

A Verified Prover Based on Ordered Resolution

Anders Schlichtkrull
DTU Compute
Technical University of Denmark
Kongens Lyngby, Denmark
andschl@dtu.dk

Jasmin Christian Blanchette
Department of Computer Science
Vrije Universiteit Amsterdam
Amsterdam, the Netherlands
j.c.blanchette@vu.nl

Dmitriy Traytel
Department of Computer Science
ETH Zürich
Zürich, Switzerland
traytel@inf.ethz.ch

Abstract

The superposition calculus, which underlies first-order theorem provers such as E, SPASS, and Vampire, combines ordered resolution and equality reasoning. As a step towards verifying modern provers, we specify, using Isabelle/HOL, a purely functional first-order ordered resolution prover and establish its soundness and refutational completeness. Methodologically, we apply stepwise refinement to obtain, from an abstract nondeterministic specification, a verified deterministic program, written in a subset of Isabelle/HOL from which we extract purely functional Standard ML code that constitutes a semidecision procedure for first-order logic.

CCS Concepts • **Theory of computation** → **Logic and verification**; *Automated reasoning*;

Keywords automatic theorem provers, proof assistants, first-order logic, stepwise refinement

ACM Reference Format:

Anders Schlichtkrull, Jasmin Christian Blanchette, and Dmitriy Traytel. 2019. A Verified Prover Based on Ordered Resolution. In *Proceedings of the 8th ACM SIGPLAN International Conference on Certified Programs and Proofs (CPP '19)*, January 14–15, 2019, Cascais, Portugal. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3293880.3294100>

1 Introduction

Automatic theorem provers based on superposition, such as E [42], SPASS [53], and Vampire [21], are often employed as backends in proof assistants and program verification tools [7, 20, 32]. Superposition is a highly successful calculus for first-order logic with equality, which generalizes both ordered resolution [2] and ordered completion [1].

Resolution operates on sets of clauses. A clause is an n -ary disjunction of literals $L_1 \vee \dots \vee L_n$ whose variables are

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CPP '19, January 14–15, 2019, Cascais, Portugal
© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-6222-1/19/01...\$15.00
<https://doi.org/10.1145/3293880.3294100>

interpreted universally. Each literal is either an atom A or its negation $\neg A$. An atom is a symbol applied to a tuple of terms—e.g., $\text{prime}(n)$. The empty clause is denoted by \perp .

Resolution works by refutation: Conceptually, the calculus proves a conjecture $\forall \bar{x}. C$ from a set of axioms \mathcal{D} by deriving \perp from $\mathcal{D} \cup \{\exists \bar{x}. \neg C\}$, indicating its unsatisfiability. As an optimization, it uses a redundancy criterion to discard tautologies, subsumed clauses, and other unnecessary clauses; for example, $p(x) \vee q(x)$ and $p(5)$ are both subsumed by $p(x)$. Compared with plain resolution, *ordered resolution* relies on an order on the atoms to further prune the search space.

Modern superposition provers are highly optimized programs that rely on sophisticated calculi, with a rich metatheory. In this paper, we propose to verify, using Isabelle/HOL [30], a purely functional prover based on ordered resolution. Although our primary interest is in metatheory per se, there are of course applications for verified provers [50].

The verification relies on stepwise refinement [55]. Four layers are connected by three refinement steps.

Our starting point, layer 1 (Section 3), is an abstract Prolog-style nondeterministic resolution prover in a highly general form, as presented by Bachmair and Ganzinger [2] and as formalized in our earlier work [39, 40]. It operates on possibly infinite sets of clauses. Its soundness and refutational completeness are inherited by the other layers.

Layer 2 (Section 4) operates on finite multisets of clauses and introduces a priority queue to ensure that inferences are performed in a fair manner, guaranteeing completeness: Given a valid conjecture, the prover will eventually derive \perp .

Layer 3 (Section 5) is a deterministic program that works on finite lists, committing to a strategy for assigning priorities to clauses. However, it is not fully executable: It abstracts over operations on atoms and employs logical specifications instead of executable functions for auxiliary notions.

Finally, layer 4 (Section 6) is a fully executable program. It provides a concrete datatype for atoms and executable definitions for all auxiliary notions, including unifiers, clause subsumption, and the order on atoms.

From layer 4, we can extract Standard ML code by invoking Isabelle's code generator [11]. The resulting prover constitutes a proof of concept: It uses an efficient calculus (layer 1) and a reasonable strategy to ensure fairness (layers 2 and 3), but depends on inefficient list-based data structures. Further refinement steps will be required to obtain a prover that is competitive with the state of the art.

The refinement steps connect vastly different levels of abstraction. The most abstract level is occupied by an infinitary logical calculus and the semantics of first-order logic. Soundness and completeness relate these two notions. At the functional programming level, soundness amounts to a safety property: Whenever the program terminates normally, its outcome is correct, whether it is a proof or a finite *saturation* witnessing unprovability. Correspondingly, refutational completeness is a liveness property: If the conjecture is valid, the program will always terminate normally. We find that, far from being academic exercises, Bachmair and Ganzinger’s framework [2] and its formalization [39, 40] adequately capture the metatheory of actual provers.

To our knowledge, our program is the first verified prover for first-order logic implementing an optimized calculus. It is also the first example of the application of refinement in a first-order context. This methodology has been used to verify SAT solvers [6, 29], which decide the satisfiability of propositional formulas, but first-order logic is semidecidable—sound and complete provers are guaranteed to terminate only for unsatisfiable (i.e., provable) clause sets. This complicates the transfer of completeness results across refinement layers.

The present work is part of the IsaFoL (Isabelle Formalization of Logic) project,¹ which aims at developing a library of results about logic and automated reasoning [3]. The Isabelle files are available in the *Archive of Formal Proofs* [38, 39] and in the IsaFoL repository.² The parts specific to the functional prover refinement amount to about 4000 lines of source text. A convenient way to study the files is to open them in Isabelle/jEdit [54], as explained in the repository’s readme file. The files were created using Isabelle version 2018, but the repositories will be updated to follow Isabelle’s evolution.

2 Atoms and Substitutions

The first three refinement layers are based on an abstract library of first-order atoms and substitutions. In the fourth and final layer, the abstract framework is instantiated with concrete datatypes and functions.

We start from IsaFoL’s library of clausal logic [6], which is parameterized by a type *'a* of logical atoms. Literals *L* are defined as an inductive datatype: *'a literal* = Pos *'a* | Neg *'a*. The type of clauses *C, D, E* is introduced as the alias *'a clause* = *'a literal multiset*, where *multiset* is the type constructor of finite multisets. Thus, the clause $\neg A \vee B$, where *A* and *B* are arbitrary atoms, is represented by the multiset {Neg *A*, Pos *B*}, and the empty clause \perp is represented by the empty multiset \emptyset . The complement operation is defined as \neg Neg *A* = Pos *A* and \neg Pos *A* = Neg *A* for any atom *A*.

In automated reasoning, it is customary to view clauses as multisets of literals rather than as sets. One reason is that

multisets behave more naturally under substitution. For example, applying $\{y \mapsto x\}$ to the two-literal clause $p(x) \vee p(y)$ results in $p(x) \vee p(x)$, which preserves the clause’s structure.

The truth value of ground (i.e., variable-free) atoms is given by a *Herbrand interpretation*: a set *I*, of type *'a set*, of all true ground atoms. The “models” predicate \models is defined as $I \models A \iff A \in I$. This definition is lifted to literals, clauses, and sets of clauses in the usual way. A set of clauses *D* is *satisfiable* if there exists an interpretation *I* such that $I \models D$.

Resolution depends on a notion of substitution and of most general unifier (MGU). These auxiliary concepts are provided by a third-party library, IsaFoR (Isabelle Formalization of Rewriting) [51]. To reduce our dependency on external libraries, we hide them behind abstract locales parameterized by a type of atoms *'a* and a type of substitutions *'s*.

We start by defining a locale *substitution_ops* that declares application (\cdot), identity (*id*), and composition (\circ):

```
locale substitution_ops =
  fixes id :: 's and o :: 's => 's => 's and dot :: 'a => 's => 'a
```

Within the locale’s scope, we introduce a number of derived concepts. Ground atoms are defined as those atoms that are left unchanged by substitutions: $\text{is_ground } A \iff \forall \sigma. A = A \cdot \sigma$. A ground substitution is a substitution whose application always results in ground atoms. Nonstrict and strict generalization are defined as

```
generalizes A B <=> exists sigma. A dot sigma = B
strictly_generalizes A B <=> generalizes A B
  and not generalizes B A
```

The operators on atoms are lifted to literals, clauses, and sets of clauses. The grounding of a clause is defined as

```
grounding_of C = {C dot sigma | is_ground sigma}
```

The operator is lifted to sets of clauses in the obvious way. Clause subsumption is defined as

```
subsumes C D <=> exists sigma. C dot sigma subseteq D
```

with *strictly_subsumes* as its strict counterpart.

The next locale, *substitution*, characterizes the operations defined by *substitution_ops*. A separate locale is necessary because we cannot interleave assumptions and definitions in a single locale. In addition, *substitution* fixes a function for renaming clauses apart (so that they do not share any variables) and a function that, given a list of atoms, constructs an atom with these as subterms:

```
locale substitution = substitution_ops +
  fixes
    renamings_apart :: 'a clause list => 's list and
    atm_of_atms :: 'a list => 'a
  assumes
    A dot id = A
    A dot (sigma o tau) = (A dot sigma) dot tau
    (forall A. A dot sigma = A dot tau) => sigma = tau
    is_ground (C dot sigma) => exists tau. is_ground tau and C dot tau = C dot sigma
```

¹<https://bitbucket.org/isafol/isafol/wiki/Home>

²https://bitbucket.org/isafol/isafol/src/master/Functional_Ordered_Resolution_Prover/

```

wf strictly_generalizes
|renamings_apart Cs| = |Cs|
ρ ∈ renamings_apart Cs ⇒ is_renaming ρ
var_disjoint (Cs · renamings_apart Cs)
atm_of_atms As · σ = atm_of_atms Bs ⇔
  map (λA. A · σ) As = Bs

```

The above definition is presented to give a flavor of our development. We refer to the Isabelle files for the exact definitions. Inside the locale, we prove further properties of the *substitution_ops* operations. Notably, we prove well-foundedness of the *strictly_subsumes* predicate based on the well-foundedness of *strictly_generalizes*, which is stated as an assumption. The *atm_of_atms* operation is used to encode a clause as a single atom in this well-foundedness proof.

Finally, a third locale, *mgu*, extends *substitution* by fixing a function $\text{mgu} :: 'a \text{ set set} \Rightarrow 's \text{ option}$ that computes an MGU σ given a set of unification constraints.

3 Bachmair and Ganzinger’s Prover

Our earlier formalization [39, 40] of a nondeterministic ordered resolution prover presented by Bachmair and Ganzinger [2] forms layer 1 of our refinement. In this paper, we restrict our focus to binary resolution, which can be implemented efficiently and forms the basis of modern provers.

The ordered resolution calculus is parameterized by a total order $>$ (“larger than”) on ground atoms. For first-order logic, the order $>$ is extended to an order $>$ on nonground atoms so that $B > A$ if and only if for all ground substitutions σ , we have $B \cdot \sigma > A \cdot \sigma$. The calculus consists of the single rule

$$\frac{C \vee A_1 \vee \dots \vee A_k \quad \neg A \vee D}{(C \vee D) \cdot \sigma}$$

where σ is the (canonical) MGU that solves the unification problem $A_1 \stackrel{?}{=} \dots \stackrel{?}{=} A_k \stackrel{?}{=} A$. In addition, each $A_i \cdot \sigma$ must be strictly $>$ -maximal with respect to the atoms in $C \cdot \sigma$ (meaning that A_i is not \leq any atom in $C \cdot \sigma$), and $A \cdot \sigma$ is $>$ -maximal with respect to the atoms in $D \cdot \sigma$. To achieve completeness, the rule must be adapted slightly to rename apart the variables occurring in different premises.

