

Formalizing Bachmair and Ganzinger’s Ordered Resolution Prover

Anders Schlichtkrull¹(✉), Jasmin Christian Blanchette^{2,3},
Dmitriy Traytel⁴, and Uwe Waldmann³

¹ DTU Compute, Technical University of Denmark, Kongens Lyngby, Denmark
andschl@dtu.dk

² Vrije Universiteit Amsterdam, The Netherlands

³ Max-Planck-Institut für Informatik, Saarland Informatics Campus, Saarbrücken, Germany

⁴ Institute of Information Security, Department of Computer Science, ETH Zürich, Switzerland

Abstract. We present a formalization of the first half of Bachmair and Ganzinger’s chapter on resolution theorem proving in Isabelle/HOL, culminating with a refutationally complete first-order prover based on ordered resolution with literal selection. We develop general infrastructure and methodology that can form the basis of completeness proofs for related calculi, including superposition. Our work clarifies several of the fine points in the chapter’s text, emphasizing the value of formal proofs in the field of automated reasoning.

1 Introduction

Much research in automated reasoning amounts to metatheoretical arguments, typically about the soundness and completeness of logical inference systems or the termination of theorem proving processes. Often the proofs contain more insights than the systems or processes themselves. For example, the superposition calculus rules [2], with their many side conditions, look rather arbitrary, whereas in the completeness proof the side conditions emerge naturally from the model construction. And yet, despite being crucial to our field, today such proofs are usually carried out without tool support beyond \TeX .

We believe proof assistants are becoming mature enough to help. In this report, we present a formalization, developed using the Isabelle/HOL system [19], of a first-order prover based on ordered resolution with literal selection. We follow Bachmair and Ganzinger’s account [3] from Chapter 2 of the *Handbook of Automated Reasoning*, which we will simply refer to as “the chapter.” Our formal development covers the refutational completeness of two resolution calculi for ground (i.e., variable-free) clauses and general infrastructure for theorem proving processes and redundancy, culminating with a completeness proof for a first-order prover expressed as transition rules operating on triples of clause sets. This material corresponds to the chapter’s first four sections.

From the perspective of automated reasoning, increased trustworthiness of the results is an obvious benefit of formal proofs. But formalizing also helps clarify arguments, by exposing and explaining difficult steps. Making theorem statements (including definitions and hypotheses) precise can be a huge gain for communicating results. Moreover, a formal proof can tell us exactly where hypotheses and lemmas are used. Once we have created a library of basic results and a methodology, we will be in a good position to study extensions and variants. Given that automatic theorem provers are in-

tegrated in modern proof assistants, there is also an undeniable thrill in applying these tools to reason about their own metatheory.

From the perspective of interactive theorem proving, formalization work constitutes a case study in the use of a proof assistant. It gives us, as developers and users of such a system, an opportunity to experiment, contribute to lemma libraries, and get inspiration for new features and improvements.

Our motivation for choosing Bachmair and Ganzinger’s chapter is manifold. The text is a standard introduction to superposition-like calculi (together with *Handbook* Chapters 7 [17] and 27 [33]). It offers perhaps the most detailed treatment of the lifting of a resolution-style calculus’s static completeness to a saturation prover’s dynamic completeness. It introduces a considerable amount of general infrastructure, including different types of inference systems (sound, reductive, counterexample-reducing, etc.), theorem proving processes, and an abstract notion of redundancy. The resolution calculus, extended with a term order and literal selection, captures most of the insights underlying ordered paramodulation and superposition, but with a simple notion of model.

The chapter’s level of rigor is uneven, as shown by the errors and imprecisions revealed by our formalization. In particular, we will see that the main completeness result does not hold, due to the improper treatment of self-inferences. Naturally, our objective is not to diminish Bachmair and Ganzinger’s outstanding achievements, which include the development of superposition; rather, it is to demonstrate that even the work of some of the most celebrated researchers in our field can benefit from formalization. Our view is that formal proofs can beneficially be used to complement and improve their informal counterparts.

This work is part of the IsaFoL (Isabelle Formalization of Logic) project,¹ which aims at developing a library of results about logical calculi used in automated reasoning. The Isabelle theory files are available in the *Archive of Formal Proofs* (AFP).² They amount to about 8000 lines of source text. Below we provide implicit hyperlinks from theory and lemma names. A better way to study the theory files, however, is to open them in Isabelle/jEdit [35], an integrated development environment for formal proof. This will ensure that logical and mathematical symbols are rendered properly (e.g., \forall instead of $\backslash\langle forall \rangle$) and let you inspect proof states. We used Isabelle version 2017, but the AFP is continuously updated to track Isabelle’s evolution. We assume the reader has some familiarity with the chapter’s content.

2 Preliminaries

Ordered resolution depends on little background metatheory. Much of it, concerning partial and total orders, well-foundedness, and finite multisets, is provided by standard Isabelle libraries. We also need literals, clauses, models, terms, and substitutions.

Clauses and Models. We reuse the same library of clauses (`Clausal_Logic.thy`) as for the verified SAT solver by Blanchette et al. [7], which is also part of IsaFoL. Atoms are represented by a type variable $'a$, which can be instantiated by arbitrary concrete

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

² https://devel.isa-afp.org/entries/Ordered_Resolution_Prover.html

types—e.g., numbers or first-order terms. A literal, of type *'a literal* (where the type constructor is written in ML-style postfix syntax), can be of the form $\text{Pos } A$ or $\text{Neg } A$, where $A :: 'a$ is an atom. The literal order $>$ (written \succ in the chapter) extends a fixed atom order $>$ by comparing polarities to break ties, with $\text{Neg } A > \text{Pos } A$. Following the chapter, a clause is a finite multiset of literals: *'a clause* = *'a literal multiset*, where *multiset* is the Isabelle type constructor of finite multisets. Thus, the clause $A \vee B$, where A and B are atoms, is identified with the multiset $\{A, B\}$; the clause $C \vee D$, where C and D are clauses, is $C \uplus D$; and the empty clause \perp is $\{\}$.

A Herbrand interpretation I is a value of type *'a set*, specifying which ground atoms are true (`Herbrand_Interpretation.thy`). The “models” operator \models is defined in the usual way on atoms, literals, clauses, sets, and multisets of clauses; for example, $I \models C \Leftrightarrow \exists L \in C. I \models L$. Satisfiability of a set or multiset of clauses N is defined by $\text{sat } N \Leftrightarrow \exists I. I \models N$.

Multisets are central to our development. Isabelle provides a multiset library, but it is much less developed than those of sets and lists. As part of IsaFoL, we have already extended it considerably and implemented further additions in a separate file (`Multiset_More.thy`). Some of these, notably a plugin for Isabelle’s simplifier to apply cancellation laws, are described in a recent paper [6, Section 3].

The main hurdle we faced concerned the multiset order. Multisets of clauses have type *'a literal multiset multiset*. The clause order is the multiset extension of the multiset extension of the literal order. In Isabelle, the multiset order was called $\#C\#$, and it relied on the element type’s $<$ operator, through Isabelle’s type class mechanism. Unfortunately, for multisets, $<$ was defined as the subset relation, so when nesting multisets (as *'a multiset multiset*), we obtained the multiset extension of the subset relation. Initially, we worked around the issue by defining an order $\#C\#\#$ on multisets of multisets, but we also saw potential for improvement. After some discussions on the Isabelle users’ mailing list, we decided to let $<$ be the multiset order and introduce the symbol $\subset\#$ for the subset relation. To avoid introducing subtle changes in the semantics of existing developments, we first renamed $<$ to $\subset\#$, freeing up $<$; then, in the next Isabelle release, we renamed $\#C\#$ to $<$. In the intermediate state, all occurrences of $<$ and of the lemmas about it were flagged as errors, easing porting. Similar changes affected the nonstrict versions of the orders (e.g., \leq) and all the lemmas about them (e.g., *add_mono*: $M \leq M' \Rightarrow N \leq N' \Rightarrow M \uplus N \leq M' \uplus N'$).

Terms and Substitutions. The `IsaFoR` (Isabelle Formalization of Rewriting) library—an inspiration for IsaFoL—contains a definition of first-order terms and results about substitutions and unification [30]. It makes sense to reuse this functionality. A practical issue is that most of `IsaFoR` is not accessible from the AFP.

Resolution depends only on basic properties of terms and atoms, such as the existence of most general unifiers (MGUs). We exploit this to keep the development parameterized by a type of atoms *'a* and an abstract type of substitutions *'s*, through Isabelle locales [4] (`Abstract_Substitution.thy`). A locale represents a module parameterized by types and terms that satisfy some assumptions. Inside the locale, we can refer to the parameters and assumptions in definitions, lemmas, and proofs. The basic operations provided by our locale are application ($\cdot :: 'a \Rightarrow 's \Rightarrow 'a$), identity ($\text{id} :: 's$), and composition ($\circ :: 's \Rightarrow 's \Rightarrow 's$), about which some assumptions are made (e.g., $A \cdot \text{id} = A$

for all atoms A). Substitution is lifted to literals, clauses, sets of clauses, and so on. Many other operations can be defined in terms of the primitives—for example:

$$\begin{aligned} \text{is_ground } A &\Leftrightarrow \forall \sigma. A = A \cdot \sigma & \text{is_renaming } \sigma &\Leftrightarrow \exists \tau. \sigma \circ \tau = \text{id} \\ \text{is_ground } \sigma &\Leftrightarrow \forall A. \text{is_ground } (A \cdot \sigma) & \text{instance_of } C D &\Leftrightarrow \exists \sigma. C \cdot \sigma = D \end{aligned}$$

MGUs are also taken as a primitive: the $\text{mgu} :: 'a \text{ set set} \Rightarrow 's \text{ option}$ operation takes a set of unification constraints, each of the form $A_1 \stackrel{?}{=} \dots \stackrel{?}{=} A_n$, and returns either an MGU or a special value (None).

To complete our formal development and ensure that our assumptions are legitimate, we instantiate the locale's parameters with IsaFoR types and operations and discharge the locale's assumptions (IsaFoR_Term.thy). This bridge is currently hosted on the IsaFoL repository, outside the AFP.

3 Refutational Inference Systems

In their Section 2.4, Bachmair and Ganzinger introduce basic conventions for refutational inference systems. In Section 3, they present two ground resolution calculi and prove them refutationally complete in Theorems 3.9 and 3.16. In Section 4.2, they introduce a notion of counterexample-reducing inference system and state Theorem 4.4 as a generalization of Theorems 3.9 and 3.16 to all such systems. For formalization, two courses of actions suggest themselves: follow the book closely and prove the three theorems separately, or focus on the most general result. We choose the latter, as being more consistent with our goal of providing a well-designed, reusable library, at the cost of widening the gap between the text and its formal companion.

