

Efficiently Checking Propositional Resolution Proofs in Isabelle/HOL

Tjark Weber¹

¹Institut für Informatik, Technische Universität München
Boltzmannstr. 3, D-85748 Garching b. München, Germany
webertj@in.tum.de

Abstract

This paper describes the integration of zChaff and MiniSat, currently two leading SAT solvers, with Isabelle/HOL. Both SAT solvers generate resolution-style proofs for (instances of) propositional tautologies. These proofs are verified by the theorem prover. The presented approach significantly improves Isabelle's performance on propositional problems, and exhibits counterexamples for unprovable conjectures. It is shown that an LCF-style theorem prover can serve as a viable proof checker even for large SAT problems. An efficient representation of the propositional problem in the theorem prover turns out to be crucial; several possible solutions are discussed.

1 Introduction

Interactive theorem provers like PVS [ORS92], HOL [GM93] or Isabelle [Pau94] traditionally support rich specification logics. Proof search and automation for these logics however is difficult, and proving a non-trivial theorem usually requires manual guidance by an expert user. Automated theorem provers on the other hand, while often designed for simpler logics, have become increasingly powerful over the past few years. New algorithms, improved heuristics and faster hardware allow interesting theorems to be proved with little or no human interaction, sometimes within seconds.

By integrating automated provers with interactive systems, we can preserve the richness of our specification logic and at the same time increase the degree of automation [Sha01]. This is an idea that goes back at least to the early nineties [KKS91]. However, to ensure that a potential bug in the automated prover does not render the whole system unsound, theorems in Isabelle, like in other LCF-style [Gor00] provers, can be derived only through a fixed set of core inference rules. Therefore it is not sufficient for the automated prover to return whether a formula is provable, but it must also generate the actual proof, expressed (or expressible) in terms of the interactive system's inference rules.

Formal verification is an important application area of interactive theorem proving. Problems in verification can often be reduced to Boolean satisfiability (SAT), and recent SAT solver advances have made this approach feasible in practice. Hence the performance of an interactive prover on propositional problems may be of significant

practical importance. In this paper we describe the integration of zChaff [MMZ⁺01] and MiniSat [ES04], two leading SAT solvers, with the Isabelle/HOL [NPW02] prover.

We have shown earlier [Web05a, Web05b] that using a SAT solver to prove theorems of propositional logic dramatically improves Isabelle’s performance on this class of formulae, even when a rather naive (and unfortunately, as we will see in Section 3, inefficient) representation of propositional problems is used. Furthermore, while Isabelle’s previous decision procedures simply fail on unprovable conjectures, SAT solvers are able to produce concrete counterexamples. In this paper we focus on recent improvements of the proof reconstruction algorithm in Isabelle/HOL, which cause a speedup by several orders of magnitude. In particular the representation of the propositional problem turns out to be crucial for performance. While the implementation in [Web05a] was still limited to relatively small SAT problems, the recent improvements now allow to check proofs with millions of resolution steps in reasonable time. This shows that, somewhat contrary to common belief, efficient proof checking in an LCF-style system is feasible.

The next section describes the integration of zChaff and MiniSat with Isabelle/HOL in more detail. In Section 3 we evaluate the performance of our approach, and report on experimental results. Related work is discussed in Section 4. Section 5 concludes this paper with some final remarks and points out directions for future research.

2 System Description

To prove a propositional tautology ϕ in the Isabelle/HOL system with the help of zChaff or MiniSat, we proceed in several steps. First ϕ is negated, and the negation is converted into an equivalent formula ϕ^* in conjunctive normal form. ϕ^* is then written to a file in DIMACS CNF format [DIM93], the standard input format supported by most SAT solvers. zChaff and MiniSat, when run on this file, return either “unsatisfiable”, or a satisfying assignment for ϕ^* .

In the latter case, the satisfying assignment is displayed to the user. The assignment constitutes a counterexample to the original (unnegated) conjecture. When the solver returns “unsatisfiable” however, things are more complicated. If we have confidence in the SAT solver, we can simply trust its result and accept ϕ as a theorem in Isabelle. The theorem is tagged with an “oracle” flag to indicate that it was proved not through Isabelle’s own inference rules, but by an external tool. In this scenario, a bug in the SAT solver could potentially allow us to derive inconsistent theorems in Isabelle/HOL.

The LCF-approach instead demands that we verify the solver’s claim of unsatisfiability within Isabelle/HOL. While this is not as simple as the validation of a satisfying assignment, the increasing complexity of SAT solvers has before raised the question of support for independent verification of their results, and in 2003 zChaff has been extended by L. Zhang and S. Malik [ZM03] to generate resolution-style proofs that can be verified by an independent checker. (This issue has also been acknowledged by the annual SAT Competition, which has introduced a special track on certified unsat answers in 2005.) More recently, a proof-logging version of MiniSat has been released [ES06], and John Matthews has extended this version to produce human-readable proofs that are easy to parse [Mat06], similar to those produced by zChaff. Hence our main task boils down to using Isabelle/HOL as an independent checker for the resolution proofs

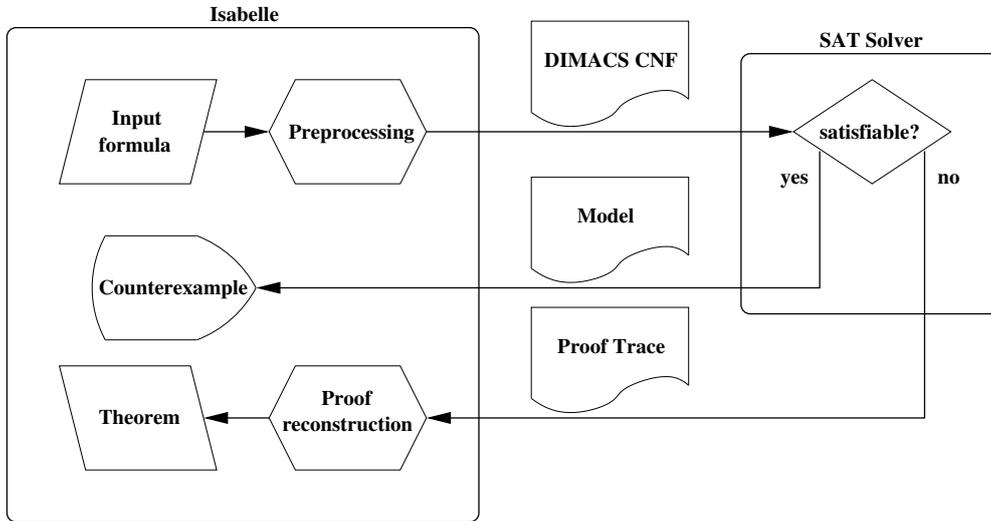


Figure 1: Isabelle – SAT System Architecture

found by zChaff and MiniSat.