A set of clauses \mathcal{D} is *saturated* if any conclusion from premises in \mathcal{D} is already in \mathcal{D} . The ordered resolution calculus is refutationally complete, meaning that any unsatisfiable saturated set of clauses necessarily contains \perp .

Resolution provers start with a finite set of initial clauses—the input problem—and successively add conclusions from premises in the set. If the inference rule is applied in a fair fashion, the set reaches saturation at the limit; if the set is unsatisfiable, this means \perp is eventually derived.

Crucially, not only do efficient provers add clauses to their working set, they also remove clauses that are deemed redundant. This requires a refined notion of saturation. We call a set of clauses \mathcal{D} *saturated up to redundancy*, written *saturated_upto* \mathcal{D} , if any inference from nonredundant clauses in \mathcal{D} yields a redundant conclusion.

Bachmair and Ganzinger’s nondeterministic first-order prover, called RP, captures the “dynamic” aspects of saturation. RP is defined as an inductive predicate \rightsquigarrow on states, which are triples $\mathcal{S} = (\mathcal{N}, \mathcal{P}, \mathcal{O})$ of *new clauses* \mathcal{N} , *processed clauses* \mathcal{P} , and *old clauses* \mathcal{O} . Initially, \mathcal{N} is the input problem, and $\mathcal{P} \cup \mathcal{O}$ is empty. Clauses can be removed if they are tautological or subsumed or after subsumption resolution has been applied. When all clauses in \mathcal{N} have been processed (either removed entirely or moved to \mathcal{P}), a clause C from \mathcal{P} can be chosen for *inference computation*: C is then moved from \mathcal{P} to \mathcal{O} , and all its conclusions with premises from the other old clauses form the new \mathcal{N} . Formally:

inductive $\rightsquigarrow :: 'a \text{ state} \Rightarrow 'a \text{ state} \Rightarrow \text{bool}$ **where**

```

Neg A ∈ C ∧ Pos A ∈ C ⇒
(N ∪ {C}, P, O) ∼1 (N, P, O)
| D ∈ P ∪ O ∧ subsumes D C ⇒
(N ∪ {C}, P, O) ∼2 (N, P, O)
| D ∈ N ∧ strictly_subsumes D C ⇒
(N, P ∪ {C}, O) ∼3 (N, P, O)
| D ∈ N ∧ strictly_subsumes D C ⇒
(N, P, O ∪ {C}) ∼4 (N, P, O)
| D ∈ P ∪ O ∧ reduces D C L ⇒
(N ∪ {C ⊔ {L}}, P, O) ∼5 (N ∪ {C}, P, O)
| D ∈ N ∧ reduces D C L ⇒
(N, P ∪ {C ⊔ {L}}, O) ∼6 (N, P ∪ {C}, O)
| D ∈ N ∧ reduces D C L ⇒
(N, P, O ∪ {C ⊔ {L}}) ∼7 (N, P ∪ {C}, O)
| (N ∪ {C}, P, O) ∼8 (N, P ∪ {C}, O)
| (∅, P ∪ {C}, O) ∼9
  (concl_of ' infers_between O C, P, O ∪ {C})

```

Subscripts on \rightsquigarrow identify the rules. The notation $f^* X$ stands for the image of X under f , *infers_between* $\mathcal{O} C$ calculates all the ordered resolution inferences whose premises are a subset of $\mathcal{O} \cup \{C\}$ that contains C , and *reduces* $D C L$ is defined as $\exists D' L' \sigma. D = D' \sqcup \{L'\} \wedge -L = L' \cdot \sigma \wedge D' \cdot \sigma \subseteq C$.

The following derivation shows that RP can diverge even on unsatisfiable clause sets:

```

({¬p(a, a), p(x, x), ¬p(f(x), y) ∨ p(x, y)}, ∅, ∅)
∼8+ (∅, {¬p(a, a), p(x, x), ¬p(f(x), y) ∨ p(x, y)}, ∅)
∼9 (∅, {¬p(a, a), p(x, x)}, {¬p(f(x), y) ∨ p(x, y)})
∼9 ({p(x, f(x))}, {¬p(a, a)}, {¬p(f(x), y) ∨ p(x, y), p(x, x)})
∼8 (∅, {¬p(a, a), p(x, f(x))}, {¬p(f(x), y) ∨ p(x, y), p(x, x)})
∼9 ({p(x, f(f(x)))}, {¬p(a, a)},
  {¬p(f(x), y) ∨ p(x, y), p(x, x), p(x, f(x))})
∼8 ...

```

We can leave $\neg p(a, a)$ in \mathcal{P} forever and always generate more clauses of the form $p(x, f^i(x))$, for increasing values of i . This emphasizes the importance of a fair strategy for selecting clauses to move from \mathcal{P} to \mathcal{O} using rule 9.

Formally, a *derivation* is a possibly infinite sequence of states $\mathcal{S}_0 \rightsquigarrow \mathcal{S}_1 \rightsquigarrow \mathcal{S}_2 \rightsquigarrow \dots$. In Isabelle, this is expressed by the codatatype of lazy lists:

codatatype $'a$ *l*list = LNil | LCons $'a$ ($'a$ *l*list)

Lazy list operation names are prefixed by an L or l to distinguish them from the corresponding operations on finite lists. For example, `lhd xs` yields `xs`'s head (if `xs ≠ LNil`), and `lnth xs i` yields the $(i + 1)$ st element of `xs` (if $i < |xs|$).

We capture the mathematical notation $\mathcal{S}_0 \rightsquigarrow \mathcal{S}_1 \rightsquigarrow \dots$ formally as chain $(\rightsquigarrow) \mathcal{S}s$, where $\mathcal{S}s$ is a lazy list of states and chain is a coinductive predicate:

coinductive chain :: ($'a \Rightarrow 'a \Rightarrow \text{bool}$) $\Rightarrow 'a$ *l*list $\Rightarrow \text{bool}$
where

chain R (LCons x LNil)
| chain R $xs \wedge R x$ (lhd xs) \Rightarrow chain R (LCons x xs)

Coinduction is used to allow infinite chains. The base case is needed to allow finite chains. Chains cannot be empty.

Another important notion is that of the limit of a sequence Xs of sets. It is defined as the set of elements that are members of all positions of Xs except for an at most finite prefix:

definition Liminf :: $'a$ set *l*list $\Rightarrow 'a$ set **where**

Liminf $Xs = \bigcup_{i < |Xs|} \bigcap_{j, i \leq j < |Xs|} \text{lnth } Xs j$

Liminf and other operators working on clause sets are lifted pointwise to states. For example, the limit of a sequence of states is defined as $\text{Liminf } \mathcal{S}s = (\text{Liminf } \mathcal{N}s, \text{Liminf } \mathcal{P}s, \text{Liminf } \mathcal{O}s)$, where $\mathcal{N}s$, $\mathcal{P}s$, and $\mathcal{O}s$ are the projections of the \mathcal{N} , \mathcal{P} , and \mathcal{O} components of $\mathcal{S}s$.

The soundness theorem states that if RP derives \perp (i.e., \emptyset) from a set of clauses, that set must be unsatisfiable:

theorem *RP_sound*:

$\emptyset \in \text{class_of } (\text{Liminf } \mathcal{S}s) \Rightarrow$
 $\neg \text{satisfiable } (\text{grounding_of } (\text{lhd } \mathcal{S}s))$

In the above, $\text{class_of } (\mathcal{N}, \mathcal{P}, \mathcal{O}) = \mathcal{N} \cup \mathcal{P} \cup \mathcal{O}$.

A stronger, finer-grained notion of soundness relates models before and after a transition:

theorem *RP_model*:

$\mathcal{S} \rightsquigarrow \mathcal{S}' \Rightarrow$
 $(I \models \text{grounding_of } \mathcal{S}' \iff I \models \text{grounding_of } \mathcal{S})$

The canonical way of expressing the unsatisfiability of a set or multiset of first-order clauses with respect to Herbrand interpretations is as the unsatisfiability of its grounding.

Completeness of the prover can only be guaranteed when its rules are executed in a fair order, such that clauses do not get stuck forever in \mathcal{N} or \mathcal{P} . Accordingly, fairness is defined as $\text{Liminf } \mathcal{N}s = \text{Liminf } \mathcal{P}s = \emptyset$. The completeness theorem states that the limit of a fair derivation $\mathcal{S}s$ is saturated:

theorem *RP_saturated_if_fair*:

$\text{fair } \mathcal{S}s \Rightarrow \text{saturated_upto } (\text{Liminf } (\text{grounding_of } \mathcal{S}s))$

In particular, if the initial problem is unsatisfiable, \perp must appear in the \mathcal{O} component of the limit of any fair derivation:

corollary *RP_complete_if_fair*:

$\text{fair } \mathcal{S}s \wedge \neg \text{satisfiable } (\text{grounding_of } (\text{lhd } \mathcal{S}s)) \Rightarrow$
 $\emptyset \in \mathcal{O_of } (\text{Liminf } \mathcal{S}s)$

4 Ensuring Fairness

The second refinement layer is the prover RP_w , which ensures fairness by assigning a *weight* to every clause and by organizing the set of processed clauses—the \mathcal{P} state component—as a priority queue, where lighter clauses are chosen first. By assigning somewhat heavier weights to newer clauses, we can guarantee that all derivations are fair.

Another necessary ingredient for fairness is that derivations must be complete. For example, the incomplete derivation consisting of the single state $(\{C\}, \emptyset, \emptyset)$ is not fair. This requirement is expressed formally as $\text{full_chain } (\rightsquigarrow_w) \mathcal{S}s$. For the rest of this section, we fix a full chain $\mathcal{S}s$ such that $\mathcal{P_of } (\text{lhd } \mathcal{S}s) = \mathcal{O_of } (\text{lhd } \mathcal{S}s) = \emptyset$.

Because each RP_w rule corresponds to an RP rule, it is straightforward to lift the soundness and completeness results from RP to RP_w . The main difficulty is to show that the priority queue ensures fairness of full derivations, which is needed to obtain an unconditional completeness theorem for RP_w , without the assumption $\text{fair } \mathcal{S}s$.

Definition. The weight of a clause C , which defines its priority in the queue, may depend both on the clause itself and on when it was generated. To reflect this, the RP_w prover represents clauses by a pair (C, i) , where i is the *timestamp*. The larger the timestamp, the newer the clause. A state is now a quadruple $\mathcal{S} = (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$, where the first three components are finite multisets and t is the timestamp to assign to the next generation of clauses. Formally, we have the following type abbreviations:

type_synonym $'a$ *w*clause = $'a$ clause \times nat

type_synonym $'a$ *w*state =

$'a$ *w*clause multiset $\times 'a$ *w*clause multiset
 $\times 'a$ *w*clause multiset \times nat

We extend the *FO_resolution_prover* locale, in which RP is defined, with a weight function that, for any given clause, is strictly monotone with respect to the timestamp, so that older copies of a clause are preferred to newer ones:

locale *weighted_FO_resolution_prover* =

FO_resolution_prover +

fixes *weight* :: $'a$ *w*clause \Rightarrow nat

assumes $i < j \Rightarrow \text{weight } (C, i) < \text{weight } (C, j)$

The weight function is otherwise arbitrary. This gives nearly unlimited freedom when selecting clauses, which is possibly the most crucial heuristic in modern provers [43]. For example, breadth-first search corresponds to the instance where $\text{weight } (C, i)$ is defined as i .