We collect the abstract hierarchy of inference systems in a single Isabelle theory file (Inference_System.thy). An inference, of type $'a \text{ inference}$, is a triple (C, D, E) that consists of a multiset of side premises C , a main premise D , and a conclusion E . An inference system, or calculus, is a possibly infinite set of inferences:

locale *inference_system* =
fixes $\Gamma :: 'a \text{ inference set}$

We use an Isabelle locale to fix, within a named context (called *inference_system*), a set Γ of inferences between clauses over atom type $'a$. Inside the locale, we define a function *infers_from* that, given a clause set N , returns the subset of Γ inferences whose premises all belong to N .

A satisfiability-preserving (or consistency-preserving) inference system enriches an inference system with an assumption. A sound system adds a different assumption.

locale *sat_preserving_inference_system* = *inference_system* +
assumes $\text{sat } N \Rightarrow \text{sat } (N \cup \text{concl_of } ' \text{infers_from } N)$
locale *sound_inference_system* = *inference_system* +
assumes $(C, D, E) \in \Gamma \Rightarrow I \models C \cup \{D\} \Rightarrow I \models E$

The notation $f ` X$ above stands for the image of the set or multiset X under function f .

Soundness is a stronger requirement than satisfiability preservation. In Isabelle, this can be expressed as a sublocale relation:

sublocale *sound_inference_system* < *sat_preserving_inference_system*

This command emits a proof goal stating that *sound_inference_system*'s assumption implies *sat_preserving_inference_system*'s. Afterwards, all the definitions and lemmas about satisfiability-preserving calculi become available about sound ones.

In reductive inference systems (*reductive_inference_system*), the conclusion of each inference is smaller than the main premise according to the clause order. A related notion, the counterexample-reducing inference systems, is specified as follows:

locale *counterex_reducing_inference_system* = *inference_system* +
fixes *I_of* :: 'a clause set \Rightarrow 'a set
assumes $\{\} \notin N \Rightarrow D \in N \Rightarrow I_of\ N \not\models D \Rightarrow$
 $(\forall C \in N. I_of\ N \not\models C \Rightarrow D \leq C) \Rightarrow$
 $\exists C \subseteq N. \exists E. I_of\ N \models C \wedge (C, D, E) \in \Gamma \wedge I_of\ N \not\models E \wedge E < D$

The “model functor” parameter *I_of* maps clause sets to candidate models. The assumption is that for any clause set *N* that does not contain $\{\}$ (i.e., \perp), if $D \in N$ is the smallest counterexample—the smallest clause in *N* that is falsified by *I_of* *N*—we can derive a smaller counterexample *E* using an inference from clauses in *N*. This property is useful because if *N* is saturated (i.e., closed under Γ inferences), we must have $E \in N$, contradicting *D*'s minimality:

theorem *saturated_model*: $\text{saturated } N \Rightarrow \{\} \notin N \Rightarrow I_of\ N \models N$

corollary *saturated_complete*: $\text{saturated } N \Rightarrow \neg \text{sat } N \Rightarrow \{\} \in N$

Bachmair and Ganzinger claim that compactness of clausal logic follows from the refutational completeness of ground resolution (Theorem 3.12), although they give no justification. Our argument relies on an inductive definition of saturation of a set of clauses: *saturate* :: 'a clause set \Rightarrow 'a clause set. Most of the work goes into proving this key lemma, by rule induction on the *saturate* function:

lemma *saturate_finite*: $C \in \text{saturate } N \Rightarrow \exists M \subseteq N. \text{finite } M \wedge C \in \text{saturate } M$

The interesting case is when $C = \perp$. We establish compactness in a locale that combines *counterex_reducing_inference_system* and *sound_inference_system*:

theorem *clausal_logic_compact*: $\neg \text{sat } N \Leftrightarrow \exists M \subseteq N. \text{finite } M \wedge \neg \text{sat } M$

To give a taste of the formalization, here is the formal proof, expressed using Isabelle's structured Isar format [34]:

```
proof
  assume  $\neg \text{sat } N$ 
  then have  $\{\} \in \text{saturate } N$ 
    using saturated_complete saturated_saturate saturate.base
    unfolding true_cls_def by meson
  then have  $\exists M \subseteq N. \text{finite } M \wedge \{\} \in \text{saturate } M$ 
    using saturate_imp_finite_subset by fastforce
  then show  $\exists M \subseteq N. \text{finite } M \wedge \neg \text{sat } M$ 
```

```

using saturate_sound by auto
next
assume  $\exists M \subseteq N. \text{finite } M \wedge \neg \text{sat } M$ 
then show  $\neg \text{sat } N$ 
by (blast intro: true_cls_mono)
qed

```

In the “implies” direction, we rely on the calculus’s refutation completeness to show that \perp belongs to saturate N , on the above key lemma to obtain a finite subset M from which \perp can be derived, and on the calculus’s soundness to conclude that M is unsatisfiable. We believe satisfiability preservation could be used instead of soundness, relying on the property that $\text{sat } N \implies \text{sat}(\text{saturate } N)$ for satisfiability-preserving calculi, but we have yet to find a good way to prove this formally.

Our compactness result is meaningful only if the locale assumptions are consistent. In the next section, we will exhibit two sound counterexample-reducing calculi that can be used to instantiate the locale and retrieve an unconditional compactness theorem.

4 Ground Resolution

A useful strategy for establishing properties of first-order calculi is to initially restrict our attention to ground calculi and then to lift the results to first-order formulas containing terms with variables. Accordingly, the chapter’s Section 3 presents two ground calculi: a simple binary resolution calculus and an ordered resolution calculus with literal selection. Both consist of a single resolution rule, with built-in positive factorization. Most of the explanations and proofs concern the simpler calculus. To avoid duplication, we factor out the candidate model construction (Ground_Resolution_Model.thy). We then define the two calculi and prove that they are sound and reduce counterexamples (Unordered_Ground_Resolution.thy, Ordered_Ground_Resolution.thy).

Candidate Models. Refutational completeness is proved by exhibiting a model for any inconsistent saturated clause set N that does not contain \perp . The model is constructed incrementally, one clause $C \in N$ at a time, starting with an empty Herbrand interpretation. The idea appears to have originated with Brand [10] and Zhang and Kapur [36].

Bachmair and Ganzinger introduce two operators to build the candidate model: I_C denotes the current interpretation before considering C , and ε_C denotes the set of (zero or one) atoms added, or *produced*, to ensure that C is satisfied. The candidate model construction is parameterized by a literal selection function S . It can be ignored by taking $S := \lambda C. \{\}$.

```

locale ground_resolution_with_selection =
fixes  $S :: 'a \text{ clause} \Rightarrow 'a \text{ clause}$ 
assumes  $S C \subseteq C$  and  $L \in S C \implies \text{is\_neg } L$ 

```

Inside the locale, we fix a clause set N , for which we try to derive a model. Then we define two operators corresponding to ε_C and I_C :

function $\text{production} :: 'a \text{ clause} \Rightarrow 'a \text{ set}$ **where**
 $\text{production } C = \{A \mid C \in N \wedge C \neq \{\} \wedge \text{Max } C = \text{Pos } A$
 $\wedge (\bigcup_{D < C} \text{production } D) \not\vdash C \wedge S C = \{\}\}$

definition $\text{interp} :: 'a \text{ clause} \Rightarrow 'a \text{ set}$ **where**
 $\text{interp } C = \bigcup_{D < C} \text{production } D$

To ensure monotonicity of the construction, any produced atom must be maximal in its clause. Moreover, productive clauses may not contain selected literals. In the chapter, ε_C and I_C are expressed in terms of each other. We simplified the definition by inlining I_C in ε_C , so that only ε_C is recursive. Since the recursive calls operate on clauses D that are smaller with respect to a well-founded order, the definition is accepted [15]. Once the operators are defined, we can fold interp 's definition in production 's equation to derive the intended mutually recursive specification as a lemma. Bachmair and Ganzinger's I^C and I_N operators are introduced as abbreviations:

$$\text{Interp } C = \text{interp } C \cup \text{production } C \quad \text{INTERP} = \bigcup_{C \in N} \text{production } C$$

We then prove a host of lemmas about these concepts. Lemma 3.4 amounts to six monotonicity properties:

lemma $\text{Interp_imp_interp}: C \leq D \Rightarrow D < D' \Rightarrow \text{Interp } D \vDash C \Rightarrow \text{interp } D' \vDash C$
lemma $\text{Interp_imp_Interp}: C \leq D \Rightarrow D \leq D' \Rightarrow \text{Interp } D \vDash C \Rightarrow \text{Interp } D' \vDash C$
lemma $\text{Interp_imp_INTERP}: C \leq D \Rightarrow \text{Interp } D \vDash C \Rightarrow \text{INTERP} \vDash C$
lemma $\text{interp_imp_interp}: C \leq D \Rightarrow D \leq D' \Rightarrow \text{interp } D \vDash C \Rightarrow \text{interp } D' \vDash C$
lemma $\text{interp_imp_Interp}: C \leq D \Rightarrow D \leq D' \Rightarrow \text{interp } D \vDash C \Rightarrow \text{Interp } D' \vDash C$
lemma $\text{interp_imp_INTERP}: C \leq D \Rightarrow \text{interp } D \vDash C \Rightarrow \text{INTERP} \vDash C$

In the chapter, the first property is wrongly stated with $D \leq D'$ instead of $D < D'$, admitting the counterexample $N = \{\{A\}\}$ and $C = D = D' = \{A\}$. Lemma 3.3, whose proof depends on monotonicity, is better proved *after* Lemma 3.4:

lemma $\text{productive_imp_INTERP}: \text{production } C \neq \{\} \Rightarrow \text{INTERP} \vDash C$

A more serious oddity is Lemma 3.7. Using our notations, it can be stated as

$$D \in N \Rightarrow C \neq D \Rightarrow (\forall D' \leq D. \text{Interp } D' \vDash C) \Rightarrow \text{interp } D \vDash D'$$

However, the last occurrence of D' is clearly wrong—the context suggests C instead. Even after this amendment, we have a counterexample, corresponding to a gap in the proof: $D = \{\}$, $C = \{\text{Pos } A\}$, and $N = \{D, C\}$. Since this “lemma” is not actually used, we can simply ignore it.

Unordered Resolution. The unordered ground resolution calculus consists of a single binary inference rule, with the side premise $C \vee A \vee \dots \vee A$, the main premise $\neg A \vee D$, and the conclusion $C \vee D$:

$$\frac{C \vee A \vee \dots \vee A \quad \neg A \vee D}{C \vee D}$$

Formally, this rule is captured by a predicate:

inductive unord_resolve :: 'a clause ⇒ 'a clause ⇒ 'a clause ⇒ bool **where**
 unord_resolve (C ⊔ replicate (n + 1) (Pos A)) ({Neg A} ⊔ D) (C ⊔ D)

Soundness is trivial to prove:

lemma unord_resolve_sound: unord_resolve C D E ⇒ I ⊨ C ⇒ I ⊨ D ⇒ I ⊨ E
using unord_resolve.cases **by** fastforce

To prove completeness, it suffices to show that the calculus reduces counterexamples. This corresponds to Theorem 3.8, except that the conclusion is strengthened slightly to match *counterex_reducing_inference_system*'s assumption:

theorem unord_resolve_counterex_reducing:
assumes {} ∉ N **and** C ∈ N **and** INTERP N ⊭ C **and**
 ∀D ∈ N. INTERP N ⊭ D ⇒ C ≤ D
obtains D E **where**
 D ∈ N **and** INTERP N ⊨ D **and** production N D ≠ {} **and**
 unord_resolve D C E **and** INTERP N ⊭ E **and** E < C

The arguments *N* to INTERP and production are necessary because we are outside the block in which *N* was fixed. This explicit dependency allows us to instantiate the locale's *I_of* :: 'a clause set ⇒ 'a set parameter with INTERP.

