problems by prover for the different encodings of λ -abstractions

Prover	Lift	Opaque lift	TF0			TH0 ⁻
			Combs	Opaque combs	Lift & combs	Native
Beagle	1064	1069	998	1026	930	
cvc5	<i>2779</i>	<i>2628</i>	<i>2838</i>	<i>2793</i>	2859	2522
E	2456	2331	2471	2444	2434	2557
iProver	2418	2320	2423	2430	2388	
Leo-III	2136	2093	2133	2152	2061	1632
Princess	1353	1369	1319	1369	1119	
Vampire	2608	2461	2602	2610	2579	2179
Z3	2249	2154	2326	2299	2325	
Zipperposition	2607	2118	2635	2590	2642	2718

■ **Table 8** Proved problems by prover for the different polymorphism-preserving encodings

Prover	FOF								TF1
	g	g?	g??	g@	t	t?	t??	t@	native
Beagle	446	560	637	514	544	650	640	487	
cvc5	1660	<i>2289</i>	<i>2552</i>	1936	<i>2092</i>	1989	2654	<i>2216</i>	
E	1451	2080	2195	1646	1800	1820	2162	1390	
ENIGMA	1366	1974	2063	1531	1570	1599	2030	1299	
iProver	1587	2016	2178	1756	2087	<i>1992</i>	2167	1368	
leanCoP	1053	1276	1455	1212	1565	1356	1451	763	
Leo-III	871	1075	1504	1076	829	890	1521	1252	503
Princess	47	344	835	136	214	416	798	454	
SPASS	755	1032	1194	935	941	859	1168	784	
Vampire	<i>1918</i>	2172	2312	<i>2032</i>	1673	1975	2294	1528	2471
Z3	1472	2054	2238	1677	1765	1692	2254	1844	
Zenon	606	634	626	605	424	564	628	403	
Zipperposition	1689	1907	2042	1660	1941	1951	2267	1631	1702

341 backends can work directly with TF0 problems but not with TH0. For Vampire, a possible
 342 explanation is that the higher-order version of the prover was not as extensively tuned as the
 343 first-order version.

344 5.5 Encoding of Types

345 Some provers do not support types at all, or they support only monomorphic types. How
 346 practical are the type encodings implemented in Sledgehammer, and is native support for
 347 types preferable when available? Table 8 presents the results for the *polymorphism-preserving*
 348 encodings—i.e., encodings that encode the rank-1 polymorphic type system in its full
 349 generality. In contrast, Table 9 presents the results for the *monomorphizing* encodings—i.e.,
 350 encodings that first heuristically instantiate type variables to eliminate polymorphism. The
 351 last column of each table offers a comparison with native TF1 or TF0.

352 Our findings are similar to those of Blanchette et al. [12]: Despite being incomplete,
 353 monomorphizing encodings outperform the polymorphism-preserving encodings. Overall,



© Martin Desharnais, Petar Vukmirović, Jasmin Blanchette, and Makarius Wenzel;
 licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:12–23:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

■ **Table 9** Proved problems by prover for the different monomorphizing encodings

Prover	FOF						TF0
	\tilde{g}	$\tilde{g}?$	$\tilde{g}??$	\tilde{t}	$\tilde{t}?$	$\tilde{t}??$	native
Beagle	473	842	901	544	887	880	1064
cvc5	2033	2479	2619	2092	2525	2687	2779
E	2133	2324	2394	1800	2303	2371	2456
ENIGMA	2065	2239	2301	1570	2220	2283	
iProver	1963	2171	2291	2087	2285	2278	2418
leanCoP	1136	1367	1487	1565	1533	1460	
Leo-III	1236	1744	1938	829	1814	1983	2136
Princess	159	923	1205	214	1010	1188	1353
SPASS	1318	1396	1550	941	1329	1533	
Vampire	2302	2426	2495	1673	2385	2484	2608
Z3	2127	2200	2226	1765	2048	2223	2249
Zenon	812	818	812	424	745	788	
Zipperposition	2360	2289	2551	1941	2455	2522	2607

the best encodings are the monomorphizing $\tilde{g}??$ and $\tilde{t}??$. But an even better option is to monomorphize the problem and use the many-sorted TF0 format directly. We also see that Leo-III’s and Zipperposition’s native polymorphism underperforms. The reason is probably that they rely on monomorphic backends. For example, Zipperposition disables its E backend when invoked on a polymorphic problem.

