

# A Modular Isabelle Framework for Verifying Saturation Provers

Sophie Touret

Research Group 1: Automation of Logic  
Max-Planck-Institut für Informatik  
Saarbrücken, Germany  
stouret@mpi-inf.mpg.de

Jasmin Blanchette

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

## Abstract

We present a formalization in Isabelle/HOL of a comprehensive framework for proving the completeness of automatic theorem provers based on resolution, superposition, or other saturation calculi. The framework helps calculus designers and prover developers derive, from the completeness of a calculus, the completeness of prover architectures implementing the calculus. It also helps derive the completeness of calculi obtained by lifting ground (i.e., variable-free) calculi. As a case study, we re-verified Bachmair and Ganzinger’s resolution prover RP to show the benefits of modularity.

## 1 Introduction

Many of the most successful automatic theorem provers today are based on saturation. These include the unit equality prover Waldmeister [18], the first-order provers E [26], SPASS [33], and Vampire [20], and the higher-order provers Leo-III [27] and Zipperposition [7]. A saturation prover starts with a problem, typically given in conjunctive normal form, and draws inferences from the problem clauses, adding the conclusions to the clause set. If the prover detects useless clauses, it may remove them. A refutation proof has been found when the prover has derived the empty clause, denoted by  $\perp$ . Which inferences are drawn and which clauses are deemed redundant depends on the logical calculus used.

Although users mostly care only about the soundness and success rate of a prover, refutational completeness is taken very seriously by researchers. Beyond its intrinsic merits, completeness is useful to optimize a calculus: From a completeness proof, calculi designers can extract precise side conditions on the inference rules, which help prune the search space. Without the guidance offered by a completeness proof,

it is easy to prune too little, too much, or both. Moreover, for decidable fragments, completeness is needed to ensure that the prover can be relied on as a decision procedure.

Completeness can be formulated in two ways:

- *Static completeness* is formulated in terms of a “saturated” clause set  $N$ —a set from which all necessary inferences have been drawn.
- *Dynamic completeness* is formulated in terms of a sequence  $(N_i)_i$  of clause sets, where each set  $N_i$  represents the prover’s state at a point in time: If  $N_0$  is unsatisfiable (i.e., the conjecture holds), then  $\perp \in N_i$  for some  $i$ .

Most saturation calculi are designed as a lifting of a ground calculus—that is, a calculus that operates on variable-free clauses. Static completeness is established on the ground level and lifted along with the calculus. This separation of concerns helps keep these technical proofs manageable.

For dynamic completeness, we distinguish between simple proving processes that operate on a single clause set and realistic provers with multiple clause sets. The *given clause procedure* [21], which is part of virtually all saturation provers, needs two clause sets, called “passive” and “active.” Clauses are initially passive and move one by one to the active set, at which point all possible inferences with active partners are drawn. This setup is reminiscent of traditional New Year celebrations in Quebec: When you arrive at the party, you must kiss all the relatives lined up before you and then join the line. In this way, everybody kisses everybody.

A crucial optimization in provers is clause subsumption. A clause  $D$ , seen as a multiset of literals, *subsumes* another clause  $C$  if  $D\sigma \subseteq C$  for some substitution  $\sigma$ ; for example,  $p(x)$  subsumes  $p(a)$  and  $p(x) \vee q(y)$ . Regrettably, much of the literature on saturation ignores subsumption and dynamic completeness. In the refinement chain leading to a verified prover, such a simplifying assumption translates into a proof gap.

Bachmair and Ganzinger’s chapter in the *Handbook of Automated Reasoning* [4] is an exception. They define a three-clause-set resolution prover, called RP, that supports clause subsumption, and they prove it dynamically complete. But their argument is highly nonmodular: They move directly from static completeness of the ground calculus to dynamic completeness of RP, awkwardly mixing issues related to

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference’17, July 2017, Washington, DC, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM. . \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

nonground lifting, subsumption, and the given clause procedure. The argument is hard to follow on paper, and even harder to formalize, as Schlichtkrull, Blanchette, Traytel, and Waldmann found out [24].

To remedy this situation, Waldmann, Tourret, Robillard, and Blanchette [32] designed a comprehensive framework that applies to a wide range of calculi, including ordered resolution [4], unifying completion [2], standard superposition [3], theory superposition [31], hierarchic superposition [5],  $\lambda$ -free superposition [8],  $\lambda$ -superposition [7], and combinatory superposition [10]. This *saturation framework* can be used to lift ground static completeness to nonground static completeness and then to dynamic completeness of minimalistic proving processes or given clause provers with support for clause subsumption.

The framework was initially developed in  $\text{\LaTeX}$ , but we also managed to formalize it in Isabelle/HOL in time to mention it in Waldmann et al. [32, Section 5]. The present paper aims at introducing this work to the interactive theorem proving community. Besides giving a thorough account of the formalization, we unveil three extensions of the framework:

- We applied the formalized framework to re-prove RP dynamically complete. The new proof is substantially simpler and shorter than that by Schlichtkrull et al.
- We formalized various basic concepts related to theorem proving that were needed to verify RP and that will likely be useful for other provers. Much of this material was generalized from Schlichtkrull et al.
- We formalized alternative, invariance-based completeness proofs for given clause procedures. These new proofs replace previous monolithic arguments and will be used, in pen-and-paper form, in a forthcoming journal submission by Waldmann et al.

We used the declarative Isar [34] in the hope that practitioners who are not Isabelle experts can follow the intermediate proof steps. We also relied extensively on locales [6], which provide parameterized modules. The vast majority of the pen-and-paper arguments survived formalization, but one serious issue was discovered and repaired.

Our work is part of the IsaFoL (Isabelle Formalization of Logic) effort, an unfunded coalition of the willing that aims at developing a library of results about logic and automated reasoning [12]. The Isabelle theory files are available in the *Archive of Formal Proofs* [11, 29]; they amount to around 5700 lines of source text.

This paper's organization largely follows Waldmann et al. Briefly, we introduce the fundamental notions that make up a calculus and related concepts such as saturation and completeness (Section 2); we cover the lifting from the ground to the nonground level and the delicate issue of clause subsumption (Section 3); we present two given clause procedures and our alternative, invariance-based proofs (Section 4); and we present our modular proof of RP's completeness (Section 5).

## 2 Calculi

The saturation framework relies on a number of basic notions that are relevant to prove a calculus complete. These include formulas, consequence relations, inference systems, and redundancy criteria (all defined in `Calculus.thy`, unless indicated otherwise).

### 2.1 Formulas

The entire framework is built on an abstract notion of formulas, generalizing clauses. Formulas are represented by a type variable  $'f$  equipped with a set `Bot` and an entailment relation  $\models$ . `Bot` is a nonempty set of patently false formulas. Typically, we take `Bot` =  $\{\perp\}$ , but it can be useful to differentiate between an active  $\perp$  and a passive  $\perp$ . Entailment must satisfy four basic properties. Formally:

```

locale consequence_relation =
  fixes
    Bot :: 'f set and
     $\models$  :: 'f set  $\Rightarrow$  'f set  $\Rightarrow$  bool
  assumes
    Bot  $\neq$   $\emptyset$  and
     $B \in$  Bot  $\Rightarrow$   $\{B\} \models N_1$  and
     $N_2 \subseteq N_1 \Rightarrow N_1 \models N_2$  and
     $(\forall C \in N_2. N_1 \models \{C\}) \Rightarrow N_1 \models N_2$  and
     $N_1 \models N_2 \Rightarrow N_2 \models N_3 \Rightarrow N_1 \models N_3$ 

```

An Isabelle locale can be seen as a record type, where fields can be types ( $'f$ ), terms (`Bot`,  $\models$ ), or properties. When we later instantiate the locale, we must supply concrete arguments for the types and terms and then discharge the proof obligations corresponding to the properties.

### 2.2 Inference Systems

An inference  $C_n, \dots, C_1 \vdash C_0$  is a recipe to derive a new formula  $C_0$  (the conclusion) from existing formulas  $C_n, \dots, C_1$  (the premises). Formally:

```

datatype 'f inference =
  Infer ('f list) 'f

```

An inference  $\iota$  is written out as  $\iota = \text{Infer } [C_n, \dots, C_1] C_0$ , with `prems_of`  $\iota = [C_n, \dots, C_1]$  and `concl_of`  $\iota = C_0$ .

