

# Three Years of Experience with Sledgehammer, a Practical Link between Automatic and Interactive Theorem Provers

Lawrence C. Paulson

Computer Laboratory  
University of Cambridge, U.K.  
lp15@cam.ac.uk and Jasmin Christian Blanchette  
Institut für Informatik  
Technische Universität München, Germany  
blanchette@in.tum.de

## Abstract

Sledgehammer is a highly successful subsystem of Isabelle/HOL that calls automatic theorem provers to assist with interactive proof construction. It requires no user configuration: it can be invoked with a single mouse gesture at any point in a proof. It automatically finds relevant lemmas from all those currently available. An unusual aspect of its architecture is its use of unsound translations, coupled with its delivery of results as Isabelle/HOL proof scripts: its output cannot be trusted, but it does not need to be trusted. Sledgehammer works well with Isar structured proofs and allows beginners to prove challenging theorems.

## 1 Introduction

Interactive theorem provers are widely used by researchers for modelling complex algorithms or systems. They typically support rich logical formalisms that include recursive functions and types. Some even make it possible to reason about a partial recursive function's domain of definition or to define a relation coinductively. The main weakness of these tools is the actual proving, which is extremely laborious.

For nearly 20 years, researchers have sought to make interactive theorem provers better at proving theorems. Much of this effort has been devoted to implementing decision procedures. Certainly decision procedures are essential, especially for arithmetic: without them, obvious identities can take hours to prove; with them, complicated facts can be proved instantly. But most problems lie beyond the scope of standard decision procedures. Automatic tools are needed that can work on any type of problem.

Automatic theorem provers (ATPs) are capable of creating long, incomprehensible chains of deduction. Many researchers have attempted to use them to support interactive theorem proving; particularly pertinent are Ahrendt et al. [1], Bezem et al. [7], Hurd [12], and Siekmann et al. [33]. But the only tool to pass the test of time is Sledgehammer [19, 28], which links Isabelle/HOL to the automatic provers E [32], SPASS [37], and Vampire [29]. Isabelle users invoke Sledgehammer routinely when undertaking difficult proofs. On a representative corpus of older proof scripts, we find that Sledgehammer can prove 38% of the goals that are not provable by the standard automatic tactics.

Sledgehammer was first released in February 2007 to users daring enough to download an Isabelle nightly build. It was announced in November 2007 as a component of Isabelle2007. This paper outlines the design goals that made Sledgehammer successful (§2) and presents an updated summary of Böhme and Nipkow's [8] empirical study (§3). It describes some of the lessons learnt in the past three years and its effect on the way Isabelle/HOL is taught (§4). Avenues for research are also discussed (§5). An earlier version of this paper was presented at the PAAR-2010 workshop [27].

## 2 Design Principles

The single most important design goal was one-click invocation. Some earlier systems required the user to gather up all the facts that could be relevant to the problem, and furthermore to reduce it to first-order form. Problem preparation could easily take hours, with no guarantee that the call to a first order theorem prover would succeed. Such a tool would be of little value to users.

The two aspects of problem preparation (translation into first-order logic and identification of relevant facts) each required a substantial research effort. The numerous choices outlined below were made on the basis of innumerable experiments that consumed many thousands of hours of processor time.

### 2.1 Translation into First-Order Logic

Most interactive theorem provers support a language much richer than that of first-order logic. Isabelle/HOL [21] supports polymorphic higher-order logic [2, 9, 10], augmented with axiomatic type classes [39].<sup>1</sup> Many user problems contain no higher-order features and might be imagined to lie within first-order logic; however, even these problems are full of typing information. Type information can take quadratic space [17] because every term must be annotated with its type, recursively, right down to the variables. Hurd [13] observed that omitting type information greatly improved the success rate of his theorem prover, Metis. This is hardly surprising, since the type information virtually buries the terms themselves. Hurd was able to omit type information because his proofs are reconstructed within HOL4 [10], which rejected any proofs that did not correspond to well-typed higher-order logic deductions.

Sledgehammer was always intended to rely on an analogous process of sound proof reconstruction, and from the outset it was clear that including complete type information would be unworkable. Completely omitting type information, although successful for HOL4, would not have worked for Isabelle because of its heavy use of type classes. We chose to include enough type information to enforce correct type class reasoning (the type class hierarchy is easily expressed using Horn clauses) but not to specify the type of every term [19, §4]. Some colleagues have expressed horror at the very idea of using unsound translations; the first author has written a lengthy exploration of the salient issues [17, §2.8].