The RP_w prover uses $'a$ *w*clause for clauses. It is defined inductively as follows:

inductive \rightsquigarrow_w :: $'a$ *w*state $\Rightarrow 'a$ *w*state $\Rightarrow \text{bool}$ **where**

Neg $A \in C \wedge \text{Pos } A \in C \Rightarrow$

$(\mathcal{N} \uplus \{(C, i)\}, \mathcal{P}, \mathcal{O}, t) \rightsquigarrow_{w1} (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$

| $D \in \text{fst } (\mathcal{P} \uplus \mathcal{O}) \wedge \text{subsumes } D C \Rightarrow$

$(\mathcal{N} + \{(C, i)\}, \mathcal{P}, \mathcal{O}, t) \rightsquigarrow_{w2} (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$

$$\begin{aligned}
& | D \in \text{fst}' \mathcal{N} \wedge C \in \text{fst}' \mathcal{P} \wedge \text{strictly_subsumes } D C \implies \\
& \quad (\mathcal{N}, \mathcal{P}, \mathcal{O}, t) \rightsquigarrow_{w3} (\mathcal{N}, \{(E, k) \in \mathcal{P}. E \neq C\}, \mathcal{O}, t) \\
& | D \in \text{fst}' \mathcal{N} \wedge \text{strictly_subsumes } D C \implies \\
& \quad (\mathcal{N}, \mathcal{P}, \mathcal{O} \uplus \{(C, i)\}, t) \rightsquigarrow_{w4} (\mathcal{N}, \mathcal{P}, \mathcal{O}, t) \\
& | D \in \text{fst}' (\mathcal{P} \uplus \mathcal{O}) \wedge \text{reduces } D C L \implies \\
& \quad (\mathcal{N} \uplus \{(C \uplus \{L\}, i)\}, \mathcal{P}, \mathcal{O}, t) \rightsquigarrow_{w5} (\mathcal{N} \uplus \{(C, i)\}, \mathcal{P}, \mathcal{O}, t) \\
& | D \in \text{fst}' \mathcal{N} \wedge \text{reduces } D C L \\
& \quad \wedge (\forall j. (C \uplus \{L\}, j) \in \mathcal{P} \implies j \leq i) \implies \\
& \quad (\mathcal{N}, \mathcal{P} \uplus \{(C \uplus \{L\}, i)\}, \mathcal{O}, t) \rightsquigarrow_{w6} (\mathcal{N}, \mathcal{P} \uplus \{(C, i)\}, \mathcal{O}, t) \\
& | D \in \text{fst}' \mathcal{N} \wedge \text{reduces } D C L \implies \\
& \quad (\mathcal{N}, \mathcal{P}, \mathcal{O} \uplus \{(C \uplus \{L\}, i)\}, t) \rightsquigarrow_{w7} (\mathcal{N}, \mathcal{P} \uplus \{(C, i)\}, \mathcal{O}, t) \\
& | (\mathcal{N} \uplus \{(C, i)\}, \mathcal{P}, \mathcal{O}, t) \rightsquigarrow_{w8} (\mathcal{N}, \mathcal{P} \uplus \{(C, i)\}, \mathcal{O}, t) \\
& | (\forall (D, j) \in \mathcal{P}. \text{weight}(C, i) \leq \text{weight}(D, j)) \wedge \\
& \quad \mathcal{N} = \text{mset_set}((\lambda D. (D, t))' \text{concl_of}' \text{infers_between} \\
& \quad (\text{set_mset}(\text{fst}' \mathcal{O})) C) \implies \\
& \quad (\emptyset, \mathcal{P} \uplus \{(C, i)\}, \mathcal{O}, t) \rightsquigarrow_{w9} \\
& \quad (\mathcal{N}, \{(D, j) \in \mathcal{P}. D \neq C\}, \mathcal{O} \uplus \{(C, i)\}, t + 1)
\end{aligned}$$

where `fst` is the function that returns the first component of a pair, `mset_set` converts a set to the multiset with exactly one copy of each element in the set, and `set_mset` converts a multiset to the set of elements in the multiset.

The most important differences with RP are in the last transition rule. This rule, which computes inferences, assigns timestamp t to each newly computed clause D and increments t . Moreover, since we want \mathcal{P} to work as a priority queue, RP_w chooses a clause C with the smallest weight.

Another difference is that RP_w uses finite multisets for representing \mathcal{N} , \mathcal{P} , and \mathcal{O} . They offer a compromise between sets in layer 1 and lists in layer 3. Finite multisets also help eliminate some unfair derivations. Finiteness guarantees that each clause in \mathcal{N} gets the opportunity to move to \mathcal{P} (and further to \mathcal{O}). Moreover, whereas the set-based RP allows idle transitions, such as $(\mathcal{N} \cup \{C\}, \mathcal{P}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P} \cup \{C\}, \mathcal{O})$ for $C \in \mathcal{N} \cap \mathcal{P}$, the use of multisets and \uplus precludes such transitions in RP_w .

Timestamps are preserved when clauses are moved between \mathcal{N} , \mathcal{P} , and \mathcal{O} . They are also preserved by reduction steps (rules 5 to 7). This works because reduction can only take place finitely many times—a k -literal clause can be reduced at most k times. Therefore, there is no risk of divergence due to an infinite chain of reductions.

Timestamps introduce a new danger. It may be the case that a clause C is in the limit if we project away the timestamps, but that no single timestamped clause (C, i) belongs to the limit because the timestamps keep changing, as in the infinite sequence $\{(C, 0)\}, \{(C, 1)\}, \{(C, 2)\}, \dots$. This could in principle arise due to subsumption, leading to derivations such as the following:

$$\begin{aligned}
(_, _ \uplus \{(C, 0)\}, _) & \rightsquigarrow_w (_, _ \uplus \{(C, 0), (C, 1)\}, _) \rightsquigarrow_w \\
(_, _ \uplus \{(C, 1)\}, _) & \rightsquigarrow_w^+ (_, _ \uplus \{(C, 1), (C, 2)\}, _) \rightsquigarrow_w \\
(_, _ \uplus \{(C, 2)\}, _) & \rightsquigarrow_w^+ \dots
\end{aligned}$$

To prevent this, the RP_w rules are formulated so that whenever they remove the earliest copy of any clause $C \in \mathcal{P}$, they

also remove all its copies from \mathcal{P} . This property is captured by the following lemma:

lemma *preserve_min_P*:

$$\begin{aligned}
S \rightsquigarrow_w S' \wedge (C, i) \in \mathcal{P}_{\text{of}} S \wedge C \in \text{fst}' \mathcal{P}_{\text{of}} S' \\
\wedge (\forall k. (C, k) \in \mathcal{P}_{\text{of}} S \implies k \geq i) \implies \\
(C, i) \in \mathcal{P}_{\text{of}} S'
\end{aligned}$$

This completes our review of RP_w . As an intermediate step towards a more concrete prover, we restrict the weight function to be a linear equation that considers both timestamps and clause sizes:

locale *weighted_FO_resolution_prover_with_size_timestamp_factors* =

$$\begin{aligned}
& \text{FO_resolution_prover} + \\
& \text{fixes } \text{size_factor} :: \text{nat} \text{ and } \text{timestamp_factor} :: \text{nat} \\
& \text{assumes } \text{timestamp_factor} > 0 \\
& \text{begin} \\
& \text{fun } \text{weight} :: 'a \text{ wclause} \implies \text{nat} \text{ where} \\
& \quad \text{weight}(C, i) = \text{size_factor} * |C| + \text{timestamp_factor} * i \\
& \text{end}
\end{aligned}$$

where $|C| = \sum_{A: A \in C \vee \neg A \in C} |A|$. It is easy to prove that this definition of `weight` is strictly monotone and hence that this locale is a sublocale of *weighted_FO_resolution_prover*. This gives us a correspondingly specialized version of RP_w that will form the basis of further refinement steps.

The idea of organizing \mathcal{P} as a priority queue is well known in the automated reasoning community. Bachmair and Ganzinger [2, p. 44] mention it in a footnote, but they require the weight to be monotone not only in the timestamp but also in the clause size, claiming that this is necessary to ensure fairness. Our proof reveals that clause size is irrelevant, even in the presence of reductions. This demonstrates how working out the details and making all assumptions explicit using a proof assistant can help clarify fine theoretical points.

Refinement Proofs. To lift the soundness and completeness results about RP to RP_w , we must first show that any possible behavior of RP_w on states of type *wstate* is a possible behavior of RP on the corresponding values of type *state*:

lemma *weighted_RP_imp_RP*:

$$S \rightsquigarrow_w S' \implies \text{state_of } S \rightsquigarrow \text{state_of } S'$$

The proof is by induction on the rules of RP_w , with one difficult case. Inference computation (rule 9) converts a set to a finite multiset using `mset_set`, which is undefined for infinite sets. Thus, we must show only a finite set of inferences may be performed from a finite clause set:

lemma *finite_ord_FO_resolution_inferences_between*:

$$\text{finite } \mathcal{D} \implies \text{finite}(\text{infers_between } \mathcal{D} C)$$

A binary resolution inference takes two premises, of the form $CAA = C \vee A_1 \vee \dots \vee A_k$ and $DA = \neg A \vee D$, and produces a conclusion $E = (C \vee D) \cdot \sigma$. It can be represented compactly by a tuple of the form (CAA, DA, AA, A, E) , where $AA = A_1 \vee \dots \vee A_k$. We must show that the set of such

tuples produced by `infers_between` is finite, assuming \mathcal{D} is finite. First, observe that the last component E of a tuple is determined by the other four. Hence it suffices to consider quadruples (CAA, DA, AA, A) . Let $\mathcal{DC} = \mathcal{D} \cup \{C\}$, and let n be the length of the longest clause in \mathcal{DC} . Moreover, let $\mathcal{A} = \bigcup_{D \in \mathcal{DC}} \text{atms_of } D$ and $\mathcal{AA} = \{\mathcal{B} \mid \text{set_mset } \mathcal{B} \subseteq \mathcal{A} \wedge |\mathcal{B}| \leq n\}$. Then all inferences between \mathcal{D} and C belong to $\mathcal{DC} \times \mathcal{DC} \times \mathcal{AA} \times \mathcal{A}$, a cartesian product of finite sets.

Soundness and Completeness Proofs. Using the refinement lemma `weighted RP_imp RP`, it is easy to lift the `RP_model` theorem (Section 3) to `RPw`:

theorem `weighted RP_model`:

$$\begin{aligned} \mathcal{S} \rightsquigarrow_w \mathcal{S}' &\implies \\ (I \models \text{grounding_of } \mathcal{S}' \iff I \models \text{grounding_of } \mathcal{S}) \end{aligned}$$

Completeness is considerably more difficult. We first show that the use of timestamps ensures that all full `RPw` derivations are fair. In principle, a full derivation could be unfair by virtue of being finite and ending in a state such as \mathcal{N} or \mathcal{P} is nonempty. However, this is impossible because a transition of rule 8 or 9 could then be taken from the last state, contradicting the hypothesis that the derivation is full. Hence, finite full derivations are necessarily fair:

lemma `fair_if_finite`:

$$\text{lfinite } \mathcal{S}s \implies \text{fair } (\text{Imap state_of } \mathcal{S}s)$$

There are two ways in which an infinite derivation $\mathcal{S}s$ in `RPw` could be unfair: A clause could get stuck forever in \mathcal{N} , or in \mathcal{P} . We show that the \mathcal{N} case is impossible by defining a measure on states that decreases with respect to the lexicographic extension of $>$ on `nat` to pairs:

abbreviation `RP_basic_measure` :: `'a wstate \implies nat2 where`

$$\begin{aligned} \text{RP_basic_measure } (\mathcal{N}, \mathcal{P}, \mathcal{O}, t) &\equiv \\ (\text{sum } ((\lambda(C, _). |C| + 1) \cdot (\mathcal{N} \uplus \mathcal{P} \uplus \mathcal{O})), |\mathcal{N}|) \end{aligned}$$

The first component of the pair is the total size of all the clauses in the state, plus 1 for each clause to ensure that empty clauses are also counted. The second component is the number of clauses in \mathcal{N} . It is easy to see why the measure is decreasing. Rule 9, inference computation, is not applicable due to our assumption that a clause remains stuck in \mathcal{N} . Rule 8, which moves a clause from \mathcal{N} to \mathcal{P} , decreases the measure's second component while leaving the first component unchanged. The other rules decrease the first component since they remove clauses or literals. Formally:

lemma `weighted RP_basic_measure_decreasing_N`:

$$\begin{aligned} \mathcal{S} \rightsquigarrow_w \mathcal{S}' \wedge (C, _) \in \mathcal{N}_of \mathcal{S} &\implies \\ (\text{RP_basic_measure } \mathcal{S}', \text{RP_basic_measure } \mathcal{S}) &\in \text{RP_basic_rel} \end{aligned}$$

where `RP_basic_rel` = `natLess <lex> natLess` and `natLess` = `{(m, n) | m < n}`.