An inference rule, with its metavariables and side conditions, generally corresponds to a set of inferences. In fact, an entire inference system, consisting of one or more inference rules, is modeled as a set of inferences `Inf`. Formally:

```

locale inference_system =
  fixes Inf :: 'f inference set

```

Inside the locale, we defined the following functions:

```

definition Inf_from :: 'f set  $\Rightarrow$  'f inference set where
  Inf_from  $N = \{\iota \in$  Inf. prems_of  $\iota \subseteq N\}$ 

```

```

definition Inf_between :: 'f set  $\Rightarrow$  'f set  $\Rightarrow$  'f inference set
where
  Inf_between  $N M =$ 
  Inf_from  $(N \cup M) - \text{Inf\_from } (N - M)$ 

```

The auxiliary function  $\text{Inf\_from } N$  returns the set of possible inferences from premises contained in  $N$ . The variant  $\text{Inf\_between } N \ M$  returns the set of possible inferences from premises in  $N \cup M$  but with at least one premise in  $M$ .

### 2.3 Calculi with Redundancy

A calculus is basically an inference system equipped with a redundancy criterion. The inference system indicates which inferences must a priori be drawn. The redundancy criterion specifies inferences that can a posteriori be omitted. It also identifies deletable formulas. Formally:

```

locale calculus =
  inference_system Inf +
  consequence_relation Bot ( $\models$ ) +
fixes
  Redl :: 'f set  $\Rightarrow$  'f inference set and
  RedF :: 'f set  $\Rightarrow$  'f set
assumes
  Redl  $N \subseteq$  Inf and
   $B \in \text{Bot} \Rightarrow N \models \{B\} \Rightarrow N - \text{Red}_F N \models \{B\}$  and
   $N \subseteq N' \Rightarrow \text{Red}_F N \subseteq \text{Red}_F N'$  and
   $N \subseteq N' \Rightarrow \text{Red}_l N \subseteq \text{Red}_l N'$  and
   $N' \subseteq \text{Red}_F N \Rightarrow \text{Red}_F N \subseteq \text{Red}_F (N - N')$  and
   $N' \subseteq \text{Red}_F N \Rightarrow \text{Red}_l N \subseteq \text{Red}_l (N - N')$  and
   $\iota \in \text{Inf} \Rightarrow \text{concl\_of } \iota \in N \Rightarrow \iota \in \text{Red}_l N$ 

```

The locale inherits  $\text{Inf}$  (and the auxiliaries  $\text{Inf\_from}$  and  $\text{Inf\_between}$ ) from  $\text{inference\_system}$  as well as  $\text{Bot}$  and  $\models$  from  $\text{consequence\_relation}$ . An inference  $\iota$  is redundant if  $\iota \in \text{Red}_l N$ ; a formula  $C$  is redundant if  $C \in \text{Red}_F N$ . The locale assumptions are needed to establish dynamic completeness. For example, the second assumption ensures that deleting a redundant formula from an unsatisfiable set preserves the set's unsatisfiability.

### 2.4 Static and Dynamic Completeness

Although dynamic completeness is the ultimate goal, static completeness is useful because it is easier to formulate and prove and, under some basic assumptions, is equivalent to dynamic completeness. Static completeness is expressed in terms of a single formula set  $N$ , which must be saturated. Saturation means that all inferences from  $N$  are redundant:

```

definition saturated :: 'f set  $\Rightarrow$  bool where
  saturated  $N \iff \text{Inf\_from } N \subseteq \text{Red}_l N$ 

```

Completeness means that any unsatisfiable saturated set must contain a patent falsehood (typically,  $\perp$ ):

```

locale statically_complete_calculus = calculus +
assumes  $B \in \text{Bot} \Rightarrow \text{saturated } N \Rightarrow N \models \{B\} \Rightarrow$ 
   $\exists B' \in \text{Bot}. B' \in N$ 

```

We abuse terminology above when we write that a set is “unsatisfiable.” Given that we only have a notion of entailment but not of modelhood, it would be more proper (but less standard) to write “inconsistent.”

Dynamic completeness more closely models saturation provers. It is expressed in terms of a finite or infinite sequence of formula sets represented by the codatatype  $'f \text{ set } \text{llist}$ . A derivation is a sequence where each pair of successive elements satisfies the following relation:

```

inductive  $\triangleright$  :: 'f set  $\Rightarrow$  'f set  $\Rightarrow$  bool where
   $M - N \subseteq \text{Red}_F N \Rightarrow M \triangleright N$ 

```

The  $\triangleright$  relation can be thought of as an abstract nondeterministic proving process. Informally, when taking a transition from a set  $M$  to another set  $N$ , a prover can add arbitrary formulas (corresponding to  $N - M$ ), and it can remove arbitrary formulas (corresponding to  $M - N$ ) as long as these are redundant w.r.t.  $N$ . Most provers would add only formulas that are entailed by  $M$ , but this is not enforced by  $\triangleright$ .

We write  $\text{chain } (\triangleright) \ Ns$  to express that the sequence  $Ns$  is a derivation:  $Ns ! 0 \triangleright Ns ! 1 \triangleright \dots$ , where the infix operator  $!$  returns the sequence element at the given index. The chain predicate is defined coinductively as in Schlichtkrull et al. [24, Section 5].

Dynamic completeness applies only to fair derivations:

```

definition fair :: 'f set llist  $\Rightarrow$  bool where
  fair  $Ns \iff$ 
   $\text{Inf\_from } (\text{Liminf } Ns) \subseteq \text{Sup } (\text{lmap } \text{Red}_l \ Ns)$ 

```

Above,  $\text{Liminf } Ns$  denotes the limit inferior  $\bigcup_i \bigcap_{j \geq i} Ns ! j$ , and  $\text{Sup } Ns$  denotes the supremum  $\bigcup_i Ns ! i$ . In accordance with the automated reasoning literature, we will refer to  $\text{Liminf } Ns$  as the *limit* of  $Ns$ .

Intuitively, a fair derivation is one where every inference with premises in the limit was made redundant at some index  $i$ , usually by adding its conclusion to  $Ns ! i$ . Using the locale assumptions about redundancy, it is easy to show that the limit of a fair derivation is saturated.

Dynamic completeness is formulated in its own locale:

```

locale dynamically_complete_calculus = calculus +
assumes
   $B \in \text{Bot} \Rightarrow \text{chain } (\triangleright) \ Ns \Rightarrow \text{fair } Ns \Rightarrow$ 
   $Ns ! 0 \models \{B\} \Rightarrow$ 
   $\exists i < \text{llength } Ns. \exists B' \in \text{Bot}. B' \in Ns ! i$ 

```

Owing to the conditions on redundancy criteria, static and dynamic completeness are equivalent:

```

sublocale statically_complete_calculus
   $\subseteq$  dynamically_complete_calculus
sublocale dynamically_complete_calculus
   $\subseteq$  statically_complete_calculus

```

The above definitions are all based on Waldmann et al., who in turn largely follow Bachmair and Ganzinger but generalize them. The literature is not entirely consistent on notions such as redundancy, saturation, and fairness. Waldmann et al. studied the variation and showed that completeness and “reduced” completeness coincide. We also formalized these proofs (in `Calculus_Variations.thy`). They

make it possible to mix and match results. For example, suppose that a calculus is proved “reducedly” statically complete in the literature. Instead of redoing the proof in a standard, nonreduced setting, we can reuse the reduced proof as is and derive dynamic completeness using the framework.

## 2.5 Intersection of Redundancy Criteria

There are various situations where a prover may need to apply several redundancy criteria in parallel and identify an inference or formula as redundant only if all the criteria agree. This can be incorporated into the framework as a single intersection criterion. One scenario where this is particularly useful is when redundancy is parameterized by a value  $q :: 'q$  that is known only at the limit. The intersection amounts to a universal quantification “for all  $q \in Q, \dots$ ” Once we know the limit, we can supply the right value for  $q$ .

Although this is not always necessary, each redundancy criterion may use its own consequence relation. Thus, we defined a `consequence_relation_family` locale parameterized by a set  $Q :: 'q \text{ set}$  and by a  $Q$ -indexed family of consequence relations  $\text{entails}_q :: 'q \Rightarrow 'f \text{ set} \Rightarrow 'f \text{ set} \Rightarrow \text{bool}$  such that each  $\text{entails}_q q$ , for  $q \in Q$ , qualifies as a consequence relation. The  $q$  subscript in the constant name is there to remind us that the family is indexed by the type  $'q$ . We will encounter this convention again later.