5.6 Portfolio

So far, we have evaluated only the performance of provers in isolation. What happens if we combine them? It turns out that a virtual portfolio consisting of all the provers in all the configurations of Tables 4 to 9, each invoked for 30 s, would prove 3508 goals. This corresponds to a success rate of 70.1%, which is clearly lower than the 76.7% figure obtained in the MaSh paper [13] on the Judgment Day benchmarks. Two possible explanations suggest themselves: The MaSh paper’s evaluation used the machine learning filters MaSh and MeSh, which are stronger than MePo, and as observed elsewhere [14, Section 6] Judgment Day might consist of relatively easy goals.

Such a virtual portfolio is not very realistic, but it does give an idea of the state of the art in automated reasoning. A more realistic alternative is to consider the *greedy sequence* of length n . The sequence is obtained by iteratively (1) taking first the (or some) best prover configuration for the goals of interest, and (2) removing all goals proved by this configuration. These two steps are repeated n times, yielding n configurations. The greedy sequence is not necessarily optimal, but it can be computed efficiently.

Table 10 presents the greedy sequence of length 16 based on our experiments. Given 480 s—or 30 s and 16 threads—we can solve 3440 goals. This shows how complementary the different provers and options are.

6 Related Work

The various evaluations of Sledgehammer, cited at the beginning of Section 4, surely constitute the most closely related work. Among these, Böhme and Nipkow’s Judgment Day study [20]



© Martin Desharnais, Petar Vukmirović, Jasmin Blanchette, and Makarius Wenzel; licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:13–23:18



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

■ **Table 10** A greedy sequence of provers and configurations

Prover	Format	Configuration	Proved problems
cvc5	TF0	lift & combs	2859
Zipperposition	TH0 ⁻	1024 facts	+230
Vampire	TF1	base	+77
veriT	SMT2	1024 facts	+60
E	TX0	base	+51
cvc5	TF0	128 facts	+38
Zipperposition	TH0 ⁻	256 facts	+24
Beagle	SMT2	base	+20
cvc5	FOF	t??	+17
Vampire	TF0	1024 facts	+13
Z3	SMT2	1024 facts	+11
E	TH0 ⁻	32 facts	+11
Vampire	TH0 ⁻	base	+9
E	TH0 ⁻	1024 facts	+9
Z3	SMT2	base	+6
Zipperposition	TH0 ⁻	base	+5
Total			3440

stands out as the most comprehensive; it considers the following aspects: success rate, running time, proof complexity, and proof minimization. But it is over a decade old. Most of the other evaluations focus on a single new Sledgehammer feature and how it improves the success rate; for example, Sultana et al. [53] evaluate encodings of higher-order features (cf. Section 5.4), and Blanchette et al. [12] evaluate type encodings (cf. Section 5.5).

Similar evaluations for other hammers or hammer-like systems include a three-prover evaluation by Kaliszyk and Urban [36], called MizAR 40, using MizAR, a three-prover evaluation by Czajka and Kaliszyk [27] using CoqHammer, an eight-prover evaluation by Filliâtre [30] using Why3, and a 19-prover evaluation by Brown et al. [25], called GRUNGE, using a custom exporter from HOL4. Among these, GRUNGE is the most similar to our effort, as it includes a large number of provers and exploits support for higher-order logic and polymorphism where available. But there are also several important differences:

- *The benchmarks presented in the GRUNGE paper include only the set of lemmas needed for a proof in HOL4, whereas our benchmarks contain heuristically selected lemmas from Isabelle’s libraries.* Including only a typically small set of lemmas makes the benchmarks easier to prove than selecting hundreds of lemmas heuristically. It is less representative of the typical hammer use case, where a proof is not available.
- *The GRUNGE benchmarks correspond to top-level lemmas or theorems in HOL4, whereas our benchmarks correspond to Isabelle goals or subgoals.* Proving a lemma is generally more difficult than proving a goal.
- *GRUNGE’s encoding of polymorphism is simple but inefficient, whereas our evaluation relies on state-of-the-art monomorphization.* One of GRUNGE’s two translation schemes [25, Section 4.3] corresponds to Sledgehammer’s t encoding, whose performance was found to be very suboptimal in Section 5.5.
- *GRUNGE uses different provers to our evaluation.* In particular, GRUNGE does not generate SMT-LIB problems, nor does it exploit the recently added support for higher-order



© Martin Desharnais, Petar Vukmirović, Jasmin Blanchette, and Makarius Wenzel; licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:14–23:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

406 logic in `cvc5`, `E`, and `Vampire`. It also misses out on the improvements to Zipperposition
407 over the past three years, which led it to win the higher-order division of the 2020 and
408 2021 editions of CASC [55, 56].

409 The top five provers on the GRUNGE benchmarks, with the number of proved problems
410 out of 12 140, are Leo-III with 7090, Vampire with 5929, CVC4 with 5709, E with 5118, and
411 HOLyHammer with 5059. The top prover is only the ninth best according to our Table 4.
412 The discrepancy is easy to explain: For GRUNGE, Leo-III was, next to Zipperposition, the
413 only system that fully supported polymorphic higher-order logic, and all the other systems
414 saw versions of the benchmarks encoded using encodings that were not primarily tuned
415 toward performance of particular provers, whereas for our evaluation we used state-of-the-art
416 encodings of polymorphic types and higher-order features.

417 As a benchmark suite, Seventeen can be used as a complement to existing libraries,
418 such as the TPTP [54], which includes both first- and higher-order problems in the TPTP
419 formats, and SMT-LIB [2], which consists of first-order problems in SMT-LIB 2 syntax.
420 Various collections of hammer-generated problems exist, such as GRUNGE, Judgment Day,
421 and MizAR 40. All these libraries are not necessarily mutually exclusive; for example,
422 many MizAR- and Sledgehammer-generated problems are in the TPTP, and Sledgehammer-
423 generated problems [47, Section 5] are also part of SMT-LIB.

424 **7 Conclusion**

425 In this paper, we evaluated 17 modern automatic provers, including SMT solvers and higher-
426 order provers, thereby superseding the aging Judgment Day study; and we evaluated several
427 aspects of a hammer, such as the encodings of types and λ -abstractions. In particular, we
428 wanted to determine whether native implementations of various features perform better than
429 encodings, which is a reasonable thing to hope for. We have now updated Sledgehammer's
430 default setup to get the best out of the provers.

431 Specifically, we found that native support for monomorphic types was beneficial, but the
432 current support for polymorphism is disappointing. Native support for features such as linear
433 arithmetic and higher-order logic helps some provers and is detrimental to others. Overall,
434 we see that simply implementing a feature in a prover is often not enough; to be useful,
435 the feature must be finely-tuned based on benchmarks. Some provers seem to misbehave
436 on Sledgehammer benchmarks, which indeed may look quite different from the problems
437 used by the provers' developers to tune their systems. Since hammers are among the most
438 useful applications of automatic provers, we contend that it is worthwhile to tune provers on
439 hammer-generated benchmarks as a complement to the TPTP and SMT-LIB.

440 **References**

- 441 1 Haniel Barbosa, Andrew Reynolds, Daniel El Ouaoui, Cesare Tinelli, and Clark W. Barrett.
442 Extending SMT solvers to higher-order logic. In Pascal Fontaine, editor, *CADE-27*, volume
443 11716 of *LNCS*, pages 35–54. Springer, 2019.
- 444 2 Clark Barrett, Pascal Fontaine, and Cesare Tinelli. The SMT-LIB standard: Version 2.6.
445 Technical report, 2017. URL: <https://www.SMT-LIB.org/>.
- 446 3 Clark W. Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanovic,
447 Tim King, Andrew Reynolds, and Cesare Tinelli. CVC4. In Ganesh Gopalakrishnan and Shaz
448 Qadeer, editors, *CAV 2011*, volume 6806 of *LNCS*, pages 171–177. Springer, 2011.