By instantiating the *sound_inference_system* and *counterex_reducing_inference_system* locales, we obtain refutational completeness (Theorem 3.9 and Corollary 3.10) and compactness of clausal logic (Theorem 3.12).

Ordered Resolution with Selection. Ordered ground resolution consists of a single rule, called ord_resolve. Like unord_resolve, it is sound and counterexample-reducing (Theorem 3.15). Moreover, it is reductive (Lemma 3.13): the conclusion is always smaller than the main premise according to the clause order. The rule is given as

$$\frac{C_1 \vee A_1 \vee \dots \vee A_1 \quad \dots \quad C_n \vee A_n \vee \dots \vee A_n \quad \neg A_1 \vee \dots \vee \neg A_n \vee D}{C_1 \vee \dots \vee C_n \vee D}$$

with multiple side conditions whose role is to prune the search space and to make the rule reductive. In Isabelle, we represent the *n* side premises by three parallel lists of length *n*: *CAs* gives the entire clauses, whereas *Cs* and *As* store the *C_i* and the *A_i* = *A_i* ∨ ... ∨ *A_i* parts separately. In addition, *As* is the list [*A₁*, ..., *A_n*]. The following inductive definition captures the rule formally:

inductive ord_resolve :: 'a clause list ⇒ 'a clause ⇒ 'a clause ⇒ bool **where**
 |CAs| = n ⇒ |Cs| = n ⇒ |As| = n ⇒ |As| = n ⇒ n ≠ 0 ⇒
 (∀i < n. CAs ! i = Cs ! i ⊔ Pos ' As ! i) ⇒ (∀i < n. As ! i ≠ {}) ⇒
 (∀i < n. ∀A ∈ As ! i. A = As ! i) ⇒ eligible As (D ⊔ Neg ' mset As) ⇒
 (∀i < n. strict_maximal_in (As ! i) (Cs ! i)) ⇒ (∀i < n. S (CAs ! i) = {}) ⇒
 ord_resolve CAs (D ⊔ Neg ' mset As) ((⋃ mset Cs) ⊔ D)

The $xs!i$ operator returns the $(i+1)$ st element of xs , and `mset` converts a list to a multiset. Before settling on the above formulation, we tried storing the n premises in a multiset, since their order is irrelevant. However, due to the permutative nature of multisets, there can be no such things as “parallel multisets”; to keep the dependencies between the C_i ’s and the A_i ’s, we must keep them in a single multiset of tuples, which is very unwieldy.

A previous version of the formalization represented each $A_i \vee \dots \vee A_i$ as a value of type $'a \times nat$ —the nat representing the number of times A_i is repeated. With this approach, the definition of `ord_resolve` did not need to state the equality of the atoms in each $As!i$. Other than that, there was nothing to win, and the approach does not work on the first-order level where atoms should be unifiable instead of equal. To achieve symmetry between the ground and first-order calculi, we went with the current approach.

Formalization revealed an error and a few ambiguities in the rule’s statement. References to $S(D)$ in the side conditions should have been to $S(\neg A_1 \vee \dots \vee \neg A_n \vee D)$. The ambiguities are discussed in our Section 8.

Soundness is a good sanity check for our definition:

lemma `ord_resolve_sound`:

$$\text{ord_resolve } CAs \ DA \ E \Rightarrow I \models \text{mset } CAs \Rightarrow I \models DA \Rightarrow I \models E$$

The proof is by case distinction: either the interpretation I contains all atoms A_i , in which case the D subclause of the main premise $\neg A_1 \vee \dots \vee \neg A_n \vee D$ must be true, or there exists an index i such that $A_i \notin I$, in which case the corresponding C_i must be true. In both cases, the conclusion $C_1 \vee \dots \vee C_n \vee D$ is true.

5 Theorem Proving Processes

In their Section 4, Bachmair and Ganzinger switch from a static to a dynamic view of saturation—from clause sets closed under inferences to theorem proving processes that start with a clause set N_0 and keep deriving new clauses until \perp is generated or no inferences are possible. Proving processes support an important optimization: redundant clauses can be deleted at any point from the clause set, and redundant inferences need not be performed at all.

A derivation performed by a proving process is a potentially infinite sequence $N_0 \triangleright N_1 \triangleright N_2 \triangleright \dots$, where \triangleright is a relation on clause sets (`Proving_Process.thy`). In Isabelle, such sequences are captured by lazy lists, a codatatype [5] generated by `LNil :: 'a llist` and `LCons :: 'a ⇒ 'a llist ⇒ 'a llist`, and equipped with `lhd` (“head”) and `ltl` (“tail”) selectors that extract `LCons`’s arguments. Unlike datatypes, codatypes allow infinite values—e.g., `LCons 0 (LCons 1 (LCons 2 ...))`. The coinductive predicate `chain` checks that its argument is a nonempty lazy list whose elements are consecutively related by a given binary predicate R :

coinductive `chain` :: $('a \Rightarrow 'a \Rightarrow bool) \Rightarrow 'a \text{ llist} \Rightarrow bool$ **where**
`chain R (LCons x LNil)`
`| chain R xs ⇒ R x (lhd xs) ⇒ chain R (LCons x xs)`

A derivation is a lazy list Ns of clause sets satisfying the chain predicate with $R = \triangleright$. Derivations depend on a redundancy criterion presented as two functions, \mathcal{R}_F and \mathcal{R}_I , that specify redundant clauses and redundant inferences, respectively:

locale *redundancy_criterion* = *inference_system* +
fixes
 $\mathcal{R}_F :: 'a \text{ clause set} \Rightarrow 'a \text{ clause set}$ **and**
 $\mathcal{R}_I :: 'a \text{ clause set} \Rightarrow 'a \text{ inference set}$
assumes
 $\mathcal{R}_I N \subseteq \Gamma$ **and**
 $N \subseteq N' \Rightarrow \mathcal{R}_F N \subseteq \mathcal{R}_F N'$ **and**
 $N \subseteq N' \Rightarrow \mathcal{R}_I N \subseteq \mathcal{R}_I N'$ **and**
 $N' \subseteq \mathcal{R}_F N \Rightarrow \mathcal{R}_F N \subseteq \mathcal{R}_F (N \setminus N')$ **and**
 $N' \subseteq \mathcal{R}_I N \Rightarrow \mathcal{R}_I N \subseteq \mathcal{R}_I (N \setminus N')$ **and**
 $\text{sat} (N \setminus \mathcal{R}_F N) \Rightarrow \text{sat } N$

By definition, a transition from M to N is possible if the only new clauses added are conclusions of inferences from M and any deleted clauses would be redundant in N :

inductive $\triangleright :: 'a \text{ clause set} \Rightarrow 'a \text{ clause set} \Rightarrow \text{bool}$ **where**
 $N \setminus M \subseteq \text{concl_of 'infers_from } M \Rightarrow M \setminus N \subseteq \mathcal{R}_F N \Rightarrow M \triangleright N$

This rule combines deduction (the addition of inferred clauses) and deletion (the removal of redundant clauses) in a single transition. The chapter keeps the two operations separated, but this is problematic, as we will see in Section 7.

A key concept to connect static and dynamic completeness is that of the set of persistent clauses, or limit of a sequence of clause sets: $N_\infty = \bigcup_i \bigcap_{j \geq i} N_j$. These are the clauses that belong to all clause sets except for at most a finite prefix of the sequence N_i . We also need the supremum of a sequence, $\bigcup_i N_i$, and of a bounded prefix, $\bigcup_{i=0}^j N_i$. We introduce these missing functions (`Lazy_List_Liminf.thy`):

definition *Liminf* :: *'a llist* \Rightarrow *'a* **where**
 $\text{Liminf } xs = \bigcup_{i < |xs|} \bigcap_{j: i \leq j < |xs|} xs ! j$
definition *Sup* :: *'a llist* \Rightarrow *'a* **where**
 $\text{Sup } xs = \bigcup_{i < |xs|} xs ! i$
definition *Sup_upto* :: *'a llist* \Rightarrow *nat* \Rightarrow *'a* **where**
 $\text{Sup_upto } xs \ j = \bigcup_{i: i < |xs| \wedge i \leq j} xs ! i$

Even though codatatypes open the door to coinductive methods, we follow whenever possible the chapter's index-based approach. When interpreting the notation $\bigcup_i \bigcap_{j \geq i} N_j$ for the case of a finite sequence of length n , it is crucial to use the right upper bounds, namely $i, j < n$. For j , it is clear that ' $< n$ ' is needed to keep N_j 's index within bounds. For i , the danger is more subtle: if $i \geq n$, then $\bigcap_{j: i \leq j < n} N_j$ collapses to the trivial infimum $\bigcap_{j \in \{\}} N_j$, i.e., the set of all clauses.

Lemma 4.2 connects the redundant clauses and inferences at the limit to those of the supremum, and the satisfiability of the limit to that of the initial clause set. Formally:

lemma *Rf_limit_eq_Rf_Sup*: chain (\triangleright) $Ns \Rightarrow \mathcal{R}_{\mathcal{F}}(\text{Liminf } Ns) = \mathcal{R}_{\mathcal{F}}(\text{Sup } Ns)$
lemma *Ri_limit_eq_Ri_Sup*: chain (\triangleright) $Ns \Rightarrow \mathcal{R}_{\mathcal{I}}(\text{Liminf } Ns) = \mathcal{R}_{\mathcal{I}}(\text{Sup } Ns)$
lemma *sat_deriv_limit_iff*: chain (\triangleright) $Ns \Rightarrow (\text{sat}(\text{Liminf } Ns) \Leftrightarrow \text{sat}(\text{Lhd } Ns))$

The proof of the last lemma relies on

lemma *deriv_sat_preserving*: chain (\triangleright) $Ns \Rightarrow \text{sat}(\text{Lhd } Ns) \Rightarrow \text{sat}(\text{Sup } Ns)$

In the chapter, this property follows “by the soundness of the inference system Γ and the compactness of clausal logic,” contradicting the claim that “we will only consider consistency-preserving inference systems” [2, Section 2.4] and not sound ones. Thanks to Isabelle, we now know that soundness is unnecessary. By compactness, it suffices to show that all finite subsets \mathcal{D} of $\bigcup_i N_i$ are satisfiable. By finiteness of \mathcal{D} , there must exist an index k such that $\mathcal{D} \subseteq \bigcup_{i=0}^k N_i$. We perform an induction on k . The base case is trivial since N_0 is assumed to be satisfiable. For the induction step, if k is beyond the end of the list, then $\bigcup_{i=0}^k N_i = \bigcup_{i=0}^{k-1} N_i$ and we can apply the induction hypothesis directly. Otherwise, we have that the set $\text{Sup_upto } Ns (k-1) \cup \text{concl_of } \text{'infers_from } (\text{Sup_upto } Ns (k-1))$ is satisfiable by the induction hypothesis and satisfiability preservation of Γ inferences. Hence, $\text{Sup_upto } Ns (k-1) \cup Ns !k$, i.e., $\text{Sup_upto } Ns k$, is satisfiable, as desired.