The next step was to define a single relation  $\models_Q$  as the intersection of the  $\text{entails}_q$ 's:

**definition**  $\models_Q :: 'f \text{ set} \Rightarrow 'f \text{ set} \Rightarrow \text{bool}$  **where**  
 $N_1 \models_Q N_2 \iff \forall q \in Q. \text{entails}_q q N_1 N_2$

We proved that  $\models_Q$  qualifies as a consequence relation.

Next, we defined a locale for intersection calculi, which work with the intersection of a  $Q$ -indexed family of redundancy criteria. The definitions (in `Intersection_Calculus.thy`) follow the same pattern as above. The locale follows:

**locale** `intersection_calculus =`  
`inference_system Inf +`  
`consequence_relation_family Bot Q entails_q +`  
**fixes**  
 $\text{Red}_{\text{I}q} :: 'q \Rightarrow 'f \text{ set} \Rightarrow 'f \text{ inference set}$  **and**  
 $\text{Red}_{\text{F}q} :: 'q \Rightarrow 'f \text{ set} \Rightarrow 'f \text{ set}$   
**assumes**  
 $Q \neq \emptyset$  **and**  
 $\forall q \in Q. \text{calculus Bot Inf (entails}_q q) (\text{Red}_{\text{I}q} q)$   
 $(\text{Red}_{\text{F}q} q)$

The redundancy criterion is defined as a pair of intersections:

**definition**  $\text{Red}_I :: 'f \text{ set} \Rightarrow 'f \text{ inference set}$  **where**  
 $\text{Red}_I N = (\bigcap q \in Q. \text{Red}_{\text{I}q} q N)$

**definition**  $\text{Red}_F :: 'f \text{ set} \Rightarrow 'f \text{ set}$  **where**  
 $\text{Red}_F N = (\bigcap q \in Q. \text{Red}_{\text{F}q} q N)$

Intersection calculi qualify as calculi with the intersections  $\text{Red}_I$  and  $\text{Red}_F$  as redundancy criterion:

**sublocale** `calculus Bot Inf ( $\models_Q$ )  $\text{Red}_I \text{Red}_F$`

To show that a calculus is statically complete, it suffices to show that for any saturated set  $N$  that does not contain a patent falsehood, there exists an index  $q \in Q$  according to which  $N$  is consistent:

**lemma** `stat_ref_comp_from_bot_in_sat`:  
 $(\forall N. \text{saturated Inf Red}_I N \wedge (\forall B \in \text{Bot}. B \notin N) \longrightarrow$   
 $\exists B \in \text{Bot}. \exists q \in Q. \neg \text{entails}_q q N \{B\}) \implies$   
`statically_complete_calculus Bot Inf ( $\models_Q$ )  $\text{Red}_I \text{Red}_F$`

## 2.6 Soundness

An inference system  $\text{Inf}$  is sound w.r.t. a consequence relation  $\models$  if every inference  $\iota \in \text{Inf}$  is sound—that is,  $\text{prems\_of } \iota \models \text{concl\_of } \iota$ . Soundness can be shown one inference rule at a time, without a sophisticated framework. Nevertheless, we formalized a few basic soundness results, as a convenience to users of the framework (in `Soundness.thy`).

The main soundness lemma, generalized from Schlichtkrull et al., states that the limit of a sound derivation is unsatisfiable if and only if the initial formula set is unsatisfiable:

**lemma** `unsat_limit_iff`:  
 $\text{chain } (\triangleright) Ns \implies \text{chain } (\models) Ns \implies$   
 $(\text{Liminf } Ns \models \text{Bot} \iff Ns ! 0 \models \text{Bot})$

For the forward direction, the proof by Schlichtkrull et al., like the one by Bachmair and Ganzinger on which it is based, requires compactness—the property that whenever  $N \models \text{Bot}$ , there exists a finite set  $M \subseteq N$  such that  $M \models \text{Bot}$ . Informally, the argument is as follows: If the limit of  $Ns$  is unsatisfiable, by compactness there exists an unsatisfiable finite subset; by the definition of limit, it takes only finitely many steps to reach a set that includes this finite subset. By transitivity of  $\models$ , the initial state  $Ns ! 0$  is unsatisfiable.

In a private communication, Waldmann supplied an almost as simple proof that dispenses with compactness. By formalizing this proof, we were able to eliminate a few bureaucratic locales and locale instantiations.

## 2.7 Standard Redundancy

Ground versions of resolution and superposition are parameterized by a well-order  $<$  that is lifted to clauses. The completeness proof is a well-founded induction over the elements of a saturated clause set  $N \not\equiv \perp$  w.r.t.  $<$ . Each clause is shown in turn to be true in a candidate model. Such inference systems are called *counterexample-reducing*. This terminology stems from an alternative formulation of induction as a proof by contradiction starting with the assumption that a minimal counterexample exists [4, Section 4.2].

An inference system  $\text{Inf}$  operating on formulas equipped with a well-order  $<$  is counterexample-reducing if the following locale requirements can be met:

**fixes**  $\text{l\_of} :: 'f \text{ set} \Rightarrow 'f \text{ set}$   
**assumes**  
 $N \cap \text{Bot} = \emptyset \implies D \in N \implies \neg \text{l\_of } N \models D \implies$

$$\begin{aligned}
& (\bigwedge C. C \in N \implies \neg \text{l\_of } N \models C \implies D \leq C) \implies \\
& \exists \iota \in \text{Inf. } \text{prems\_of } \iota \neq [] \wedge \text{main\_prem\_of } \iota = D \\
& \quad \wedge \text{side\_prems\_of } \iota \subseteq N \\
& \quad \wedge \text{l\_of } N \models \text{side\_prems\_of } \iota \\
& \quad \wedge \neg \text{l\_of } N \models \text{concl\_of } \iota \wedge \text{concl\_of } \iota < D
\end{aligned}$$

Observe that the rightmost, or *main*, premise is distinguished from the remaining *side* premises. Since the framework provides no notion of interpretations, we represent an interpretation as a formula set  $M$  and write  $M \models N$  to mean “ $M$  is a model of  $N$ .” Given a formula set, the  $\text{l\_of}$  operator yields a candidate model. The locale assumption states that if  $N$  contains no patent falsehoods and  $D$  is a minimal counterexample formula in  $N$ , then there exists an inference from  $N$  that produces an even smaller counterexample.

It is not hard to do nonsensical things using a proof assistant. Adversarial users might define  $\text{l\_of } N$  as  $\{\perp\}$  and use this to show that *any* calculus is counterexample-reducing. This would not take them very far, however, because they could never connect that definition of  $\text{l\_of}$  with a concrete notion of model suitable for their logic.

Counterexample-reducing inference systems are compatible with the *standard redundancy criterion*, defined as follows (in `Standard_Redundancy_Criterion.thy`):

**definition**  $\text{Red}_l :: 'f \text{ set} \implies 'f \text{ inference set}$  **where**

$$\begin{aligned}
\text{Red}_l N &= \{\iota \in \text{Inf. } \exists M \subseteq N. \\
& \quad M \cup \text{side\_prems\_of } \iota \models \{\text{concl\_of } \iota\} \\
& \quad \wedge \forall D \in M. D < \text{main\_prem\_of } \iota\}
\end{aligned}$$

**definition**  $\text{Red}_F :: 'f \text{ set} \implies 'f \text{ set}$  **where**

$$\text{Red}_F N = \{C. \exists M \subseteq N. M \models \{C\} \wedge \forall D \in M. D < C\}$$

For inference systems  $\text{Inf}$  that reduce counterexamples, we obtain a completeness theorem: If  $N$  is saturated w.r.t. the standard redundancy criterion and does not contain a patent falsehood, then  $\text{l\_of } N$  is a model of  $N$ .

## 2.8 Clauses

Most saturation calculi work with clauses. In Isabelle, we reuse the representation of clauses as finite multisets of literals introduced by Schlichtkrull et al. [24, Section 2.3]. The empty clause  $\perp$  is represented by  $\emptyset$ , but we will continue to write  $\perp$ . We take  $\text{Bot} = \{\perp\}$ . Literals over atoms of type  $'a$  take the form  $\text{Pos } A$  or  $\text{Neg } A$ , where  $A :: 'a$ .