© Martin Desharnais, Petar Vukmirović, Jasmin Blanchette, and Makarius Wenzel;
licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:15–23:18

Leibniz International Proceedings in Informatics



LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

- 449 4 Clark W. Barrett, Leonardo de Moura, and Aaron Stump. SMT-COMP: satisfiability modulo
450 theories competition. In Kousha Etessami and Sriram K. Rajamani, editors, *CAV 2005*, volume
451 3576 of *LNCS*, pages 20–23. Springer, 2005.
- 452 5 Peter Baumgartner, Joshua Bax, and Uwe Waldmann. Beagle—a hierarchic superposition
453 theorem prover. In Amy P. Felty and Aart Middeldorp, editors, *CADE-25*, volume 9195 of
454 *LNCS*, pages 367–377. Springer, 2015.
- 455 6 Alexander Bentkamp, Jasmin Blanchette, Simon Cruanes, and Uwe Waldmann. Superposition
456 for lambda-free higher-order logic. *Log. Meth. Comput. Sci.*, 17(2):1:1–1:38, 2021.
- 457 7 Alexander Bentkamp, Jasmin Blanchette, Sophie Touret, and Petar Vukmirović. Superposition
458 for full higher-order logic. In André Platzer and Geoff Sutcliffe, editors, *CADE-28*, volume
459 12699 of *LNCS*, pages 396–412. Springer, 2021.
- 460 8 Alexander Bentkamp, Jasmin Blanchette, Sophie Touret, Petar Vukmirovic, and Uwe Wald-
461 mann. Superposition with lambdas. *J. Autom. Reason.*, 65(7):893–940, 2021.
- 462 9 Ahmed Bhayat and Giles Reger. A combinator-based superposition calculus for higher-order
463 logic. In Nicolas Peltier and Viorica Sofronie-Stokkermans, editors, *IJCAR 2020, Part I*,
464 volume 12166 of *LNCS*, pages 278–296. Springer, 2020.
- 465 10 Jasmin Blanchette. Hammering away: A user’s guide to Sledgehammer for Isabelle/HOL,
466 2021. [https://isabelle.in.tum.de/website-Isabelle2021-1/dist/Isabelle2021-1/doc/
467 sledgehammer.pdf](https://isabelle.in.tum.de/website-Isabelle2021-1/dist/Isabelle2021-1/doc/sledgehammer.pdf).
- 468 11 Jasmin Christian Blanchette, Sascha Böhme, and Lawrence C. Paulson. Extending Sledge-
469 hammer with SMT solvers. *J. Autom. Reason.*, 51(1):109–128, 2013.
- 470 12 Jasmin Christian Blanchette, Sascha Böhme, Andrei Popescu, and Nicholas Smallbone. En-
471 coding monomorphic and polymorphic types. *Log. Meth. Comput. Sci.*, 12(4), 2016.