Next, we show that the limit is saturated, under some assumptions and for a relaxed notion of saturation. A clause set N is saturated up to redundancy if all inferences from nonredundant clauses in N are redundant:

definition *saturated_upto* :: 'a clause set \Rightarrow bool **where**
saturated_upto $N \Leftrightarrow \text{infers_from } (N \setminus \mathcal{R}_{\mathcal{F}} N) \subseteq \mathcal{R}_{\mathcal{I}} N$

The limit is assured to be saturated for fair derivations—derivations in which no inferences from nonredundant persisting clauses are delayed indefinitely:

definition *fair_cls_seq* :: 'a clause set llist \Rightarrow bool **where**
fair_cls_seq $Ns \Leftrightarrow \text{let } N' = \text{Liminf } Ns \setminus \mathcal{R}_{\mathcal{F}}(\text{Liminf } Ns) \text{ in}$
 $\text{concl_of } \text{'infers_from } N' \setminus \mathcal{R}_{\mathcal{I}} N' \subseteq \text{Sup } Ns \cup \mathcal{R}_{\mathcal{F}}(\text{Sup } Ns)$

The criterion must also be effective, which is expressed by a locale:

locale *effective_redundancy_criterion* = *redundancy_criterion* +
assumes $\gamma \in \Gamma \Rightarrow \text{concl_of } \gamma \in N \cup \mathcal{R}_{\mathcal{F}} N \Rightarrow \gamma \in \mathcal{R}_{\mathcal{I}} N$

In a locale that combines *sat_preserving_inference_system* and *effective_redundancy_criterion*, we have Theorem 4.3:

theorem *fair_derive_saturated_upto*:
chain (\triangleright) $Ns \Rightarrow \text{fair_cls_seq } Ns \Rightarrow \text{saturated_upto } (\text{Liminf } Ns)$

It is easy to show that the trivial criterion defined by $\mathcal{R}_{\mathcal{F}} N = \{\}$ and $\mathcal{R}_{\mathcal{I}} N = \{\gamma \in \Gamma \mid \text{concl_of } \gamma \in N\}$ satisfies the requirements on a *effective_redundancy_criterion*. A more useful instance is the standard redundancy criterion, which depends on a counter-example-reducing inference system Γ (*Standard_Redundancy.thy*):

definition *R_F* :: 'a clause set \Rightarrow 'a clause set **where**
 $\mathcal{R}_{\mathcal{F}} N = \{C \mid \exists D \subseteq N. (\forall I. I \models D \Rightarrow I \models C) \wedge (\forall D \in \mathcal{D}. D < C)\}$

definition $\mathcal{R}_I :: 'a \text{ clause set} \Rightarrow 'a \text{ inference set}$ **where**

$$\mathcal{R}_I N = \{\gamma \in \Gamma \mid \exists \mathcal{D} \subseteq N. (\forall I. I \models \mathcal{D} \uplus \text{side_prems_of } \gamma \Rightarrow I \models \text{concl_of } \gamma) \wedge (\forall D \in \mathcal{D}. D < \text{main_prem_of } \gamma)\}$$

Standard redundancy qualifies as an *effective_redundancy_criterion*. In the chapter, this is stated as Theorems 4.7 and 4.8, which depend on two auxiliary properties, Lemmas 4.5 and 4.6. The main theorem, Theorem 4.9, is that counterexample-reducing calculi are refutationally complete also under the application of standard redundancy:

theorem *saturated_upto_complete*: $\text{saturated_upto } N \Rightarrow (\neg \text{sat } N \Leftrightarrow \{\} \in N)$

The informal proof of Lemma 4.6 applies Lemma 4.5 in a seemingly impossible way, confusing redundant clauses and redundant inferences and exploiting properties that appear only in the first lemma’s proof. Our solution is to generalize the core argument into the following lemma and apply it to prove Lemmas 4.5 and 4.6:

lemma *wlog_non_Rf*:

$$(\exists \mathcal{D} \subseteq N. (\forall I. I \models \mathcal{D} \uplus \mathcal{C} \Rightarrow I \models E) \wedge (\forall D' \in \mathcal{D}. D' < D)) \Rightarrow \exists \mathcal{D} \subseteq N \setminus \mathcal{R}_F N. (\forall I. I \models \mathcal{D} \uplus \mathcal{C} \Rightarrow I \models E) \wedge (\forall D' \in \mathcal{D}. D' < D)$$

Incidentally, the informal proof of Theorem 4.9 also needlessly invokes Lemma 4.5.

Finally, given a redundancy criterion $(\mathcal{R}_F, \mathcal{R}_I)$ for Γ , its standard extension for $\Gamma' \supseteq \Gamma$ is defined as $(\mathcal{R}_F, \mathcal{R}'_I)$, where $\mathcal{R}'_I N = \mathcal{R}_I N \cup (\Gamma' \setminus \Gamma)$ (*Proving_Process.thy*). The standard extension is itself a redundancy criterion and it preserves effectiveness, saturation up to redundancy, and fairness. In Isabelle, this can be expressed by leaving the locales and using the locale predicates—explicit predicates named after the locales and parameterized by the locale arguments:

lemma *standard_redundancy_criterion_extension*:

$$\Gamma \subseteq \Gamma' \Rightarrow \text{redundancy_criterion } \Gamma \mathcal{R}_F \mathcal{R}_I \Rightarrow \text{redundancy_criterion } \Gamma' \mathcal{R}_F \mathcal{R}'_I$$

lemma *standard_redundancy_criterion_extension_effective*:

$$\Gamma \subseteq \Gamma' \Rightarrow \text{effective_redundancy_criterion } \Gamma \mathcal{R}_F \mathcal{R}_I \Rightarrow \text{effective_redundancy_criterion } \Gamma' \mathcal{R}_F \mathcal{R}'_I$$

lemma *standard_redundancy_criterion_extension_saturated_upto_iff*:

$$\Gamma \subseteq \Gamma' \Rightarrow \text{redundancy_criterion } \Gamma \mathcal{R}_F \mathcal{R}_I \Rightarrow (\text{redundancy_criterion.saturated_upto } \Gamma \mathcal{R}_F \mathcal{R}_I N \Leftrightarrow \text{redundancy_criterion.saturated_upto } \Gamma' \mathcal{R}_F \mathcal{R}'_I N)$$

lemma *standard_redundancy_criterion_extension_fair_iff*:

$$\Gamma \subseteq \Gamma' \Rightarrow \text{effective_redundancy_criterion } \Gamma \mathcal{R}_F \mathcal{R}_I \Rightarrow (\text{effective_redundancy_criterion.fair_class_seq } \Gamma' \mathcal{R}_F \mathcal{R}'_I Ns \Leftrightarrow \text{effective_redundancy_criterion.fair_class_seq } \Gamma \mathcal{R}_F \mathcal{R}_I Ns)$$

6 First-Order Resolution

The chapter’s Section 4.3 presents a first-order version of the ordered resolution rule and a first-order prover, RP, based on that rule. The first step towards lifting the completeness of ground resolution is to show that we can lift individual ground resolution inferences (*F0_Ordered_Resolution.thy*).

Inference Rule. First-order ordered resolution consists of a single rule. In the chapter, ground and first-order resolution are both called $O_{\Sigma}^>$. In the formalization, we also let the rules share the same name, but since they exist in separate locales the system generates qualified names which make this unambiguous: Isabelle generates the name `ground_resolution_with_selection.ord_resolve`, which refers to ground resolution, and `FO_resolution.ordered_resolve`, which refers to first-order resolution. If the user is in doubt at any time, the system can always tell which one is meant.

The rule is given as

$$\frac{C_1 \vee A_{11} \vee \dots \vee A_{1k_1} \quad \dots \quad C_n \vee A_{n1} \vee \dots \vee A_{nk_n} \quad \neg A_1 \vee \dots \vee \neg A_n \vee D}{C_1 \cdot \sigma \vee \dots \vee C_n \cdot \sigma \vee D \cdot \sigma}$$

where σ is the (canonical) MGU that solves all unification problems $A_{i1} \stackrel{?}{=} \dots \stackrel{?}{=} A_{ik_i} \stackrel{?}{=} A_i$, for $1 \leq i \leq n$. As expected, the rule has several side conditions. The Isabelle representation of this rule is based on that of its ground counterpart, generalized to apply σ :

inductive `ord_resolve` :: *'a clause list* \Rightarrow *'a clause* \Rightarrow *'s* \Rightarrow *'a clause* \Rightarrow *bool* **where**
 $|CA_s| = n \Rightarrow |Cs| = n \Rightarrow |As| = n \Rightarrow |As| = n \Rightarrow n \neq 0 \Rightarrow$
 $(\forall i < n. CA_s ! i = Cs ! i \uplus Pos \cdot As ! i) \Rightarrow (\forall i < n. As ! i \neq \{\}) \Rightarrow$
 $Some \sigma = mgu (set_mset \cdot set (map2 add_mset As As)) \Rightarrow$
 $eligible \sigma As (D \uplus Neg \cdot mset As) \Rightarrow$
 $(\forall i < n. strict_maximal_in (As ! i \cdot \sigma) (Cs ! i \cdot \sigma)) \Rightarrow (\forall i < n. S (CA_s ! i) = \{\}) \Rightarrow$
`ord_resolve` $CA_s (D \uplus Neg \cdot mset As) \sigma (((\bigcup mset Cs) \uplus D) \cdot \sigma)$

The rule as stated is incomplete; for example, $p(x)$ and $\neg p(f(x))$ cannot be resolved because x and $f(x)$ are not unifiable. Such issues arise when the same variable names appear in different premises. In the chapter, the authors circumvent this issue by stating, “We also implicitly assume that different premises and the conclusion have no variables in common; variables are renamed if necessary.” For the formalization, we first considered enforcing the invariant that all derived clauses use mutually disjoint variables, but this does not help when a clause is repeated in an inference’s premises. An example is the inference

$$\frac{p(x) \quad p(y) \quad \neg p(a) \vee \neg p(b)}{\perp}$$

where $p(x)$ and $p(y)$ are the same clause up to renaming. Instead, we rely on a predicate `ord_resolve_rename`, based on `ord_resolve`, that standardizes the premises apart. The renaming is performed by a function called `renamings_apart` :: *'a clause list* \Rightarrow *'s list* that, given a list of clauses, returns a list of corresponding substitutions to apply. This function is part of the abstract interface for terms and substitutions (which we presented in Section 2) and is implemented using `lsaFoR`.