As a sanity check and as a convenience to users, we specialized various framework concepts to clauses (in `Clausal_Calculus.thy`). In particular, we proved that entailment of ground clauses is compact and tailored the notion of counterexample-reducing inference system as follows:

**fixes**  $\text{J\_of} :: 'a \text{ clause set} \implies 'a \text{ set}$

**assumes**

$$\perp \notin N \implies D \in N \implies \neg \text{J\_of } N \models D \implies$$

$$\begin{aligned}
& (\bigwedge C. C \in N \implies \neg \text{J\_of } N \models C \implies D \leq C) \implies \\
& \exists \iota \in \text{Inf. } \text{prems\_of } \iota \neq [] \wedge \text{main\_prem\_of } \iota = D \\
& \quad \wedge \text{side\_prems\_of } \iota \subseteq N \\
& \quad \wedge \text{J\_of } N \models \text{side\_prems\_of } \iota \\
& \quad \wedge \neg \text{J\_of } N \models \text{concl\_of } \iota \wedge \text{concl\_of } \iota < D
\end{aligned}$$

For this locale, we were able to represent a candidate model more conventionally as a set of true atoms and to use the modelhood relation  $\models$  instead of entailment  $\models$  of clauses.

We connected the two formulations of counterexample reduction as a convenience for end users. The abstract locale was instantiated by taking

$$\text{l\_of } N = \{\{\text{if } A \in \text{J\_of } N \text{ then Pos } A \text{ else Neg } A\} \mid A :: 'a\}$$

We retrieved the following completeness theorem tuned for clauses:  $\text{saturated } N \implies \perp \notin N \implies \text{J\_of } N \models N$ .

## 3 Lifting Ground Calculi

The saturation framework supports several ways to lift ground calculi to the nonground level (in `Lifting_to_Non_Ground_Calculi.thy` and `Labeled_Lifting_to_Non_Ground_Calculi.thy`). We start with the folklore approach and introduce features incrementally until the result is flexible enough to support the calculi of realistic provers.

### 3.1 Standard Lifting

The standard lifting is based on a pair of grounding functions  $\mathcal{G}_F$  and  $\mathcal{G}_l$  that link nonground formulas (of some type  $'f$ ) to sets of ground formulas (of some type  $'g$ ) and similarly for nonground and ground inferences. Using these functions, it connects a nonground inference system to a ground calculus:

**locale** `standard_lifting =`  
`inference_system InfF +`  
`calculus BotG InfG ( $\models$ ) RedlG RedFG +`

**fixes**  
 $\mathcal{G}_F :: 'f \implies 'g \text{ set}$  **and**  
 $\mathcal{G}_l :: 'f \text{ inference} \implies 'g \text{ inference set option}$

**assumes**  
 $\text{Bot}_F \neq \emptyset$  **and**  
 $B \in \text{Bot}_F \implies \mathcal{G}_F B \neq \emptyset$  **and**  
 $B \in \text{Bot}_F \implies \mathcal{G}_F B \subseteq \text{Bot}_G$  **and**  
 $\mathcal{G}_F C \cap \text{Bot}_G \neq \emptyset \longrightarrow C \in \text{Bot}_F$  **and**  
 $\iota \in \text{Inf}_F \implies \mathcal{G}_l \iota \neq \text{None} \implies$   
`the`  $(\mathcal{G}_l \iota) \subseteq \text{Red}_{lG} (\mathcal{G}_F (\text{concl\_of } \iota))$

An instance of the locale for first-order logic would map the nonground formula  $p(x)$  to a Herbrand set such as  $\{p(a), p(b), \dots, p(f(a)), p(f(b)), \dots\}$ .

Entailment of sets of nonground formulas is defined via grounding:  $N_1 \models_G N_2 \iff \mathcal{G}_{Fset} N_1 \models \mathcal{G}_{Fset} N_2$ , where  $\mathcal{G}_{Fset} :: 'fset \implies 'gset$  is the lifting of  $\mathcal{G}_F$  to sets. This relation is called Herbrand entailment. For first-order logic, it coincides with the familiar Tarski entailment when  $N_2 = \{\perp\}$ , allowing us to use  $\models_G$  in a refutational setting.

Based on these notions, we extended the nonground inference system  $\text{Inf}_F$  with a lifted redundancy criterion to form a nonground calculus:

**definition**  $\text{Red}_I :: 'f \text{ set} \Rightarrow 'f \text{ inference set}$  **where**

$$\begin{aligned} \text{Red}_I N = \{ & \iota \in \text{Inf}_F. \\ & (\mathcal{G}_I \iota \neq \text{None} \wedge \text{the} (\mathcal{G}_I \iota) \subseteq \text{Red}_{IG} (\mathcal{G}_{Fset} N)) \\ & \vee (\mathcal{G}_I \iota = \text{None} \\ & \wedge \mathcal{G}_F (\text{concl\_of } \iota) \subseteq \mathcal{G}_{Fset} N \cup \text{Red}_{FG} (\mathcal{G}_{Fset} N)) \} \end{aligned}$$

**definition**  $\text{Red}_F :: 'f \text{ set} \Rightarrow 'f \text{ set}$  **where**

$$\text{Red}_F N = \{ C. \forall D \in \mathcal{G}_F C. D \in \text{Red}_{FG} (\mathcal{G}_{Fset} N) \}$$

**sublocale** calculus  $\text{Bot}_F \text{Inf}_F (\models_{\mathcal{G}}) \text{Red}_I \text{Red}_F$

This brings us to our first main result. If the ground calculus is statically complete and is overapproximated by the nonground calculus, then the nonground calculus is statically complete as well:

**theorem** *stat\_ref\_comp\_to\_non\_ground:*

**assumes**

$$\begin{aligned} & \text{statically\_complete\_calculus Bot}_G \text{Inf}_G (\models) \\ & \text{Red}_{IG} \text{Red}_{FG} \text{ and} \\ & \wedge N. \text{saturated } N \Rightarrow \text{ground\_Inf\_redundant } N \end{aligned}$$

**shows**

$$\begin{aligned} & \text{statically\_complete\_calculus Bot}_F \text{Inf}_F (\models_{\mathcal{G}}) \\ & \text{Red}_I \text{Red}_F \end{aligned}$$

### 3.2 Adding Tiebreaker Orders

The completeness proofs for standard superposition and many other calculi are organized as a standard lifting. However, the lifted redundancy criterion cannot be used to justify the use of clause subsumption in provers. This points to a gap between theory and practice.

Waldmann et al. generalized standard lifting to bridge this gap. They use a well-founded order, or more generally a family of orders, which they call a tiebreaker order. The tiebreaker can be instantiated with subsumption. In a given clause procedure, it can also be used to make a passive formula redundant w.r.t. the active version of the same formula.

The tiebreaker lifting is defined as an extension of standard lifting with a family  $\text{Prec}_{FG}$  of well-founded orders:

**locale** *tiebreaker\_lifting* =

$$\begin{aligned} & \text{empty\_ord: standard\_lifting Bot}_F \text{Inf}_F \text{Bot}_G \text{Inf}_G \\ & (\models) \text{Red}_{IG} \text{Red}_{FG} \mathcal{G}_F \mathcal{G}_I + \end{aligned}$$

**fixes**  $\text{Prec}_{FG} :: 'g \Rightarrow 'f \Rightarrow 'f \Rightarrow \text{bool}$

**assumes** *minimal\_element* ( $\text{Prec}_{FG} g$ ) UNIV

Redundancy of inferences  $\text{Red}_I$  is as for the standard lifting, but redundancy of formulas can use the tiebreaker:

**definition**  $\text{Red}_F :: 'f \text{ set} \Rightarrow 'f \text{ set}$  **where**

$$\begin{aligned} \text{Red}_F N = \{ & C. \forall D \in \mathcal{G}_F C. D \in \text{Red}_{FG} (\mathcal{G}_{Fset} N) \\ & \vee \exists E \in N. \text{Prec}_{FG} D E C \wedge D \in \mathcal{G}_F E \} \end{aligned}$$

Based on this definition, we proved that *tiebreaker\_lifting* produces a nonground calculus.

Waldmann et al. claim that *standard\_lifting* produces a nonground calculus but point to the more general *tiebreaker\_*

lifting for a proof. To follow this structure, we interpreted *standard\_lifting* as an instance of *tiebreaker\_lifting* by taking  $\text{Prec}_{FG} := \lambda g C C'$ . False, the empty tiebreaker family.