Higher-order problems posed special difficulties. We never expected first-order theorem provers to perform deep higher-order reasoning, but merely hoped to automate proofs where the higher-order steps were trivial. We examined several methods of translating higher-order problems into first-order logic, allowing truth values to be the values of terms and curried functions to take varying numbers of arguments [17]. We eventually adopted a translation based on the one that we used for first-order logic, modified to introduce higher-order mechanisms (such as an “apply operator” for function values) only when absolutely necessary. We thereby eliminated our original distinction between first-order and higher-order problems. A higher-order feature within a problem affects the translation locally, yielding a smooth transition from purely first-order to heavily higher-order problems.

We also experimented with two methods of eliminating  $\lambda$ -abstractions in terms: by translating them into combinator form or by declaring equivalent functions. We ultimately opted for a naive translation scheme based on the combinators S, K, I, B, and C: more sophisticated schemes delivered no additional benefits. Unfortunately, our experience suggests that Sledgehammer is seldom successful on problems containing higher-order elements. Integration with a genuine higher-order automatic theorem prover, such as LEO-II [5] and Satallax [3], seems necessary. This would pose interesting problems for proof reconstruction: LEO-II’s approach is to reduce higher-order problems to first-order ones by repeatedly applying specialised inference rules and then calling first-order ATPs. A LEO-II proof will therefore

---

<sup>1</sup>Isabelle [26] is a generic theorem prover, based on a logical framework [23]. Isabelle/HOL is the specialisation of Isabelle for higher-order logic.

consist of a string of higher-order steps followed by a first-order proof. The latter part we know how to do; the crucial challenge is to devise a reliable way of emulating the higher-order steps within Isabelle.

Arithmetic remains an issue. A purely arithmetic problem can be solved using decision procedures, but what about problems that combine arithmetic with a significant amount of logic? In principle, Sledgehammer could solve such problems with the help of an ATP that combined arithmetic and logical reasoning, analogous to LEO-II's approach to higher-order logic. Current SMT solvers are probably of little value, because they do not handle quantified formulas well. But progress in that field is extremely rapid, and soon this option could become attractive.

## 2.2 Relevance Filtering

Our initial goals for Sledgehammer were modest: to improve upon Isabelle's built-in automatic tools, using only the lemma libraries used by those tools. There were two libraries, one consisting of facts useful for forward and backward chaining, the other consisting of rewriting rules for simplification. Each library contained hundreds of lemmas. We discovered that automatic theorem provers could solve only trivial problems in the presence so many extraneous facts; by developing a lightweight, symbol-based relevance filter, we greatly improved the success rate [18]. Users would still have to identify relevant facts that did not belong to these lemma libraries.

Tobias Nipkow made the crucial suggestion to dispense with the lemma libraries, substituting the full collection of Isabelle theorems (around 10,000). This idea offered the enticing prospect that any relevant existing theorem, however obscure, could be located. We thought this goal to be unrealistic; it seemed to have too much in common with McAllester's Ontic system [14]. Ontic was intended to be able to prove mathematical results using known results that it identified automatically, and it seems fair to say that this objective was too ambitious. But Sledgehammer would only be part of the system rather than all-encompassing, and it could take advantage of 20 years of increasing hardware performance.

We were able to scale up the relevance filter to cope with the 20-fold increase in the number of facts to process. However, it relies on ad hoc heuristics that sometimes deliver poor results. Briefly, it assigns a score to every available theorem based upon how many constants that theorem shares with the conjecture; this process iterates to include theorems relevant to those just accepted, but with a decay factor to ensure termination. The constants are weighted to give unusual ones greater significance. The relevance filter copes best when the conjecture contains some unusual constants; if all the constants are common, it is unable to discriminate among the hundreds of facts that are picked up. The relevance filter is also memoryless: it has no information about how many times a particular fact has been used in a proof, and it cannot learn.

It would obviously be preferable for the automatic theorem provers themselves to perform relevance filtering. Or we could use a sophisticated system based on machine learning, such as Josef Urban's MaLAREa [36], where successful proofs provide information to guide other proofs. Unfortunately, any such approach will fail given Sledgehammer's use of unsound translations. In unpublished work by Urban, MaLAREa easily proved the full Sledgehammer test suite by identifying an inconsistency in the translated lemma library; once MaLAREa had found the inconsistency in one proof, it easily found it in all the others. Sledgehammer is successful only because its relevance filter generally selects too few lemmas to produce an inconsistent axiom set, even with the unsound translations.

To accomplish better relevance filtering, we must decide whether to adopt a general first-order approach or to build a sophisticated relevance filter directly into Isabelle. The former approach could take advantage of the efforts of the entire ATP community, but it would have to be good enough to cope with soundly translated (and presumably enormous) formulas. The latter approach would avoid translation issues, but it would impose the entire effort onto a few Isabelle developers.