Like for the ground case, it is important to establish soundness. We prove that any ground instance of the rule `ord_resolve` is sound:

lemma `ord_resolve_ground_inst_sound`:
`ord_resolve` $CA_s DA As As \sigma E \Rightarrow I \models mset CA_s \cdot \sigma \cdot \eta \Rightarrow I \models DA \cdot \sigma \cdot \eta \Rightarrow$
`is_ground_subst` $\eta \Rightarrow I \models E \cdot \eta$

Likewise, ground instances of `ord_resolve_rename` are sound. It then follows that the rules `ord_resolve` and `ord_resolve_rename` are sound:

lemma *ord_resolve_rename_sound*:
 $\text{ord_resolve_rename } CAs\ DA\ \mathcal{A}s\ As\ \sigma\ E \Rightarrow$
 $(\forall \sigma. \text{is_ground_subst } \sigma \Rightarrow I \models (\text{mset } CAs + \{DA\}) \cdot \sigma) \Rightarrow$
 $\text{is_ground_subst } \eta \Rightarrow I \models E \cdot \eta$

Lifting Lemma. To lift ground inferences to the first-order level, we consider a set of clauses M and introduce an adjusted version S_M of the selection function S . This new selection function depends on both S and M and works in such a way that any ground instance inherits the selection of at least one of the nonground clauses of which it is an instance. This property is captured formally as

lemma *S_M_grounding_of_cls*:
 $C \in \text{grounding_of } M \Rightarrow$
 $\exists D \in M. \exists \sigma. C = D \cdot \sigma \wedge S_M C = S D \cdot \sigma \wedge \text{is_ground_subst } \sigma$

where `grounding_of` M is the set of ground instances of a set of clauses M .

The lifting lemma, Lemma 4.12, states that whenever there exists a ground inference of E from clauses belonging to `grounding_of` M , there exists a (potentially) more general inference from clauses belonging to M :

lemma *ord_resolve_rename_lifting*:
 $(\forall \rho C. \text{is_renaming } \rho \Rightarrow S (C \cdot \rho) = S C \cdot \rho) \Rightarrow$
 $\text{ord_resolve } S_M CAs\ DA\ \mathcal{A}s\ As\ \sigma\ E \Rightarrow$
 $\{DA\} \cup \text{set } CAs \subseteq \text{grounding_of } M \Rightarrow$
 $\exists \eta_s\ \eta\ \theta\ CA_{s_0}\ DA_0\ \mathcal{A}_{s_0}\ A_{s_0}\ E_0\ \tau.$
 $\text{ord_resolve_rename } S CA_{s_0}\ DA_0\ \mathcal{A}_{s_0}\ A_{s_0}\ \tau\ E_0$
 $\wedge CA_{s_0} \cdot \eta_s = CAs \wedge DA_0 \cdot \eta = DA \wedge E_0 \cdot \theta = E \wedge \{DA_0\} \cup \text{set } CA_{s_0} \subseteq M$

The informal proof of this lemma consists of two sentences spanning four lines of text. In Isabelle, these two sentences translate to 75 lines and 400 lines, respectively, excluding auxiliary lemmas. Our proof involves six steps:

1. Obtain a list of first-order clauses CA_{s_0} and a first-order clause DA_0 that belong to M and that generalize CAs and DA with substitutions η_s and η , respectively.
2. Choose atoms \mathcal{A}_{s_0} and A_{s_0} in the first-order clauses on which to resolve.
3. Standardize CA_{s_0} and DA_0 apart, yielding CA'_{s_0} and DA'_0 .
4. Obtain the MGU τ of the literals on which to resolve.
5. Show that ordered resolution on CA'_{s_0} and DA'_0 with τ as MGU is applicable.
6. Show that the resulting resolvent E_0 generalizes E with substitution θ .

In step 1, suitable clauses must be chosen so that $S (CA_{s_0} ! i)$ generalizes $S_M (CAs ! i)$, for $0 \leq i < n$, and $S DA_0$ generalizes $S_M DA$. By definition of S_M , this is always possible. In step 2, we choose the literals to resolve upon in the first-order inference depending on the selection on the ground inference. If some literals are selected in DA ,

we define As_0 as the selected literals in DA_0 , such that $(As_0 ! i) \cdot \eta = As ! i$ for each i . Otherwise, As must be a singleton list containing some atom A , and we define As_0 as the singleton list consisting of an arbitrary $A_0 \in DA_0$ such that $A_0 \cdot \eta = A$. Step 3 may seem straightforward until one realizes that renaming variables can influence selection. To rule this out, our lemma assumes stability under renaming: $S(C \cdot \rho) = S C \cdot \rho$ for any renaming substitution ρ and clause C . This requirement seems natural, but it is not mentioned in the chapter.

The above choices allow us to perform steps 4 to 6. In the chapter, the authors assume that the obtained CA_{s_0} and DA_0 are standardized apart from each other as well as their conclusion E_0 . This means that they can obtain a single ground substitution μ that connect CA_{s_0} , DA_0 , E_0 to CA_s , DA , E . By contrast, we provide separate substitutions η_s , η , θ for the different side premises, the main premise, and the conclusion.

7 A First-Order Prover

Modern resolution provers interleave inference steps with steps that delete or reduce (simplify) clauses. In their Section 4.3, Bachmair and Ganzinger introduce the non-deterministic abstract prover RP that works on triples of clause sets and that generalizes the Otter-style and DISCOUNT-style loops [12, 16]. RP's core rule, called inference computation, performs first-order ordered resolution as described above; the other rules delete or reduce clauses or move them between clause sets. We formalize RP and prove it complete assuming a fair strategy (`F0_Ordered_Resolution_Prover.thy`).

Abstract First-Order Prover. The RP prover is a relation \rightsquigarrow on states of the form $(\mathcal{N}, \mathcal{P}, \mathcal{O})$, where \mathcal{N} is the set of *new clauses*, \mathcal{P} is the set of *processed clauses*, and \mathcal{O} is the set of *old clauses*. RP's formal definition is very close to the original formulation:

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

- Neg $A \in C \Rightarrow$ Pos $A \in C \Rightarrow (\mathcal{N} \cup \{C\}, \mathcal{P}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P}, \mathcal{O})$
- | $D \in \mathcal{P} \cup \mathcal{O} \Rightarrow$ subsumes $DC \Rightarrow (\mathcal{N} \cup \{C\}, \mathcal{P}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P}, \mathcal{O})$
- | $D \in \mathcal{N} \Rightarrow$ strictly_subsumes $DC \Rightarrow (\mathcal{N}, \mathcal{P} \cup \{C\}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P}, \mathcal{O})$
- | $D \in \mathcal{N} \Rightarrow$ strictly_subsumes $DC \Rightarrow (\mathcal{N}, \mathcal{P}, \mathcal{O} \cup \{C\}) \rightsquigarrow (\mathcal{N}, \mathcal{P}, \mathcal{O})$
- | $D \in \mathcal{P} \cup \mathcal{O} \Rightarrow$ reduces $DCL \Rightarrow (\mathcal{N} \cup \{C \uplus \{L\}\}, \mathcal{P}, \mathcal{O}) \rightsquigarrow (\mathcal{N} \cup \{C\}, \mathcal{P}, \mathcal{O})$
- | $D \in \mathcal{N} \Rightarrow$ reduces $DCL \Rightarrow (\mathcal{N}, \mathcal{P} \cup \{C \uplus \{L\}\}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P} \cup \{C\}, \mathcal{O})$
- | $D \in \mathcal{N} \Rightarrow$ reduces $DCL \Rightarrow (\mathcal{N}, \mathcal{P}, \mathcal{O} \cup \{C \uplus \{L\}\}) \rightsquigarrow (\mathcal{N}, \mathcal{P} \cup \{C\}, \mathcal{O})$
- | $(\mathcal{N} \cup \{C\}, \mathcal{P}, \mathcal{O}) \rightsquigarrow (\mathcal{N}, \mathcal{P} \cup \{C\}, \mathcal{O})$
- | $(\{\}, \mathcal{P} \cup \{C\}, \mathcal{O}) \rightsquigarrow (\text{concl_of } ' \text{infers_between } \mathcal{O} C, \mathcal{P}, \mathcal{O} \cup \{C\})$

The rules correspond, respectively, to tautology deletion, forward subsumption, backward subsumption in \mathcal{P} and \mathcal{O} , forward reduction, backward reduction in \mathcal{P} and \mathcal{O} , clause processing, and inference computation.

Initially, \mathcal{N} consists of the problem clauses and the other two sets are empty. Clauses in \mathcal{N} are reduced using $\mathcal{P} \cup \mathcal{O}$, or even deleted if they are tautological or subsumed by $\mathcal{P} \cup \mathcal{O}$; conversely, \mathcal{N} can be used for reducing or subsuming clauses in $\mathcal{P} \cup \mathcal{O}$. Clauses eventually move from \mathcal{N} to \mathcal{P} , one at a time. As soon as \mathcal{N} is empty, a clause from \mathcal{P} is selected to move to \mathcal{O} . Then all possible resolution inferences between this given clause and the clauses in \mathcal{O} are computed and put in \mathcal{N} , closing the loop.

The subsumption and reduction rules depend on the following predicates:

$$\begin{aligned} \text{subsumes } DC &\Leftrightarrow \exists \sigma. D \cdot \sigma \subseteq C \\ \text{strictly_subsumes } DC &\Leftrightarrow \text{subsumes } DC \wedge \neg \text{subsumes } CD \\ \text{reduces } DCL &\Leftrightarrow \exists D' L' \sigma. D = D' \uplus \{L'\} \wedge -L = L' \cdot \sigma \wedge D' \cdot \sigma \subseteq C \end{aligned}$$

The definition of the set $\text{infers_between } \mathcal{O} C$, on which inference computation depends, is more subtle. In the chapter, the set of inferences between C and \mathcal{O} consists of all inferences from $\mathcal{O} \cup \{C\}$ that have C as *exactly one* of their premises. This, however, leads to an incomplete prover, because it ignores inferences that need multiple copies of C . For example, assuming a maximal selection function (one that always returns all negative literals), the resolution inference

$$\frac{p \quad p \quad \neg p \vee \neg p}{\perp}$$

is possible. Yet if the clause $\neg p \vee \neg p$ reaches \mathcal{O} earlier than p , the inference would not be performed. This counterexample requires ternary resolution, but there also exists a more complicated one for binary resolution, where both premises are the same clause. Consider the clause set containing

$$(1) \ q(a, c, b) \quad (2) \ \neg q(x, y, z) \vee q(y, z, x) \quad (3) \ \neg q(b, a, c)$$

and an order $>$ on atoms such that $q(c, b, a) > q(b, a, c) > q(a, c, b)$. Inferences between (1) and (2) or between (2) and (3) are impossible due to order restrictions. The only possible inference involves two copies of (2):

$$\frac{\neg q(x, y, z) \vee q(y, z, x) \quad \neg q(x', y', z') \vee q(y', z', x')}{\neg q(x, y, z) \vee q(z, x, y)}$$

From the conclusion, we derive $\neg q(a, c, b)$ by (3) and \perp by (1). This incompleteness is a severe flaw, although it is probably just an oversight. Fortunately, it can easily be repaired by defining $\text{infers_between } \mathcal{O} C$ as $\{(C, D, E) \in \Gamma \mid \mathcal{C} \cup \{D\} \subseteq \mathcal{O} \cup \{C\} \wedge C \in \mathcal{C} \cup \{D\}\}$.