Both solvers store their proof in a text file that is read in by Isabelle, and the individual resolution steps are replayed in Isabelle/HOL. Section 2.1 briefly describes the necessary preprocessing of the input formula, and details of the proof reconstruction are explained in Section 2.2. The overall system architecture is shown in Figure 1.

2.1 Preprocessing

Isabelle/HOL offers higher-order logic (on top of Isabelle’s meta logic), whereas most SAT solvers only support formulae of propositional logic in conjunctive normal form (CNF). Therefore the (negated) input formula ϕ must be preprocessed before it can be passed to the solver.

Two different CNF conversions are currently implemented: a naive encoding that may cause an exponential blowup of the formula, and a Tseitin-style encoding [Tse83] that may introduce (existentially quantified) auxiliary Boolean variables, cf. [Gor01]. The technical details can be found in [Web05a]. More sophisticated CNF conversions, e.g. from [NRW98], could be implemented as well. The main focus of our work however is on efficient proof reconstruction, less on transformations of the input formula: the benchmark problems used for evaluation in Section 3 are already given in CNF anyway.

Note that it is not sufficient to convert ϕ into an equivalent formula ϕ^* in CNF. Rather, we have to *prove* this equivalence inside Isabelle/HOL. The result is not a single formula, but a theorem of the form $\vdash \phi = \phi^*$. The fact that our CNF transformation must be proof-producing leaves some potential for optimization. One could implement a non proof-producing (and therefore much faster) version of the same CNF transformation, and use it for preprocessing instead. Application of the proof-producing version would then be necessary only if the SAT solver has in fact shown a formula to be unsatisfiable. The total runtime on provable formulae would increase slightly, as the CNF

transformation needed to be done twice – first without, later with proofs. Preprocessing times for unprovable formulae however should improve.

ϕ^* is written to a file in DIMACS CNF format, and the SAT solver is invoked on this input file.

2.2 Proof Reconstruction

When zChaff and MiniSat return “unsatisfiable”, they also generate a resolution-style proof of unsatisfiability and store the proof in a text file [ZM03, Mat06]. While the precise format of this file differs between the solvers, the essential proof technique is the same. Both SAT solvers use propositional resolution to derive new clauses from existing ones:

$$\frac{P \vee x \quad Q \vee \neg x}{P \vee Q}$$

It is well-known that this single inference rule is sound and complete for propositional logic. A set of clauses is unsatisfiable iff the empty clause is derivable via resolution. (For the purpose of proof reconstruction, we are only interested in the proof returned by the SAT solver, not in the techniques and heuristics that the solver uses internally to find this proof. Therefore the integration of zChaff and MiniSat is quite similar – minor differences in their proof trace format aside –, and further SAT solvers capable of generating resolution-style proofs could be integrated with Isabelle in the exact same manner.)

We assign a unique identifier – simply a non-negative integer, starting with 0 – to each clause of the original CNF formula. Further clauses derived by resolution are assigned identifiers by the solver. Often we are not interested in the clause obtained by resolving just two existing clauses, but only in the result of a whole resolution *chain*, where two clauses are resolved, the result is resolved with yet another clause, and so on. Consequently, we define an ML [MTHM97] type of propositional resolution proofs as a pair whose first component is a table mapping integers (to be interpreted as the identifiers of clauses derived by resolution) to lists of integers (to be interpreted as the identifiers of previously derived clauses that are part of the defining resolution chain). The second component of the proof is just the identifier of the empty clause.

```
type proof = int list Inttab.table * int
```

This type is merely intended as an internal format to store the information contained in a resolution proof. There are many restrictions on valid proofs that are not enforced by this type. For example, it does not ensure that its second component indeed denotes the empty clause, that every resolution step is legal, or that there are no circular dependencies between derived clauses. It is only important that every resolution proof can be represented as a value of type `proof`, not conversely. The proof returned by zChaff or MiniSat is translated into this internal format, and passed to the actual proof reconstruction algorithm. This algorithm will either generate an Isabelle/HOL theorem, or fail in case the proof is invalid (which should not happen unless the SAT solver contains a bug).

2.2.1 zChaff's Proof Trace Format

The format of the proof trace generated by zChaff has not been documented before. Therefore a detailed description of the format and its interpretation, although not the main focus of this paper, seems in order.

The proof file generated by zChaff consists of three sections, the first two of which are optional (but present in any non-trivial proof). The first section defines clauses derived from the original problem by resolution. A typical line would be “CL: 7 <= 2 3 0”, meaning that a new clause was derived by resolving clauses 2 and 3, and resolving the result with clause 0. In this example, the new clause is assigned the identifier 7, which may then be used in further lines of the proof file. Clauses of the original CNF formula are implicitly assigned identifiers starting from 0, in the order they are given in the DIMACS file. When converting zChaff's proof into our internal format, the clause identifiers in a CL line can immediately be added to the table which constitutes the proof's first component, with the new identifier as the key, and the list of resolvents as the associated value.

The second section of the proof file records variable assignments that are implied by the first section, and by other variable assignments. As an example, consider “VAR: 3 L: 2 V: 0 A: 1 Lits: 4 7”. This line states that variable 3 must be false (i.e. its value must be 0; zChaff uses “V: 1” for variables that must be true) at decision level 2, the *antecedent* being clause 1. The antecedent is a clause in which every literal except for the one containing the assigned variable must evaluate to false because of variable assignments at lower decision levels (or because the antecedent is already a unit clause). The antecedent's literals are given explicitly by zChaff, using an encoding that multiplies each variable by 2 and adds 1 for negative literals. Hence “Lits: 4 7” corresponds to the clause $x_2 \vee \neg x_3$. Our internal proof format does not allow to record variable assignments directly, but we can translate them by observing that they correspond to unit clauses. For each variable assignment in zChaff's trace, a new clause identifier is generated (using the number of clauses derived in the first section as a basis, and the variable itself as offset) and added as a key to the proof's table. The associated list of resolvents contains the antecedent, and is otherwise obtained from the explicitly given literals: for each literal's variable (except for the one that is being assigned), a similar unit clause must have been added to the table before; its identifier computed according to the same formula. We ignore both the value and the level information in zChaff's trace. The former is implicit in the derived unit clause (which contains the variable either positively or negatively), and the latter is implicit in the overall proof structure.

The last section of the proof file consists of just one line which specifies the *conflict clause*, a clause which has only false literals: e.g. “CONF: 3 == 4 6”. Literals are encoded in the same way as in the second section, so clause 3 would be $x_2 \vee x_3$ in this case. We translate this line into our internal proof format by generating a new clause identifier i which is added to the proof's table, with the conflict clause itself and the unit clause for each of the conflict clause's variables as associated resolvents. Finally, we set the proof's second component to i .