## 2.3 Proof Reconstruction

Isabelle subscribes to the LCF philosophy: all proofs ultimately reduce to primitives executed by a logical kernel. Isabelle users would not trust a tool that uncritically accepted proofs from an external source. But Sledgehammer had an even stronger design objective: to deliver Isabelle proofs in source form. We envisaged that Sledgehammer runs would demand substantial computational resources; if somebody used Sledgehammer many times while constructing a proof, would it be feasible to run that proof again, perhaps to modify it using a laptop while at a conference? To be useful, Sledgehammer would have to return a piece of proof script that could be executed cheaply.

### 2.3.1 Reconstruction of the Resolution Proof

The original plan was to emulate the inference rules of automatic theorem provers directly within Isabelle. We should have known better: Hurd [12] had noticed that the proofs delivered by Gandalf [35] were not detailed and explicit enough. We made the same discovery with SPASS and, despite considerable efforts, were only able to reconstruct a handful of proofs [19].

We came up with a new plan: to use a general theorem prover, Metis, to reconstruct each proof step. Metis was designed to be interfaced with LCF-style interactive theorem provers, specifically HOL4. Integrating it with Isabelle’s proof kernel required significant effort [28]. Metis then became available to Isabelle users, and it turned out to be capable of reconstructing proof steps easily. The output of Sledgehammer was now a list of calls to Metis, each of which proved a clause. While the output is primarily designed for replaying proofs, it also has a pedagogical value: unlike Isabelle’s automatic tactics, which are black boxes, the proofs delivered by Sledgehammer can be inspected and understood.

Consider the theorem “ $\text{length } (tl\ xs) \leq \text{length } xs$ ”, which states that the tail of a list (the list from which we remove its first element, or the empty list if the list is empty) is shorter than or of equal length as the original list. The proof produced by Vampire, expressed in Isabelle’s structured Isar format, looks as follows:

```

proof neg_clausify
  assume “ $\neg \text{length } (tl\ xs) \leq \text{length } xs$ ”
  hence “ $\text{drop } (\text{length } xs) (tl\ xs) \neq []$ ” by (metis drop_eq_Nil)
  hence “ $tl (\text{drop } (\text{length } xs) xs) \neq []$ ” by (metis drop_tl)
  hence “ $\forall u. xs @ u \neq xs \vee tl\ u \neq []$ ” by (metis append_eq_conv_conj)
  hence “ $tl [] \neq []$ ” by (metis append_Nil2)
  thus “False” by (metis tl.simps(1))
qed

```

The *neg\_clausify* method transforms the Isabelle conjecture into negated clause form, ensuring that it has the same shape as the corresponding ATP conjecture. The negation of the clause is introduced by the **assume** keyword, and a series of intermediate facts introduced by **hence** lead to a contradiction.

This approach was inspired by the Otterfier proof transformation service [40]. Resolution proofs should ideally be translated to natural, intuitive Isabelle proofs. The best-known prior work on translating resolution proofs is TRAMP [16]; its applicability to Sledgehammer is unexplored.

Preliminary work has commenced at Munich to see to what extent resolution proofs can be transformed into intelligible proofs. The first step is to transform the proof into a direct proof by applying contraposition repeatedly and introducing case splits where appropriate. For example, the proof above is transformed into

```

proof –
  have “ $tl [] = []$ ” by (metis tl.simps(1))

```

```

hence “ $\exists u. xs @ u = xs \wedge tl\ u = []$ ” by (metis append_Nil2)
hence “ $tl\ (drop\ (length\ xs)\ xs) = []$ ” by (metis append_eq_conv_conj)
hence “ $drop\ (length\ xs)\ (tl\ xs) = []$ ” by (metis drop_tl)
thus “ $length\ (tl\ xs) \leq length\ xs$ ” by (metis drop_eq_Nil)
qed

```

For most Isabelle users, the direct proof is much easier to understand and maintain.

### 2.3.2 One-Line Reconstruction, or ATPs as Relevance Filters

Having Metis available made possible an entirely different approach to proof reconstruction: to throw the proof away and allow Metis to find its own proof, using the lemmas that took part in the original resolution proof [28]. For example, the Isar proofs from §2.3.1 reduce to simply

```

by (metis append_Nil2 append_eq_conv_conj drop_eq_Nil drop_tl tl.simps(1))

```

With this approach, Sledgehammer became merely a lemma finder, one that used automatic theorem provers merely as relevance filters. But this approach was generally effective and had the great advantage that each Sledgehammer call now delivered a one-line result, rather than a lengthy and often incomprehensible proof script in which all formulas were in clause form. At least one ATP implementer expressed disbelief that his system could be used merely as a relevance filter, but this approach allows any ATP to be used provided it returns the list of axioms used in its proof.