Projection to Theorem Proving Process. On the first-order level, a derivation can be expressed as a lazy list $\mathcal{S}s$ of states, or as three parallel lazy lists $\mathcal{N}s, \mathcal{P}s, \mathcal{O}s$. The limit state of a derivation $\mathcal{S}s$ is defined as $\text{Liminf } \mathcal{S}s = (\text{Liminf } \mathcal{N}s, \text{Liminf } \mathcal{P}s, \text{Liminf } \mathcal{O}s)$, where Liminf on the right-hand side is as in Section 5.

Bachmair and Ganzinger use the completeness of ground resolution to prove RP complete. The first step is to show that first-order derivations can be projected down to theorem proving processes on the ground level. This corresponds to Lemma 4.10. Adapted to our conventions, its statement is as follows:

If $\mathcal{S} \rightsquigarrow \mathcal{S}'$, then $\text{grounding_of } \mathcal{S} \triangleright^* \text{grounding_of } \mathcal{S}'$, with \triangleright based on some extension of ordered resolution with selection function S and the standard redundancy criterion $(\mathcal{R}_F, \mathcal{R}_I)$.

This raises some questions: (1) Exactly which instance of the calculus are we extending? (2) Which calculus extension should we use? (3) How can we repair the mismatch between \triangleright^* in the lemma statement and \triangleright where the lemma is invoked?

Regarding question (1), it is not clear which selection function to use. Is the function the same S as in the definition of RP or is it arbitrary? It takes a close inspection of the proof of Lemma 4.13, where Lemma 4.10 is invoked, to find out that the selection function used there is $S_{\text{Liminf } \mathcal{O}_S}$.

Regarding question (2), the phrase “some extension” is cryptic. It suggests an existential reading, and from the context it would appear that a standard extension (Section 5) is meant. However, neither the lemma’s proof nor the context where it is invoked supplies the desired existential witness. A further subtlety is that the witness should be independent of S and S' , so that transitions can be joined to form a single theorem proving derivation. Our approach is to let \triangleright be the standard extension for the proof system consisting of all *sound* derivations: $\Gamma = \{(\mathcal{C}, D, E) \mid \forall I. I \models \mathcal{C} \cup \{D\} \implies I \models E\}$. This also eliminates the need for Bachmair and Ganzinger’s subsumption resolution rule, a special calculus rule that is, from what we understand, implicitly used in the proof of Lemma 4.10 for the subcases associated with RP’s reduction rules.

As for question (3), when the lemma is invoked, it is used to join transitions together to whole theorem proving processes. That requires these transitions to be of \triangleright – not \triangleright^* . The need for \triangleright^* instead of \triangleright arises because one of the cases requires a combination of deduction and deletion, which Bachmair and Ganzinger model as separate transitions. By merging the two transitions (Section 5), we avoid the issue altogether and can use \triangleright in the formal counterpart of Lemma 4.10.

With these issues resolved, we can prove Lemma 4.10. In Section 6 we established that ground instances of the resolution rule are sound. Since our ground proof system consists of all sound inference rules we can reuse that lemma in proving the inference computation case. We prove Lemma 4.10 for single steps and extend it to entire derivations:

lemma *RP_ground_derive*: $S \rightsquigarrow S' \implies \text{grounding_of } S \triangleright \text{grounding_of } S'$

lemma *RP_ground_derive_chain*:

$\text{chain } (\rightsquigarrow) Ss \implies \text{chain } (\triangleright) (\text{lmap } \text{grounding_of } Ss)$

The `lmap` function applies its first argument elementwise to its second argument.

Fairness and Clause Movement. From a given initial state $(\mathcal{N}_0, \{\}, \{\})$, many derivations are possible, reflecting RP’s nondeterminism. In some derivations, we could leave a crucial clause in \mathcal{N} or \mathcal{P} without ever reducing it or moving it to \mathcal{O} , and then fail to derive \perp even if \mathcal{N}_0 is unsatisfiable. For this reason, refutational completeness is guaranteed only for fair derivations. These are defined as derivations such that $\text{Liminf } \mathcal{N}_s = \text{Liminf } \mathcal{P}_s = \{\}$, guaranteeing that no clause will stay forever in \mathcal{N} or \mathcal{P} .

Fairness is expressed by the `fair_state_seq` predicate, which is distinct from the `fair_cls_seq` predicate presented in Section 5. In particular, Theorem 4.3 is used in neither the informal nor the formal proof, and appears to play a purely pedagogic role in the chapter. For the rest of this section, we fix a lazy list of states Ss , and its projections \mathcal{N}_s , \mathcal{P}_s , and \mathcal{O}_s , such that $\text{chain } (\rightsquigarrow) Ss$, `fair_state_seq` Ss , and $\text{lhd } \mathcal{O}_s = \{\}$.

Thanks to fairness, any nonredundant clause C in Ss ’s projection to the ground level eventually ends up in \mathcal{O} and stays there. This is proved informally as Lemma 4.11, but again there are some difficulties. The vagueness concerning the selection function can be resolved as for Lemma 4.10, but there is another, deeper flaw.

Bachmair and Ganzinger’s proof idea is as follows. By hypothesis, the ground clause C must be an instance of a first-order clause D in $\mathcal{N}s!j \cup \mathcal{P}s!j \cup \mathcal{O}s!j$ for some index j . If $C \in \mathcal{N}s!j$, then by nonredundancy of C , fairness of the derivation, and Lemma 4.10, there must exist a clause D' that generalizes C in $\mathcal{P}s!l \cup \mathcal{O}s!l$ for some $l > j$. By a similar argument, if D' belongs to $\mathcal{P}s!l$, it will be in $\mathcal{O}s!l'$ for some $l' > l$, and finally in all $\mathcal{O}s!k$ with $k \geq l'$. The flaw is that backward subsumption can delete D' without moving it to \mathcal{O} . The subsumer clause would then be a strictly more general version of D' (and of the ground clause C).

Our solution is to choose D , and consequently D' , such that it is minimal, with respect to subsumption, among the clauses that generalize C in the derivation. This works because strict subsumption is well founded—which we also proved, by reduction to a well-foundedness result about the strict generalization relation on first-order terms, included in IsaFoR [14, Section 2]. By minimality, D' cannot be deleted by backward subsumption. This line of reasoning allows us to prove Lemma 4.11, where \mathcal{O}_{of} extracts the \mathcal{O} component of a state:

lemma *fair_imp_Liminf_minus_Rf_subset_ground_Liminf_state*:
 $Gs = \text{Imap grounding_of } \mathcal{S}s \Rightarrow$
 $\text{Liminf } Gs - \mathcal{R}_{\mathcal{F}}(\text{Liminf } Gs) \subseteq \text{grounding_of } (\mathcal{O}_{\text{of}}(\text{Liminf } \mathcal{S}s))$

In the formalization of the above proof, we do not prove $l > j$ and $l' > l$. While they guide human intuition, they are not necessary to prove the lemma.

Soundness and Completeness. The chapter’s main result is Theorem 4.13, which states that, for fair derivations, the prover is sound and complete. Soundness follows from Lemma 4.2 (*sat_deriv_Liminf_iff*) and is, perhaps not surprisingly, independent of whether the derivation is fair.

theorem *RP_sound*:
 $\{\} \in \text{cls_of } (\text{Liminf } \mathcal{S}s) \Rightarrow \neg \text{sat } (\text{grounding_of } (\text{lhd } \mathcal{S}s))$

Because we have brought Lemmas 4.10, 4.11, and 4.12 into a suitable shape, completeness is not difficult to formalize:

theorem *RP_saturated_if_fair*:
 $\text{saturated_upto } (\text{Liminf } (\text{Imap grounding_of } \mathcal{S}s))$
corollary *RP_complete_if_fair*:
 $\neg \text{sat } (\text{grounding_of } (\text{lhd } \mathcal{S}s)) \Rightarrow \{\} \in \mathcal{O}_{\text{of}}(\text{Liminf } \mathcal{S}s)$

A crucial point that is not clear from the text is that we must always use the selection function S on the first-order level and $S_{\text{Liminf } \mathcal{O}s}$ on the ground level. Another noteworthy part of the proof is the passage “ $\text{Liminf } Gs$ (and hence $\text{Liminf } \mathcal{S}s$) contains the empty clause” (using our notations). Obviously, if $\text{grounding_of } (\text{Liminf } \mathcal{S}s)$ contains \perp , then $\text{Liminf } \mathcal{S}s$ must as well. However, the authors do not explain the step from $\text{Liminf } Gs$, the limit of the grounding, to $\text{grounding_of } (\text{Liminf } \mathcal{S}s)$, the grounding of the limit. Fortunately, by Lemma 4.11, the latter contains all the nonredundant clauses of the former, and the empty clause is nonredundant. Hence the informal argument is fundamentally correct. For the other direction, which is used in the soundness proof, we can prove that the former includes the latter.

8 Errors and Imprecisions Discovered in the Chapter

In the previous sections, we discussed several mathematical errors and imprecisions, of various severity levels, in Bachmair and Ganzinger’s chapter. We also found lemmas that are stated but not explicitly applied afterwards. In this section, we list our findings exhaustively for reference.

Let us start with the errors and imprecisions. We ignore infelicities that are not mathematical in nature, such as typos and \LaTeX macros gone wrong (e.g., “by the defn[candidate model]candidate model for N ” on page 34); for such errors, careful reading is a more effective antidote than formalization. We also ignore minor ambiguities, such as whether the clause $C \vee A \vee \dots \vee A$ may contain zero occurrences of A , if they can be resolved easily by appealing to the context and the reader’s common sense.