What if a clause C is stuck in \mathcal{P} ? Lemma `preserve_min_P` states that in any step, either all copies of C are removed or the one with the lowest timestamp is kept. Hence, C 's

timestamp will either remain stable or decrease over time. Since $>$ is well founded on natural numbers, eventually a fixed i will be reached and will belong to the limit:

lemma `persistent_wclause_in_P_if_persistent_clause`:

$$\begin{aligned} C \in \text{Liminf } (\text{Imap } \mathcal{P}_of (\text{Imap state_of } \mathcal{S}s)) &\implies \\ \exists i. (C, i) \in \text{Liminf } (\text{Imap } (\text{set_mset} \circ \mathcal{P}_of) \mathcal{S}s) \end{aligned}$$

Again, we define a measure, but it must also decrease when inferences are computed and new clauses appear in \mathcal{N} . (In this case, `RP_basic_measure` may increase.) Our new measure is parameterized by a predicate p that can be used to filter out undesirable clauses:

abbreviation `RP_filtered_measure` ::

$$\begin{aligned} ('a \text{ wclause} \implies \text{bool}) \implies 'a \text{ wstate} \implies \text{nat}^3 \text{ where} \\ \text{RP_filtered_measure } p (\mathcal{N}, \mathcal{P}, \mathcal{O}, t) &\equiv \\ (\text{sum } ((\lambda(C, _). |C| + 1) \cdot \{\mathcal{D}i \in \mathcal{N} \uplus \mathcal{P} \uplus \mathcal{O} \mid p \mathcal{D}i\}), \\ |\{\mathcal{D}i \in \mathcal{N} \mid p \mathcal{D}i\}|, |\{\mathcal{D}i \in \mathcal{P} \mid p \mathcal{D}i\}|) \end{aligned}$$

Notice that the case `RP_filtered_measure` $(\lambda_ . \text{True})$ essentially amounts to `RP_basic_measure`. In the formalization, we use `RP_filtered_measure` $(\lambda_ . \text{True})$ to avoid duplication.

Suppose the clause C that is stuck in \mathcal{P} has weight w in the limit, and suppose that a clause D is moved from \mathcal{P} to \mathcal{O} by rule 9. That clause's weight must be at most w ; otherwise, it would not have been preferred to C . Thus, infinite derivations necessarily consist of segments each consisting of finitely many applications of rules other than rule 9 followed by an application of rule 9: $(\rightsquigarrow_{w1-8}^* \circ \rightsquigarrow_{w9})^\omega$. Since each application of rule 9 increases the t component of the state, eventually we reach a state in which $t > w$. As a consequence of strict monotonicity of weight, any clauses generated by inference computation from that point on will have weights above C 's, and if C remains stuck, then so must these clauses. Thus, we can ignore these clauses altogether, by using $\lambda(C, i). i \leq w$ as the filter p . We adapt the corresponding relation to consider the extra argument:

abbreviation `RP_filtered_rel` :: $(\text{nat}^3)^2 \text{ set where}$

$$\begin{aligned} \text{RP_filtered_rel} &\equiv \\ \text{natLess} <\text{lex}> \text{natLess} <\text{lex}> \text{natLess} \end{aligned}$$

The measure `RP_filtered_measure` $(\lambda(_, i). i \leq w)$ decreases for steps occurring between inference computations and for all steps once we have reached a state where $t > w$ (at which point all inference computations are blocked by C). To obtain a measure that also decreases on inference computation, we add a component $w + 1 - t$ to the measure. We also add a component `RP_basic_measure` \mathcal{S} to ensure that the measure decreases when a clause (C, i) such that $i > w$ is simplified. This yields the combined measure

abbreviation `RP_combined_measure` ::

$$\begin{aligned} \text{nat} \implies 'a \text{ wstate} \implies \text{nat} \times \text{nat}^3 \times \text{nat}^3 \text{ where} \\ \text{RP_combined_measure } w \mathcal{S} &\equiv \\ (w + 1 - \text{t_of } \mathcal{S}, \\ \text{RP_filtered_measure } (\lambda(_, i). i \leq w) \mathcal{S}, \\ \text{RP_basic_measure } \mathcal{S}) \end{aligned}$$

This measure is indeed decreasing with respect to a left-to-right lexicographic order:

lemma *weighted_RP_basic_measure_decreasing_P*:

$$\begin{aligned} S \rightsquigarrow_w S' \wedge Ci \in \mathcal{P}_{\text{of}} S \implies \\ (\text{RP_combined_measure}(\text{weight } Ci) S', \\ \text{RP_combined_measure}(\text{weight } Ci) S) \\ \in \text{natLess} \langle \text{lex} \rangle \text{RP_filtered_rel} \langle \text{lex} \rangle \text{RP_basic_rel} \end{aligned}$$

By combining the two lemmas *weighted_RP_basic_measure_decreasing_N* and *weighted_RP_basic_measure_decreasing_P*, we can prove all derivations starting with $\mathcal{P} = \mathcal{O} = \emptyset$ fair:

theorem *weighted_RP_fair*: fair (lmap state_of Ss)

Since all derivations are fair and RP_w derivations correspond to RP derivations, it is trivial to lift RP's saturation and completeness theorems:

corollary *weighted_RP_saturated*:

$$\text{saturated_upto}(\text{Liminf}(\text{lmap grounding_of } Ss))$$

corollary *weighted_RP_complete*:

$$\begin{aligned} \neg \text{satisfiable}(\text{grounding_of}(\text{lhd } Ss)) \implies \\ \emptyset \in \mathcal{O}_{\text{of}}(\text{Liminf}(\text{lmap state_of } Ss)) \end{aligned}$$

5 Eliminating Nondeterminism

The third refinement layer defines a functional program RP_d that embodies a specific rule application strategy, thereby eliminating RP_w 's nondeterminism. Clauses are represented by lists, and multisets of clauses by lists of lists.

Definition. Our prover corresponds roughly to the following pseudocode:

```
function  $\text{RP}_d(\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$ 
  repeat forever
    if  $\perp \in \mathcal{P} \uplus \mathcal{O}$  then
      return  $\mathcal{P} \uplus \mathcal{O}$ 
    else if  $N = P = \emptyset$  then
      return  $\mathcal{O}$ 
    else if  $N = \emptyset$  then
      let  $C$  be a minimal-weight clause in  $\mathcal{P}$ ;
       $\mathcal{N} :=$  conclusions of all inferences from  $\mathcal{O} \uplus \{C\}$ 
        involving  $C$ , with timestamp  $t$ ;
      move  $C$  from  $\mathcal{P}$  to  $\mathcal{O}$ ;
       $t := t + 1$ 
    else
      remove an arbitrary clause  $C$  from  $\mathcal{N}$ ;
      reduce  $C$  using  $\mathcal{P} \uplus \mathcal{O}$ ;
      if  $C = \perp$  then
        return  $\{\perp\}$ 
      else if  $C$  is neither a tautology nor subsumed by
        a clause in  $\mathcal{P} \uplus \mathcal{O}$  then
        reduce  $\mathcal{P}$  using  $C$ ;
        reduce  $\mathcal{O}$  using  $C$ , moving any reduced
          clauses from  $\mathcal{O}$  to  $\mathcal{P}$ ;
        remove all clauses from  $\mathcal{P}$  and  $\mathcal{O}$  that are
          strictly subsumed by  $C$ ;
        add  $C$  to  $\mathcal{P}$ 
```

The function should be invoked with \mathcal{N} as the input problem, $\mathcal{P} = \mathcal{O} = \emptyset$, and an arbitrary timestamp t (e.g., 0). The loop is loosely modeled after Vampire's proof procedure [52].

In Isabelle, the list-based representations compel us to introduce the following type abbreviations:

type_synonym 'a lclause = 'a literal list

type_synonym 'a dclause = 'a lclause \times nat

type_synonym 'a dstate =

'a dclause list \times 'a dclause list \times 'a dclause list \times nat

The prover is defined inside a locale that inherits *weighted_FO_resolution_prover_with_size_timestamp_factors*. The core function, $\text{RP}_d_{\text{step}}$, performs a single iteration of the main loop. Here is the definition, excluding auxiliary functions:

fun $\text{RP}_d_{\text{step}} :: 'a \text{dstate} \Rightarrow 'a \text{dstate}$ **where**

$\text{RP}_d_{\text{step}}(\mathcal{N}, \mathcal{P}, \mathcal{O}, t) =$

if $\exists Ci \in \mathcal{P} @ \mathcal{O}. \text{fst } Ci = []$ then

([], [], remdups $\mathcal{P} @ \mathcal{O}, t + |\text{remdups } \mathcal{P}|)$

else case \mathcal{N} of

[] \Rightarrow

(case \mathcal{P} of

[] $\Rightarrow (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$

| $P_0 \# \mathcal{P}' \Rightarrow$

let

(C, i) = select_min_weight_clause $P_0 \mathcal{P}'$;

$\mathcal{N} = \text{map}(\lambda D. (D, t))(\text{remdups}$

(resolve_rename $C C @ \text{concat}(\text{map}$

(resolve_rename_both_ways $C \circ \text{fst } \mathcal{O})))$;

$\mathcal{P} = \text{filter}(\lambda(D, j). \text{mset } D \neq \text{mset } C) \mathcal{P}$;

$\mathcal{O} = (C, i) \# \mathcal{O}$;

$t = t + 1$

in

($\mathcal{N}, \mathcal{P}, \mathcal{O}, t$)

| (C, i) $\# \mathcal{N} \Rightarrow$

let $C = \text{reduce}(\text{map } \text{fst}(\mathcal{P} @ \mathcal{O})) [] C$ in

if $C = []$ then

([], [], [([], i)], $t + 1$)

else if is_tautology C

$\vee \text{subsume}(\text{map } \text{fst}(\mathcal{P} @ \mathcal{O})) C$ then

($\mathcal{N}, \mathcal{P}, \mathcal{O}, t$)

else let

$\mathcal{P} = \text{reduce_all } C \mathcal{P}$;

(back_to_ \mathcal{P} , \mathcal{O}) = reduce_all2 $C \mathcal{O}$;

$\mathcal{P} = \text{back_to_}\mathcal{P} @ \mathcal{P}$;

$\mathcal{O} = \text{filter}((\neg) \circ \text{strictly_subsume } C \circ \text{fst}) \mathcal{O}$;

$\mathcal{P} = \text{filter}((\neg) \circ \text{strictly_subsume } C \circ \text{fst}) \mathcal{P}$;

$\mathcal{P} = (C, i) \# \mathcal{P}$

in

($\mathcal{N}, \mathcal{P}, \mathcal{O}, t$)

The # operator abbreviates the Cons constructor, and @ is the append operator.

The existential quantifier above is unproblematic because it ranges over a finite set, but some of the auxiliary functions use infinite quantification. Notably, subsumption of a clause

D by another clause C is defined as $\exists\sigma. C \cdot \sigma \subseteq D$ (Section 2), where σ ranges over substitutions. Nonexecutable constructs are acceptable if we know that we can replace them by equivalent executable constructs further down the refinement chain; for example, an implementation of subsumption can compute a witness σ using matching, instead of blindly enumerating all possible substitutions.