- 472 13 Jasmin Christian Blanchette, David Greenaway, Cezary Kaliszyk, Daniel Kühlwein, and Josef
473 Urban. A learning-based fact selector for Isabelle/HOL. *J. Autom. Reason.*, 57(3):219–244,
474 2016.
- 475 14 Jasmin Christian Blanchette, Maximilian Haslbeck, Daniel Matichuk, and Tobias Nipkow.
476 Mining the archive of formal proofs. In Manfred Kerber, editor, *CICM 2015*, volume 9150 of
477 *LNCS*, pages 1–15. Springer, 2015.
- 478 15 Jasmin Christian Blanchette, Cezary Kaliszyk, Lawrence C. Paulson, and Josef Urban. Ham-
479 mering towards QED. *J. Formaliz. Reason.*, 9(1):101–148, 2016.
- 480 16 Jasmin Christian Blanchette, Nicolas Peltier, and Simon Robillard. Superposition with
481 datatypes and codatatypes. In Didier Galmiche, Stephan Schulz, and Roberto Sebastiani,
482 editors, *IJCAR 2018*, volume 10900 of *LNCS*, pages 370–387. Springer, 2018.
- 483 17 Jasmin Christian Blanchette, Andrei Popescu, Daniel Wand, and Christoph Weidenbach.
484 More SPASS with Isabelle: Superposition with hard sorts and configurable simplification. In
485 Lennart Beringer and Amy Felty, editors, *ITP 2012*, volume 7406 of *LNCS*, pages 345–360.
486 Springer, 2012.
- 487 18 François Bobot and Andrei Paskevich. Expressing polymorphic types in a many-sorted
488 language. In Cesare Tinelli and Viorica Sofronie-Stokkermans, editors, *FroCoS 2011*, volume
489 6989 of *LNCS*, pages 87–102. Springer, 2011.
- 490 19 Sascha Böhme. *Proving Theorems of Higher-Order Logic with SMT Solvers*. PhD thesis,
491 Technische Universität München, 2012.
- 492 20 Sascha Böhme and Tobias Nipkow. Sledgehammer: Judgement Day. In Jürgen Giesl and
493 Reiner Hähnle, editors, *IJCAR 2010*, volume 6173 of *LNCS*, pages 107–121. Springer, 2010.
- 494 21 Sascha Böhme and Tjark Weber. Fast LCF-style proof reconstruction for Z3. In Matt
495 Kaufmann and Lawrence C. Paulson, editors, *ITP 2010*, volume 6172 of *LNCS*, pages 179–194.
496 Springer, 2010.
- 497 22 Richard Bonichon, David Delahaye, and Damien Doligez. Zenon: An extensible automated
498 theorem prover producing checkable proofs. In Nachum Dershowitz and Andrei Voronkov,
499 editors, *LPAR 2007*, volume 4790 of *LNCS*, pages 151–165. Springer, 2007.



