

Automatic Proof and Disproof in Isabelle/HOL

Jasmin Blanchette, Lukas Bulwahn,
Tobias Nipkow

Fakultät für Informatik
TU München

- ① Introduction
- ② Isabelle's Standard Proof Methods
- ③ Sledgehammer
- ④ Quickcheck: Counterexamples by Testing
- ⑤ Nitpick: Counterexamples by SAT Solving

- 1 Introduction
- 2 Isabelle's Standard Proof Methods
- 3 Sledgehammer
- 4 Quickcheck: Counterexamples by Testing
- 5 Nitpick: Counterexamples by SAT Solving

A tale of two worlds

<i>FOL</i>		<i>HOL</i>
$f(s, t)$		$f\ s\ t, f\ s, \lambda x.t$
<i>Otter</i> (1987)		<i>Isabelle</i> (1986)

They did not talk to each other
because they spoke different languages.

This is the tale of how these two worlds
began to understand and boost each other.

Isabelle

- is an interactive theorem prover
- that has always embraced automation
- but without sacrificing soundness:

All proofs
must ultimately go through the Isabelle kernel

This is the *LCF principle* (Robin Milner).

Two decades of Isabelle development

1990s Basic proof automation

Our own proof search in ML:

simplifier, automatic provers, arithmetic

2000s Embrace external tools

Let them do the proof search,

but don't trust them:

ATPs (FOL provers)

SMT solvers

SAT solvers

Programming languages

- 1 Introduction
- 2 Isabelle's Standard Proof Methods**
- 3 Sledgehammer
- 4 Quickcheck: Counterexamples by Testing
- 5 Nitpick: Counterexamples by SAT Solving

Simplifier

N.

- First and higher-order equations (λ)
- Conditional equations
- Contextual simplification
- Special solvers (eg orderings)
- Arithmetic
- Case splitting (triggered by `if` and `case`)
- Large library of default equations

Isabelle's workhorse

The power of Isabelle's internal automated proof methods

- relies on large sets of **default rules**
- that are **user-extensible** (`[simp]`)
- and **tuned over time**.

Tableaux prover

Paulson

- Based on lean TAP (Beckert & Posegga)
- Generic
- User-extensible by `intro` and `elim` rules
- Proof search in ML,
proof checking via Isabelle kernel
- Works well for pure logic and set theory
- Does not know anything about equality

Isabelle Demo

- ① Introduction
- ② Isabelle's Standard Proof Methods
- ③ Sledgehammer**
- ④ Quickcheck: Counterexamples by Testing
- ⑤ Nitpick: Counterexamples by SAT Solving

Sledgehammer

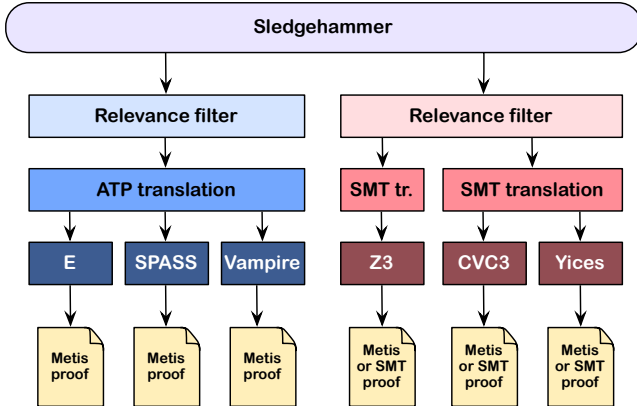
Paulson et al.

- Connects Isabelle with ATPs and SMT solvers
E, SPASS, Vampire, CVC3, Yices, Z3, ...
- One-click invocation:
 - Users don't need to select facts
 - ... or ensure the problem is first-order
- Exploits local parallelism, remote servers

Sledgehammer: Demo



Sledgehammer: Architecture



Sledgehammer: Fact selection

Meng & Paulson

Provers perform poorly given 1000s of facts

A lightweight, symbol-based filter greatly improves the success rate

Number of facts is optimized for each prover

Sledgehammer: Translation

Meng & Paulson Bl., Böhme & Smallbone

Source: higher-order, polymorphism + type classes

Target: first-order, untyped/simply-typed

① Firstorderize

- SK combinators, λ -lifting
- Explicit application operator

② Encode types

- Monomorphize
- ... or encode polymorphism

Sledgehammer: Reconstruction

Paulson & Susanto Böhme & Weber

Four approaches (the 4 Rs):

- A. Re-find using Metis
- B. Rerun external prover
- C. Recheck stored proof
- D. Recast into Isar proof

A. Re-find using Metis

lemma `length (tl xs) ≤ length xs`

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

Usually fast and reliable

Metis sometimes too slow (5% loss on avg)

B. Rerun external prover

```
lemma length (tl xs) ≤ length xs  
by (smt append_Nil2 append_eq_conv_conj  
      drop_eq_Nil drop_tl tl.simps(1))
```

Reinvokes the SMT solver each time!