Waldmann et al. showed that the tiebreaker order can be ignored when proving static completeness and recovered for dynamic completeness. Using the above construction, it was easy to prove formally:

**theorem** *static\_to\_dynamic:*

$$\begin{aligned} & \text{statically\_complete\_calculus Bot}_F \text{Inf}_F (\models_{\mathcal{G}}) \\ & \text{Red}_I \text{empty\_ord.Red}_F \longleftrightarrow \\ & \text{dynamically\_complete\_calculus Bot}_F \text{Inf}_F (\models_{\mathcal{G}}) \\ & \text{Red}_I \text{Red}_F \end{aligned}$$

The theorem holds because the definitions of saturation, fairness, and completeness depend on  $\text{Red}_I$  but not on  $\text{Red}_F$ , and because static and dynamic completeness are equivalent.

### 3.3 Intersection of Liftings

Ground calculi defined via the intersection of redundancy criteria can be lifted using the same general approach. The locale declaration is as follows:

**locale** *lifting\_intersection* =

$$\begin{aligned} & \text{inference\_system Inf}_F + \\ & \text{inference\_system\_family } Q \text{Inf}_{Gq} + \\ & \text{consequence\_relation\_family Bot}_G Q \text{entails}_q + \end{aligned}$$

**fixes**

$\text{Bot}_F :: 'f \text{ set}$  **and**

$\mathcal{G}_{Fq} :: 'q \Rightarrow 'f \Rightarrow 'g \text{ set}$  **and**

$\mathcal{G}_{Iq} :: 'q \Rightarrow 'f \text{ inference} \Rightarrow 'g \text{ inference set option}$  **and**

$\text{Prec}_{FG} :: 'g \Rightarrow 'f \Rightarrow 'f \Rightarrow \text{bool}$

**assumes**

$$\begin{aligned} & \forall q \in Q. \text{tiebreaker\_lifting Bot}_F \text{Inf}_F \text{Bot}_G \\ & (\text{entails}_q q) (\text{Inf}_{Gq} q) (\text{Red}_{Iq} q) (\text{Red}_{Fq} q) \\ & (\mathcal{G}_{Fq} q) (\mathcal{G}_{Iq} q) \text{Prec}_{FG} \end{aligned}$$

For each index, we defined redundancy with and without tiebreaker order ( $\text{Red}_{IGq}$  and either  $\text{Red}_{FGq}$  or  $\text{Red}_{FGq\emptyset}$ ) and entailment ( $\text{entails}_{Gq}$ ) such that the intersection of the corresponding nonground calculi are also calculi:

**sublocale** *intersection\_calculus*  $\text{Bot}_F \text{Inf}_F Q \text{entails}_{Gq}$   
 $\text{Red}_{IGq} \text{Red}_{FGq}$

**sublocale** *empty\_ord:* *intersection\_calculus*  $\text{Bot}_F \text{Inf}_F$   
 $Q \text{entails}_{Gq} \text{Red}_{IGq} \text{Red}_{FGq\emptyset}$

These commands introduce intersected notions of entailment ( $\models_{\mathcal{G}}$ ) and redundancies ( $\text{Red}_{IG}$  and either  $\text{Red}_{FG}$  or  $\text{empty\_ord.Red}_F$ ).

As in *standard\_lifting*, the static completeness of the lifted intersection follows from that of the members of the ground family, and as in *tiebreaker\_lifting*, the tiebreaker order can be ignored for static completeness:

**theorem** *stat\_eq\_dyn\_ref\_comp\_fam\_inter:*

$$\begin{aligned} & \text{statically\_complete\_calculus Bot}_F \text{Inf}_F (\models_{\mathcal{G}}) \\ & \text{Red}_I \text{empty\_ord.Red}_F \longleftrightarrow \end{aligned}$$

dynamically\_complete\_calculus Bot<sub>F</sub> Inf<sub>F</sub> ( $\models_{\cap G}$ )  
Red<sub>I</sub> Red<sub>F</sub>

Although the proofs are theoretically easy and close to their informal counterparts, their formalization was complicated by the stack of locales and required tedious definition unfolding sequences and **interpret** commands to instantiate locales in the middle of proofs.

### 3.4 Adding Labels

The last kind of lifting attaches labels to formulas. These labels can be used to partition a set of formulas into its passive and active subsets. More fine-grained partitions are also possible. The labels have no logical meaning.

The locale extends `tiebreaker_lifting` with a set of labeled inferences  $\text{Inf}_{\text{FL}}$  such that whenever an  $\text{Inf}_{\text{F}}$  inference is possible, an  $\text{Inf}_{\text{FL}}$  is also possible, and vice versa:

**locale** `labeled_tiebreaking_lifting` =  
*no\_labels*: `tiebreaker_lifting BotF InfF BotG ( $\models$ )  
 InfG RedIG RedFG  $\mathcal{G}_{\text{F}}$   $\mathcal{G}_{\text{I}}$   $\text{Prec}_{\text{Fg}}$  +`  
**fixes**  $\text{Inf}_{\text{FL}} :: ('f \times 'l)$  *inference set*  
**assumes**  
 $\iota_{\text{F}} \in \text{Inf}_{\text{F}} \implies \text{length } Ls = \text{length } (\text{prems\_of } \iota_{\text{F}}) \implies$   
 $\exists L_0. \text{Infer } (\text{zip } (\text{prems\_of } \iota_{\text{F}}) Ls)$   
 $(\text{concl\_of } \iota_{\text{F}}, L_0) \in \text{Inf}_{\text{FL}}$  **and**  
 $\iota_{\text{FL}} \in \text{Inf}_{\text{FL}} \implies \text{Infer } (\text{map } \text{fst } (\text{prems\_of } \iota_{\text{FL}}))$   
 $(\text{fst } (\text{concl\_of } \iota_{\text{FL}})) \in \text{Inf}_{\text{F}}$

$\text{Bot}_{\text{FL}} :: ('f \times 'l)$  *set* is defined as  $\text{Bot}_{\text{F}} \times \text{UNIV}$ . The grounding functions ( $\mathcal{G}_{\text{IL}}$  and  $\mathcal{G}_{\text{FL}}$ ) and the entailment relation ( $\models_{\mathcal{G}_{\text{L}}}$ ) simply ignore labels. The main result lets us derive the static completeness of the labeled calculus from the static completeness of the base calculus:

**lemma** `stat_ref_comp_to_labeled_sta_ref_comp`:  
`statically_complete_calculus BotF InfF ( $\models_{\mathcal{G}}$ )  
 no_labels.RedI no_labels.RedF  $\implies$   
 statically_complete_calculus BotFL InfFL ( $\models_{\mathcal{G}_{\text{L}}}$ )  
 RedI RedF`

Labels can also be added to calculi that work with families of redundancy criteria. Since completeness results are independent of the tiebreaker order, we only considered the empty tiebreaker order:

**locale** `labeled_lifting_intersection` =  
*no\_labels*: `lifting_intersection InfF BotG Q InfGq  
 entailsq RedIq RedFq BotF  $\mathcal{G}_{\text{Fq}}$   $\mathcal{G}_{\text{Iq}}$   
 $(\lambda g CL CL'. \text{False}) +$   
fixes  $\text{Inf}_{\text{FL}} :: ('f \times 'l)$  inference set`

We proved that static completeness with labels can be reduced to static completeness without labels.

## 4 Prover Architectures

Dynamic completeness is formulated in terms of an abstract single-formula-set prover  $\triangleright$  (Section 2.4). In practice, saturation provers implement the given clause procedure described

in the introduction. This simplifies bookkeeping: Instead of having to keep track of all possible  $n$ -ary inferences, actual provers only need to track active formulas.

Waldmann et al. describe two versions of the given clause procedure: an eager version called GC and a lazy version called LGC. LGC strictly generalizes GC, but because GC is sufficient for most applications and simpler to use, we formalized both separately (in `Given-Clause-Architectures.thy` and `Given-Clause-Architectures-Revisited.thy`).