2.2.2 MiniSat’s Proof Trace Format

The proof-logging version of MiniSat originally generated proof traces in a rather compact (and again undocumented) binary format, for which we have not implemented a parser. John Matthews [Mat06] however has extended this version with the ability to produce readable proof traces in ASCII format, similar to those produced by zChaff. We describe the precise proof trace format, and its translation into our internal proof format.

MiniSat’s proof traces, unlike zChaff’s, are not divided into sections. They contain four different types of statements: “R” to reference original clauses, “C” for clauses derived via resolution, “D” to delete clauses that are not needed anymore, and “X” to indicate the end of proof. Aside from “X”, which must appear exactly once and at the end of the proof trace, the other statements may appear in any number and (almost) any order.

MiniSat does not implicitly assign identifiers to clauses in the original CNF formula. Instead, “R” statements, e.g. “R 0 <= -1 3 4”, are used to establish clause identifiers. This particular line introduces a clause identifier 0 for the clause $\neg x_1 \vee x_3 \vee x_4$, which must have been one of the original clauses in this example. (Note that MiniSat, unlike zChaff, uses the DIMACS encoding of literals in its proof trace.) Since our internal proof format uses different identifiers for the original clauses, the translation of MiniSat’s proof trace into the internal format becomes parameterized by a renaming \mathcal{R} of clause identifiers. An “R” statement does not affect the proof itself, but it extends the renaming. The given literals are used to look up the identifier of the corresponding original clause, and the clause identifier introduced by the “R” statement is mapped to the clause’s original (internal) identifier.

New clauses are derived from existing clauses via resolution chains. A typical line would be “C 7 <= 2 5 3 4 0”, meaning that a new clause with identifier 7 was derived by resolving clauses 2 and 3 (with x_5 as the pivot variable), and resolving the result with clause 0 (with x_4 as the pivot variable). In zChaff’s notation, this would correspond to “CL: 7 <= 2 3 0”. We add this line to the proof’s table just like for zChaff, but with one difference: MiniSat’s clause identifiers cannot be used directly. Instead, we generate a new internal clause identifier for this line, extend the renaming \mathcal{R} by mapping MiniSat’s clause identifier (7 in this example) to the newly generated identifier, and apply \mathcal{R} to the identifiers of resolvents as well.

Clauses that are not needed anymore can be indicated by a “D” statement, followed by a clause identifier. Currently such statements are ignored. Making beneficial use of them would require not only a modified proof format, but also a different algorithm for proof reconstruction.

Finally a line like “X 0 42” indicates the end of proof. The numbers are the minimum and maximum, respectively, identifiers of clauses used in the proof. We ignore the first identifier (which is usually 0 anyway), and use the second identifier, mapped from MiniSat’s identifier scheme to our internal one by applying \mathcal{R} , as the identifier of the empty clause, i.e. as the proof’s second component.

There is one significant difference between MiniSat’s and zChaff’s proof traces that should have become apparent from the foregoing description. MiniSat, unlike zChaff, records the pivot variable for each resolution step in its trace, i.e. the variable that occurs

positively in one clause partaking in the resolution, and negatively in the other. This information is redundant, as the pivot variable can always be determined from those two clauses: If two clauses containing more than one variable both positively and negatively were to be resolved, the resulting clause would be tautological, i.e. contain a variable and its negation. Both zChaff and MiniSat are smart enough not to derive such tautological clauses in the first place. We have decided to ignore the pivot information in MiniSat’s traces, since proof reconstruction for zChaff requires the pivot variable to be determined anyway, and using MiniSat’s pivot data would need a modified internal proof format. This however leaves some potential for optimization wrt. replaying MiniSat proofs.

2.2.3 Proof Reconstruction

We now come to the core of this paper. The task of proof reconstruction is to derive False from the original clauses, using information from a value of type `proof` (which represents a resolution proof found by a SAT solver). This can be done in various ways. In particular the precise representation of the problem as an Isabelle/HOL theorem (or a collection of Isabelle/HOL theorems) turns out to be crucial for performance.

Naive HOL Representation In an early implementation [Web05a], the whole problem was represented as a single theorem $\vdash (\phi^* \implies \text{False}) \implies (\phi^* \implies \text{False})$, where ϕ^* was completely encoded in HOL as a conjunction of disjunctions. Step by step, this theorem was then modified to reduce the antecedent $\phi^* \implies \text{False}$ to True, which would eventually prove $\vdash \phi^* \implies \text{False}$.

This was extremely inefficient for two reasons. First, every resolution step required manipulation of the whole (possibly huge) problem at once. Second, and just as important, SAT solvers treat clauses as *sets* of literals, making implicit use of associativity, commutativity and idempotence of disjunction. Likewise, CNF formulae are treated as sets of clauses, making implicit use of the same properties for conjunction. The encoding in HOL however required numerous explicit rewrites (with theorems like $\vdash (P \vee Q) = (Q \vee P)$) to reorder clauses and literals before each resolution step.

Separate Clauses Representation A better representation of the CNF formula was discussed in [FMM⁺06]. In order to understand it, we need to look at the ML datatype of theorems that Isabelle uses internally. Every theorem encodes a *sequent* $\Gamma \vdash \phi$, where ϕ is a single formula, and Γ is a finite set of formulae (implemented as an ordered list of terms, although this detail doesn’t matter to us). The intended interpretation is that ϕ holds when every formula in Γ is assumed as a *hypothesis*. So far we have only considered theorems where $\Gamma = \emptyset$, written $\vdash \phi$ for short. This was motivated by the normal user-level view on theorems in Isabelle, where assumptions are encoded using implications \implies , rather than hypotheses. Isabelle’s inference kernel however provides rules that let us convert between hypotheses and implications as we like:

$$\frac{}{\{\phi\} \vdash \phi} \text{Assume} \quad \frac{\Gamma \vdash \psi}{\Gamma \setminus \phi \vdash \phi \implies \psi} \text{impI} \quad \frac{\Gamma \vdash \phi \implies \psi \quad \Gamma' \vdash \phi}{\Gamma \cup \Gamma' \vdash \psi} \text{impE}$$

Let us use $\llbracket A_1; \dots; A_n \rrbracket \implies B$ as a short hand for $A_1 \implies \dots \implies A_n \implies B$ (with implication associating to the right). In [FMM⁺06], each clause $p_1 \vee \dots \vee p_n$ is encoded

as an implication $\overline{p_1} \Longrightarrow \dots \Longrightarrow \overline{p_n} \Longrightarrow \text{False}$ (where $\overline{p_i}$ denotes the negation normal form of $\neg p_i$, for $1 \leq i \leq n$), and turned into a separate theorem

$$\{p_1 \vee \dots \vee p_n\} \vdash \llbracket \overline{p_1}; \dots; \overline{p_n} \rrbracket \Longrightarrow \text{False}.$$