Metis sometimes fails to reconstruct the result of a Sledgehammer call within a reasonable time. Reconstruction necessarily fails if the resolution proof does not correspond to a well-typed Isabelle proof (recall that, normally, Sledgehammer omits most type information when translating Isabelle formulas into first-order logic). This type of failure could be eliminated by using sound translations, but the overall success rate would actually decrease considerably. The seasoned user eventually learns to recognise unsound proofs—certain lemmas always seem to be mentioned. The number of such proofs is reduced by removing certain lemmas from the scope of Sledgehammer, both manually and automatically. Machine learning could perhaps be used here to determine which lemmas are harmful.

### 2.3.3 Proof Minimisation

Sometimes Metis is simply not powerful enough to prove a theorem that has already been proved by a more powerful system, despite being given a small list of axioms. ATPs frequently use many more axioms than are necessary. Sledgehammer’s minimisation tool takes a set of axioms returned by a given ATP and repeatedly calls the same ATP with subsets of those axioms in order to find a minimal set. Reducing the number of axioms improves Metis’s speed and success rate, while also removing superfluous clutter from the proof scripts.

ATPs themselves could return proofs using a minimum of axioms or, alternatively, proofs of a minimum length. Vampire’s well-known limited resource strategy [30], although designed to cope with limited processor time, could probably be modified to minimise proofs efficiently.

## 2.4 Parallelism

Parallelism was another design objective, both to exploit the abundance of cheap processing power and so that users would not have to wait. Sledgehammer was intended to run in the background; Isabelle would continue to respond to commands, and users could keep working. This idea has turned out to be somewhat misconceived: thinking is difficult, and users typically wait for Sledgehammer to return,

hoping to get a proof for free. We hope that eventually Sledgehammer will be configured to run spontaneously, without even the need for a mouse click. Then users will simply work and occasionally be delighted to have solutions displayed for them. Such a configuration would require a machine with enough processing power to support several ATP executions without becoming sluggish. An agent-based implementation of similar ideas, using a blackboard architecture, has for some time been part of the  $\Omega$ MEGA system [6, 33].

The parallel invocation of different theorem provers is invaluable. Böhme and Nipkow [8] have demonstrated that running three different theorem provers (E, SPASS, and Vampire) for five seconds solves as many problems as running the best theorem prover (Vampire) for a full two minutes. It would be better to employ even more theorem provers. We have undertaken informal, unpublished experiments involving many other systems.

- Gandalf [35] shows great potential, but unfortunately it does not output useful proofs; one cannot easily identify which axioms have taken part in the proof. A simple source code modification to improve the legibility of proofs would allow Gandalf to make useful contributions. Unfortunately, we were unable to identify the necessary changes. Gandalf has been found to be unsound,<sup>2</sup> but a small percentage of incorrect (and hence unreconstructable) proofs would be tolerable.
- SInE, the Sumo Inference Engine [11], is a wrapper around E that is designed to cope with large axiom bases. We pass it more facts than can be handled by the other ATPs, and it sometimes surprises us with original proofs. In the current experimental setup, it is invoked remotely via SystemOnTPTP [34] in parallel with Vampire.
- People sometimes suggest that we include Prover9 [15]. In our experiments, Prover9 performed poorly on the large problems generated by Sledgehammer. It could be effective in conjunction with an advanced and selective relevance filter.
- We could also run multiple instances of a theorem prover with different heuristics. This is not necessary with Vampire, which attempts a variety of heuristics in separate time slices. It could be particularly effective with E, but designing suitable heuristics requires highly specialised skills.

### 3 Evaluation

In their “Judgement Day” study, Böhme and Nipkow [8] evaluated Sledgehammer with E, SPASS, and Vampire on 1240 provable proof goals arising in seven representative Isabelle theories:

<i>Arrow</i>	Arrow’s impossibility theorem
<i>NS</i>	Needham–Schroeder shared-key protocol
<i>Hoare</i>	Completeness of Hoare logic with procedures
<i>Jinja</i>	Type soundness of a subset of Java
<i>SN</i>	Strong normalisation of the typed $\lambda$ -calculus with de Bruijn indices
<i>FTA</i>	Fundamental theorem of algebra
<i>FFT</i>	Fast Fourier transform

Sledgehammer has been developed further since they ran their experiments. In particular, it now communicates with ATPs using full first-order logic instead of clause form, adds SInE to the collection of ATPs, and employs the latest versions of SPASS, Vampire, and Metis. To account for these changes, we ran the Judgement Day benchmark suite on the same hardware as Böhme and Nipkow but with the latest version of Sledgehammer and of the Isabelle theories.

<sup>2</sup>See <http://www.cs.miami.edu/~tptp/TPTP/BustedAsUnsound.html>.