- 500 23 Thomas Bouton, Diego Caminha Barbosa De Oliveira, David Déharbe, and Pascal Fontaine.
501 veriT: An open, trustable and efficient SMT-solver. In Renate A. Schmidt, editor, *CADE-22*,
502 volume 5663 of *LNCS*, pages 151–156. Springer, 2009.
- 503 24 Chad E. Brown. Satallax: An automatic higher-order prover. In Bernhard Gramlich, Dale
504 Miller, and Uli Sattler, editors, *IJCAR 2012*, volume 7364 of *LNCS*, pages 111–117. Springer,
505 2012.
- 506 25 Chad E. Brown, Thibault Gauthier, Cezary Kaliszyk, Geoff Sutcliffe, and Josef Urban.
507 GRUNGE: A grand unified ATP challenge. In Pascal Fontaine, editor, *CADE-27*, volume
508 11716 of *LNCS*, pages 123–141. Springer, 2019.
- 509 26 Lukas Bulwahn. The new Quickcheck for Isabelle—random, exhaustive and symbolic testing
510 under one roof. In Chris Hawblitzel and Dale Miller, editors, *CPP 2012*, volume 7679 of *LNCS*,
511 pages 92–108. Springer, 2012.
- 512 27 Łukasz Czajka and Cezary Kaliszyk. Hammer for Coq: Automation for dependent type theory,
513 2018.
- 514 28 Leonardo Mendonça de Moura and Nikolaž Bjørner. Z3: An efficient SMT solver. In C. R.
515 Ramakrishnan and Jakob Rehof, editors, *TACAS 2008*, volume 4963 of *LNCS*, pages 337–340.
516 Springer, 2008.
- 517 29 André Duarte and Konstantin Korovin. Implementing superposition in iProver (system
518 description). In Nicolas Peltier and Viorica Sofronie-Stokkermans, editors, *IJCAR 2020, Part*
519 *II*, volume 12167 of *LNCS*, pages 388–397. Springer, 2020.
- 520 30 Jean-Christophe Filiâtre. One logic to use them all. In Maria Paola Bonacina, editor,
521 *CADE-24*, volume 7898 of *LNCS*, pages 1–20. Springer, 2013.
- 522 31 Krystof Hoder and Andrei Voronkov. Sine qua non for large theory reasoning. In Nikolaj
523 Bjørner and Viorica Sofronie-Stokkermans, editors, *CADE-23*, volume 6803 of *LNCS*, pages
524 299–314. Springer, 2011.
- 525 32 Edvard K. Holden and Konstantin Korovin. Heterogeneous heuristic optimisation and schedul-
526 ing for first-order theorem proving. In Fairouz Kamareddine and Claudio Sacerdoti Coen,
527 editors, *CICM 2021*, volume 12833 of *LNCS*, pages 107–123. Springer, 2021.
- 528 33 R. J. M. Hughes. Super-combinators: A new implementation method for applicative languages.
529 In *LFP 1982*, pages 1–10. ACM Press, 1982.
- 530 34 Jan Jakubův and Josef Urban. ENIGMA: efficient learning-based inference guiding machine. In
531 Herman Geuvers, Matthew England, Osman Hasan, Florian Rabe, and Olaf Teschke, editors,
532 *CICM 2017*, volume 10383 of *LNCS*, pages 292–302. Springer, 2017.
- 533 35 Cezary Kaliszyk and Josef Urban. HOL(y)Hammer: Online ATP service for HOL Light. *Math.*
534 *Comput. Sci.*, 9(1):5–22, 2015.
- 535 36 Cezary Kaliszyk and Josef Urban. MizAR 40 for Mizar 40. *J. Autom. Reason.*, 55(3):245–256,
536 2015.
- 537 37 Konstantin Korovin. iProver—an instantiation-based theorem prover for first-order logic
538 (system description). In Alessandro Armando, Peter Baumgartner, and Gilles Dowek, editors,
539 *IJCAR 2008*, volume 5195 of *LNCS*, pages 292–298. Springer, 2008.
- 540 38 Laura Kovács and Andrei Voronkov. First-order theorem proving and Vampire. In Natasha
541 Sharygina and Helmut Veith, editors, *CAV 2013*, volume 8044 of *LNCS*, pages 1–35. Springer,
542 2013.
- 543 39 Fredrik Lindblad. A focused sequent calculus for higher-order logic. In Stéphane Demri,
544 Deepak Kapur, and Christoph Weidenbach, editors, *IJCAR 2014*, volume 8562 of *LNCS*, pages
545 61–75. Springer, 2014.
- 546 40 Fredrik Lindblad and Marcin Benke. A tool for automated theorem proving in Agda. In
547 Jean-Christophe Filiâtre, Christine Paulin-Mohring, and Benjamin Werner, editors, *TYPES*
548 *2004*, volume 3839 of *LNCS*, pages 154–169. Springer, 2004.
- 549 41 Jia Meng and Lawrence C. Paulson. Translating higher-order clauses to first-order clauses. *J.*
550 *Autom. Reason.*, 40(1):35–60, 2008.



