



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, 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,
Zurich, 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 paper, we present a formalization, developed using the Isabelle/HOL system [16], of a first-order prover based on ordered resolution with literal selection. We follow Bachmair and Ganzinger’s account [3] from Chap. 2 of the *Handbook of Automated Reasoning*, which we will simply refer to as “the chapter.” Our formal development covers the refutationally 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 integrated 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* Chaps. 7 [14] and 27 [26]). 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. 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 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. The Isabelle 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 names. A better way to study the theory files, however, is to open them in Isabelle/jEdit [28]. We used Isabelle version 2017, but the AFP is continuously updated to track Isabelle’s evolution. Due to lack of space, we assume the reader has some familiarity with the chapter’s content. An extended version of this paper is available as a technical report [21].

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

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

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 use the same library of clauses ([Clausal_Logic.thy](#)) as for the verified SAT solver by Blanchette et al. [6], which is also part of IsaFoL. Atoms are represented by a type variable $'a$, which can be instantiated by arbitrary concrete 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 $>$ extends a fixed atom order $>$ by comparing polarities to break ties, with $\text{Neg } A > \text{Pos } A$. 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 $\{\}$. The clause order is the multiset extension of the literal order.

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 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 [7, Sect. 3].

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 [23]. 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., $A \cdot \text{id} = A$). 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, $\text{is_ground } A \Leftrightarrow \forall \sigma. A = A \cdot \sigma$.

To complete our development and ensure that our assumptions are legitimate, we instantiate the locale’s parameters with `IsaFoR` types and operations and

discharge its assumptions ([IsaFoR_Term.thy](#)). This bridge is currently hosted outside the AFP.

3 Refutational Inference Systems

In their Sect. 2.4, Bachmair and Ganzinger introduce basic conventions for refutational inference systems. In Sect. 3, they present two ground resolution calculi and prove them refutationally complete in Theorems 3.9 and 3.16. In Sect. 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 the goal of providing a well-designed, reusable library.

We collect the abstract hierarchy of inference systems in a single Isabelle theory file ([Inference.System.thy](#)). An inference, of type *'a inference*, is a triple (\mathcal{C}, D, E) that consists of a multiset of side premises \mathcal{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 (*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 the inference system locale with an assumption, whereas sound systems are characterized by a different assumption:

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

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:

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 *l_of* :: 'a clause set \Rightarrow 'a set
assumes $\{\} \notin N \Rightarrow D \in N \Rightarrow \text{l_of } N \not\models D \Rightarrow$
 $(\forall C \in N. \text{l_of } N \not\models C \Rightarrow D \leq C) \Rightarrow$
 $\exists C \subseteq N. \exists E. \text{l_of } N \models C \wedge (C, D, E) \in \Gamma \wedge \text{l_of } N \not\models E \wedge E < D$

The “model functor” parameter *l_of* maps clause sets to candidate models. The assumption is that for any set N that does not contain $\{\}$ (i.e., \perp), if $D \in N$ is the smallest counterexample—the smallest clause in N falsified by $\text{l_of } N$ —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), we must have $E \in N$, violating D 's minimality:

theorem *saturated_model*: $\text{saturated } N \Rightarrow \{\} \notin N \Rightarrow \text{l_of } N \models N$
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*: $\neg \text{sat } N \Leftrightarrow \exists M \subseteq N. \text{finite } M \wedge \neg \text{sat } M$

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 Sect. 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 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 [29].

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 $S :: 'a \text{ clause} \Rightarrow 'a \text{ clause}$. We also fix a clause set N . Then we define two operators corresponding to ε_C and I_C :

function `production` $:: 'a \text{ clause} \Rightarrow 'a \text{ set}$ **where**
`production` $C = \{A \mid C \in N \wedge C \neq \{\}\} \wedge \text{Max } C = \text{Pos } A$
 $\wedge \left(\bigcup_{D < C} \text{production } D \right) \not\equiv C \wedge S \ C = \{\}$
definition `interp` $:: 'a \text{ clause} \Rightarrow 'a \text{ set}$ **where**
`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. Bachmair and Ganzinger's I^C and I_N operators are introduced as abbreviations: `Interp` $C = \text{interp } C \cup \text{production } C$ and `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, including these:

lemma `interp_imp_Interp`:

$$C \leq D \Rightarrow D \leq D' \Rightarrow \text{interp } D \models C \Rightarrow \text{Interp } D' \models C$$

lemma `Interp_imp_INTERP`: $C \leq D \Rightarrow \text{Interp } D \models C \Rightarrow \text{INTERP} \models C$