- One of Lemma 3.4’s claims is that if clause C is true in I^D , then C is also true in $I_{D'}$, where $C \preceq D \preceq D'$. This does not hold if $C = D = D'$ and C is productive. Similarly, the first sentence of the proof is wrong if $D = D'$ and D is productive: “First, observe that $I_D \subseteq I^D \subseteq I_{D'} \subseteq I^{D'} \subseteq I_N$, whenever $D' \succeq D$.”
- The last occurrence of D' in the statement of Lemma 3.7 should be changed to C . In addition, it is not clear whether the phrase “another clause C ” implies that $C \neq D$, but the counterexample we gave in Section 4 works in both cases. Correspondingly, in the proof, the case distinction is incomplete, as can be seen by specializing the proof for the counterexample.
- In the chapter’s Figure 2, in Section 3, the selection function is wrongly applied: references to $S(D)$ should be changed to $S(\neg A_1 \vee \dots \vee \neg A_n \vee D)$. Moreover, in condition (iii), it is not clear with respect to which clause the “selected atom” must be considered, the two candidates being $S(\neg A_1 \vee \dots \vee \neg A_n \vee D)$ and $S(C_i \vee A_i \vee \dots \vee A_j)$. We assume the latter is meant.
- The notation $\bigcup_i \bigcap_{j \geq i} N_j$ used in the chapter’s Section 4.1 only partially specifies the range of i and j if N_j is a finite sequence. Clearly, j must be bounded by the length of the list, but it is less obvious that i also needs a bound, to avoid the inner intersection to expand to be the set of all clauses for indices i beyond the list’s end.
- Soundness is required in the chapter’s Section 4.1, even though it is claimed in Section 2.4 that only consistency-preserving inference systems will be considered.
- In the proof of Theorem 4.3, the case where $\gamma \in \mathcal{R}_{\mathcal{I}}(N_\infty \setminus \mathcal{R}_{\mathcal{F}}(N_\infty))$ is not covered.
- In Section 4.2, the phrase “side premises that are true in N ” must be understood as meaning that the side premises both belong to N and are true in I_N .
- Lemma 4.5 states the basic properties of the redundant clause operator $\mathcal{R}_{\mathcal{F}}$ (monotonicity and independence). Lemma 4.6 states the corresponding properties of the redundant inference operator $\mathcal{R}_{\mathcal{I}}$. As justification for Lemma 4.6, the authors tell us that “the proof uses Lemma 4.5,” but redundant inferences are a more general concept than redundant clauses, and we see no way to bridge the gap.
- Similarly, in the proof Theorem 4.9, the application of Lemma 4.5 does not fit. What is needed is a generalization of Lemma 4.6.
- In condition (ii) of Figure 4, Section 4.2, $A_{ii}\sigma$ should be changed to $A_{ij}\sigma$.

- In n th side premise of Figure 4, Section 4.2, A_{1n} should be changed to A_n .
- Section 4.3 states “Subsumption defines a well-founded ordering on clauses.” A simple counterexample is an infinite sequence repeating some clause. A correct statement would instead be “Proper subsumption defines a well-founded ordering on clauses.”
- In Lemma 4.10 it is not clear which selection function is used. When the lemma is applied in the proofs of Lemma 4.11 and Theorem 4.13, it has to be $S_{\mathcal{O}_\infty}$.
- In Lemma 4.10 $G(\mathcal{S})$ and $G(\mathcal{S}')$ are related by \triangleright^* , but \triangleright is needed in the proofs of Lemma 4.11 and Lemma 4.13 since then derivations in RP, which are possibly infinite, can be projected to theorem proving processes. However $G(\mathcal{S}) \triangleright G(\mathcal{S}')$ does not hold in one of the cases since a combination of deduction and deletion is required. A solution is to change the definition of \triangleright to allow such combinations.
- In Lemma 4.10 it is not clear that the extension used should be the same between any considered pair of states. Otherwise, the lemma cannot be used to project derivations in RP to theorem proving processes.
- In Lemma 4.11 it is not clear which selection function is used. When the lemma is applied in the proofs of Theorem 4.13, it has to be $S_{\mathcal{O}_\infty}$.
- A step in the proof of Lemma 4.11 considers a clause $D \in \mathcal{P}_l$ which has a nonredundant instance C . It is claimed that when D is removed from \mathcal{P} , another clause D' with C as instance appears in some \mathcal{O}'_l . That, however, does not follow if D was removed by backward subsumption. The problem can be resolved by choosing D as minimal, with respect to subsumption, among the clauses that generalize C in the derivation. This can be done since proper subsumption is well founded.
- In Lemma 4.11, a very minor inconsistency is that the described first-order derivation is indexed from 1 instead of 0.
- In the proof of Theorem 4.13, the conclusion of Lemma 4.11 is stated as $N_\infty \setminus \mathcal{R}(N_\infty) \subseteq \mathcal{O}_\infty$, but it should have been $N_\infty \setminus \mathcal{R}(N_\infty) \subseteq G(\mathcal{O}_\infty)$. Furthermore, when lemma 4.11 was first stated the conclusion was $N_\infty \setminus \mathcal{R}_{\mathcal{F}}(N_\infty) \subseteq G(\mathcal{S}_\infty)$. The two are by fairness equivalent, but we find $N_\infty \setminus \mathcal{R}(N_\infty) \subseteq G(\mathcal{O}_\infty)$ more intuitive since it more clearly expresses that all nonredundant clauses grow old.

Chief among the factors that make the chapter hard to follow is that many lemmas are stated (and usually proved) but not referenced later. We already mentioned the unfortunate Lemma 3.7. Section 4 contains several other specimens:

- Theorem 4.3 (*fair_derive_saturated_upto*) states a completeness theorem for fair derivations. However, in Section 4.3, fairness is defined differently, and neither the text nor the formalization applies this theorem.
- For the same reason, the property stated in the next-to-last sentence of Section 4.1 (*standard_redundancy_criterion_extension_fair_iff*), which lifts fairness with respect to $(\mathcal{R}_{\mathcal{F}}, \mathcal{R}_{\mathcal{I}})$ to a standard extension $(\mathcal{R}_{\mathcal{F}}, \mathcal{R}'_{\mathcal{I}})$, is not needed later.
- Lemma 4.2 (*sat_deriv_Liminf_iff, Ri_Sup_subset_Ri_Liminf*) is not referenced in the text, but we need it to prove Theorem 4.13 (*fair_state_seq_complete*).

- Lemma 4.2 (*Rf_Sup_subset_Rf_Liminf*) is not referenced in the text, but we need it to prove Lemma 4.11 (*fair_imp_Liminf_minus_Rf_subset_ground_Liminf_state*).
- Lemma 4.6 (*saturated_upto_complete_if*) is not referenced in the text, but we need it to prove Lemma 4.10 (*resolution_prover_ground_derivation*), Lemma 4.11 (*fair_imp_Liminf_minus_Rf_subset_ground_Liminf_state*), and Theorem 4.13 (*fair_state_seq_complete*).
- Theorem 4.8 (*Ri_effective*) is not referenced in the text, but we need it to prove Theorem 4.13 (*fair_state_seq_complete*).
- Theorem 4.9 (*saturated_upto_complete*) is invoked implicitly in the next-to-last sentence in the proof of Theorem 4.13 (*fair_state_seq_complete*).
- The sentence “In that case, if the derivation is fair with respect to inferences in Γ the derivation is also fair with respect to inferences in Γ' , and vice versa” on page 38 of Section 4.1 (*redundancy_criterion_standard_extension_saturated_upto_iff*) is not referenced in the text, but we need it to prove Theorem 4.13 (*fair_state_seq_complete*).

9 Discussion

Bachmair and Ganzinger cover a lot of ground in a few pages. We found much of the material straightforward to formalize: it took us about two weeks to reach their Section 4.3, which introduces the RP prover and establishes its refutational completeness. By contrast, we needed months to fully understand and formalize that section. While the *Handbook* chapter succeeds at conveying the key ideas at the propositional level, the lack of rigor makes it difficult to develop a deep understanding of ordered resolution proving on first-order clauses.

There are several reasons why Section 4.3 did not lend itself easily to a formalization. The proofs often depend on lemmas and theorems from previous sections without explicitly mentioning them. The lemmas and proofs do not quite fit together. And while the general idea of the proofs stands up, they have many confusing flaws that must be repaired. Our methodology involved the following steps: (1) rewrite the informal proofs to a handwritten pseudo-Isabelle; (2) fill in the gaps, emphasizing which lemmas are used where; (3) turn the pseudo-Isabelle to real Isabelle, but with **sorry** placeholders for the proofs; and (4) replace the **sorry**s with proofs. Progress was not always linear. As we worked on each step, more than once we discovered an earlier mistake.

The formalization helps us answer questions such as, “Is effectiveness of ordered resolution (Lemma 3.13) actually needed, and if so, where?” It also allows us to track definitions and hypotheses precisely, so that we always know the scope and meaning of every definition, lemma, or theorem. If a hypothesis appears too strong or superfluous, we can try to rephrase or eliminate it; the proof assistant tells us where the proof breaks. If a definition is changed, the proof assistant tells us where proofs of the related lemmas break. In the best case, the proofs do not break at all since the automation of the proof assistant is flexible enough to still prove them. This happened, for example, when we changed the definition of \triangleright to combine deduction and deletion.

Starting from RP, we could refine it to obtain an efficient imperative implementation, following the lines of Fleury, Blanchette, and Lammich’s verified SAT solver

with the two-watched-literals optimization [13]. However, this would probably involve a huge amount of work. To increase provers' trustworthiness, a more practical approach is to have them generate detailed proofs that record all inferences leading to the empty clause [25, 29]. Such output can be independently checked by verified programs or reconstructed using a proof assistant's inference kernel. This is the approach implemented in Sledgehammer [8], which integrates automatic provers in Isabelle. Formalized metatheory could in principle be used to deduce a formula's satisfiability from a finite saturation.

We found Isabelle/HOL eminently suitable to this kind of formalization work. Its logic—classical simple type theory extended with polymorphism, type classes, and the axiom of choice—balances expressiveness and automatability. We nowhere felt the need for dependent types. We benefited from many features of the system, including co-datatypes [5], Isabelle/jEdit [35], the Isar proof language [34], locales [4], and Sledgehammer [8]. It is perhaps indicative of the maturity of theorem proving technology that most of the issues we encountered were unrelated to Isabelle. The main challenge was to understand the informal proof well enough to design suitable locale hierarchies and state the definitions and lemmas precisely, and correctly.

10 Related Work

Formalizing the metatheory of logic and deduction is an enticing proposition for many researchers in interactive theorem proving. In this section, we briefly review some of the main related work, without claim to exhaustiveness. Two recent, independent developments are particularly pertinent.

Peltier [22] proved static refutational completeness of a variant of the superposition calculus in Isabelle/HOL. Since superposition generalizes ordered resolution, his result subsumes our static completeness theorem. On the other hand, he did not formalize a prover or dynamic completeness, nor did he attempt to develop general infrastructure. It would be interesting to extend his formal development to obtain a verified superposition prover. We could also consider calculus extensions such as polymorphism [11, 32], type classes [32], and AVATAR [31]. Two significant differences between Peltier's work and ours is that he represents clauses as sets instead of multisets (to exploit Isabelle's better proof automation for sets) and that he relies, for terms and unification, on an example theory file included in Isabelle (`Unification.thy`) instead of `IsaFoR`.

Hirokawa et al. [14] formalized, also in Isabelle/HOL, an abstract Knuth–Bendix completion procedure as well as ordered (unfailing) completion, a method developed by Bachmair, Ganzinger, and Plaisted [1]. Superposition combines ordered resolution (to reason about clauses) and ordered completion (to reason about equality). There are many similarities between their formalization and ours, which is unsurprising given that both follow work by Bachmair and Ganzinger; for example, they need to reason about the limit of fair infinite sequences of sets of equations and rewrite rules to establish completeness.