### 4.1 Basic Setup

Much of the setup between GC and LGC can be shared, in a locale called `given_clause_basis` that extends `labeled_lifting_intersection`. The locale assumes that the unlabeled version of the intersection calculus is complete:

**assumes**  
`statically_complete_calculus BotF InfF ( $\models_{\cap G}$ )  
 no_labels.RedI no_labels.RedF`

It introduces an equivalence relation  $\doteq$  and an order  $\prec$  on unlabeled formulas that must be compatible with each other and with the grounding function:

**assumes**  
 $C_1 \doteq D_1 \implies C_2 \doteq D_2 \implies C_1 \prec C_2 \implies D_1 \prec D_2$  **and**  
 $q \in Q \implies C_1 \doteq C_2 \implies \mathcal{G}_{\text{Fq}} q C_1 \implies \mathcal{G}_{\text{Fq}} q C_2$  **and**  
 $q \in Q \implies C_2 \prec C_1 \implies \mathcal{G}_{\text{Fq}} q C_1 \implies \mathcal{G}_{\text{Fq}} q C_2$

It also introduces an order  $\sqsubseteq_{\text{L}}$  on labels and a distinguished label `active` that must be minimal w.r.t.  $\sqsubseteq_{\text{L}}$ . Both orders are assumed to be well founded. Inside the locale, we used these relations to define a tiebreaker order  $\sqsubset$  on labeled formulas and showed it to be well founded:

**definition**  $\sqsubset :: ('f \times 'l) \Rightarrow ('f \times 'l) \Rightarrow \text{bool}$  **where**  
 $CL_1 \sqsubset CL_2 \iff \text{fst } CL_1 \prec \text{fst } CL_2$   
 $\vee (\text{fst } CL_1 \doteq \text{fst } CL_2 \wedge \text{snd } CL_1 \sqsubseteq_{\text{L}} \text{snd } CL_2)$

**lemma** `wf_prec_FL`: `minimal_element` ( $\sqsubset$ ) UNIV

This order can be used to strengthen the redundancy criterion of the labeled intersection calculus. Below, `RedI` and `RedF` refer to the strengthened criterion.

Finally, we present two lemmas that make explicit the relation between the labeled and unlabeled redundancy criteria:

**lemma** `labeled_red_inf_eq_red_inf`:  
 $\iota \in \text{Inf}_{\text{FL}} \implies \iota \in \text{Red}_{\text{I}}$   $\iff$   
 $\text{to\_F } \iota \in \text{no\_labels.Red}_{\text{I}} (\text{fst } \iota)$

**lemma** `red_labeled_clauses`:  
 $C \in \text{no\_labels.Red}_{\text{F}} (\text{fst } \iota) \vee (\exists C' \in \text{fst } \iota. C' \prec C)$   
 $\vee (\exists (C', L') \in N. L' \sqsubseteq_{\text{L}} L \wedge C' \prec C) \implies$   
 $(C, L) \in \text{Red}_{\text{F}} N$

Above, `to_F` denotes a function that erases all labels, and the infix operator  $\iota$  denotes the image of a set under a function.

Although the first lemma is obvious when we compare the definitions of `RedI` and `no_labels.RedI`, the corresponding

Isabelle proof is over 100 lines long. This is an unfortunate consequence of the deep nesting of locales. The identity of notions is hidden under multiple layers of names, and Isabelle provides no automation that can unfold these definitions automatically so as to reveal such connections. We conjecture that proof assistants based on computational type theories such as Agda [14], Coq [9], Lean [15], and Matita [1] would cope more gracefully with these definitional equalities.

The `given_clause_basis` locale also ensures that inferences never produce active formulas. Only the procedure itself is allowed to make a formula active.

When we developed this locale and the two locales based on it, we did not immediately try to instantiate them. This led to an amusing incident when we attempted to verify RP. The `given_clause_basis` locale axioms turned out to be flawed. First, we had stated that  $\succ$  must be well founded when we meant  $\prec$ . Strangely enough, we had successfully proved  $\sqsubset$  well founded, even though it is based on  $\prec$ .

How could that be? The answer is that we had relied on very powerful proof automation: Sledgehammer [22]. That tool had discovered *another* flawed assumption and exploited it. This time, we had misspelled a constant's name. In Isabelle, this creates an implicitly universal variable. This convention, which certainly saves a lot of typing and is considered by some as an advantage of Isabelle over its rivals, had led to an inconsistent assumption.

## 4.2 The Eager Given Clause Procedure

The eager given clause procedure GC is described as a transition system that operates on labeled formula sets. This corresponds naturally to an inductive predicate  $\rightsquigarrow_{GC}$ :

**inductive**  $\rightsquigarrow_{GC} :: ('f \times 'l) \text{ set} \Rightarrow ('f \times 'l) \text{ set} \Rightarrow \text{bool}$   
**where**  
*process*:  $N_1 = N \cup M \Rightarrow N_2 = N \cup M' \Rightarrow$   
 $M \subseteq \text{Red}_F(N \cup M') \Rightarrow \text{active\_subset } M' = \emptyset \Rightarrow$   
 $N_1 \rightsquigarrow_{GC} N_2$   
*infer*:  $N_1 = N \cup \{(C, L)\} \Rightarrow$   
 $N_2 = N \cup \{(C, \text{active})\} \cup M \Rightarrow L \neq \text{active} \Rightarrow$   
 $\text{active\_subset } M = \emptyset \Rightarrow$   
 $\text{no\_labels.Inf\_between } (\text{fst } ^\text{c} \text{ active\_subset } N) \{C\}$   
 $\subseteq \text{no\_labels.Red}_1(\text{fst } ^\text{c} (N \cup (C, \text{active}) \cup M)) \Rightarrow$   
 $N_1 \rightsquigarrow_{GC} N_2$

The *process* rule can be used to add arbitrary passive formulas  $M'$  or to remove redundant formulas  $M$ . For example, the simplification of  $p(2 + 2)$  to  $p(4)$  would be modeled by taking  $M := \{p(2 + 2)\}$  and  $M' := \{p(4)\}$ . The *infer* rule selects a passive formula  $C$ —called the given formula or given clause—and makes it active. In addition, it draws all inferences between  $C$  and all the active formulas; more precisely, the rule requires the addition of enough passive formulas  $M$  to make these inferences redundant. The locale assumes that the inference system  $\text{Inf}_F$  contains no nullary inferences.

The main benefit of using GC is that it makes it easier to draw inferences fairly. We only need to ensure that the prover eventually makes every passive formula active.

Via a refinement step, we showed that  $\rightsquigarrow_{GC}$ -transitions correspond to  $\triangleright$ -transitions. This was straightforward. Then we needed to connect GC's notion of fairness with the fair predicate on  $\triangleright$ -derivations (Section 2.4). Fairness for GC means that no formulas are passive at the limit, starting from a state in which no formulas are active:

**lemma** *gc\_fair*:

**assumes**

$\text{chain } (\rightsquigarrow_{GC}) \text{ } Ns$  **and**

$\text{active\_subset } (Ns ! 0) = \emptyset$  **and**

$\text{passive\_subset } (\text{Liminf } Ns) = \emptyset$

**shows**  $\text{fair } Ns$

The proof by Waldmann is 13 lines long, but we needed 197 lines of Isabelle. The explanation for this poor De Bruijn factor is that the proof considers the maximal element of a set of integers, each representing a step in the derivation that is the first to satisfy a complex property for one of the premises of a given inference. The previous sentence already hints at the number of hidden reasoning steps necessary to obtain this particular element, although it is fairly easy to write the notion mathematically in a way that enables readers to build an adequate mental representation of the object. In Isabelle, however, every intermediate object in this construction must be shown to exist.

For each premise, the existence of the particular step of the derivation that verifies the property must be proved. This alone accounts for almost half of the proof. Then a set is built of those integers that must be proved nonempty and finite, before the maximal element can finally be obtained. The other half of the proof is dominated by the need to reason by elimination to identify that this particular step must be an infer step. This is obvious on paper by inspection of the transition rules but requires detailed reasoning in Isabelle.

We were not satisfied with this proof. A proof that is tedious to write is also tedious to read, and hides its insights from the reader. The situation reminded us of a 700-line Isabelle formalization by Fleury of an 8-line pen-and-paper proof by Weidenbach [13, Section 4.3]. The solution in that case was to reason by invariance.

Can invariance help prove *gc\_fair* in a more comprehensible way? It is well established that the given clause procedure satisfies the following key invariant:

All inferences from active formulas are redundant w.r.t. the current set of formulas.

Initially, all formulas are passive, so the invariant holds vacuously. Whenever a given formula is made active, all inferences with active partners are made redundant, thereby maintaining the invariant. And if all formulas are active at



the limit, then this means that all inferences from the entire limit  $N$  are redundant. Hence the derivation is fair and consequently  $N$  is saturated.

Paradoxically, although the invariant above is widely recognized as an important property of the given clause procedure, which is itself central in automated reasoning, we are not aware of any proofs in the literature that rely on it. Instead, the invariant is presented as a guide to intuition, if at all. For example, Bachmair and Ganzinger do not mention the invariant in their *Handbook* chapter.