When running the four ATPs in parallel for 120 seconds, followed by Metis with a 30-second time limit, Sledgehammer now solves 52% of the goals (compared with 48% in Böhme and Nipkow). The table below gives the success rates for each ATP and theory.

	<i>Arrow</i>	<i>NS</i>	<i>Hoare</i>	<i>Jinja</i>	<i>SN</i>	<i>FTA</i>	<i>FFT</i>	Avg.
SInE 0.4	18%	22%	43%	31%	61%	53%	17%	40%
E 1.0	19%	39%	45%	33%	66%	57%	17%	44%
SPASS 3.7	30%	35%	43%	32%	59%	58%	17%	44%
Vampire 1.0	36%	40%	50%	35%	63%	60%	17%	47%
Together	43%	45%	54%	41%	68%	65%	26%	52%

About one third of the goals are “trivial” in the sense that they can be solved by standard Isabelle tactics invoked with no arguments. For these, the overall success rate is 84%. Regrettably, for the nontrivial goals, which users are especially keen on seeing solved by Sledgehammer, the overall success rate is much lower at 38%.<sup>3</sup>

Proof reconstruction using Metis loses about 10% of ATP proofs, partly because some of the ATP proofs are unsound in a typed setting, but mostly because Metis times out. Proof minimisation reduces the required number of facts by about one third, thereby helping 20% of the failed Metis proofs to succeed. It would be interesting to measure the effect of reconstructing Isar proof texts when the one-line Metis calls time out; regrettably, we currently lack the infrastructure to do so.

## 4 Sledgehammer and Teaching

Sledgehammer was not designed specifically as an aid to novices. Experienced users have come to rely on it. But Sledgehammer seems to offer the greatest benefits to the least experienced users. It has certainly transformed the way Isabelle is taught. There are two reasons for this:

- Because it identifies relevant facts, users no longer need to memorise lemma libraries.
- Because it works in harmony with Isar structured proofs, users no longer need to learn many low-level tactics.

Demonstrations of interactive theorem provers necessarily involve deception. The implementers naturally want to show off their system in the best possible light, so they present examples that look more difficult than they really are. Typically they define some recursive functions and prove properties using obvious inductions followed by some sort of automatic tactic. The audience will be duly impressed, and some among them will decide to adopt that tool as the basis for their Ph.D. research. Too late, they encounter the crucial issue:

*What do I do when the automatic tactics fail?*

Typically, the answer is that one must write incomprehensible scripts that invoke a plethora of obscure commands. These generally include tactics to manipulate the set of assumptions in natural deduction or a sequent calculus. There may be tactics to transform an assumption by applying a rewrite rule or theorem, to create a case split from an assumption, or to substitute user-supplied terms into theorems.

In Isabelle, the simple combination of structured proofs and Sledgehammer takes the user surprisingly far. This is not the place to give a detailed tutorial on Isar structured proofs [20, 38]. In brief,

<sup>3</sup>Böhme and Nipkow reported a success rate of 34% for nontrivial goals. Their result cannot be directly compared with ours because of methodological differences: to establish triviality, they relied on a syntactic criterion and ignored two of the theories, as opposed to actually running the Isabelle tactics on all proof goals.

they support natural deduction through local scopes that can introduce assumptions (using the keyword **assume**) as well as local variables and definitions. Moreover, while traditional tactic scripts contain only commands, a structured proof explicitly states the assumptions and goals. When proving a proposition, users can state intermediate properties that they believe to be helpful. If they understand the problem well enough to propose some intermediate properties, then all they have to do is state a progression of properties in small enough steps for Sledgehammer to be able to prove each one.

In the example below, taken a measure theory development, two intermediate facts (introduced by **hence**) assist in the proof of the conclusion (introduced by **thus**):

```

proof –
  assume  $x$ : “ $x \in \text{lambda\_system } M \ f$ ”
  hence “ $x \subseteq \text{space } M$ ”
    by (metis sets_into_space lambda_system_sets)
  hence “ $\text{space } M - (\text{space } M - x) = x$ ”
    by (metis double_diff equalityE)
  thus “ $\text{space } M - x \in \text{lambda\_system } M \ f$ ” using  $x$ 
    by (force simp add: lambda_system_def)
qed

```

Each of the intermediate facts is proved by a call to Metis that was generated using Sledgehammer. While the example features a linear progression of facts, Isar proofs can also be nested to any depth.