Lemma 3.3, whose proof depends on monotonicity, is better proved *after* 3.4:

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

A more serious oddity is Lemma 3.7. Using our notations, it can be stated as $D \in N \Rightarrow C \neq D \Rightarrow (\forall D' < D. \text{Interp } D' \models C) \Rightarrow \text{interp } D \models D'$. 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 $C \vee D$. Formally, this rule is captured by a predicate:

inductive `unord_resolve` $:: 'a \text{ clause} \Rightarrow 'a \text{ clause} \Rightarrow 'a \text{ clause} \Rightarrow \text{bool}$ **where**
`unord_resolve` $(C \uplus \text{replicate } (n + 1) (\text{Pos } A)) (\{\text{Neg } A\} \uplus D) (C \uplus D)$

To prove completeness, it suffices to show that the calculus reduces counterexamples (Theorem 3.8). 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, `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 $A_i = A_i \vee \dots \vee A_i$ parts separately. In addition, As is the list $[A_1, \dots, A_n]$. The following inductive definition captures the rule formally:

inductive `ord_resolve` :: ‘*a clause list* \Rightarrow ‘*a clause* \Rightarrow ‘*a clause* \Rightarrow *bool*

where

$$\begin{aligned} & |CAs| = n \Rightarrow |Cs| = n \Rightarrow |As| = n \Rightarrow |As| = n \Rightarrow n \neq 0 \Rightarrow \\ & (\forall i < n. CAs ! i = Cs ! i \uplus \text{Pos } 'As ! i) \Rightarrow (\forall i < n. As ! i \neq \{\}) \Rightarrow \\ & (\forall i < n. \forall A \in As ! i. A = As ! i) \Rightarrow \text{eligible } As (D \uplus \text{Neg } 'mset As) \Rightarrow \\ & (\forall i < n. \text{strict_max_in } (As ! i) (Cs ! i)) \Rightarrow (\forall i < n. S (CAs ! i) = \{\}) \Rightarrow \\ & \text{ord_resolve } CAs (D \uplus \text{Neg } 'mset As) ((\bigcup \text{mset } Cs) \uplus D) \end{aligned}$$

The $xs ! i$ operator returns the $(i + 1)$ st element of xs , and `mset` converts a list to a multiset. Initially, 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”; the alternative, a single multiset of tuples, is very unwieldy.

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 technical report [21, Appendix A].

5 Theorem Proving Processes

In their Sect. 4, Bachmair and Ganzinger switch 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 no inferences are possible. Redundant clauses can be deleted at any point, and redundant inferences need not be performed.

A derivation performed by a proving process is a possibly infinite sequence $N_0 \triangleright N_1 \triangleright N_2 \triangleright \dots$, where \triangleright relates 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* \Rightarrow ‘*a llist* \Rightarrow ‘*a llist*, and equipped with `lhd` (“head”)

and `ltl` (“tail”) selectors that extract `LCons`’s arguments. The coinductive predicate `chain` checks that its argument is a nonempty lazy list whose elements are linked by a binary predicate `R`:

```
coinductive chain :: ('a ⇒ 'a ⇒ bool) ⇒ 'a llist ⇒ 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, $\mathcal{R}_{\mathcal{F}}$ and $\mathcal{R}_{\mathcal{I}}$:

```
locale redundancy_criterion = inference_system +
fixes  $\mathcal{R}_{\mathcal{F}} :: 'a \text{ clause set} \Rightarrow 'a \text{ clause set}$  and
   $\mathcal{R}_{\mathcal{I}} :: 'a \text{ clause set} \Rightarrow 'a \text{ inference set}$ 
assumes  $\mathcal{R}_{\mathcal{I}} N \subseteq \Gamma$  and  $\text{sat } (N \setminus \mathcal{R}_{\mathcal{F}} N) \Rightarrow \text{sat } N$  and
   $N \subseteq N' \Rightarrow \mathcal{R}_{\mathcal{F}} N \subseteq \mathcal{R}_{\mathcal{F}} N' \wedge \mathcal{R}_{\mathcal{I}} N \subseteq \mathcal{R}_{\mathcal{I}} N'$  and
   $N' \subseteq \mathcal{R}_{\mathcal{F}} N \Rightarrow \mathcal{R}_{\mathcal{F}} N \subseteq \mathcal{R}_{\mathcal{F}} (N \setminus N') \wedge \mathcal{R}_{\mathcal{I}} N \subseteq \mathcal{R}_{\mathcal{I}} (N \setminus N')$ 
```

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 :

$$M \triangleright N \Leftrightarrow N \setminus M \subseteq \text{concl_of } ' \text{infers_from } M \wedge M \setminus N \subseteq \mathcal{R}_{\mathcal{F}} N$$

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 Sect. 7.

A key concept to connect static and dynamic completeness is that of the set of persistent clauses, or limit: $N_{\infty} = \bigcup_i \bigcap_{j \geq i} N_j$. 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, $\bigcup_i N_i$. We introduce these missing functions ([Lazy_List.Liminf.thy](#)):

$$\text{Liminf } xs = \bigcup_{i < |xs|} \bigcap_{j: i \leq j < |xs|} xs ! j \quad \text{Sup } xs = \bigcup_{i < |xs|} xs ! i$$

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 i , the danger is 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_Sup: chain ( $\triangleright$ )  $Ns \Rightarrow \mathcal{R}_{\mathcal{F}} (\text{Liminf } Ns) = \mathcal{R}_{\mathcal{F}} (\text{Sup } Ns)$ 
lemma Ri_limit_Sup: chain ( $\triangleright$ )  $Ns \Rightarrow \mathcal{R}_{\mathcal{I}} (\text{Liminf } Ns) = \mathcal{R}_{\mathcal{I}} (\text{Sup } Ns)$ 
lemma sat_limit_iff: chain ( $\triangleright$ )  $Ns \Rightarrow (\text{sat } (\text{Liminf } Ns) \Leftrightarrow \text{sat } (\text{lhd } Ns))$ 
```


In the chapter, the proof relies on “the soundness of the inference system,” contradicting the claim that “we will only consider consistency-preserving inference systems” [3, Sect. 2.4]. Thanks to Isabelle, we now know that soundness is unnecessary.

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:

$$\text{saturated_upto } N \Leftrightarrow \text{infers_from } (N \setminus \mathcal{R}_{\mathcal{F}} N) \subseteq \mathcal{R}_{\mathcal{I}} N$$

The limit is saturated for fair derivations, defined by $\text{fair_class_seq } Ns \Leftrightarrow \text{concl_of } \text{‘infers_from } N' \setminus \mathcal{R}_{\mathcal{I}} N' \subseteq \text{Sup } Ns \cup \mathcal{R}_{\mathcal{F}} (\text{Sup } Ns)$ with $N' = \text{Liminf } Ns \setminus \mathcal{R}_{\mathcal{F}} (\text{Liminf } Ns)$. The criterion must also be effective, meaning $\gamma \in \Gamma \Rightarrow \text{concl_of } \gamma \in N \cup \mathcal{R}_{\mathcal{F}} N \Rightarrow \gamma \in \mathcal{R}_{\mathcal{I}} N$. Under these assumptions, we have Theorem 4.3:

theorem *fair_derive_saturated_upto:*

$$\text{chain } (\triangleright) Ns \Rightarrow \text{fair_class_seq } Ns \Rightarrow \text{saturated_upto } (\text{Liminf } Ns)$$

The standard redundancy criterion is an instance of the framework. It relies on a counterexample-reducing inference system Γ ([Standard_Redundancy.thy](#)):

$$\mathcal{R}_{\mathcal{F}} N = \{C \mid \exists \mathcal{D} \subseteq N. (\forall I. I \models \mathcal{D} \Rightarrow I \models C) \wedge \forall D' \in \mathcal{D}. D' < C\}$$

$$\mathcal{R}_{\mathcal{I}} N = \{(C, D, E) \in \Gamma \mid \exists \mathcal{D} \subseteq N. (\forall I. I \models \mathcal{D} \uplus C \Rightarrow I \models E) \wedge \forall D' \in \mathcal{D}. D' < D\}$$

Standard redundancy qualifies as *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 result, 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 a lemma and apply it to prove Lemmas 4.5 and 4.6. Incidentally, the informal proof of Theorem 4.9 also needlessly invokes Lemma 4.5.