C. Recheck stored proof

lemma `length (tl xs) ≤ length xs`

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

Fast

No need for SMT solver for replay

Fragile

D. Recast into Isar proof

lemma $\text{length } (\text{tl } xs) \leq \text{length } xs$

proof –

have $\text{tl } [] = []$ **by** (metis tl.simps(1))

hence $\exists u. xs @ u = xs \wedge \text{tl } u = []$ **by** (metis append_Nil2)

hence $\text{tl } (\text{drop } (\text{length } xs) xs) = []$ **by** (metis append_eq_conv_conj)

hence $\text{drop } (\text{length } xs) (\text{tl } xs) = []$ **by** (metis drop_tl)

thus $\text{length } (\text{tl } xs) \leq \text{length } xs$ **by** (metis drop_eq_Nil)

qed

Fast, self-explanatory

Experimental, bulky

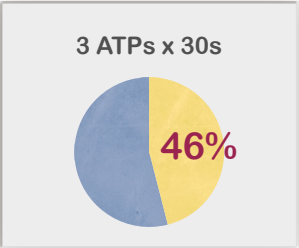
Sledgehammer: Judgment Day

Böhme & N. Bl., Böhme & Paulson

- 1240 goals arising in 7 older theories
Arrow, FFT, FTA, Hoare, Jinja, NS, SN
- In 2010: E, SPASS, Vampire (5 to 120 s)
 $ESV \times 5s \approx V \times 120s$
- In 2011: Also E-SInE, CVC3, Yices, Z3 (30 s)
 $Z3 > V$
- In 2012: Tighter integration with SPASS
SPASS most successful backend (by a small margin)

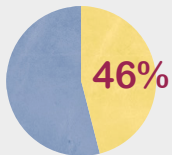
2010

2010

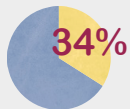


2010

3 ATPs x 30s

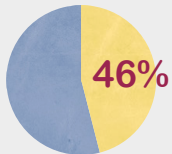


3 ATPs x 30 s
nontrivial goals

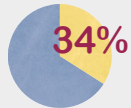


2010

3 ATPs x 30s

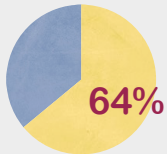


3 ATPs x 30 s
nontrivial goals

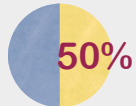


2012

(4 ATPs + 3 SMTs) x 30s



(4 ATPs + 3 SMTs) x 30s
nontrivial goals



Sledgehammer & Teaching

Paulson

Old way: Low-level tactics + lemma libraries

New way: Isar + Sledgehammer + simp etc.

lemma *blah*

sorry

proof –

have *blah*₀ sorry**by** (metis foo bar)

hence *blah*₁ sorry**by** metis

hence *blah*₂ sorry**by** auto

thus *blah* sorry**by** (metis baz)

qed

Sledgehammer: Success story

Guttman, Struth & Weber

Developed large Isabelle/HOL repository of algebras for modeling imperative programs

(Kleene Algebra, Hoare logic, . . . , \approx 1000 lemmas)

Intricate refinement and termination theorems

Surprise: Sledgehammer and Z3 automate algebraic proofs at textbook level!

“The integration of ATP, SMT, and Nitpick is for our purposes **very very helpful.**” — G. Struth

Theorem proving and testing

*Testing can show only the presence of errors,
but not their absence. (Dijkstra)*

*Testing cannot prove theorems,
but it can refute conjectures!*

Two facts of life:

- 95% of all conjectured theorems are wrong.
- Theorem proving is an expensive debugging technique.

Theorem provers need counterexample finders!

- ① Introduction
- ② Isabelle's Standard Proof Methods
- ③ Sledgehammer
- ④ Quickcheck: Counterexamples by Testing**
- ⑤ Nitpick: Counterexamples by SAT Solving

Quickcheck

Berghofer & N. Bul.

- Adds lightweight validation by testing
- Motivated by Haskell's QuickCheck
- Employs Isabelle's code generator
- Quick response time
- No-click invocation:
 automatic after parsing a proposition
 (well, at least in ProofGeneral)

Quickcheck: Demo

Quickcheck

Berghofer & N. Bul.