Isar also supports calculational reasoning [4]. A chain of reasoning steps, connected by familiar relations such as  $=$ ,  $\leq$ , and  $<$ , can be written with separate proofs for each step of the calculation. Once again, if the user can see the intermediate stages of the transformation, then the proof of each step can easily be found. The example below illustrates this type of reasoning:

```

proof –
   $\vdots$ 
  have “ $f(u \cap (x \cap y)) + f(u - x \cap y) = (f(u \cap (x \cap y)) + f(u \cap y - x)) + f(u - y)$ ”
    by (metis class_semiring.add_a ey)
  also have “ $\dots = (f((u \cap y) \cap x) + f(u \cap y - x)) + f(u - y)$ ”
    by (metis Int_commute Int_left_commute)
  also have “ $\dots = f(u \cap y) + f(u - y)$ ” using  $fx \text{ Int } y \ u$ 
    by auto
  also have “ $\dots = fu$ ”
    by (metis fy u)
  finally show “ $f(u \cap (x \cap y)) + f(u - x \cap y) = f \ u$ ” .
qed

```

Top down proof development is greatly assisted by a trivial Isar feature: the ability to omit proofs. Where a proof is required, the user may simply insert the word **sorry**. Isabelle then regards the theorem as proved.<sup>4</sup> The user can then check that the newly introduced proposition indeed suffices to prove the next proposition in the development. A difficult proof can develop as a series of propositions, each initially “proved” using **sorry** but eventually using either Sledgehammer, an automatic tactic, or a nested proof development of the same form. Progress in such a proof can be measured in terms of the difficulty of the propositions that lack real proofs. Although we can never be certain that a proof development can be completed until the very end, the ability to write **sorry** in place of a proof reduces the risk of discovering that a lemma is useless only after spending weeks proving it.

<sup>4</sup>The existence of **sorry** does not compromise Isabelle’s soundness, because it is only permitted during interactive sessions. A theory file containing an occurrence of **sorry** may not be imported by another theory.



In January 2010, as part of its new M.Phil. programme, the University of Cambridge offered a lecture course on Isabelle [22]. The course materials included almost no information about the low-level tactics that had been the mainstay of Isabelle proofs for nearly 20 years. Only two of the 12 lectures were devoted to Isar structured proofs, and they took a novel approach: rather than proceeding methodically through the Isar fundamentals, the lectures presented the outer skeleton of a proof, with crucial sections replaced by **sorry**. They described the idea of trying to eliminate each **sorry** using either Sledgehammer or some automatic tactic. Practical work submitted by the students later demonstrated that several of them had learnt how to write complex, well-structured proofs. We were happy to reassure them that submitting work generated largely by Sledgehammer was by no means cheating!

The first author has taught Isabelle on a number of occasions, starting in the mid-1990s. Sledgehammer is obviously not the only thing to have changed in that period. The introduction of Isar, continual improvements to the automated reasoners and counterexample generators, and 15 years of Moore’s Law have transformed the user experience. Interactive theorem proving has never been practical because it required far too much effort, even from highly specialised experts. For the first time, we can envisage the day when interactive theorem proving becomes straightforward enough to be adopted on a large scale.

## 5 Conclusions

Sledgehammer has been available for over three years, and in that time it has become an essential part of the Isabelle user’s workflow. It is the only interface between an interactive theorem prover and automatic ones to achieve such popularity with users. It has transformed the way beginners perceive Isabelle.

Sledgehammer has a number of limitations which we hope to address. The relevance filter is primitive, but an improved one will have to be part of Sledgehammer itself as long as unsound translations are used. Unsound translations can be used safely because Sledgehammer does not trust the proofs that it receives from ATPs but merely uses them as hints to generate Isabelle proof scripts; proofs that violate Isabelle’s typing rules are eliminated at this stage.

Sledgehammer’s performance on higher-order problems is unimpressive, and given the inherent difficulty of performing higher-order reasoning using first-order theorem provers, the way forward is to integrate Sledgehammer with higher-order automatic theorem provers, such as LEO-II [5] and Satallax [3]. Proof reconstruction would benefit from new ideas, especially so that it can deliver natural, intuitive proofs. Finally, we hope to generalise Sledgehammer so that it can also cope with Isabelle/ZF, the Zermelo–Fraenkel set-theory instantiation of Isabelle [24, 25].

## Acknowledgements

Christoph Benzmlüller commented on the version of this paper presented at PAAR-2010. Mark Summerfield suggested several textual improvements to a later draft.

The Cambridge team that developed Sledgehammer included Jia Meng, Claire Quigley, and Kong Woei Susanto. Sledgehammer has since been refined, improved, and tested by the Isabelle team in Munich—notably Sascha Böhme, Fabian Immler, Philipp Meyer, Tobias Nipkow, and Markus Wenzel. Throughout the development, we benefited from Joe Hurd’s expert help with Metis.

The research was supported by the Engineering and Physical Sciences Research Council (*Automation for Interactive Proof*, grant number GR/S57198/01) and also partly by the Deutsche Forschungsgemeinschaft (*Quis Custodiet*, grant number Ni 491/11-1).