This allows resolution to operate on comparatively small objects, and resolving two clauses $\Gamma \vdash \llbracket p_1; \dots; p_n \rrbracket \Longrightarrow \text{False}$ and $\Gamma' \vdash \llbracket q_1; \dots; q_m \rrbracket \Longrightarrow \text{False}$, where $\neg p_i = q_j$ for some i and j , essentially becomes an application of the cut rule. The first clause is rewritten to $\Gamma \vdash \llbracket p_1; \dots; p_{i-1}; p_{i+1}; \dots; p_n \rrbracket \Longrightarrow \neg p_i$. A derived Isabelle tactic then performs the cut to obtain

$$\Gamma \cup \Gamma' \vdash \llbracket q_1; \dots; q_{j-1}; p_1; \dots; p_{i-1}; p_{i+1}; \dots; p_n; q_{j+1}; \dots; q_m \rrbracket \Longrightarrow \text{False}$$

from the two clauses. Note that this representation, while breaking apart the given clauses into separate theorems allows us to view the CNF formula as a set of clauses, still does not allow us to view each individual clause as a set of literals. Some reordering of literals is necessary before cuts can be performed, and after each cut, duplicate literals have to be removed from the result.

Sequent Representation We can further exploit the fact that Isabelle's inference kernel treats a theorem's hypotheses as a set of formulae, by encoding each clause using hypotheses *only*. Consider the following representation of a clause $p_1 \vee \dots \vee p_n$ as an Isabelle/HOL theorem:

$$\{p_1 \vee \dots \vee p_n, \overline{p_1}, \dots, \overline{p_n}\} \vdash \text{False}.$$

Resolving two clauses $p_1 \vee \dots \vee p_n$ and $q_1 \vee \dots \vee q_m$, where $\neg p_i = q_j$, now starts with two applications of the impI rule to obtain theorems

$$\{p_1 \vee \dots \vee p_n, \overline{p_1}, \dots, \overline{p_{i-1}}, \overline{p_{i+1}}, \dots, \overline{p_n}\} \vdash \neg p_i \Longrightarrow \text{False}$$

and

$$\{q_1 \vee \dots \vee q_m, \overline{q_1}, \dots, \overline{q_{j-1}}, \overline{q_{j+1}}, \dots, \overline{q_m}\} \vdash p_i \Longrightarrow \text{False}.$$

We then instantiate a previously proven theorem

$$\vdash (P \Longrightarrow \text{False}) \Longrightarrow (\neg P \Longrightarrow \text{False}) \Longrightarrow \text{False}$$

(where P is an arbitrary proposition) with p_i for P . Instantiation is another basic operation provided by Isabelle's inference kernel. Finally two applications of impE yield

$$\{p_1 \vee \dots \vee p_n, \overline{p_1}, \dots, \overline{p_{i-1}}, \overline{p_{i+1}}, \dots, \overline{p_n}\} \cup \{q_1 \vee \dots \vee q_m, \overline{q_1}, \dots, \overline{q_{j-1}}, \overline{q_{j+1}}, \dots, \overline{q_m}\} \vdash \text{False}.$$

This approach requires no explicit reordering of literals anymore. Furthermore, duplicate literals do not need to be eliminated after resolution. This is all handled by the inference kernel now; the sequent representation is as close to a SAT solver's view of clauses as sets of literals as possible in Isabelle. With this representation, we do not rely on derived tactics anymore to perform resolution, but we can give a precise description of the implementation in terms of (five, as we see above) applications of core inference rules.

CNF Sequent Representation The sequent representation has the disadvantage that each clause contains itself as a hypothesis. Since hypotheses are accumulated during resolution, this leads to larger and larger sets of hypotheses, which will eventually contain every clause used in the resolution proof. Forming the union of these sets takes the kernel a significant amount of time. It is therefore faster to use a slightly different clause representation, where each clause contains the whole CNF formula ϕ^* as a hypothesis. Let $\phi^* \equiv \bigwedge_{i=1}^k C_i$, where k is the number of clauses. Using the Assume rule, we obtain a theorem $\{\bigwedge_{i=1}^k C_i\} \vdash \bigwedge_{i=1}^k C_i$. Repeated elimination of conjunction (with the help of two theorems, namely $\vdash P \wedge Q \implies P$ and $\vdash P \wedge Q \implies Q$) yields a list of theorems $\{\bigwedge_{i=1}^k C_i\} \vdash C_1, \dots, \{\bigwedge_{i=1}^k C_i\} \vdash C_k$. Each of these theorems is then converted into the sequent form described above, with literals as hypotheses and False as the theorem’s conclusion. This representation increases preprocessing times slightly, but throughout the entire proof, the set of hypotheses for each clause now consists of $\bigwedge_{i=1}^k C_i$ and the clause’s literals only. It is therefore much smaller than before, which speeds up resolution. Furthermore, memory requirements do *not* increase: the term $\bigwedge_{i=1}^k C_i$ needs to be kept in memory only once, and can be shared between different clauses. This can also be exploited when the union of hypotheses is formed (assuming that the inference kernel and the underlying ML system support it): a simple pointer comparison is sufficient to determine that both theorems contain $\bigwedge_{i=1}^k C_i$ as a hypothesis (and hence that the resulting theorem needs to contain it only once); no lengthy term traversal is required.

We should mention that this representation of clauses, despite its superior practical performance, has a small downside. The resulting theorem always has *every* given clause as a premise, while the theorem produced by the sequent representation only has those clauses as premises that were actually used in the proof. If the logically stronger theorem is needed, it can be obtained by analyzing the resolution proof to identify the used clauses beforehand, and filtering out the unused ones *before* proof reconstruction.

We still need to determine the pivot literal (i.e. p_i and $\neg p_i$ in the above example) before resolving two clauses. This could be done by directly comparing the hypotheses of the two clauses, and searching for a term that occurs both positively and negatively. It turns out to be slightly faster however (and also more robust, since we make fewer assumptions about the actual implementation of hypotheses in Isabelle) to use our own data structure. With each clause, we associate a table that maps integers – one for each literal in the clause – to the Isabelle term representation of a literal. The table is an inverse of the mapping from literals to integers that was constructed for translation into DIMACS format, but restricted to the literals that actually occur in a clause. Positive integers are mapped to positive literals (atoms), and negative integers are mapped to negative literals (negated atoms). This way term negation simply corresponds to integer negation. The table associated with the result of a resolution step is the union of the two tables that were associated with the resolvents, but with the entry for p_i ($\neg p_i$, respectively) removed from the table associated with the first (second, respectively) clause.

Another optimization, related not to the representation of individual clauses, but to the overall proof structure, is perhaps more obvious and has been present in our implementations since the beginning. zChaff and MiniSat, during proof search, may generate many clauses that are ultimately not needed to derive the empty clause. Instead

of replaying the whole proof trace in chronological order, we perform “backwards” proof reconstruction, starting with the identifier of the empty clause, and recursively proving the required resolvents using depth-first search.