The main program is a tail-recursive function that repeatedly calls $\text{RP}_d\text{-step}$ until a final state $([], [], \mathcal{O}, t)$ is reached, at which point it returns the set \mathcal{O} stripped of its timestamps:

partial_function (*option*)

$\text{RP}_d :: 'a \text{ dstate} \Rightarrow 'a \text{ lclause list option}$

where

$\text{RP}_d \mathcal{S} = \text{if is_final } \mathcal{S} \text{ then Some (map fst } (\mathcal{O}\text{-of } \mathcal{S}))$
 $\text{else RP}_d (\text{RP}_d\text{-step } \mathcal{S})$

Since the recursion may diverge, we cannot introduce the function using the **fun** command [22]. Instead, we use **partial_function** (*option*) [23], which puts the computation in an option monad. The function's result is of the form $\text{Some } R$ if the recursion terminates and None otherwise.

Refinement Proofs. Using refinement, we connect the $\text{RP}_d\text{-step}$ function to the RP_w predicate. $\text{RP}_d\text{-step}$ has a coarser granularity than RP_w : A single invocation on a nonfinal state \mathcal{S} can amount to a chain of RP_w transitions. This is captured by the following (weak-)refinement property:

lemma nonfinal_deterministic_RP_step:

$\neg \text{is_final } \mathcal{S} \Rightarrow$
 $\text{wstate_of } \mathcal{S} \rightsquigarrow_w^+ \text{wstate_of } (\text{RP}_d\text{-step } \mathcal{S})$

where wstate_of converts RP_d states to RP_w states. The entire proof, including key lemmas, is about 1300 lines long. It follows the case distinctions present in $\text{RP}_d\text{-step}$'s definition:

case $\exists Ci \in \mathcal{P} @ \mathcal{O}. \text{fst } Ci = []$:

By induction on $|\text{remdups } \mathcal{P}|$ (where remdups removes duplicates), there must exist a derivation of the form

$\text{wstate_of } (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w2}^* \text{wstate_of } ([], \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w9} \text{wstate_of } (\mathcal{N}', \mathcal{P}', (C, i) \# \mathcal{O}, t + 1)$
 $\rightsquigarrow_w^* \text{wstate_of } ([], [], \mathcal{O}', t + |\text{remdups } \mathcal{P}'|)$

for $\mathcal{P}' = \text{filter } (\lambda(D, j). \text{mset } D \neq \text{mset } C) \mathcal{P}$, $\mathcal{O}' = \text{remdups } \mathcal{P}' @ \mathcal{O}$, and suitable \mathcal{N}' and $(C, i) \in \mathcal{P}$. The last step is justified by the induction hypothesis.

case $\mathcal{N} = \mathcal{P} = []$:

Contradiction with the assumption that $(\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$ is a nonfinal state.

case $\mathcal{N} = []$:

It suffices to show that the transition

$\text{wstate_of } ([], \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w9} \text{wstate_of } (\mathcal{N}', \mathcal{P}', (C, i) \# \mathcal{O}, t + 1)$

is possible, where $(C, i) \in \mathcal{P}$ is a minimal-weight clause, $\mathcal{N}' = \text{map } (\lambda D. (D, t)) (\text{remdups } (\text{resolve_rename } C C$

@ concat (map (resolve_rename_both_ways $C \circ \text{fst } \mathcal{O})))$, and $\mathcal{P}' = \text{filter } (\lambda(D, j). \text{mset } D \neq \text{mset } C) \mathcal{P}$. The main proof obligation is that \mathcal{N}' , converted to multisets, equals the multiset $\text{mset_set } ((\lambda D. (D, t)) \text{ ` } \text{concl_of } \text{ ` } \text{infers_between } (\text{set_mset } (\text{fst } \text{ ` } \mathcal{O})) C)$ specified in rule \rightsquigarrow_{w9} . The distance between the functional program and its mathematical specification is at its greatest here. The proof is tedious but straightforward.

otherwise:

Let $C' = \text{reduce } (\text{map fst } \mathcal{P} @ \text{map fst } \mathcal{O}) [] C$. If $C' = []$, then

$\text{wstate_of } ((C, i) \# \mathcal{N}', \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w5}^* \text{wstate_of } ([], i) \# \mathcal{N}', \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w3}^* \text{wstate_of } ([], i) \# \mathcal{N}', [], \mathcal{O}, t)$
 $\rightsquigarrow_{w4}^* \text{wstate_of } ([], i) \# \mathcal{N}', [], [], t)$
 $\rightsquigarrow_{w8} \text{wstate_of } (\mathcal{N}', [([], i)], [], t)$
 $\rightsquigarrow_{w2}^* \text{wstate_of } ([], [([], i)], [], t)$
 $\rightsquigarrow_{w9} \text{wstate_of } ([], [], [([], i)], t)$

Otherwise, if $\text{is_tautology } C' \vee \text{subsume } (\text{map fst } (\mathcal{P} @ \mathcal{O})) C'$, then

$\text{wstate_of } ((C, i) \# \mathcal{N}, \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w5}^* \text{wstate_of } ((C', i) \# \mathcal{N}, \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w1,2} \text{wstate_of } (\mathcal{N}, \mathcal{P}, \mathcal{O}, t)$

Otherwise:

$\text{wstate_of } ((C, i) \# \mathcal{N}', \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w5}^* \text{wstate_of } ((C', i) \# \mathcal{N}', \mathcal{P}, \mathcal{O}, t)$
 $\rightsquigarrow_{w6}^* \text{wstate_of } ((C', i) \# \mathcal{N}', \mathcal{P}', \mathcal{O}, t)$
 $\rightsquigarrow_{w7}^* \text{wstate_of } ((C', i) \# \mathcal{N}', \text{back_to_} \mathcal{P} @ \mathcal{P}', \mathcal{O}', t)$
 $\rightsquigarrow_{w4}^* \text{wstate_of } ((C', i) \# \mathcal{N}', \text{back_to_} \mathcal{P} @ \mathcal{P}', \mathcal{O}'', t)$
 $\rightsquigarrow_{w3}^* \text{wstate_of } ((C', i) \# \mathcal{N}', \mathcal{P}'', \mathcal{O}'', t)$
 $\rightsquigarrow_{w8} \text{wstate_of } (\mathcal{N}', (C', i) \# \mathcal{P}'', \mathcal{O}'', t)$

for suitable lists \mathcal{P}' , $\text{back_to_} \mathcal{P}$, \mathcal{O}' , \mathcal{O}'' , and \mathcal{P}'' .

Soundness and Completeness Proofs. Let $\mathcal{S}_0 = (\mathcal{N}_0, [], [], t_0)$ be an arbitrary initial state. Soundness means that whenever $\text{RP}_d \mathcal{S}_0$ terminates with some clause set R , then R is a saturation that satisfies the same models as \mathcal{N}_0 . Moreover, if \mathcal{N}_0 is unsatisfiable, then R contains \perp , providing a simple syntactic check for unsatisfiability. Completeness means that divergence is possible only if \mathcal{N}_0 is satisfiable. For satisfiable clause sets \mathcal{N}_0 , both termination and divergence are possible.

To lift soundness and completeness results from RP_w to RP_d , we first define \mathcal{S}_s as a full chain of nontrivial RP_d steps starting from \mathcal{S}_0 . We let $\mathcal{S}_s = \text{derivation_from } \mathcal{S}_0$, with

primcorec $\text{derivation_from} :: 'a \text{ dstate} \Rightarrow 'a \text{ dstate llist}$
where

$\text{derivation_from } \mathcal{S} = \text{LCons } \mathcal{S} \text{ (if is_final } \mathcal{S} \text{ then LNil}$
 $\text{else derivation_from } (\text{RP}_d\text{-step } \mathcal{S}))$

Based on \mathcal{S}_s , we let $w\mathcal{S}_s = \text{Imap wstate_of } \mathcal{S}_s$ and note that $w\mathcal{S}_s$ is a full chain of “big” \rightsquigarrow_w^+ steps. Using a lemma that will be proved below, we obtain a full chain $s\mathcal{S}_s$ of “small”

\rightsquigarrow_w steps. This chain satisfies the conditions postulated on Ss in Section 4, allowing us to lift the results presented there.

The soundness results are proved in a nameless locale, or *context*, that assumes termination of RP_d :

fixes $R :: 'a$ lclause list
assumes $RP_d \mathcal{S}_0 = \text{Some } R$

The definition of RP_d , using **partial_function**, gives us an induction rule restricted to the case where RP_d terminates (i.e., returns a *Some* value). This rule can be used to prove that Ss and hence wSs and $sswSs$ are finite sequences.

Soundness takes the form of a pair of theorems that lift *weighted_RP_model* and *weighted_RP_saturated*:

theorem deterministic_RP_model:
 $I \models \text{grounding_of } \mathcal{N}_0 \iff I \models \text{grounding_of } R$

theorem deterministic_RP_saturated:
 $\text{saturated_upto } (\text{grounding_of } R)$

In most applications, all that matters is the satisfiability status of the set \mathcal{N}_0 . It can be retrieved syntactically:

corollary deterministic_RP_refutation:
 $\neg \text{satisfiable } (\text{grounding_of } \mathcal{N}_0) \iff \emptyset \in R$

Completeness is proved in a separate context that assumes nontermination: $RP_d \mathcal{S}_0 = \text{None}$. The strongest result we prove is that this assumption implies the satisfiability of \mathcal{N}_0 :

theorem deterministic_RP_complete:
 $\text{satisfiable } (\text{grounding_of } \mathcal{N}_0)$

The proof is by contradiction:

Assume that $\neg \text{satisfiable } (\text{grounding_of } \mathcal{N}_0)$. Hence, by *weighted_RP_complete* we have $\emptyset \in \mathcal{O}_{\text{of}} \text{ } sswSs$. It is easy to show that $sswSs$'s limit is a subset of wSs 's limit; hence $\emptyset \in \mathcal{O}_{\text{of}} wSs$. This implies the existence of a natural number k such that $\emptyset \in \mathcal{O}_{\text{of}} (\text{Inth } wSs \ k)$. Hence $\emptyset \in \mathcal{O}_{\text{of}} (RP_d\text{-step}^k \ \mathcal{S}_0)$. However, by an induction on k , we can show that RP_d must terminate after at most k iterations, contradicting the assumption that RP_d diverges.

A Coinductive Puzzle. A single “big” step of the deterministic prover RP_d may correspond to many “small” steps of the weighted prover RP_w . To transfer the results from RP_w to RP_d , we must expand the big steps. The core of the expansion is an abstract property of chains and transitive closure:

Let R be a relation and xs a chain of R^+ transitions. There exists a chain of R transitions that embeds xs —i.e., that contains all elements of xs in the same order and with only finitely many elements inserted between each pair of consecutive elements of xs .

On finite chains, this property can be proved by straightforward induction. But the completeness proof must also consider infinite chains. Coinduction and corecursion up-to techniques are useful for such tasks.

The desired property is stated formally as follows:

lemma chain_tranclp_imp_exists_chain:
 $\text{chain } R^+ \ xs \implies$
 $\exists ys. \text{chain } R \ ys \wedge xs \sqsubseteq ys \wedge \text{lhd } xs = \text{lhd } ys$
 $\wedge \text{llast } xs = \text{llast } ys$

where the embedding \sqsubseteq of lazy lists is defined coinductively using $++$, which prepends a finite list to a lazy list:

coinductive $\sqsubseteq :: 'a$ llist $\implies 'a$ llist $\implies \text{bool}$ **where**
 $\text{lfinite } xs \implies \text{LNil} \sqsubseteq xs$
 $| \ xs \sqsubseteq ys \implies \text{LCons } x \ xs \sqsubseteq zs ++ \text{LCons } x \ ys$
fun $++ :: 'a$ list $\implies 'a$ llist $\implies 'a$ llist **where**
 $[] ++ xs = xs$
 $| (z \# zs) ++ xs = \text{LCons } z \ (zs ++ xs)$

The definition of \sqsubseteq ensures that infinite lazy lists only embed other infinite lazy lists, but not the finite ones: $xs \sqsubseteq ys \implies (\text{lfinite } xs \iff \text{lfinite } ys)$. The unguarded calls to llast may seem worrying, but the function is conveniently defined to always return the same unspecified element for infinite lists.