In Isabelle, we defined the invariant as follows:

**definition**  $gc\_invar :: ('f \times 'l) set llist \Rightarrow enat \Rightarrow bool$

**where**

$gc\_invar\ Ns\ i \iff$   
 $Inf\_from\ (active\_subset\ (Liminf\_upto\ Ns\ i))$   
 $\subseteq\ Sup\_upto\ (lmap\ Red_{lG}\ Ns)\ i$

The predicate  $gc\_invar\ Ns\ i$  determines whether the invariant holds at index  $i$  (which can be  $\infty$ ). If  $k$  is within bounds,  $Liminf\_upto\ Ns\ k$  is defined as the bounded limit inferior  $\bigcup_i \bigcap_{j=i}^k Ns\ !\ j$ , which is equal to the limit if  $k = \infty$  and to  $Ns\ !\ k$  otherwise. Similarly,  $Sup\_upto\ Ns\ k$  is the bounded supremum  $\bigcup_i^k Ns\ !\ i$ .

The invariant holds for all derivations with a reasonable initial state, and it extends to the limit:

**lemma**  $gc\_invar\_gc$ :

$chain\ (\sim_{GC})\ Ns \implies active\_subset\ (Ns\ !\ 0) = \emptyset \implies$   
 $i < llength\ Ns \implies gc\_invar\ Ns\ i$

**lemma**  $gc\_invar\_infinity$ :

$Ns \neq [] \implies (\forall i < llength\ Ns. gc\_invar\ Ns\ i) \implies$   
 $gc\_invar\ Ns\ \infty$

Assuming that all formulas eventually become active, the invariant at the limit simplifies to  $Inf\_from\ (Liminf\ Ns) \subseteq Sup\ (lmap\ Red_{lG}\ Ns)$ . This corresponds exactly to the definition of fairness, which is what we want.

The new formal proof is not much shorter in terms of number of lines (166 vs. 197), but it is much more compact horizontally and about half the size in bytes (7 vs. 13 kB). Written up as an informal proof, however, it is longer than the original; specifying the invariant and stating intermediate lemmas takes up vertical space. Even so, the new proof is more lucid, and we plan to include in the journal version of Waldmann et al.

### 4.3 The Given Clause Procedure with Delayed Inferences

The lazy given clause procedure LGC decouples inference scheduling and computation. Each inference  $\iota$  is first added to a “to do” set  $T$  of scheduled inferences. At some later point,  $\iota$  is computed, meaning that it is removed from  $T$  and its conclusion is added to the set  $N$  of labeled formulas. Alternatively,  $\iota$  might be identified as redundant and deleted,

a technique with the unfortunate name of “orphan deletion.” The formal definition of LGC follows:

**inductive**  $\sim_{LGC} :: 'f\ inference\ set \times ('f \times 'l) set \Rightarrow$   
 $'f\ inference\ set \times ('f \times 'l) set \Rightarrow bool$

**where**

$process: N_1 = N \cup M \implies N_2 = N \cup M' \implies$   
 $M \subseteq Red_F\ (N \cup M') \implies active\_subset\ M' = \emptyset \implies$   
 $(T, N_1) \sim_{LGC}\ (T, N_2)$   
 $| schedule\_infer: T_2 = T_1 \cup T' \implies N_1 = N \cup \{(C, L)\} \implies$   
 $N_2 = N \cup \{(C, active)\} \implies L \neq active \implies$   
 $T' = no\_labels.Inf\_between$   
 $(fst\ 'active\_subset\ N)\ \{C\} \implies$   
 $(T_1, N_1) \sim_{LGC}\ (T_2, N_2)$   
 $| compute\_infer: T_1 = T_2 \cup \{\iota\} \implies N_2 = N_1 \cup M \implies$   
 $active\_subset\ M = \emptyset \implies$   
 $\iota \in no\_labels.Red_1\ (fst\ '(N_1 \cup M)) \implies$   
 $(T_1, N_1) \sim_{LGC}\ (T_2, N_2)$   
 $| delete\_orphans: T_1 = T_2 \cup T' \implies T' \cap$   
 $no\_labels.Inf\_from\ (fst\ 'active\_subset\ N) = \emptyset \implies$   
 $(T_1, N) \sim_{LGC}\ (T_2, N)$

LGC has several benefits beyond orphan deletion. The “to do” set  $T$  can be used to support nullary inferences, or axioms. It can be used to represent infinite streams of inferences, which might for example result from higher-order unification (as in Zipperposition). Finally, inferences can often be stored more compactly than their conclusions, so LGC can be used to enable a low-level optimization (as in Waldmeister).

The key result is to show that  $\sim_{LGC}$ -derivations are fair under some reasonable assumptions:

**lemma**  $lgc\_fair$ :

**assumes**

$chain\ (\sim_{LGC})\ TNs\ \mathbf{and}$   
 $active\_subset\ (snd\ (TNs\ !\ 0)) = \emptyset\ \mathbf{and}$   
 $passive\_subset\ (Liminf\ (lmap\ snd\ TNs)) = \emptyset\ \mathbf{and}$   
 $\forall \iota \in Inf_F. prems\_of\ \iota = [] \longrightarrow \iota \in fst\ (TNs\ !\ 0)\ \mathbf{and}$   
 $Liminf\ (lmap\ fst\ TNs) = \emptyset$

**shows**  $fair\ (lmap\ snd\ TNs)$

Compared with  $\sim_{GC}$ , an extra condition is added to ensure that all premise-free inferences are scheduled up front. The proof has the same structure as that of  $gc\_fair$ , but the reasoning steps are more complex due to the extra rules and the presence of scheduled inferences. We needed 326 lines.

Can invariance help re-prove  $lgc\_fair$  in the same way that it helped us with  $gc\_fair$ ? The lazy given clause procedure is little discussed in the literature, but it is easy to come up with a suitable generalization of the previous invariant:

All inferences from active formulas are either scheduled in the “to do” set  $T$  or redundant w.r.t. the formula set  $N$ .

Initially, all formulas are passive, and all nullary inferences are in  $T$ , so the invariant holds. Whenever a given formula is made active, all inferences with active partners are scheduled,

thereby maintaining the invariant. Removing an inference from  $T$  and computing it makes it redundant, preserving the invariant. Orphan deletion also preserves the invariant, because by definition an orphan is not an inference from active formulas. And if  $T$  is empty and all formulas are active at the limit, then all inferences from the limit  $N$  are redundant, as required by fairness.

We defined the invariant as follows in Isabelle:

**definition** `lgc_invar` ::

*(f inference set × (f × l) set) llist ⇒ enat ⇒ bool*

**where**

`lgc_invar`  $TNs$   $i$   $\longleftrightarrow$

`Inf_from` (active\_subset

(`Liminf_upto` (`lmap snd`  $TNs$ )  $i$ ))

$\subseteq$  `from_F` ' `Liminf_upto` (`lmap fst`  $TNs$ )  $i$ )

$\cup$  `Sup_upto` (`lmap` (`RedIG`  $\circ$  `snd`)  $TNs$ )  $i$

The formal proof by invariance resembles that for GC but with more cases to consider. Compared with the monolithic proof, we have a small improvement in number of lines (301 vs. 326) and a more substantial improvement in byte size (13 vs. 22 kB). For both GC and LGC, we believe the gain in lucidity is much higher than the numbers suggest.

## 5 Application to a Resolution Prover

A framework is not a framework until it is used. Waldmann et al. [32, Example 29] applied their pen-and-paper framework to Bachmair and Ganzinger's resolution prover RP. To validate our formalized framework, a natural option was to apply it to RP as well, or rather RP's formalization by Schlichtkrull et al. We lifted ground static completeness explicitly and re-proved the dynamic completeness of RP. This example uses all the main components of the saturation framework and can serve as a blueprint for verifying other provers.

### 5.1 The Resolution Prover

Ordered resolution with selection, the calculus underlying RP, works on first-order clauses without equality (e.g.,  $p(x) \vee \neg q(b, f(a))$ ). It is parameterized by a well-order  $<$  on ground atoms, which is lifted to literals and clauses, and by a selection function  $S$  that maps each clause to one of its subclauses. The ground version of the calculus consists of a single  $n$ -ary inference rule:

$$\frac{(C_i \vee A_i \vee \dots \vee A_i)_{i=1}^n \quad \neg A_1 \vee \dots \vee \neg A_n \vee D}{C_1 \vee \dots \vee C_n \vee D}$$

Side conditions restrict the applicability of the rule, based on the well-order and the selection function. The order prunes the search space by identifying clauses that stray away from the goal ( $\perp$ ). The selection function guides the search.