The literature contains many other formalized completeness proofs, mostly for inference systems of theoretical interest. Early work was carried out in the 1980s and 1990s, notably by Shankar [27] and Persson [23]. Some of our own efforts are also re-

lated: completeness of unordered resolution using semantic trees by Schlichtkrull [26]; completeness of a Gentzen system following the Beth–Hintikka style and soundness of a cyclic proof system for first-order logic with inductive definitions by Blanchette, Popescu, and Traytel [9]; and completeness of a SAT solver based on CDCL (conflict-driven clause learning) by Blanchette, Fleury, and Weidenbach [7]. We refer to our earlier papers [7, 9, 26] for further discussions of related work. Beyond completeness, Gödel’s first incompleteness theorem has been formalized in Nqthm by Shankar [28], in Coq by O’Connor [20], in HOL Light by Harrison (in unpublished work), and in Isabelle/HOL by Paulson [21]. Paulson additionally proved Gödel’s second incompleteness theorem.

11 Conclusion

We presented a formal proof that captures the core of Bachmair and Ganzinger’s *Handbook* chapter on resolution theorem proving. For all its idiosyncrasies, the chapter withstood the test of formalization, once we had added self-inferences to the RP prover. Given that the text is a basic building block of automated reasoning (as confirmed by its placement as Chapter 2), we believe there is value in clarifying its mathematical content for the next generations of researchers. We hope that our work will be useful to the editors of a future revision of the *Handbook*.

Formalization of the metatheory of logical calculi is one of the many connections between automatic and interactive theorem proving. We expect to see wider adoption of proof assistants by researchers in automated reasoning, as a convenient way to develop metatheory. By building formal libraries of standard results, we aim to make it easier to formalize state-of-the-art research as it emerges. We also see potential uses of formal proofs in teaching automated reasoning, inspired by the use of proof assistants in courses on the semantics of programming languages [18, 24].

Acknowledgment. Christoph Weidenbach discussed Bachmair and Ganzinger’s chapter with us on many occasions and generously hosted Schlichtkrull at the Max-Planck-Institut in Saarbrücken. Christian Sternagel and René Thiemann answered our questions about IsaFoR. Mathias Fleury, Florian Haftmann, and Tobias Nipkow helped enrich and reorganize Isabelle’s library of multisets. Mathias Fleury, Robert Lewis, Mark Summerfield, and Sophie Turret suggested many textual improvements.

Blanchette was partly supported by the Deutsche Forschungsgemeinschaft (DFG) project Hardening the Hammer (grant NI 491/14-1). He also received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 713999, Matryoshka). Traytel was partly supported by the DFG program Program and Model Analysis (PUMA, doctorate program 1480).

References

- [1] Bachmair, L., Dershowitz, N., Plaisted, D.A.: Completion without failure. In: Ait-Kaci, H., Nivat, M. (eds.) *Rewriting Techniques—Resolution of Equations in Algebraic Structures*, vol. 2, pp. 1–30. Academic Press (1989)

- [2] Bachmair, L., Ganzinger, H.: Rewrite-based equational theorem proving with selection and simplification. *J. Log. Comput.* 4(3), 217–247 (1994)
- [3] Bachmair, L., Ganzinger, H.: Resolution theorem proving. In: Robinson, A., Voronkov, A. (eds.) *Handbook of Automated Reasoning*, vol. I, pp. 19–99. Elsevier and MIT Press (2001)
- [4] Ballarín, C.: Locales: A module system for mathematical theories. *J. Autom. Reason.* 52(2), 123–153 (2014)
- [5] Biendarra, J., Blanchette, J.C., Bouzy, A., Desharnais, M., Fleury, M., Hölzl, J., Kuncar, O., Lochbihler, A., Meier, F., Panny, L., Popescu, A., Sternagel, C., Thiemann, R., Traytel, D.: Foundational (co)datatypes and (co)recursion for higher-order logic. In: Dixon, C., Finger, M. (eds.) *FroCoS 2017. LNCS*, vol. 10483, pp. 3–21. Springer (2017)
- [6] Blanchette, J.C., Fleury, M., Traytel, D.: Nested multisets, hereditary multisets, and syntactic ordinals in Isabelle/HOL. In: Miller, D. (ed.) *FSCD 2017. LIPIcs*, vol. 84, pp. 11:1–11:18. Schloss Dagstuhl—Leibniz-Zentrum für Informatik (2017)
- [7] Blanchette, J.C., Fleury, M., Weidenbach, C.: A verified SAT solver framework with learn, forget, restart, and incrementality. In: Olivetti, N., Tiwari, A. (eds.) *IJCAR 2016. LNCS*, vol. 9706, pp. 25–44. Springer (2016)
- [8] Blanchette, J.C., Kaliszyk, C., Paulson, L.C., Urban, J.: Hammering towards QED. *J. Formaliz. Reason.* 9(1), 101–148 (2016)
- [9] Blanchette, J.C., Popescu, A., Traytel, D.: Soundness and completeness proofs by inductive methods. *J. Autom. Reason.* 58(1), 149–179 (2017)
- [10] Brand, D.: Proving theorems with the modification method. *SIAM J. Comput.* 4(4), 412–430 (1975)
- [11] Cruanes, S.: Logtk: A logic toolkit for automated reasoning and its implementation. In: Schulz, S., de Moura, L., Konev, B. (eds.) *PAAR-2014. EPiC Series in Computing*, vol. 31, pp. 39–49. EasyChair (2014)
- [12] Denzinger, J., Kronenburg, M., Schulz, S.: DISCOUNT—A distributed and learning equational prover. *J. Autom. Reason.* 18(2), 189–198 (1997)
- [13] Fleury, M., Blanchette, J.C., Lammich, P.: A verified SAT solver with watched literals using Imperative HOL. In: Andronick, J., Felty, A.P. (eds.) *CPP 2018*. pp. 158–171. ACM (2018)
- [14] Hirokawa, N., Middeldorp, A., Sternagel, C., Winkler, S.: Infinite runs in abstract completion. In: Miller, D. (ed.) *FSCD 2017. LIPIcs*, vol. 84, pp. 19:1–19:16. Schloss Dagstuhl—Leibniz-Zentrum für Informatik (2017)
- [15] Krauss, A.: Partial recursive functions in higher-order logic. In: Furbach, U., Shankar, N. (eds.) *IJCAR 2006. LNCS*, vol. 4130, pp. 589–603. Springer (2006)
- [16] McCune, W.: OTTER 2.0. In: Stickel, M.E. (ed.) *10th International Conference on Automated Deduction, Kaiserslautern, FRG, July 24-27, 1990, Proceedings. LNCS*, vol. 449, pp. 663–664. Springer (1990), https://doi.org/10.1007/3-540-52885-7_131
- [17] Nieuwenhuis, R., Rubio, A.: Paramodulation-based theorem proving. In: Robinson, A., Voronkov, A. (eds.) *Handbook of Automated Reasoning*, vol. I, pp. 371–443. Elsevier and MIT Press (2001)
- [18] Nipkow, T.: Teaching semantics with a proof assistant: No more LSD trip proofs. In: Kunčák, V., Rybalchenko, A. (eds.) *VMCAI 2012. LNCS*, vol. 7148, pp. 24–38. Springer (2012)
- [19] Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL: A Proof Assistant for Higher-Order Logic. *LNCS*, vol. 2283. Springer (2002)
- [20] O’Connor, R.: Essential incompleteness of arithmetic verified by Coq. In: Hurd, J., Melham, T.F. (eds.) *TPHOLs 2005. LNCS*, vol. 3603, pp. 245–260. Springer (2005)
- [21] Paulson, L.C.: A machine-assisted proof of Gödel’s incompleteness theorems for the theory of hereditarily finite sets. *Rew. Symb. Logic* 7(3), 484–498 (2014)
- [22] Peltier, N.: A variant of the superposition calculus. *Archive of Formal Proofs 2016* (2016), <https://www.isa-afp.org/entries/SuperCalc.shtml>

- [23] Persson, H.: Constructive completeness of intuitionistic predicate logic—A formalisation in type theory (1996)
- [24] Pierce, B.C.: Lambda, the ultimate TA: Using a proof assistant to teach programming language foundations. In: Hutton, G., Tolmach, A.P. (eds.) ICFP 2009. pp. 121–122. ACM (2009)
- [25] Reger, G., Suda, M.: Checkable proofs for first-order theorem proving. In: Reger, G., Traytel, D. (eds.) ARCADE 2017. EPiC Series in Computing, vol. 51, pp. 55–63. EasyChair (2017)
- [26] Schlichtkrull, A.: Formalization of the resolution calculus for first-order logic. In: Blanchette, J.C., Merz, S. (eds.) ITP 2016. LNCS, vol. 9807, pp. 341–357. Springer (2016)
- [27] Shankar, N.: Towards mechanical metamathematics. *J. Autom. Reason.* 1(4), 407–434 (1985)
- [28] Shankar, N.: *Metamathematics, Machines, and Gödel’s Proof*, Cambridge Tracts in Theoretical Computer Science, vol. 38. Cambridge University Press (1994)
- [29] Sutcliffe, G., Zimmer, J., Schulz, S.: TSTP data-exchange formats for automated theorem proving tools. In: Zhang, W., Sorge, V. (eds.) Distributed Constraint Problem Solving and Reasoning in Multi-Agent Systems. *Frontiers in Artificial Intelligence and Applications*, vol. 112, pp. 201–215. IOS Press (2004)
- [30] Thiemann, R., Sternagel, C.: Certification of termination proofs using CeTA. In: Berghofer, S., Nipkow, T., Urban, C., Wenzel, M. (eds.) TPHOLs 2009. LNCS, vol. 5674, pp. 452–468. Springer (2009)
- [31] Voronkov, A.: AVATAR: The architecture for first-order theorem provers. In: Biere, A., Bloem, R. (eds.) CAV 2014. LNCS, vol. 8559, pp. 696–710. Springer (2014)
- [32] Wand, D.: Polymorphic+typeclass superposition. In: Schulz, S., de Moura, L., Konev, B. (eds.) PAAR-2014. EPiC Series in Computing, vol. 31, pp. 105–119. EasyChair (2014)
- [33] Weidenbach, C.: Combining superposition, sorts and splitting. In: Robinson, A., Voronkov, A. (eds.) *Handbook of Automated Reasoning*, vol. II, pp. 1965–2013. Elsevier and MIT Press (2001)
- [34] Wenzel, M.: Isabelle/Isar—A generic framework for human-readable proof documents. In: Matuszewski, R., Zalewska, A. (eds.) *From Insight to Proof: Festschrift in Honour of Andrzej Trybulec*, *Studies in Logic, Grammar, and Rhetoric*, vol. 10(23). University of Białystok (2007)
- [35] Wenzel, M.: Isabelle/jEdit—A prover IDE within the PIDE framework. In: Jeuring, J., Campbell, J.A., Carette, J., Reis, G.D., Sojka, P., Wenzel, M., Sorge, V. (eds.) CICM 2012. LNCS, vol. 7362, pp. 468–471. Springer (2012)
- [36] Zhang, H., Kapur, D.: First-order theorem proving using conditional rewrite rules. In: Lusk, E.L., Overbeek, R.A. (eds.) CADE-9. LNCS, vol. 310, pp. 1–20. Springer (1988)