While some clauses may not be needed at all, others may be used multiple times in the resolution proof. It would be highly inefficient to prove these clauses over and over again. Therefore all clauses proved are stored in an array, which is allocated at the beginning of proof reconstruction (with a size big enough to possibly hold all clauses derived during the proof). Initially, this array only contains clauses present in the original CNF formula, still in their original format as a disjunction of literals. Whenever an original clause is used as a resolvent, it is converted into the sequent format described above. (Note that this avoids converting original clauses that are not used in the proof at all.) The converted clause, along with its literal table, is stored in the array instead of the original (unconverted) clause. Each clause obtained as the result of a resolution chain is stored as well. Reusing a previously proved clause merely causes an array lookup.

For this reason, it could be beneficial to analyze the resolution chains in more detail: sometimes very similar chains occur in a proof, differing only in a clause or two. Common parts of resolution chains could be stored as additional lemmas (which only need to be derived once), thereby reducing the total number of resolution steps. A detailed evaluation of this idea is beyond the scope of this paper.

2.3 A Simple Example

In this section we illustrate the proof reconstruction using a small example. Consider the following input formula

$$\phi \equiv (\neg x_1 \vee x_2) \wedge (\neg x_2 \vee \neg x_3) \wedge (x_1 \vee x_2) \wedge (\neg x_2 \vee x_3).$$

Since ϕ is already in conjunctive normal form, preprocessing simply yields the theorem $\vdash \phi = \phi$. The corresponding DIMACS CNF file, aside from its header, contains one line for each clause in ϕ :

```
-1 2 0
-2 -3 0
1 2 0
-2 3 0
```

zChaff and MiniSat easily detect that this problem is unsatisfiable. zChaff creates a text file with the following data:

```
CL: 4 <= 2 0
VAR: 2 L: 0 V: 1 A: 4 Lits: 4
VAR: 3 L: 1 V: 0 A: 1 Lits: 5 7
CONF: 3 == 5 6
```

We see that first a new clause, with identifier 4, is derived by resolving clause 2, $x_1 \vee x_2$, with clause 0, $\neg x_1 \vee x_2$. The pivot variable which occurs both positively (in clause 2) and negatively (in clause 0) is x_1 ; this variable is eliminated by resolution.

$$\begin{array}{c}
\frac{\neg x_2 \vee x_3}{x_3} \quad \frac{\frac{x_1 \vee x_2}{x_2} \quad \frac{\neg x_1 \vee x_2}{x_2}}{\perp} \quad \frac{\neg x_2 \vee \neg x_3}{\neg x_3} \quad \frac{\frac{x_1 \vee x_2}{x_2} \quad \frac{\neg x_1 \vee x_2}{x_2}}{\perp} \\
\hline
\perp
\end{array}$$

Figure 2: Resolution Proof found by zChaff

Now the value of x_2 (VAR: 2) can be deduced from clause 4 (A: 4). x_2 must be true (V: 1). Clause 4 contains only one literal (Lits: 4), namely x_2 (since $4 \div 2 = 2$), occurring positively (since $4 \bmod 2 = 0$). This decision is made at level 0 (L: 0), before any decision at higher levels.

Likewise, the value of x_3 can then be deduced from clause 1, $\neg x_2 \vee \neg x_3$. x_3 must be false (V: 0).

Finally clause 3 is our conflict clause. It contains two literals, $\neg x_2$ (since $5 \div 2 = 2$, $5 \bmod 2 = 1$) and x_3 (since $6 \div 2 = 3$, $6 \bmod 2 = 0$). But we already know that both literals must be false, so this clause is not satisfiable.

In Isabelle, the resolution proof corresponding to zChaff's proof trace is constructed backwards from the conflict clause. A tree-like representation of the proof is shown in Figure 2. Note that information concerning the level of decisions, the actual value of variables, or the literals that occur in a clause is redundant in the sense that it is not needed by Isabelle to validate zChaff's proof. The clause x_2 , although used twice in the proof, is derived only once during resolution (and reused the second time), saving one resolution step in this little example.

The proof trace produced by MiniSat for the same problem happens to encode a different resolution proof:

```

R 0 <= -1 2
R 1 <= -2 -3
R 2 <= 1 2
R 3 <= -2 3
C 4 <= 3 3 1
C 5 <= 0 2 4
C 6 <= 2 2 4
C 7 <= 5 1 6
X 0 7

```

The first four lines introduce clause identifiers for all four clauses in the original problem, in their original order as well (effectively making the renaming \mathcal{R} from MiniSat's clause identifiers to internal clause identifiers the identity in this case). The next four lines define four new clauses (one clause per line), derived by resolution. Clause 4 is the result of resolving clause 3 ($\neg x_2 \vee x_3$) with clause 1 ($\neg x_2 \vee \neg x_3$), where x_3 is used as pivot literal. Hence clause 4 is equal to $\neg x_2$. Likewise, clause 5 is the result of resolving clauses 0 and 4, and clause 6 is obtained by resolving clauses 2 and 4. Finally resolving clauses 5 and 6 yields the empty clause, which is assigned clause identifier 7. The proof is shown in Figure 3. Again one resolution step is saved in the Isabelle implementation because clause $\neg x_2$ is proved only once.

$$\begin{array}{c}
\frac{\neg x_1 \vee x_2}{\neg x_1} \quad \frac{\frac{\neg x_2 \vee x_3 \quad \neg x_2 \vee \neg x_3}{\neg x_2}}{\neg x_1} \quad \frac{x_1 \vee x_2}{x_1} \quad \frac{\frac{\neg x_2 \vee x_3 \quad \neg x_2 \vee \neg x_3}{\neg x_2}}{x_1} \\
\hline
\perp
\end{array}$$

Figure 3: Resolution Proof found by MiniSat

Problem Representation	Proof Reconstruction (zChaff)
Naive HOL	726.5
Separate Clauses	7.8
Sequent	1.2
CNF Sequent	0.5

Table 1: Runtimes (in seconds) for MSC007-1.008

3 Evaluation

In [Web05a], we compared the performance of our approach, using the naive HOL problem representation, to that of Isabelle’s existing automatic proof procedures on all 42 problems contained in version 2.6.0 of the TPTP library [SS98] that have a representation in propositional logic. The problems were negated, so that unsatisfiable problems became provable. The benchmarks were run on a machine with a 3 GHz Intel Xeon CPU and 1 GB of main memory.

19 of these 42 problems are rather easy, and were solved in less than a second each by both the existing procedures and the SAT solver approach. On the remaining 23 problems, zChaff proved to be clearly superior to Isabelle’s built-in proof procedures. zChaff solved all problems in less than a second, and proof reconstruction in Isabelle/HOL took a few seconds at most for all but one problem: with the naive HOL representation, the proof for problem MSC007-1.008 was reconstructed in just over 12 minutes.