To prove the lemma above, we instantiate the existential quantifier by the following corecursively defined witness:

corec $\text{wit} :: ('a \implies 'a \implies \text{bool}) \implies 'a$ llist $\implies 'a$ llist **where**
 $\text{wit } R \ xs = (\text{case } xs \text{ of}$
 $\text{LCons } x \ (\text{LCons } y \ ys) \implies$
 $\text{LCons } x \ (\text{pick } R \ x \ y ++ \text{wit } R \ (\text{LCons } y \ ys))$
 $| _ \implies xs)$

Here $\text{pick } R \ x \ y$ returns an arbitrary finite list of R -related states connecting the R^+ -related x and y . Its definition is $\text{pick } R \ x \ y = (\text{SOME } zs. \text{chain } R \ (\text{lhist_of } (x \# zs @ [y])))$, where lhist_of converts finite lists into lazy lists and SOME is Hilbert's choice operator. Thus, pick satisfies the characteristic property $R^+ \ x \ y \implies \text{chain } R \ (\text{lhist_of } (x \# \text{pick } R \ x \ y @ [y]))$. The nonexecutability entailed by the use of Hilbert choice is unproblematic because the wit function is used only in the proofs and not in the prover's code.

The definition of wit is not primitively corecursive. Although there is a guarding LCons constructor, the corecursive call occurs under $++$, which makes the productivity of this function nontrivial. This syntactic structure of the definition is called *corecursive up to* $++$. Ultimately, wit is productive because $++$ does not remove any LCons constructors from its second arguments. A slightly weaker requirement, called *friendliness*, is supported by Isabelle's **corec** command [5]. For the above definition to be accepted $++$ must be registered as a “friend.” This involves a one-line proof.

The four conjuncts in *chain_tranclp_imp_exists_chain* are discharged in turn under the assumption $\text{chain } R^+ \ xs$. In order of increasing difficulty: $\text{lhd } (\text{wit } R \ xs) = \text{lhd } xs$ follows by simple rewriting. Next, $\text{llast } (\text{wit } R \ xs) = \text{llast } xs$ requires an induction in the case of finite chains xs . For any infinite chain zs , $\text{llast } zs$ is defined as a fixed unspecified $'a$ value. The properties $xs \sqsubseteq \text{wit } R \ xs$ and $\text{chain } R \ (\text{wit } R \ xs)$ require a coinduction on \sqsubseteq and chain , respectively. In keeping with the

definition, plain coinduction on \sqsubseteq and chain does not suffice, and we must use coinduction up to $\dashv\vdash$ on \sqsubseteq and chain.

The property *chain_tranclp_imp_exists_chain* easily extends to full chains.

6 Obtaining Executable Code

Our deterministic prover RP_d is already quite close to being an executable program. The fourth refinement, the prover RP_x , adds the missing ingredients: a concrete representation of terms and an executable algorithm for clause subsumption.

First-Order Terms. We instantiate our abstract notion of atom using a particularly comprehensive formalization of terms developed as part of the IsaFoR library [51]. This rewriting-independent part of IsaFoR has recently moved to the *Archive of Formal Proofs* [47].

IsaFoR terms are defined as the following datatype:

datatype $(f, 'v)$ *term* = Var $'v$ | Fun $'f$ $((f, 'v)$ *term list*)

To simplify notation, in this paper we fix $'f = 'v = \text{nat}$ and abbreviate $(f, 'v)$ *term* by *term*. In the formalization, polymorphic types are used whenever possible. IsaFoR also defines the standard monadic term substitution $\cdot :: \text{term} \Rightarrow ('v \Rightarrow \text{term}) \Rightarrow \text{term}$ and a unify $:: (\text{term} \times \text{term}) \text{ list} \Rightarrow \text{lsubst} \Rightarrow \text{lsubst}$ function, where $\text{lsubst} = ('v \times \text{term}) \text{ list}$ is the list-based representation of a finite substitution. The function unify computes the MGU for a list of unification constraints that is compatible with a given substitution. IsaFoR includes a wealth of theorems, including the correctness of unify and the well-foundedness of strict term generalization, defined as $(\exists \sigma. s \cdot \sigma = t) \wedge (\nexists \sigma. t \cdot \sigma = s)$.

This infrastructure allows us to conveniently instantiate our locales *substitution_ops*, *substitution*, and *mgu*. We instantiate the type $'a$ of atoms with *term* and the type $'s$ of substitutions with $'v \Rightarrow \text{term}$ and the constants \cdot , *id*, \circ , and *atm_of_atms* with \cdot , Var, $\lambda \sigma \tau x. \sigma x \cdot \tau$, and (arbitrarily) Fun 0. For the MGU computation, there is a slight type mismatch: IsaFoR offers a list-based unifier, whereas our locale requires the type *term set set* $\Rightarrow ('v \Rightarrow \text{term}) \text{ option}$. It is easy to translate a finite set of finite sets of terms into a finite list of pairs of constraints. To be executable, the translation requires us to sort the terms belonging to set with respect to an arbitrary (but executable) linear order.

Only the function *renamings_apart* was not present in IsaFoR. We supply a definition:

fun *renamings_apart* $:: \text{term clause list} \Rightarrow ('v \Rightarrow \text{term}) \text{ list}$
where

renamings_apart [] = []
 | *renamings_apart* $(C \# Cs) =$
 let $\sigma_s = \text{renamings_apart } Cs$ in
 $(\lambda v. v + \max (\{0\} \cup \text{vars_clause_list } (Cs \cdot \sigma_s)) + 1) \# \sigma_s$

where *vars_clause_list* $:: \text{term clause list} \Rightarrow 'v \text{ set}$ returns the variables contained in a list of clauses. The creation of fresh variable names relies on $'v = \text{nat}$.

Finally, the *FO_resolution_prover* locale requires that the type of atoms supports two comparison operators: a well-order $>$ and a comparison $>$ that is stable under substitution (i.e., $B > A \Rightarrow B \cdot \sigma > A \cdot \sigma$). Moreover, $>$ and $>$ must coincide on ground atoms. We instantiate $>$ with the Knuth-Bendix order [19] on terms, provided by IsaFoR [46]. This order is executable, stable under substitution, well founded, and total on ground terms. The well-order $>$, which must be total on *all* terms, is then defined as an arbitrary extension of a partial well-founded order $>$ to a well-order, using Hilbert choice. This makes $>$ nonexecutable, but this is acceptable since it is $>$, not $>$, that is used in the prover's code.

Clause Subsumption. The second hurdle concerns clause subsumption. Its mathematical definition, $\text{subsumes } C D \iff \exists \sigma. C \cdot \sigma \subseteq D$, involves an infinite quantification.

The problem of deciding whether such a substitution exists is NP-complete [18]. We start with the following naive code. In contrast to the mathematical definition, which operates on multisets of literals, our function operates on lists:

```
fun subsumes_list  $:: \text{term literal list} \Rightarrow$   

 $\text{term literal list} \Rightarrow ('v \Rightarrow \text{term option}) \Rightarrow \text{bool}$   

where  

    subsumes_list []  $Ks \sigma = \text{True}$   

    | subsumes_list  $(L \# Ls) Ks \sigma =$   

 $(\exists K \in \text{set } Ks. \text{is\_pos } K = \text{is\_pos } L \wedge$   

 $\text{case match\_term\_list } [(\text{atm\_of } L, \text{atm\_of } K)] \sigma \text{ of}$   

 $\text{None} \Rightarrow \text{False}$   

 $| \text{Some } \rho \Rightarrow \text{subsumes\_list } Ls (\text{remove1 } K Ks) \rho)$ 
```

In the Cons case, we must consider all possible matching literals for L from Ks compatible with the substitution σ . The bounded existential quantification that expresses this non-determinism can be executed by iterating over the finite list Ks . The functions *is_pos* and *atm_of* are the discriminator and selector for literals. The function *match_term_list* is provided by IsaFoR. It attempts to extend a given substitution into Some matcher for a list of matching constraints, given as term pairs. If the extension is impossible, *match_term_list* returns None. This substitution-passing style is typical of purely functional implementations of matching.

It is easy to prove that the above executable function implements clause subsumption: $\text{subsumes } (\text{mset } Ls) (\text{mset } Ks) = \text{subsumes_list } Ls Ks (\lambda x. \text{None})$, where *mset* converts lists to multisets by forgetting the order of the elements. After the registration of this equation, Isabelle's code generator will rewrite any code that contains the nonexecutable left-hand side to use the executable right-hand side instead.

Clause subsumption is a hot spot in a resolution prover [41]. Following Tammet [49], we implement a heuristic that often reduces the number of calls to *match_term_list*, which is linear in the size of the input terms, by first performing a simpler, imprecise comparison. For example, terms with different root symbols will never match, and these can be compared in

constant time. Similarly, literals with opposite polarities cannot match. We sort our (list-represented) clauses with respect to a literal quasi-order (i.e., a transitive and reflexive relation) leq_lit such that $\text{leq_lit } L K$ only if $\text{is_pos } L = \text{is_pos } K$ and $\text{match_term_list } [(\text{atm_of } L, \text{atm_of } K)] \sigma = \text{Some } \rho$ for some σ and ρ . Any quasi-order satisfying this property can be used in a refinement of subsumes_list to remove too small literals (with respect to leq_lit), as highlighted in gray below:

```

fun subsumes_list' :: term literal list  $\Rightarrow$ 
  term literal list  $\Rightarrow$  ('v  $\Rightarrow$  term option)  $\Rightarrow$  bool
where
  subsumes_list' [] Ks  $\sigma =$  True
| subsumes_list' (L # Ls) Ks  $\sigma =$ 
  let Ks = filter (leq_lit L) Ks in
  ( $\exists K \in \text{set } Ks. \text{is\_pos } K = \text{is\_pos } L \wedge$ 
   case match_term_list [(atm_of L, atm_of K)]  $\sigma$  of
     None  $\Rightarrow$  False
   | Some  $\rho \Rightarrow$  subsumes_list' Ls (remove1 K Ks)  $\rho$ )

```

The lemma $\text{subsumes_list } Ls Ks \rho = \text{subsumes_list}'$ (sort $\text{leq_lit } Ls$) $Ks \rho$ allows the code generator to refine the original version. In RP_x , we let leq_lit be a quasi-order that (1) considers negative literals smaller than positive ones; (2) considers variables smaller than nonvariables; and (3) sorts atoms according to a total order on their root symbols.

This refinement is a local optimization: It requires us to explicitly sort one of the input clauses. A potentially more efficient refinement would be to ensure that all clauses in the prover's state are sorted with respect to leq_lit . Sorting Ls at each invocation of subsumption could then be avoided, and filtering Ks could be performed more efficiently. However, maintaining the invariant would require changes throughout the prover's code.

The End Result. Finally, Isabelle can export our prover to Standard ML, Haskell, OCaml, or Scala. The command

```
export_code  $\text{RP}_x$  in SML module_name RP
```

generates an ML module containing the implementation of our prover in about 1000 lines of code, including dependencies. The generated module exports the function $\text{RP}_x : (\text{term literal list} * \text{nat}) \text{list} \rightarrow \text{bool}$. The input is the \mathcal{N} component of an initial state, which consists of pairs of clauses and arbitrary timestamps (e.g., 0).

Even though in Isabelle we have proved that for any unsatisfiable input RP_x will terminate and return False, the code generator guarantees only partial correctness of its output: If the ML program terminates on the ML input generated from the Isabelle term t and evaluates to the Boolean result b , the proposition $\text{RP}_x t = b$ is provable in Isabelle; by soundness, b indicates the satisfiability of the input clause set. There is recent work towards providing stronger guarantees and reducing the generator's trusted code base [14].