## References

- [1] Wolfgang Ahrendt, Bernhard Beckert, Reiner Hähnle, Wolfram Menzel, Wolfgang Reif, Gerhard Schellhorn, and Peter H. Schmitt. Integrating automated and interactive theorem proving. In Wolfgang Bibel and Peter H. Schmitt, editors, *Automated Deduction—A Basis for Applications*, volume II. Systems and Implementation Techniques, pages 97–116. Kluwer Academic Publishers, 1998.
- [2] Peter B. Andrews. *An Introduction to Mathematical Logic and Type Theory: To Truth Through Proof (2nd Ed.)*. Springer, 2002. Applied Logic 27.
- [3] Julian Backes and Chad E. Brown. Analytic tableaux for higher-order logic with choice. In Jürgen Giesl and Reiner Hähnle, editors, *Automated Reasoning—5th International Joint Conference, IJCAR 2010*, LNAI 6173, pages 76–90. Springer, 2010.
- [4] Gertrud Bauer and Markus Wenzel. Calculational reasoning revisited (an Isabelle/Isar experience). In Richard J. Boulton and Paul B. Jackson, editors, *Theorem Proving in Higher Order Logics: TPHOLS 2001*, LNCS 2152, pages 75–90. Springer, 2001. Online at <http://link.springer.de/link/service/series/0558/tocs/t2152.htm>.
- [5] Christoph Benz Müller, Lawrence C. Paulson, Frank Theiss, and Arnaud Fietzke. LEO-II—A cooperative automatic theorem prover for higher-order logic. In Alessandro Armando, Peter Baumgartner, and Gilles Dowek, editors, *Automated Reasoning—4th International Joint Conference, IJCAR 2008*, LNAI 5195, pages 162–170. Springer, 2008.
- [6] Christoph Benz Müller and Volker Sorge. OANTS—An open approach at combining interactive and automated theorem proving. In Manfred Kerber and Michael Kohlhase, editors, *Symbolic Computation and Automated Reasoning*, pages 81–97. A. K. Peters, 2000.
- [7] Marc Bezem, Dimitri Hendriks, and Hans de Nivelle. Automatic proof construction in type theory using resolution. *Journal of Automated Reasoning*, 29(3–4):253–275, 2002.
- [8] Sascha Böhme and Tobias Nipkow. Sledgehammer: Judgement day. In Jürgen Giesl and Reiner Hähnle, editors, *Automated Reasoning (IJCAR 2010)*, LNCS 6173, pages 107–121. Springer, 2010.
- [9] Alonzo Church. A formulation of the simple theory of types. *Journal of Symbolic Logic*, 5:56–58, 1940.
- [10] M. J. C. Gordon and T. F. Melham. *Introduction to HOL: A Theorem Proving Environment for Higher Order Logic*. Cambridge University Press, 1993.
- [11] Kryštof Hoder. SInE (Sumo Inference Engine). <http://www.cs.man.ac.uk/~hoderk/sine/>.
- [12] Joe Hurd. Integrating Gandalf and HOL. In Yves Bertot, Gilles Dowek, André Hirschowitz, Christine Paulin, and Laurent Théry, editors, *Theorem Proving in Higher Order Logics: TPHOLS '99*, LNCS 1690, pages 311–321. Springer, 1999.
- [13] Joe Hurd. First-order proof tactics in higher-order logic theorem provers. In Myla Archer, Ben Di Vito, and César Muñoz, editors, *Design and Application of Strategies/Tactics in Higher Order Logics*, number NASA/CP-2003-212448 in NASA Technical Reports, pages 56–68, September 2003.
- [14] David McAllester. Ontic: A knowledge representation system for mathematics. In Ewing Lusk and Ross Overbeek, editors, *9th International Conference on Automated Deduction*, LNCS 310, pages 742–743. Springer, 1988.
- [15] William McCune. Prover9 and Mace4. <http://www.cs.unm.edu/~mccune/prover9/>.
- [16] Andreas Meier. TRAMP: Transformation of machine-found proofs into natural deduction proofs at the assertion level (system description). In David McAllester, editor, *Automated Deduction—CADE-17 International Conference*, LNAI 1831, pages 460–464. Springer, 2000.
- [17] Jia Meng and Lawrence C. Paulson. Translating higher-order clauses to first-order clauses. *Journal of Automated Reasoning*, 40(1):35–60, 2008.
- [18] Jia Meng and Lawrence C. Paulson. Lightweight relevance filtering for machine-generated resolution problems. *Journal of Applied Logic*, 7(1):41–57, 2009.
- [19] Jia Meng, Claire Quigley, and Lawrence C. Paulson. Automation for interactive proof: First prototype. *Information and Computation*, 204(10):1575–1596, 2006.