To give an impression of the effect that the different clause representations discussed in Section 2.2.3 have on performance, Table 1 shows the different times required to prove problem MSC007-1.008. The proof found by zChaff for this problem has 8,705 resolution steps. MiniSat finds a proof with 40,790 resolution steps for the same problem, which is reconstructed in about 3.8 seconds total with the sequent representation, and in 1.1 seconds total with the CNF sequent representation. The times to prove the other problems from the TPTP library have decreased in a similar fashion and are well below one second each now.

This enables us to evaluate the performance on some significantly larger problems, taken from the SATLIB library [HS00]. These problems do not only push Isabelle’s inference kernel to its limits, but also other parts of the prover. While the smaller TPTP problems were converted to Isabelle’s input syntax by a Perl script, this approach turns out to be infeasible for the larger SATLIB problems. The Perl script still works fine, but Isabelle’s parser (which was mainly intended for small, hand-crafted terms) is unable to parse the resulting theory files, which are several megabytes large, in reasonable time. Also, the prover’s user interface is unable to display the resulting formulae. We have therefore implemented our own little parser, which builds ML terms directly from

Problem	Variables	Clauses	zChaff (s)	Proof (s)	Resolutions	Total (s)
c7552mul.miter	11282	69529	73	70	252200	145
6pipe	15800	394739	167	321	268808	512
6pipe_6_ooo	17064	545612	308	2575	870345	3179
7pipe	23910	751118	495	1132	357136	1768

Table 2: Runtimes (in seconds) for SATLIB problems, zChaff

Problem	Variables	Clauses	MiniSat (s)	Proof (s)	Resolutions	Total (s)
c7552mul.miter	11282	69529	25	49	908231	106
6pipe	15800	394739	x	—	—	—
6pipe_6_ooo	17064	545612	x	—	—	—
7pipe	23910	751118	x	—	—	—

Table 3: Runtimes (in seconds) for SATLIB problems, MiniSat

DIMACS files, and we work entirely at the system’s ML level, avoiding the usual user interface, to prove unsatisfiability.

Statistics for four unsatisfiable SATLIB problems (chosen among those that were used to evaluate zChaff’s performance in [ZM03]) are shown in Tables 2 and 3, for zChaff and MiniSat respectively. The first column shows the time in seconds that it takes the SAT solver to find a proof of unsatisfiability. The second column, “Proof”, shows the time in seconds required to replay the proof’s resolutions steps in Isabelle/HOL, using the CNF sequent representation of clauses. The third column shows the number of resolution steps performed during proof replay. The last column, “Total”, finally shows the total time to prove the problem unsatisfiable in Isabelle, including SAT solving time, proof replay, parsing of input and output files, and any other intermediate pre- and postprocessing. These timings were obtained on an AMD Athlon 64 X2 Dual Core Processor 3800+ with 4 GB of main memory. An **x** indicates that the solver ran out of memory, or that the proof trace file exceeded a size of 2 GB. Needless to say that none of these problems can be solved automatically by Isabelle’s built-in proof procedures.

It seems that proof checking in Isabelle/HOL, despite all optimizations that we have implemented, is sometimes about an order of magnitude slower than proof verification with an external checker written in C++ [ZM03]. From Table 2 we conclude that proving unsatisfiability in Isabelle/HOL is by a factor of roughly 2 to 10 slower than using zChaff alone. This additional overhead was to be expected: it is the price that we have to pay for using Isabelle’s LCF-style kernel, which is not geared towards propositional logic. However, we also see that proof reconstruction in Isabelle scales quite well with our latest implementation, and that it remains feasible even for large SAT problems.

Comparing the runtimes for problem c7552mul.miter on the proofs found by zChaff and MiniSat, we see that the time taken to reconstruct a proof does not solely depend on the number of resolutions steps. In particular our algorithm for resolving two clauses, as described in Section 2.2.3, is linear in the length (i.e. number of literals) of those clauses. The average length of a clause is about 31.0 for the MiniSat proof, and about 98.6 for the proof found by zChaff. This explains why the zChaff proof, despite its smaller number of resolution steps, takes longer to reconstruct.

4 Related Work

Michael Gordon has implemented *HolSatLib* [Gor01], a library which is now part of the HOL 4 theorem prover. This library provides functions to convert HOL 4 terms into CNF, and to analyze them using a SAT solver. In the case of unsatisfiability however, the user only has the option to trust the external solver. No proof reconstruction takes place, “since there is no efficient way to check for unsatisfiability using pure Hol98 theorem proving” [Gor01]. A bug in the SAT solver could ultimately lead to an inconsistency in HOL 4.

Perhaps closer related to our work is the integration of automated first-order provers, in the context of Isabelle recently further explored by Joe Hurd [Hur99, Hur02], Jia Meng [Men03], and Lawrence Paulson [MP04, MP06]. Proofs found by the automated system are either verified by the interactive prover immediately [Hur99], or translated into a proof script that can be executed later [MP04]. Also Andreas Meier’s TRAMP system [Mei00] transforms the output of various automated first-order provers into natural deduction proofs. The main focus of their work however is on the necessary translation from the interactive prover’s specification language to first-order logic. In contrast our approach is so far restricted to instances of *propositional* tautologies, but we have focused on performance (rather than on difficult translation issues), and we use a SAT solver, rather than a first-order prover. Other work on combining proof and model search includes [dNM06].

A custom-built SAT solver has been integrated with the CVC Lite system [BB04] by Clark Barrett et al. [BBD03]. While this solver produces proofs that can be checked independently, our work shows that it is possible to integrate existing, highly efficient solvers with an LCF-style prover: the information provided by recent versions of zChaff and MiniSat is sufficient to produce a proof object in a theorem prover, no custom-built solver is necessary.

An earlier version of this work was presented in [Web05a], and improved by Alwen Tiu et al. [FMM⁺06]. Furthermore Hasan Amjad [Amj06b] has recently integrated a proof-generating version of the MiniSat solver with HOL 4 in a similar fashion. In this paper we have discussed our most recent implementation, which is based on a novel clause representation and constitutes a significant performance improvement when compared to earlier work.

5 Conclusions and Future Work

The SAT solver approach dramatically outperforms the automatic proof procedures that were previously available in Isabelle/HOL. With the help of zChaff or MiniSat, many formulae that were previously out of the scope of Isabelle’s built-in tactics can now be proved – or refuted – automatically, often within seconds. Isabelle’s applicability as a tool for formal verification, where large propositional problems occur in practice, has thereby improved considerably.

Furthermore, using the data structures and optimizations described in this paper, proof reconstruction for propositional logic scales quite well even to large SAT problems and proofs with several hundred thousand resolution steps. The additional confidence gained by using an LCF-style prover to check the proof obviously comes at a price (in

terms of runtime), but it's not nearly as expensive as one might have expected after earlier implementations.