- Covers different testing approaches
 - Random and exhaustive testing
 - Smart test data generators
 - Narrowing-based testing
- Creates test data generators automatically

Test generators for datatypes

Fast iteration over the large number of tests using continuation-passing-style programming:

For datatype α $list = Nil \mid Cons \alpha (\alpha list)$

we create a test function for property P :

$test_{\alpha list} P =$

$P Nil$ and also

$test_{\alpha} (\lambda x. test_{\alpha list} (\lambda xs. P (Cons x xs)))$

Test generators for predicates

Testing propositions with preconditions

$\text{distinct } xs \implies \text{distinct } (\text{remove1 } x \text{ } xs)$

Problem:

Exhaustive testing creates useless test data

Solution:

Use precondition's definition for smarter generator

Test generators for predicates

From the definition:

distinct Nil = True

distinct (Cons x xs) = $(x \notin xs \wedge \text{distinct } xs)$

we create a test function for property P :

test-distinct _{α list} $P =$

P Nil andalso

test _{α} ($\lambda x.$ test-distinct _{α list} ($\lambda xs.$

if $x \notin xs$ then P (Cons x xs) else True))

Non-distinct lists are never generated

Test generators for predicates

Construct generators using data flow analysis:

- 1 Transform predicates to system of horn clauses
 $x \notin xs \implies \text{distinct } xs \implies \textit{distinct} (\text{Cons } x \ xs)$
- 2 Perform data flow analysis:
which variables can be computed,
which variables must be generated?
- 3 Synthesize test data generator

Narrowing-based testing

- Symbolic execution with demand-driven refinement:
 - Test cases can contain variables
 - If execution cannot proceed, variables are instantiated, again by symbolic terms
- Pays off if large search spaces can be discarded
distinct (`Cons 1 (Cons 1 x)`) is false for every x
No further instantiations for x

Implementations of narrowing

- Programming language with native narrowing currently still too slow
- Lazy execution with outer refinement loop results in many recomputations, but fast

Limitations

Quickcheck only checks *executable* specifications:

- No equality on functions with infinite domain
- No axiomatic specifications

- ① Introduction
- ② Isabelle's Standard Proof Methods
- ③ Sledgehammer
- ④ Quickcheck: Counterexamples by Testing
- ⑤ Nitpick: Counterexamples by SAT Solving

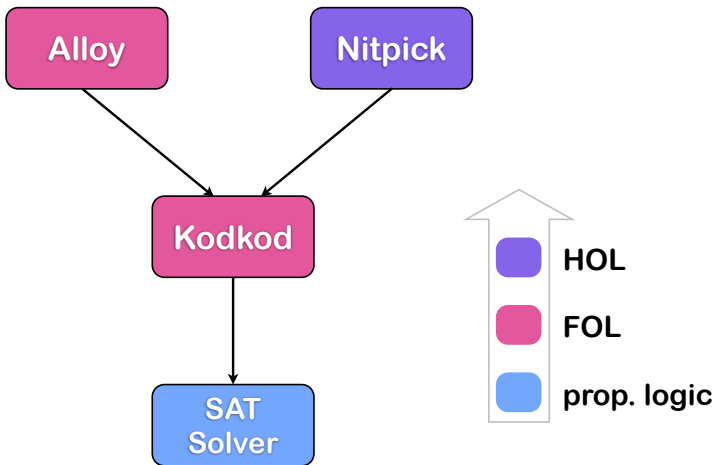
Finite model finder

Based on SAT via Kodkod (Alloy's backend)

Soundly approximates infinite types

Nitpick: Demo

Nitpick: Architecture



Nitpick: Basic translation

For fixed finite cardinalities (1, 2, 3, ..., 10)

First-order:

$$\tau_1 \rightarrow \cdots \rightarrow \tau_n \rightarrow \mathit{bool} \quad \mapsto \quad A_1 \times \cdots \times A_n$$

$$\tau_1 \rightarrow \cdots \rightarrow \tau_n \rightarrow \tau \quad \mapsto \quad A_1 \times \cdots \times A_n \times A \\ + \text{constraint}$$

Higher-order args of type $\sigma \rightarrow \tau \quad \mapsto$

$$\underbrace{A \times \cdots \times A}_{|\sigma| \text{ times}}$$

Nitpick: Datatypes

Soundly approximated by finite sets (3-valued logic)

Efficient axiomatization:

Subterm-closed substructures (Kuncak & Jackson)

Examples

nat: $\{0, \text{Suc } 0, \text{Suc } (\text{Suc } 0)\}$

α list: $\{[], [a_1], [a_2], [a_2, a_1]\}$

Motto: Let the SAT solver spin!

(and trust Kodkod's symmetry breaking)

Nitpick: Inductive predicates

p is the least solution to $p = F(p)$ for some F

Naive idea: Take $p = F(p)$ as p 's specification!

Unsound in general, but:

- Sound if recursion $p = F(p)$ is well-founded
- Sound for negative occurrences of p

Otherwise: Unroll! (cf. Biere, Cimatti, Clarke & Zhu)

$$p_0 = (\lambda x. \text{False}) \quad p_{i+1} = F(p_i)$$

Nitpick: Success stories

Algebraic methods (Guttman, Struth & Weber)

C++ memory model (Bl., Weber, Batty, Owens & Sarkar)

Soundness bugs in TPS and LEO-II

Typical fan mail:

“Last night I got stuck on a goal I was sure was a theorem. After 5–10 minutes I gave Nitpick a try, and within a few secs it had found a splendid counterexample—despite the mess of locales and type classes in the context!”