Empirical Evaluation. To measure the gap with the state of the art, we compare our prover's performance with that

of three other provers on a benchmark suite. TPTP (Thousands of Problems for Theorem Provers) [48] is the de facto standard library for benchmarking automatic provers. We extended RP_x with the trusted TPTP parser from Metis [15]. We benchmarked E 2.1, Vampire 4.2.2, Metis 2.4, and RP_x on 1000 randomly selected equality-free problems from the TPTP's FOF (first-order formulas) and CNF (first-order formulas in conjunctive normal form) categories. We converted all FOF problems to CNF using E's clausifier. Each prover was run on each problem for 60 s on an Intel Core i9-7900X (3.3 GHz 10-Core) with 128 GB of RAM.

The results are summarized in the following table, showing for each prover how many unsatisfiable and satisfiable problems were solved and how many seconds were needed on average by each prover on the problems that were solved by all four:

	Vampire	E	Metis	RP_x
Unsatisfiable	675	635	436	331
Satisfiable	158	135	91	22
Average time (s)	0.032	0.014	0.637	3.126

The detailed results of the evaluation are available online, together with instructions for reproducing them.³

As expected, RP_x is not competitive. A prover's performance comes from its calculus, its heuristics, and its indexing data structures. RP_x employs an excellent calculus but mediocre heuristics and data structures. Better performance could be achieved by working on these last two aspects. Heuristics are often easy to verify, because their input-output specifications are permissive, but formalizing optimized data structures can be very laborious [10].

Nevertheless, sometimes the calculus is all that matters. Benchmark MSC015 from the TPTP library is a particularly challenging family Φ_n of first-order problems, each consisting of the following $n + 2$ clauses:

$$\begin{aligned}
 & \neg p(b, \dots, b) \quad p(a, \dots, a) \\
 & \neg p(a, b, \dots, b) \vee p(b, a, \dots, a) \\
 & \neg p(x_1, a, b, \dots, b) \vee p(x_1, b, a, \dots, a) \\
 & \quad \vdots \\
 & \neg p(x_1, \dots, x_{n-2}, a, b) \vee p(x_1, \dots, x_{n-2}, b, a) \\
 & \neg p(x_1, \dots, x_{n-2}, x_{n-1}, a) \vee p(x_1, \dots, x_{n-2}, x_{n-1}, b)
 \end{aligned}$$

A comment in the benchmark warns us that back in 2007, no prover could solve the Φ_{23} within an hour. Even in 2018, only one prover solves Φ_{22} within 300 s, and four provers solve Φ_{20} within 300 s. RP_x solves Φ_{20} in 100 s and Φ_{22} in 200 s. Presumably, the reason for this success is that RP_x fortuitously chooses an instance of the Knuth–Bendix order that is well suited to this benchmark.

³http://matryoshka.gforge.inria.fr/pubs/fun_rp_data.tar.gz

7 Discussion and Related Work

We found Bachmair and Ganzinger’s [2] chapter and its formalization [39, 40] suitable as a starting point for a verified prover. Nonetheless, we faced some difficulties, notably concerning the identification of suitable refinement layers. We developed layers 2, 3, and 4 largely in parallel, with each of the authors working on a separate layer. Bringing layer 2 into a state such that it both ensures fairness and could be refined further by layer 3 required several iterations.

Stepwise refinement helped us achieve separation of concerns: fairness, determinism, and executability were achieved successively. Another strength of this methodology is that it allows us to prove results at a high level of abstraction; for example, fairness is established at layer 2 already and is inherited by subsequent layers. The main difficulty with refinement is that some nontrivial machinery is necessary to lift results from one layer to the next. We believe the gain in modularity makes this worthwhile.

It took us quite some time to design a suitable measure to prove the fairness of the layer 2 prover RP_w . Our solution amounts to advancing to a state carrying a suitably high timestamp and filtering out all overly heavy clauses. Initially, our proof consisted of two steps—advancing and filtering—each with its own measure. This proof gave us the assurance that RP_w was fair, but we found that combining the measures is both more succinct and more intelligible.

Our main objective was not to reach **qed** as quickly as possible but rather to investigate how to harness a modern proof assistant to formalize the metatheory of automatic theorem provers. We found Isabelle suitable for this verification task. The Isar proof language allows us to state key intermediate steps, as in a paper proof. Standard tactics, including Isabelle’s simplifier, can be used to discharge proof obligations. The Sledgehammer tool [32] employs superposition provers and SMT (satisfiability modulo theories) solvers to swiftly identify which lemmas can be used to prove a goal; standard Isabelle tactics are then used to reconstruct the proof. Isabelle’s support for coinductive methods, including the **coinductive**, **codatatype**, and **corec** commands, helps us reason about infinite processes. Locales are a useful abstraction for defining the refinement layers. And Isabelle’s libraries, the *Archive of Formal Proofs*, and IsaFoR certainly saved us months of labor.

The *Archive* also includes a refinement framework [25], which has been used in a separate effort to connect the imperative code of an efficient SAT solver to an abstract calculus [6]. The framework is helpful in a variety of situations, including when the refinement relation between a concrete and an abstract data representation is not a function. But since converting a list to a multiset (between our levels 3 and 2) or a multiset to a set (between levels 2 and 1) is a function, we did not see a need to employ it. Moreover, the framework is currently not designed for refining semidecision procedures,

as acknowledged privately by its developer. We conjecture that its support for separation logic could be useful if we were to refine the prover further to obtain imperative code.

Thanks to the verification, we can trust to a very high extent that our ordered resolution prover is sound and complete. To make the prover’s performance competitive with E, SPASS, and Vampire, we would need to extend the current work along two axes. First, we should use superposition, together with its extensive simplification machinery, as the base calculus. A good starting point would be to apply our methodology to Peltier’s [33] formalization of superposition. Given that a large part of a modern superposition prover’s code consists of heuristics, which are easy to verify, the full verification of a competitive superposition prover appears to be a realistic objective for a forthcoming Ph.D. thesis. Second, the refinement chain should be continued to cover optimized algorithms and data structures. These could be specified by refining layer 4 further, along the lines of Fleury et al.’s [10] refinement of an imperative SAT solver.

In computer science, a metatheory may inspire an implementation, or vice versa, but the connection is seldom made explicit. By formalizing the metatheory, the implementation, and their connection, we can demonstrate not only the implementation’s correctness but also the metatheory’s adequacy for describing potential implementations. In particular, we have now confirmed that Bachmair and Ganzinger [2] accurately describe the abstract principles of an executable functional prover (with a few exceptions [40]), even though they provide few details beyond layer 1.

We built the prover on our earlier formalization [39, 40] of ordered resolution. Related efforts developed using Isabelle/HOL include Peltier’s [33] formalization of superposition and Schlichtkrull’s [37] formalization of unordered resolution. These developments cover only logical calculi; had we started with any of them, the first step would have been to define an abstract prover in the style of layer 1 and prove basic properties about it. Another related effort is Hirokawa et al.’s [13] formalization of ordered completion, which (like ordered resolution) can be regarded as a special case of superposition.

Formalizing a theorem proving tool using a theorem proving tool is a thrilling (if self-referential) prospect for many researchers. An early result is Ridge and Margetson’s [35] verified first-order prover, based on a sequent calculus for first-order logic without full first-order terms but only variables. Kumar et al. [24] formalized the soundness of a proof assistant for higher-order logic. Jensen et al. [16] verified the soundness of a kernel for a proof assistant for first-order logic that includes a tableau prover. There are several verified SAT solvers [6, 27–29, 31, 44]. SAT being a decidable problem, termination has been proved for most solvers. First-order logic, on the other hand, is semidecidable.

A pragmatic approach to combining the efficiency of unverified code with the trustworthiness of verified code involves checking certificates produced by reasoning tools—

e.g., proofs produced by SAT solvers [9, 26]. Researchers from the first-order theorem proving community are now advocating this approach for their systems [34]. An ad hoc version of this approach is used in Sledgehammer and HOLy-Hammer to reconstruct proofs found by external automatic provers [4, 17].

8 Conclusion

Starting from an abstract description of an ordered resolution prover [39, 40], we verified, through a refinement chain, a purely functional prover that uses lists as its main data structure. The resulting program is interesting in its own right and could be refined further to obtain an implementation that is competitive with the state of the art.

Stepwise refinement is a keystone of our methodology, and we found it adequate. Each refinement step cleanly isolates concerns, yielding intelligible proof obligations. Refinement also helped us identify an unnecessary assumption in Bachmair and Ganzinger's [2] chapter and clarify the argument. Lifting results from one layer to another required some thought, especially the completeness results, which correspond to liveness properties.

Having now established a methodology and built basic formal libraries, we expect that verifying other provers, using Isabelle or other systems, will be substantially easier. Because it is based on Bachmair and Ganzinger's framework, our approach generally applies to all saturation-based provers, with or without redundancy. This includes resolution, paramodulation, ordered rewriting, superposition, and variants thereof, covering many of the most successful provers for equational [8, 12], first-order [21, 42, 53], and higher-order logic [45].

Acknowledgments

Johannes Hölzl gave us some useful advice on how to specify and reason about possibly nonterminating functions in Isabelle/HOL. Alexander Bentkamp, Mathias Fleury, Andreas Halkjær From, Carsten Fuhs, Peter Lammich, Mark Summerfield, Jørgen Villadsen, and the anonymous reviewers suggested many textual improvements. Eugene Kotelnikov, Stephan Schulz, and Geoff Sutcliffe graciously answered our questions about Vampire, E, and TPTP.

Schlichtkrull has received funding from a Ph.D. scholarship in the Algorithms, Logic and Graphs section of DTU Compute and from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement No. 700321, LIGHTest). Blanchette has received funding from the ERC under the European Union's Horizon 2020 research and innovation program (grant agreement No. 713999, Matryoshka). An earlier version of this paper was included as a chapter of Schlichtkrull's Ph.D. thesis [36] with the same authors' list.