- [20] Tobias Nipkow. A tutorial introduction to structured Isar proofs. <http://isabelle.in.tum.de/dist/Isabelle/doc/isar-overview.pdf>.
- [21] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Springer, 2002. LNCS 2283.
- [22] Lawrence C. Paulson. Interactive formal verification. <http://www.cl.cam.ac.uk/teaching/0910/L21/>. Lecture course materials.
- [23] Lawrence C. Paulson. The foundation of a generic theorem prover. *Journal of Automated Reasoning*, 5(3):363–397, 1989.
- [24] Lawrence C. Paulson. Set theory for verification: I. From foundations to functions. *Journal of Automated Reasoning*, 11(3):353–389, 1993.
- [25] Lawrence C. Paulson. Set theory for verification: II. Induction and recursion. *Journal of Automated Reasoning*, 15(2):167–215, 1995.
- [26] Lawrence C. Paulson. Tool support for logics of programs. In Manfred Broy, editor, *Mathematical Methods in Program Development: Summer School Marktoberdorf 1996*, NATO ASI Series F, pages 461–498. Springer, 1997.
- [27] Lawrence C. Paulson. Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In Boris Konev, Renate Schmidt, and Stephan Schulz, editors, *PAAR-2010: Workshop on Practical Aspects of Automated Reasoning*, 2010.
- [28] Lawrence C. Paulson and Kong Woei Susanto. Source-level proof reconstruction for interactive theorem proving. In Klaus Schneider and Jens Brandt, editors, *Theorem Proving in Higher Order Logics: TPHOLS 2007*, LNCS 4732, pages 232–245. Springer, 2007.
- [29] Alexander Riazanov and Andrei Voronkov. Vampire 1.1 (system description). In Rajeev Goré, Alexander Leitsch and Tobias Nipkow, editors, *Automated Reasoning—First International Joint Conference, IJCAR 2001*, LNAI 2083, pages 376–380. Springer, 2001.
- [30] Alexandre Riazanov and Andrei Voronkov. Limited resource strategy in resolution theorem proving. *Journal of Symbolic Computation*, 36(1–2):101–115, 2003.
- [31] Alan Robinson and Andrei Voronkov, editors. *Handbook of Automated Reasoning*. Elsevier Science, 2001.
- [32] Stephan Schulz. System description: E 0.81. In David Basin and Michaël Rusinowitch, editors, *Automated Reasoning—Second International Joint Conference, IJCAR 2004*, LNAI 3097, pages 223–228. Springer, 2004.
- [33] Jörg Siekmann, Christoph Benzmüller, Armin Fiedler, Andreas Meier, Immanuel Normann and Martin Pollet. Proof development with  $\Omega$ MEGA: The irrationality of  $\sqrt{2}$ . In Fairouz Kamareddine, editor, *Thirty Five Years of Automating Mathematics*, pages 271–314. Kluwer Academic Publishers, 2003.
- [34] Geoff Sutcliffe. System description: SystemOnTPTP. In David McAllester, editor, *Automated Deduction—CADE-17 International Conference*, LNAI 1831, pages 406–410. Springer, 2000.
- [35] Tanel Tammet. Gandalf. *Journal of Automated Reasoning*, 18(2):199–204, 1997.
- [36] Josef Urban. MaLAREa: A metasystem for automated reasoning in large theories. In Geoff Sutcliffe, Josef Urban, and Stephan Schulz, editors, *ESARLT 2007: Empirically Successful Automated Reasoning in Large Theories*, volume 257 of *CEUR Workshop Proceedings*, 2007.
- [37] Christoph Weidenbach. Combining superposition, sorts and splitting. In Robinson and Voronkov [31], chapter 27, pages 1965–2013.
- [38] Makarius Wenzel. Isabelle/Isar—A generic framework for human-readable proof documents. In Roman Matuszewski and Anna Zalewska, editors, *From Insight to Proof—Festschrift in Honour of Andrzej Trybulec*. University of Białystok, 2007. Studies in Logic, Grammar, and Rhetoric 10(23).
- [39] Markus Wenzel. Type classes and overloading in higher-order logic. In Elsa L. Gunter and Amy Felty, editors, *Theorem Proving in Higher Order Logics: TPHOLS '97*, LNCS 1275, pages 307–322. Springer, 1997.
- [40] Jürgen Zimmer, Andreas Meier, Geoff Sutcliffe, and Yuan Zhan. Integrated proof transformation services. In Christoph Benzmüller and Wolfgang Windsteiger, editors, *Workshop on Computer-Supported Mathematical Theory Development, IJCAR 2004*, ENTCS, 2004.