- 551 42 Jia Meng and Lawrence C. Paulson. Lightweight relevance filtering for machine-generated
552 resolution problems. *J. Applied Logic*, 7(1):41–57, 2009.
- 553 43 Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL: A Proof Assistant*
554 *for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002.
- 555 44 Jens Otten. leanCoP 2.0 and ileanCoP 1.2: High performance lean theorem proving in classical
556 and intuitionistic logic (system descriptions). In Alessandro Armando, Peter Baumgartner,
557 and Gilles Dowek, editors, *IJCAR 2008*, volume 5195 of *LNCS*, pages 283–291. Springer, 2008.
- 558 45 Lawrence C. Paulson and Jasmin Christian Blanchette. Three years of experience with
559 Sledgehammer, a practical link between automatic and interactive theorem provers. In Geoff
560 Sutcliffe, Stephan Schulz, and Eugenia Ternovska, editors, *IWIL-2010*, volume 2 of *EPiC*,
561 pages 1–11. EasyChair, 2012.
- 562 46 Lawrence C. Paulson and Kong Woei Susanto. Source-level proof reconstruction for interactive
563 theorem proving. In Klaus Schneider and Jens Brandt, editors, *TPHOLs 2007*, volume 4732
564 of *LNCS*, pages 232–245. Springer, 2007.
- 565 47 Andrew Reynolds and Jasmin Christian Blanchette. A decision procedure for (co)datatypes in
566 SMT solvers. *J. Autom. Reason.*, 58(3):341–362, 2017.
- 567 48 Philipp Rümmer. A constraint sequent calculus for first-order logic with linear integer
568 arithmetic. In Iliano Cervesato, Helmut Veith, and Andrei Voronkov, editors, *LPAR 2008*,
569 volume 5330 of *LNCS*, pages 274–289. Springer, 2008.
- 570 49 Stephan Schulz, Simon Cruanes, and Petar Vukmirovic. Faster, higher, stronger: E 2.3. In
571 Pascal Fontaine, editor, *CADE-27*, volume 11716 of *LNCS*, pages 495–507. Springer, 2019.
- 572 50 Hans-Jörg Schurr, Mathias Fleury, and Martin Desharnais. Reliable reconstruction of fine-
573 grained proofs in a proof assistant. In André Platzer and Geoff Sutcliffe, editors, *CADE-28*,
574 volume 12699 of *LNCS*, pages 450–467. Springer, 2021.
- 575 51 Alexander Steen and Christoph Benzmüller. Extensional higher-order paramodulation in
576 Leo-III. *J. Autom. Reason.*, 65(6):775–807, 2021.
- 577 52 Aaron Stump, Geoff Sutcliffe, and Cesare Tinelli. StarExec: A cross-community infrastructure
578 for logic solving. In Stéphane Demri, Deepak Kapur, and Christoph Weidenbach, editors,
579 *IJCAR 2014*, volume 8562 of *LNCS*, pages 367–373. Springer, 2014.
- 580 53 Nik Sultana, Jasmin Christian Blanchette, and Lawrence C. Paulson. LEO-II and Satallax on
581 the Sledgehammer test bench. *J. Applied Logic*, 11(1):91–102, 2013.
- 582 54 Geoff Sutcliffe. The TPTP problem library and associated infrastructure—from CNF to TH0,
583 TPTP v6.4.0. *J. Autom. Reason.*, 59(4):483–502, 2017.
- 584 55 Geoff Sutcliffe. The 10th IJCAR automated theorem proving system competition—CASC-J10.
585 *AI Commun.*, 34(2):163–177, 2021.
- 586 56 Geoff Sutcliffe and Martin Desharnais. The CADE-28 automated theorem proving system
587 competition—CASC-28. *AI Communications*, 34(4):259–276, 2022.
- 588 57 David A. Turner. A new implementation technique for applicative languages. *Softw. Pract.*
589 *Exper.*, 9:31–49, 1979.
- 590 58 Josef Urban, Piotr Rudnicki, and Geoff Sutcliffe. ATP and presentation service for Mizar
591 formalizations. *J. Autom. Reason.*, 50(2):229–241, 2013.
- 592 59 Petar Vukmirović, Alexander Bentkamp, Jasmin Blanchette, Simon Cruanes, Visa Nummelin,
593 and Sophie Turrett. Making higher-order superposition work. In André Platzer and Geoff
594 Sutcliffe, editors, *CADE-28*, volume 12699 of *LNCS*, pages 415–432. Springer, 2021.
- 595 60 Petar Vukmirović, Jasmin Christian Blanchette, Simon Cruanes, and Stephan Schulz. Extend-
596 ing a brainiac prover to lambda-free higher-order logic. In Tomas Vojnar and Lijun Zhang,
597 editors, *TACAS 2019*, volume 11427 of *LNCS*, pages 192–210. Springer, 2019.
- 598 61 Christoph Weidenbach, Dilyana Dimova, Arnaud Fietzke, Rohit Kumar, Martin Suda, and
599 Patrick Wischniewski. SPASS version 3.5. In Renate A. Schmidt, editor, *CADE-22*, volume
600 5663 of *LNCS*, pages 140–145. Springer, 2009.