References

- [1] Leo Bachmair, Nachum Dershowitz, and David A. Plaisted. 1989. Completion without Failure. In *Rewriting Techniques—resolution of Equations in Algebraic Structures*, Hassan Aït-Kaci and Maurice Nivat (Eds.). Vol. 2. Academic Press, 1–30.
- [2] Leo Bachmair and Harald Ganzinger. 2001. Resolution Theorem Proving. In *Handbook of Automated Reasoning*, Alan Robinson and Andrei Voronkov (Eds.). Vol. I. Elsevier and MIT Press, 19–99.
- [3] Jasmin Christian Blanchette. 2019. Formalizing the Metatheory of Logical Calculi and Automatic Provers in Isabelle/HOL (Invited Paper). In *CPP 2019*, Assia Mahboubi and Magnus O. Myreen (Eds.). ACM.
- [4] Jasmin Christian Blanchette, Sascha Böhme, Mathias Fleury, Steffen Juilf Smolka, and Albert Steckermeier. 2016. Semi-intelligible Isar Proofs from Machine-Generated Proofs. *J. Autom. Reasoning* 56, 2 (2016), 155–200.
- [5] Jasmin Christian Blanchette, Aymeric Bouzy, Andreas Lochbihler, Andrei Popescu, and Dmitriy Traytel. 2017. Friends with Benefits: Implementing Corecursion in Foundational Proof Assistants. In *ESOP 2017*, Hongseok Yang (Ed.). LNCS, Vol. 10201. Springer, 111–140.
- [6] Jasmin Christian Blanchette, Mathias Fleury, Peter Lammich, and Christoph Weidenbach. 2018. A Verified SAT Solver Framework with Learn, Forget, Restart, and Incrementality. *J. Autom. Reasoning* 61, 1–4 (2018), 333–365.
- [7] François Bobot, Jean-Christophe Filliâtre, Claude Marché, and Andrei Paskevich. 2011. Why3: Shepherd Your Herd of Provers. In *Boogie 2011*, K. Rustan M. Leino and Michał Moskal (Eds.). 53–64.
- [8] Koen Claessen and Nicholas Smallbone. 2018. Efficient Encodings of First-Order Horn Formulas in Equational Logic. In *IJCAR 2018*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). LNCS, Vol. 10900. Springer, 388–404.
- [9] Luís Cruz-Filipe, Marijn J. H. Heule, Warren A. Hunt Jr., Matt Kaufmann, and Peter Schneider-Kamp. 2017. Efficient Certified RAT Verification. In *CADE-26*, Leonardo de Moura (Ed.). LNCS, Vol. 10395. Springer, 220–236.
- [10] Mathias Fleury, Jasmin Christian Blanchette, and Peter Lammich. 2018. A Verified SAT Solver with Watched Literals using Imperative HOL. In *CPP 2018*, June Andronick and Amy P. Felty (Eds.). ACM, 158–171.
- [11] Florian Haftmann and Tobias Nipkow. 2010. Code Generation via Higher-Order Rewrite Systems. In *FLOPS 2010*, Matthias Blume, Naoki Kobayashi, and Germán Vidal (Eds.). LNCS, Vol. 6009. Springer, 103–117.
- [12] Thomas Hillenbrand, Arnim Buch, Roland Vogt, and Bernd Löchner. 1997. WALDMEISTER—High-Performance Equational Deduction. *J. Autom. Reasoning* 18, 2 (1997), 265–270.
- [13] Nao Hirokawa, Aart Middeldorp, Christian Sternagel, and Sarah Winkler. 2017. Infinite Runs in Abstract Completion. In *FSCD 2017*, Dale Miller (Ed.). LIPIcs, Vol. 84. Schloss Dagstuhl—Leibniz-Zentrum für Informatik, 19:1–19:16.
- [14] Lars Hupel and Tobias Nipkow. 2018. A Verified Compiler from Isabelle/HOL to CakeML. In *ESOP 2018*, Amal Ahmed (Ed.). LNCS, Vol. 10801. Springer, 999–1026.
- [15] Joe Hurd. 2003. First-Order Proof Tactics in Higher-Order Logic Theorem Provers. In *Design and Application of Strategies/Tactics in Higher Order Logics (STRATA)*, Myla Archer, Ben Di Vito, and César Muñoz (Eds.). 56–68.
- [16] Alexander Birch Jensen, John Bruntse Larsen, Anders Schlichtkrull, and Jørgen Villadsen. 2018. Programming and Verifying a Declarative First-Order Prover in Isabelle/HOL. *AI Commun.* 31, 3 (2018), 281–299.
- [17] Cezary Kaliszyk and Josef Urban. 2013. PROCH: Proof Reconstruction for HOL Light. In *CADE-24*, Maria Paola Bonacina (Ed.). LNCS, Vol. 7898. Springer, 267–273.
- [18] Deepak Kapur and Paliath Narendran. 1986. NP-Completeness of the Set Unification and Matching Problems. In *CADE-8*, Jörg H. Siekmann

- (Ed.). LNCS, Vol. 230. Springer, 489–495.
- [19] Donald E. Knuth and Peter B. Bendix. 1970. Simple Word Problems in Universal Algebras. In *Computational Problems in Abstract Algebra*, John Leech (Ed.). Pergamon Press, 263–297.
- [20] Laura Kovács and Andrei Voronkov. 2009. Finding Loop Invariants for Programs over Arrays using a Theorem Prover. In *SYNASC 2009*, Stephen M. Watt, Viorel Negru, Tetsuo Ida, Tudor Jebelean, Dana Petcu, and Daniela Zaharie (Eds.). IEEE Computer Society, 10.
- [21] Laura Kovács and Andrei Voronkov. 2013. First-Order Theorem Proving and Vampire. In *CAV 2013*, Natasha Sharygina and Helmut Veith (Eds.). LNCS, Vol. 8044. Springer, 1–35.
- [22] Alexander Krauss. 2006. Partial Recursive Functions in Higher-Order Logic. In *IJCAR 2006*, Ulrich Furbach and Natarajan Shankar (Eds.). LNCS, Vol. 4130. Springer, 589–603.
- [23] Alexander Krauss. 2010. Recursive Definitions of Monadic Functions. *EPTCS* 43 (2010), 1–13.
- [24] Ramana Kumar, Rob Arthan, Magnus O. Myreen, and Scott Owens. 2016. Self-Formalisation of Higher-Order Logic: Semantics, Soundness, and a Verified Implementation. *J. Autom. Reasoning* 56, 3 (2016), 221–259.
- [25] Peter Lammich. 2013. Automatic Data Refinement. In *ITP 2013*, Sandrine Blazy, Christine Paulin-Mohring, and David Pichardie (Eds.). LNCS, Vol. 7998. Springer, 84–99.
- [26] Peter Lammich. 2017. The GRAT Tool Chain—Efficient (UN)SAT Certificate Checking with Formal Correctness Guarantees. In *SAT 2017*, Serge Gaspers and Toby Walsh (Eds.). LNCS, Vol. 10491. Springer, 457–463.
- [27] Stephane Lescuyer. 2011. *Formalizing and Implementing a Reflexive Tactic for Automated Deduction in Coq*. Ph.D. Dissertation. Université Paris-Sud.
- [28] Filip Marić. 2008. Formal Verification of Modern SAT Solvers. *Archive of Formal Proofs* (2008). Formal Proof Development. <http://isa-afp.org/entries/SATSolverVerification.html>.
- [29] Filip Marić. 2010. Formal Verification of a Modern SAT Solver by Shallow Embedding into Isabelle/HOL. *Theoret. Comput. Sci.* 411, 50 (2010), 4333–4356.
- [30] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. 2002. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. LNCS, Vol. 2283. Springer.
- [31] Duckki Oe, Aaron Stump, Corey Oliver, and Kevin Clancy. 2012. versat: A Verified Modern SAT Solver. In *VMCAI 2012*, Viktor Kuncak and Andrey Rybalchenko (Eds.). LNCS, Vol. 7148. Springer, 363–378.
- [32] Lawrence C. Paulson and Jasmin Christian Blanchette. 2012. Three Years of Experience with Sledgehammer, a Practical Link Between Automatic and Interactive Theorem Provers. In *IWIL-2010*, Geoff Sutcliffe, Stephan Schulz, and Eugenia Ternovska (Eds.). EPiC Series in Computing, Vol. 2. EasyChair, 1–11.
- [33] Nicolas Peltier. 2016. A Variant of the Superposition Calculus. *Archive of Formal Proofs* (2016). Formal Proof Development. <http://isa-afp.org/entries/SuperCalc.html>.
- [34] Giles Reger and Martin Suda. 2017. Checkable Proofs for First-Order Theorem Proving. In *ARCADE 2017*, Giles Reger and Dmitriy Traytel (Eds.). EPiC Series in Computing, Vol. 51. EasyChair, 55–63.
- [35] Tom Ridge and James Margetson. 2005. A Mechanically Verified, Sound and Complete Theorem Prover for First Order Logic. In *TPHOLS 2005*, Joe Hurd and Tom Melham (Eds.). LNCS, Vol. 3603. Springer, 294–309.
- [36] Anders Schlichtkrull. 2018. *Formalization of Logic in the Isabelle Proof Assistant*. Ph.D. Dissertation. Technical University of Denmark.
- [37] Anders Schlichtkrull. 2018. Formalization of the Resolution Calculus for First-Order Logic. *J. Autom. Reasoning* 61, 1–4 (2018), 455–484.
- [38] Anders Schlichtkrull, Jasmin Christian Blanchette, and Dmitriy Traytel. 2018. A Verified Functional Implementation of Bachmair and Ganzinger’s Ordered Resolution Prover. *Archive of Formal Proofs* (2018). Formal Proof Development. http://isa-afp.org/entries/Functional_Ordered_Resolution_Prover.html.
- [39] Anders Schlichtkrull, Jasmin Christian Blanchette, Dmitriy Traytel, and Uwe Waldmann. 2018. Formalization of Bachmair and Ganzinger’s Ordered Resolution Prover. *Archive of Formal Proofs* (2018). Formal Proof Development. http://isa-afp.org/entries/Ordered_Resolution_Prover.html.
- [40] Anders Schlichtkrull, Jasmin Christian Blanchette, Dmitriy Traytel, and Uwe Waldmann. 2018. Formalizing Bachmair and Ganzinger’s Ordered Resolution Prover. In *IJCAR 2018*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). LNCS, Vol. 10900. Springer, 89–107.
- [41] Stephan Schulz. 2013. Simple and Efficient Clause Subsumption with Feature Vector Indexing. In *Automated Reasoning and Mathematics—Essays in Memory of William W. McCune*, Maria Paola Bonacina and Mark E. Stickel (Eds.). LNCS, Vol. 7788. Springer, 45–67.
- [42] Stephan Schulz. 2013. System Description: E 1.8. In *LPAR-19*, Ken McMillan, Aart Middeldorp, and Andrei Voronkov (Eds.). LNCS, Vol. 8312. Springer, 735–743.
- [43] Stephan Schulz and Martin Möhrmann. 2016. Performance of Clause Selection Heuristics for Saturation-Based Theorem Proving. In *IJCAR 2016*, Nicola Olivetti and Ashish Tiwari (Eds.). LNCS, Vol. 9706. Springer, 330–345.
- [44] Natarajan Shankar and Marc Vaucher. 2011. The Mechanical Verification of a DPLL-Based Satisfiability Solver. *Electr. Notes Theor. Comput. Sci.* 269 (2011), 3–17. LSF 2010.
- [45] Alexander Steen and Christoph Benzmüller. 2018. The Higher-Order Prover Leo-III. In *IJCAR 2018*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). LNCS, Vol. 10900. Springer, 108–116.
- [46] Christian Sternagel and René Thiemann. 2013. Formalizing Knuth-Bendix Orders and Knuth-Bendix Completion. In *RTA 2013*, Femke van Raamsdonk (Ed.). LIPICs, Vol. 21. Schloss Dagstuhl—Leibniz-Zentrum für Informatik, 287–302.
- [47] Christian Sternagel and René Thiemann. 2018. First-Order Terms. *Archive of Formal Proofs* (2018). Formal Proof Development. http://isa-afp.org/entries/First_Order_Terms.html.
- [48] Geoff Sutcliffe. 2017. The TPTP Problem Library and Associated Infrastructure: From CNF to TH0, TPTP v6.4.0. *J. Autom. Reasoning* 59, 4 (2017), 483–502.
- [49] Tanel Tammet. 1998. Towards Efficient Subsumption. In *CADE-15*, Claude Kirchner and Hélène Kirchner (Eds.). LNCS, Vol. 1421. Springer, 427–441.
- [50] René Thiemann. 2018. Extending a Verified Simplex Algorithm. In *IWIL-2018*, Geoff Sutcliffe, Stephan Schulz, and Eugenia Ternovska (Eds.).
- [51] René Thiemann and Christian Sternagel. 2009. Certification of Termination Proofs using CeTA. In *TPHOLS 2009*, Stefan Berghofer, Tobias Nipkow, Christian Urban, and Makarius Wenzel (Eds.). LNCS, Vol. 5674. Springer, 452–468.
- [52] Andrei Voronkov. 2014. AVATAR: The Architecture for First-Order Theorem Provers. In *CAV 2014*, Armin Biere and Roderick Bloem (Eds.). LNCS, Vol. 8559. Springer, 696–710.
- [53] Christoph Weidenbach, Dilyana Dimova, Arnaud Fietzke, Rohit Kumar, Martin Suda, and Patrick Wischniewski. 2009. SPASS Version 3.5. In *CADE-22*, Renate A. Schmidt (Ed.). LNCS, Vol. 5663. Springer, 140–145.
- [54] Makarius Wenzel. 2012. Isabelle/jEdit—a Prover IDE within the PIDE Framework. In *CICM 2012*, Johan Jeuring, John A. Campbell, Jacques Carette, Gabriel Dos Reis, Petr Sojka, Makarius Wenzel, and Volker Sorge (Eds.). LNCS, Vol. 7362. Springer, 468–471.
- [55] Niklaus Wirth. 1971. Program Development by Stepwise Refinement. *Commun. ACM* 14, 4 (1971).