Finally, given a redundancy criterion $(\mathcal{R}_{\mathcal{F}}, \mathcal{R}_{\mathcal{I}})$ for Γ , its standard extension for $\Gamma' \supseteq \Gamma$ is defined as $(\mathcal{R}_{\mathcal{F}}, \mathcal{R}'_{\mathcal{I}})$, where $\mathcal{R}'_{\mathcal{I}} N = \mathcal{R}_{\mathcal{I}} N \cup (\Gamma' \setminus \Gamma)$ ([Proving_Process.thy](#)). The standard extension preserves effectiveness, saturation up to redundancy, and fairness.

6 First-Order Resolution

The chapter’s Sect. 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 ([FO_Ordered_Resolution.thy](#)).

Inference Rule. First-order ordered resolution consists of the single rule

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

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

$$\begin{aligned} &|CAs| = n \Rightarrow |Cs| = n \Rightarrow |As| = n \Rightarrow |As| = n \Rightarrow n \neq 0 \Rightarrow \\ &(\forall i < n. CAs!i = Cs!i \uplus \text{Pos} \cdot As!i) \Rightarrow (\forall i < n. As!i \neq \{\}) \Rightarrow \\ &\text{Some } \sigma = \text{mgu} (\text{set_mset} \cdot \text{set} (\text{map2} \text{add_mset} As As)) \Rightarrow \\ &\text{eligible } \sigma As (D \uplus \text{Neg} \cdot \text{mset} As) \Rightarrow \\ &(\forall i < n. \text{strict_max_in} (As!i \cdot \sigma) (Cs!i \cdot \sigma)) \Rightarrow (\forall i < n. S (CAs!i) = \{\}) \Rightarrow \\ &\text{ord_resolve } CAs (D \uplus \text{Neg} \cdot \text{mset} As) \sigma (((\bigcup \text{mset } Cs) \uplus D) \cdot \sigma) \end{aligned}$$

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. 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. 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 Sect. 2) and is implemented using `IsaFoR`.

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 $S_M_grounding_of_cls$:

$$C \in \text{grounding_of } M \Rightarrow \exists D \in M. \exists \sigma. C = D \cdot \sigma \wedge S_M C = S D \cdot \sigma$$

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 (possibly) more general inference from clauses belonging to M :

lemma *ord_resolve_rename_lifting*:

$$\begin{aligned}
& (\forall \rho \ C. \text{is_renaming } \rho \implies S (C \cdot \rho) = S C \cdot \rho) \implies \\
& \text{ord_resolve } S_M \ CAs \ DA \ As \ As \ \sigma \ E \implies \\
& \{DA\} \cup \text{set } CAs \subseteq \text{grounding_of } M \implies \\
& \exists \eta s \ \eta \ \theta \ CAs_0 \ DA_0 \ As_0 \ As_0 \ E_0 \ \tau. \\
& \text{ord_resolve_rename } S \ CAs_0 \ DA_0 \ As_0 \ As_0 \ \tau \ E_0 \wedge \\
& CAs_0 \cdot \eta s = CAs \wedge DA_0 \cdot \eta = DA \wedge E_0 \cdot \theta = E \wedge \{DA_0\} \cup \text{set } CAs_0 \subseteq M
\end{aligned}$$

The informal proof of this lemma consists of two sentences spanning four lines of text. In Isabelle, these two sentences translate to 250 lines and 400 lines, respectively, excluding auxiliary lemmas. Our proof involves six steps:

1. Obtain a list of first-order clauses CAs_0 and a first-order clause DA_0 that belong to M and that generalize CAs and DA with substitutions ηs and η , respectively.
2. Choose atoms As_0 and As_0 in the first-order clauses on which to resolve.
3. Standardize CAs_0 and DA_0 apart, yielding CAs'_0 and DA'_0 .
4. Obtain the MGU τ of the literals on which to resolve.
5. Show that ordered resolution on CAs'_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 $S (CAs_0 ! i)$ generalizes $S_M (CAs ! i)$, for $0 \leq i < n$, and $S DA_0$ generalizes $S_M DA$. By the 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 in principle 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 CAs_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 CAs_0 , DA_0 , E_0 to CAs , DA , E . By contrast, we provide separate substitutions ηs , 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 Sect. 4.3, Bachmair and Ganzinger introduce the nondeterministic abstract prover RP that works on triples of clause sets and that generalizes the Otter-style and DISCOUNT-style loops. 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 ([FO_Ordered_Resolution_Prover.thy](#)).