While improving the performance of our implementation, we confirmed an almost self-evident truth: use profiling to see which functions take a lot of time, and focus on improving them – this is where the greatest benefits lie. This was an iterative process. A better implementation would allow us to tackle larger SAT problems, which in turn would uncover new performance bottlenecks. More importantly, we discovered some inefficiencies in the implementation of the Isabelle kernel. (Instantiating a theorem with a term, for example, was linear in the size of the term, rather than in constant time.) These inefficiencies played no important role as long as the kernel only had to deal with relatively small terms, but in our context, where formulae sometimes consist of millions of literals, they turned out to have a negative impact on performance. Subsequently the kernel implementation was modified, and these inefficiencies were removed.

Tuning an implementation to the extent presented here requires a great deal of familiarity with the underlying theorem prover. Nevertheless our results are applicable beyond Isabelle/HOL. Other interactive provers for higher-order logic, e.g. HOL 4 and HOL-Light, use very similar data structures to represent their theorems. Hasan Amjad has confirmed that the CNF sequent representation works equally well in these provers [Amj06b].

We have already mentioned some possible directions for future work. There is probably not very much potential left to optimize the implementation of resolution itself at this point. However, to further improve the performance of proof reconstruction, it could be beneficial to analyze the resolution proof found by the SAT solver in more detail. Merging similar resolution chains may reduce the overall number of resolutions required, and re-sorting resolutions may help to derive shorter clauses during the proof, which should improve the performance of individual resolution steps. Some preliminary results along these lines are reported in [Amj06a].

The approach presented in this paper has applications beyond propositional reasoning. The decision problem for richer logics (or fragments thereof) can be reduced to SAT [ABC⁺02, Str02, MS05, RH06]. Consequently, proof reconstruction for propositional logic can serve as a foundation for proof reconstruction for other logics. Based on our work, only a proof-generating implementation of the reduction is needed to integrate the more powerful, yet SAT-based decision procedure with an LCF-style theorem prover. This has already been used to integrate haRVey, a Satisfiability Modulo Theories (SMT) prover, with Isabelle [Hur06]. haRVey, like other SMT systems, uses various decision procedures (e.g. congruence closure for uninterpreted functions) on top of a SAT solver.

Acknowledgments. The author would like to thank Markus Wenzel for several good ideas and his extensive help with tuning the implementation, and Hasan Amjad for more good ideas and the challenge that he posed. The clause representations used in this paper were suggested to me by various people, including Hasan Amjad, John Harrison, John Matthews, and Markus Wenzel.

References

- [ABC⁺02] G. Audemard, P. Bertoli, A. Cimatti, A. Kornilowicz, and R. Sebastiani. A SAT based approach for solving formulas over Boolean and linear mathematical propositions. In Andrei Voronkov, editor, *Proceedings of the 18th International Conference on Automated Deduction (CADE-18)*, volume 2392 of *Lecture Notes in Artificial Intelligence*, pages 195–210, Copenhagen, Denmark, July 2002. Springer.
- [Amj06a] Hasan Amjad. Compressing propositional refutations. In Stephan Merz and Tobias Nipkow, editors, *Sixth International Workshop on Automated Verification of Critical Systems (AVOCS '06) – Preliminary Proceedings*, pages 7–18, 2006.
- [Amj06b] Hasan Amjad. A HOL/MiniSat interface. Personal communication, September 2006.
- [BB04] Clark Barrett and Sergey Berezin. CVC Lite: A new implementation of the cooperating validity checker. In *Proceedings of the 16th International Conference on Computer Aided Verification (CAV 2004)*, Boston, Massachusetts, USA, July 2004.
- [BBD03] Clark Barrett, Sergey Berezin, and David L. Dill. A proof-producing Boolean search engine. In *Proceedings of the Workshop on Pragmatics of Decision Procedures in Automated Reasoning (PDPAR 2003)*, Miami, Florida, USA, July 2003.
- [DIM93] DIMACS satisfiability suggested format, 1993. Available online at <ftp://dimacs.rutgers.edu/pub/challenge/satisfiability/doc>.
- [dNM06] Hans de Nivelle and Jia Meng. Geometric resolution: A proof procedure based on finite model search. In Ulrich Furbach and Natarajan Shankar, editors, *Automated Reasoning – Third International Joint Conference, IJ-CAR 2006, Seattle, WA, USA, August 2006, Proceedings*, volume 4130 of *Lecture Notes in Artificial Intelligence*, pages 303–317, 2006.
- [ES04] Niklas Eén and Niklas Sörensson. An extensible SAT-solver. In Enrico Giunchiglia and Armando Tacchella, editors, *Theory and Applications of Satisfiability Testing, 6th International Conference, SAT 2003. Santa Margherita Ligure, Italy, May 5-8, 2003 Selected Revised Papers*, volume 2919 of *Lecture Notes in Computer Science*, pages 502–518. Springer, 2004.
- [ES06] Niklas Eén and Niklas Sörensson. MiniSat-p-v1.14 – A proof-logging version of MiniSat, September 2006. Available online at <http://www.cs.chalmers.se/Cs/Research/FormalMethods/MiniSat/>.
- [FMM⁺06] Pascal Fontaine, Jean-Yves Marion, Stephan Merz, Leonor Prensa Nieto, and Alwen Tiu. Expressiveness + automation + soundness: Towards combining SMT solvers and interactive proof assistants. In Holger Hermanns

and Jens Palsberg, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 12th International Conference, TACAS 2006 Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2006, Vienna, Austria, March 25 - April 2, 2006, Proceedings*, volume 3920 of *Lecture Notes in Computer Science*, pages 167–181. Springer, 2006.

- [GM93] M. J. C. Gordon and T. F. Melham, editors. *Introduction to HOL: A theorem proving environment for higher order logic*. Cambridge University Press, 1993.
- [Gor00] M. J. C. Gordon. From LCF to HOL: A short history. In G. Plotkin, Colin P. Stirling, and Mads Tofte, editors, *Proof, Language, and Interaction*. MIT Press, 2000.
- [Gor01] M. J. C. Gordon. HolSatLib documentation, version 1.0b, June 2001. Available online at <http://www.cl.cam.ac.uk/~mjcg/HolSatLib/HolSatLib.html>.
- [HS00] Holger H. Hoos and Thomas Stütze. SATLIB: An online resource for research on SAT. In Ian Gent, Hans van Maaren, and Toby Walsh, editors, *SAT 2000*, pages 283–292. IOS Press, 2000. Available online at <http://www.satlib.org/>.
- [Hur99] 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, 12th International Conference, TPHOLS '99*, volume 1690 of *Lecture Notes in Computer Science*, pages 311–321, Nice, France, September 1999. Springer.
- [Hur02] Joe Hurd. An LCF-style interface between HOL and first-order logic. In Andrei Voronkov, editor, *Proceedings of the 18th International Conference on Automated Deduction (CADE-18)*, volume 2392 of *Lecture Notes in Artificial Intelligence*, pages 134–138, Copenhagen, Denmark, July 2002. Springer.
- [Hur06] Clément Hurlin. Proof reconstruction for first-order logic and set-theoretical constructions. In Stephan Merz and Tobias Nipkow, editors, *Sixth International Workshop on Automated Verification of Critical Systems (AVOCS '06) – Preliminary Proceedings*, pages 157–162, 2006.
- [KKS91] R. Kumar, T. Kropf, and K. Schneider. Integrating a first-order automatic prover in the HOL environment. In M. Archer, J. J. Joyce, K. N. Levitt, and P. J. Windley, editors, *Proceedings of the 1991 International Workshop on the HOL Theorem Proving System and its Applications*, pages 170–176, Davis, California, USA, August 1991. IEEE Computer Society Press, 1992.
- [Mat06] John Matthews. ASCII proof traces for MiniSat. Personal communication, August 2006.