Abstract First-Order Prover. The RP prover is a relation \sim 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 $\sim :: 'a \text{ state} \Rightarrow 'a \text{ state} \Rightarrow \text{bool}$ **where**

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

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:

$$\text{subsumes } D C \Leftrightarrow \exists \sigma. D \cdot \sigma \subseteq C$$

$$\text{strictly_subsumes } D C \Leftrightarrow \text{subsumes } D C \wedge \neg \text{subsumes } C D$$

$$\text{reduces } D C L \Leftrightarrow \exists D' L' \sigma. D = D' \uplus \{L'\} \wedge -L = L' \cdot \sigma \wedge D' \cdot \sigma \subseteq C$$

The definition of the set `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, the resolution inference

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

is possible. Yet if the clause $\neg \text{p} \vee \neg \text{p}$ 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

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

and an order $>$ on atoms such that $\mathbf{q}(c, b, a) > \mathbf{q}(b, a, c) > \mathbf{q}(a, c, b)$. 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):

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

From the conclusion, we derive $\neg \mathbf{q}(a, c, b)$ by (3) and \perp by (1). This incompleteness is a severe flaw, although it is probably just an oversight.

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 Sect. 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 $(\mathcal{R}_{\mathcal{F}}, \mathcal{R}_{\mathcal{I}})$.

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 (Sect. 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 \mathcal{S} and \mathcal{S}' , so that transitions can be joined to form a single theorem proving derivation. Our approach is to let \triangleright be the extension 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), 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 (Sect. 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 for single steps and extend it to entire derivations:

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

lemma *RP_ground_derive_chain*:

chain $(\rightsquigarrow) \mathcal{S}s \implies \text{chain } (\triangleright) (\text{lmap grounding_of } \mathcal{S}s)$

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 Sect. 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 $\mathcal{S}s$, and its projections $\mathcal{N}s$, $\mathcal{P}s$, and $\mathcal{O}s$, such that `chain` $(\rightsquigarrow) \mathcal{S}s$, `fair_state_seq` $\mathcal{S}s$, and `lhd` $\mathcal{O}s = \{\}$.

Thanks to fairness, any nonredundant clause C in $\mathcal{S}s$ ’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` [13, Sect. 2]. By minimality, D' cannot be deleted by backward subsumption. This line of reasoning allows us to prove Lemma 4.11, where `O_of` extracts the \mathcal{O} component of a state:

lemma *fair_imp_Liminf_minus_Rf_subset_ground_Liminf_state*:
 $G_s = \text{Imap_grounding_of } \mathcal{S}_s \implies$
 $\text{Liminf } G_s - \mathcal{R}_{\mathcal{F}}(\text{Liminf } G_s) \subseteq \text{grounding_of } (\mathcal{O}.\text{of } (\text{Liminf } \mathcal{S}_s))$

Completeness. Once we have brought Lemmas 4.10, 4.11, and 4.12 into a suitable shape, the main completeness result, Theorem 4.13, 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*:
 $\neg \text{sat } (\text{grounding_of } (\text{lhd } \mathcal{S}_s)) \implies \{\} \in \mathcal{O}.\text{of } (\text{Liminf } \mathcal{S}_s)$

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 “Liminf G_s (and hence Liminf \mathcal{S}_s) contains the empty clause” (using our notations). Obviously, if $\text{grounding_of } (\text{Liminf } \mathcal{S}_s)$ contains \perp , then Liminf \mathcal{S}_s must as well. However, the authors do not explain the step from Liminf G_s , 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.

8 Discussion and Related Work

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 Sect. 4.3, which introduces the RP prover. 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 Sect. 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 into 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.

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 [12]. 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. Such output can be independently 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—based on classical simple type theory—balances expressiveness and automatability. We benefited from many features of the system, including codatatypes [5], Isabelle/jEdit [28], the Isar proof language [27], 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.