- [Mei00] Andreas Meier. TRAMP: Transformation of machine-found proofs into natural deduction proofs at the assertion level. In David A. McAllester, editor, *Automated Deduction – CADE-17, 17th International Conference on Automated Deduction, Pittsburgh, PA, USA, June 17-20, 2000, Proceedings*, volume 1831 of *Lecture Notes in Artificial Intelligence*, pages 460–464. Springer, 2000.
- [Men03] Jia Meng. Integration of interactive and automatic provers. In Manuel Carro and Jesus Correias, editors, *Second CologNet Workshop on Implementation Technology for Computational Logic Systems, FME 2003*, September 2003.
- [MMZ⁺01] M. Moskewicz, C. Madigan, Y. Zhao, L. Zhang, and S. Malik. Chaff: Engineering an efficient SAT solver. In *Proceedings of the 38th Design Automation Conference*, Las Vegas, June 2001.
- [MP04] Jia Meng and Lawrence C. Paulson. Experiments on supporting interactive proof using resolution. In David Basin and Michaël Rusinowitch, editors, *Automated Reasoning: Second International Joint Conference, IJCAR 2004, Cork, Ireland, July 4-8, 2004, Proceedings*, volume 3097 of *Lecture Notes in Artificial Intelligence*, pages 372–384. Springer, 2004.
- [MP06] Jia Meng and Lawrence C. Paulson. Translating higher-order problems to first-order clauses. In Geoff Sutcliffe, Renate Schmidt, and Stephan Schulz, editors, *ESCoR: Empirically Successful Computerized Reasoning*, volume 192 of *CEUR Workshop Proceedings*, pages 70–80, 2006.
- [MS05] Andreas Meier and Volker Sorge. Applying SAT solving in classification of finite algebras. *Journal of Automated Reasoning*, 35(1–3):201–235, October 2005.
- [MTHM97] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML - Revised*. MIT Press, May 1997.
- [NPW02] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL – A Proof Assistant for Higher-Order Logic*, volume 2283 of *Lecture Notes in Computer Science*. Springer, 2002.
- [NRW98] Andreas Nonnengart, Georg Rock, and Christoph Weidenbach. On generating small clause normal forms. In Claude Kirchner and Hélène Kirchner, editors, *Automated Deduction – CADE-15, 15th International Conference on Automated Deduction, Lindau, Germany, July 5-10, 1998, Proceedings*, volume 1421 of *Lecture Notes in Computer Science*, pages 397–411. Springer, 1998.
- [ORS92] S. Owre, J. M. Rushby, and N. Shankar. PVS: A prototype verification system. In Deepak Kapur, editor, *11th International Conference on Automated Deduction (CADE)*, volume 607 of *Lecture Notes in Artificial Intelligence*, pages 748–752, Saratoga, NY, June 1992. Springer.

- [Pau94] Lawrence C. Paulson. *Isabelle: A Generic Theorem Prover*, volume 828 of *Lecture Notes in Computer Science*. Springer, 1994.
- [RH06] Erik Reeber and Warren A. Hunt, Jr. A SAT-based decision procedure for the subclass of unrollable list formulas in ACL2 (SULFA). In Ulrich Furbach and Natarajan Shankar, editors, *Automated Reasoning – Third International Joint Conference, IJCAR 2006, Seattle, WA, USA, August 2006, Proceedings*, volume 4130 of *Lecture Notes in Artificial Intelligence*, pages 453–467, 2006.
- [Sha01] Natarajan Shankar. Using decision procedures with a higher-order logic. In Richard J. Boulton and Paul B. Jackson, editors, *Theorem Proving in Higher Order Logics, 14th International Conference, TPHOLs 2001, Edinburgh, Scotland, UK, September 3-6, 2001, Proceedings*, volume 2152 of *Lecture Notes in Computer Science*, pages 5–26. Springer, 2001.
- [SS98] Geoff Sutcliffe and Christian Suttner. The TPTP problem library: CNF release v1.2.1. *Journal of Automated Reasoning*, 21(2):177–203, 1998. Available online at <http://www.cs.miami.edu/~tptp/>.
- [Str02] Ofer Strichman. On solving Presburger and linear arithmetic with SAT. In M. D. Aagaard and J. W. O’Leary, editors, *Formal Methods in Computer-Aided Design: 4th International Conference, FMCAD 2002, Portland, OR, USA, November 6-8, 2002, Proceedings*, volume 2517 of *Lecture Notes in Computer Science*, pages 160–169. Springer, 2002.
- [Tse83] G. S. Tseitin. On the complexity of derivation in propositional calculus. In J. Siekmann and G. Wrightson, editors, *Automation Of Reasoning: Classical Papers On Computational Logic, Vol. II, 1967-1970*, pages 466–483. Springer, 1983. Also in *Structures in Constructive Mathematics and Mathematical Logic Part II*, ed. A. O. Slisenko, 1968, pp. 115–125.
- [Web05a] Tjark Weber. Integrating a SAT solver with an LCF-style theorem prover. In A. Armando and A. Cimatti, editors, *Proceedings of PDPAR’05 – Third International Workshop on Pragmatical Aspects of Decision Procedures in Automated Reasoning*, Edinburgh, UK, July 2005.
- [Web05b] Tjark Weber. Using a SAT solver as a fast decision procedure for propositional logic in an LCF-style theorem prover. In Joe Hurd, Edward Smith, and Ashish Darbari, editors, *Theorem Proving in Higher Order Logics – 18th International Conference, TPHOLs 2005, Oxford, UK, August 2005, Emerging Trends Proceedings*, pages 180–189, Oxford, UK, August 2005. Oxford University Computing Laboratory, Programming Research Group. Research Report PRG-RR-05-02.
- [ZM03] Lintao Zhang and Sharad Malik. Validating SAT solvers using an independent resolution-based checker: Practical implementations and other applications. In *Design, Automation and Test in Europe (DATE 2003)*, pages 10880–10885. IEEE Computer Society, 2003.