Formalizing the metatheory of logic and deduction is an enticing proposition for many researchers. Two recent, independent developments are particularly pertinent. Peltier [17] proved static completeness of a variant of the superposition calculus in Isabelle/HOL. Since superposition generalizes ordered resolution, his result subsumes our static completeness theorem. 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, 25], type classes [25], and AVATAR [24]. Hirokawa et al. [13] formalized, also in Isabelle/HOL, an abstract Knuth–Bendix completion procedure as well as ordered (unfailing) completion [1]. Superposition combines ordered resolution (to reason about clauses) and ordered completion (to reason about equality).

The literature contains many other formalized completeness proofs. Early work was carried out by Shankar [22] and Persson [18]. Some of our own efforts are also related: completeness of unordered resolution using semantic trees by Schlichtkrull [20]; completeness of a Gentzen system by Blanchette, Popescu, and Traytel [9]; and completeness of CDCL by Blanchette, Fleury, Lammich, and Weidenbach [6]. We refer to our earlier papers for further discussions of related work.

9 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, 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 [15, 19].

Acknowledgment. Christoph Weidenbach discussed Bachmair and Ganzinger's chapter with us on many occasions and 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 multiset library. Mathias Fleury, Robert Lewis, Mark Summerfield, Sophie Tourret, and the anonymous reviewers 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: Aït-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., et al.: Foundational (co)datatypes and (co)recursion for higher-order logic. In: Dixon, C., Finger, M. (eds.) *FroCoS 2017. LNCS (LNAI)*, vol. 10483, pp. 3–21. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66167-4_1
6. Blanchette, J.C., Fleury, M., Lammich, P., Weidenbach, C.: A verified SAT solver framework with learn, forget, restart, and incrementality. *J. Autom. Reason.* **61**(3), 333–366
7. 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)

8. Blanchette, J.C., Kaliszky, C., Paulson, L.C., Urban, J.: Hammering towards QED. *J. Formal. Reason.* **9**(1), 101–148 (2016)
9. Blanchette, J.C., Popescu, A., Traytel, D.: Soundness and completeness proofs by coinductive 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. 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)
13. 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)
14. 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)
15. Nipkow, T.: Teaching Semantics with a proof assistant: no more LSD trip proofs. In: Kuncak, V., Rybalchenko, A. (eds.) VMCAI 2012. LNCS, vol. 7148, pp. 24–38. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-27940-9_3
16. Nipkow, T., Wenzel, M., Paulson, L.C. (eds.): Isabelle/HOL: A Proof Assistant for Higher-Order Logic. LNCS, vol. 2283. Springer, Heidelberg (2002). <https://doi.org/10.1007/3-540-45949-9>
17. Peltier, N.: A variant of the superposition calculus. *Archive of Formal Proofs* 2016 (2016). <https://www.isa-afp.org/entries/SuperCalc.shtml>
18. Persson, H.: Constructive Completeness of Intuitionistic Predicate Logic—a Formalisation in Type Theory. Licentiate thesis, Chalmers tekniska högskola and Göteborgs universitet (1996)
19. 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)
20. Schlichtkrull, A.: Formalization of the resolution calculus for first-order logic. *J. Autom. Reason.* **61**(4), 455–484
21. Schlichtkrull, A., Blanchette, J.C., Traytel, D., Waldmann, U.: Formalizing Bachmair and Ganzinger’s ordered resolution prover (technical report). Technical report (2018). http://matryoshka.gforge.inria.fr/pubs/rp_report.pdf
22. Shankar, N.: Towards mechanical metamathematics. *J. Autom. Reason.* **1**(4), 407–434 (1985)
23. 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, Heidelberg (2009). https://doi.org/10.1007/978-3-642-03359-9_31
24. 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, Cham (2014). https://doi.org/10.1007/978-3-319-08867-9_46
25. 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)

26. 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)
27. 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, no. 23, University of Białystok (2007)
28. Wenzel, M.: Isabelle/jEdit—a prover IDE within the PIDE framework. In: Jeuring, J., Campbell, J.A., Carette, J., Dos Reis, G., Sojka, P., Wenzel, M., Sorge, V. (eds.) *CICM 2012. LNCS (LNAI)*, vol. 7362, pp. 468–471. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-31374-5_38
29. Zhang, H., Kapur, D.: First-order theorem proving using conditional rewrite rules. In: Lusk, E., Overbeek, R. (eds.) *CADE 1988. LNCS*, vol. 310, pp. 1–20. Springer, Heidelberg (1988). <https://doi.org/10.1007/BFb0012820>