The nonground version of the inference rule is

$$\frac{(C_i \vee A_{1i} \vee \dots \vee A_{ik_i})_{i=1}^n \quad \neg A_1 \vee \dots \vee \neg A_n \vee D}{(C_1 \vee \dots \vee C_n \vee D)\sigma}$$

where  $\sigma$  is a most general simultaneous solution of all unification problems  $A_{i1} \stackrel{?}{=} \dots \stackrel{?}{=} A_{ik_i} \stackrel{?}{=} A_i$ , where  $1 \leq i \leq n$ . As in the ground case, further side conditions apply.

RP is defined as a transition system  $\rightsquigarrow_{RP}$  over states of the form  $\mathcal{S} = (N, P, O)$ , where  $N$  contains the “new” clauses,  $P$  contains the “processed” clauses, and  $O$  contains the “old” clauses. RP implements a variant of the given clause procedure, with  $N \cup P$  as the passive set and  $O$  as the active set. RP is modeled naturally as an inductive predicate in Isabelle, with nine introduction rules:

**inductive**  $\rightsquigarrow_{RP}$  :: 'a state  $\Rightarrow$  'a state  $\Rightarrow$  bool **where**

- {Neg  $A$ , Pos  $A$ }  $\subseteq C \Rightarrow (N \cup \{C\}, P, O) \rightsquigarrow_{RP} (N, P, O)$
- |  $D \in P \cup O \Rightarrow$  subsumes  $D C \Rightarrow$   
( $N \cup \{C\}, P, O) \rightsquigarrow_{RP} (N, P, O)$
- |  $D \in N \Rightarrow$  strict\_subsumes  $D C \Rightarrow$   
( $N, P \cup \{C\}, O) \rightsquigarrow_{RP} (N, P, O)$
- |  $D \in N \Rightarrow$  strict\_subsumes  $D C \Rightarrow$   
( $N, P, O \cup \{C\}) \rightsquigarrow_{RP} (N, P, O)$
- |  $D \in P \cup O \Rightarrow$  reduces  $D C L \Rightarrow$   
( $N \cup \{C \uplus \{L\}\}, P, O) \rightsquigarrow_{RP} (N \cup \{C\}, P, O)$
- |  $D \in N \Rightarrow$  reduces  $D C L \Rightarrow$   
( $N, P \cup \{C \uplus \{L\}\}, O) \rightsquigarrow_{RP} (N, P \cup \{C\}, O)$
- |  $D \in N \Rightarrow$  reduces  $D C L \Rightarrow$   
( $N, P, O \cup \{C \uplus \{L\}\}) \rightsquigarrow_{RP} (N, P \cup \{C\}, O)$
- | ( $N \cup \{C\}, P, O) \rightsquigarrow_{RP} (N, P \cup \{C\}, O)$
- | ( $\emptyset, P \cup \{C\}, O) \rightsquigarrow_{RP}$   
(concl\_of ' `Inf_between`  $O \{C\}, P, O \cup \{C\}$ )

The first seven rules perform optional clause processing steps: tautology elimination, clause subsumption, and clause reduction (i.e., removal of needless literals). The eighth rule moves a clause  $C$  from  $N$  to  $P$ , and the ninth rule moves  $C$  further from  $P$  to  $O$  and computes all possible inferences between  $C$  and partners from  $O$ .

A sequence of states  $\mathcal{S}s$  is called fair, written fair  $\mathcal{S}s$ , if  $N\_of$  (`Liminf`  $\mathcal{S}s$ ) =  $P\_of$  (`Liminf`  $\mathcal{S}s$ ) =  $\emptyset$ . The completeness theorem for RP can be stated as follows:

**theorem** `RP_complete_if_fair`:

**assumes**

chain ( $\rightsquigarrow_{RP}$ )  $\mathcal{S}s$  **and**

$O\_of$  ( $\mathcal{S}s ! 0$ ) =  $\emptyset$  **and**

fair  $\mathcal{S}s$  **and**

$\forall I. \neg I \models \mathcal{G}_S (\mathcal{S}s ! 0)$

**shows**  $\perp \in O\_of$  (`Liminf`  $\mathcal{S}s$ )

Informally, if we start in a state where all clauses are passive (i.e.,  $O = \emptyset$ ) and make sure that all clauses are active (i.e.,  $N = P = \emptyset$ ) at the limit, and the initial clause set is unsatisfiable, then  $\perp$  will belong to the limit. This, in turn, means that  $\perp$  is derived after finitely many transitions—ideally before the prover runs out of time and the user, of patience.

### 5.2 The Original Completeness Proof

The Isabelle completeness proof by Schlichtkrull et al. (in `FO_Ordered_Resolution.thy` and `FO_Ordered_Resolution_`

Prover.thy), based on Bachmair and Ganzinger, consists of the following sequence of steps:

1. An  $\rightsquigarrow_{\text{RP}}$ -transition from  $\mathcal{S}$  to  $\mathcal{S}'$  corresponds to a  $\triangleright$ -transition from  $\mathcal{G}_{\mathcal{S}} \mathcal{S}$  to  $\mathcal{G}_{\mathcal{S}'} \mathcal{S}'$ , where  $\mathcal{G}_{\mathcal{S}}$  returns, for a given state, the set consisting of the  $\mathcal{G}_{\text{F}}$ -groundings of all the clauses in the state.
2. Let  $\mathcal{S}$ s be a fair RP derivation. Then any nonredundant clause  $C$  in  $\mathcal{S}$ s's ground projection eventually reaches  $O$  and stays there.
3. Moreover, if  $O$  is initially empty, then the limit of the  $\mathcal{G}_{\mathcal{S}}$ -grounding of  $\mathcal{S}$ s is saturated.
4. As a corollary, since ground resolution is complete, if the initial clause set is unsatisfiable, then  $\perp$  occurs in the limit (grounded or not).

It is difficult to relate to these proof steps. Even after the imprecisions in Bachmair and Ganzinger's proofs have been identified and clarified [24, Appendix A], the proof is hard to understand and remember. And yet, it is not long. In the *Handbook*, it spans four pages. In Isabelle, it amounts to around 2500 lines, but this includes additional minor results. About half of the lines are dedicated to lifting ground inferences and coping with  $\alpha$ -equivalence—issues that are only alluded to in the *Handbook*.

### 5.3 The New Completeness Proof

The new completeness proof (in FO\_Ordered\_Resolution\_Prover\_Revisited.thy) has a more abstract structure. It exploits the concepts presented in Sections 2 to 4. It also consists of four steps:

1. Ground resolution is statically complete.
2. Hence nonground resolution is statically complete.
3. Hence a given clause prover GC based on resolution is dynamically complete.
4. Thus, via refinement, RP is dynamically complete.

The four steps are identified by letter codes: G (ground), F (first-order), FL (F with labels), and RP.

For step 1, it sufficed to show that ground resolution  $G\_Inf$  is a clausal counterexample-reducing inference system (Sections 2.7 and 2.8). The core of the proof is a lemma already provided by Schlichtkrull et al.

For step 2, we defined nonground resolution as an inference system  $F\_Inf$ . This could be done in one line in terms of the definition by Schlichtkrull et al. We also defined grounding functions:  $\mathcal{G}_{\text{F}}$  on clauses and  $\mathcal{G}_{\text{I}}$  on inferences. The grounding of a clause  $C$  consists of all ground clauses of the form  $C\sigma$ ; the grounding of an inference  $C_n, \dots, C_1 \vdash C_0$  consists of all the inferences  $C_n\sigma_n, \dots, C_1\sigma_1 \vdash C_0\sigma_0 \in G\_Inf$ .

The main difficulty we encountered concerned the selection function. Recall that resolution is parameterized by a function  $S$  that restricts inferences. Inconveniently,  $S$  is generally not stable under substitution—it might behave differently on  $C\sigma$  and  $C$ . Based on the limit set  $M$ , we can define

an appropriate selection function  $S_M$  for the ground level as a modification of  $S$ . However, this definition is circular: The limit  $M$  depends on the nonground calculus with its selection function  $S$ , which in turn is lifted from a ground calculus with its selection function  $S_M$ , which depends on  $M$ .

Remarkably, this circularity had escaped the attention of Waldmann et al. as they worked in  $\text{\LaTeX}$ . It is all too easy to dismiss selection as an orthogonal concern. We noticed the issue only as we formalized RP. Fortunately, it was not too late to adapt the pen-and-paper framework.

The notion of intersection of redundancy criteria (Section 2.5) was designed to break the circularity. Instead of being lifted from  $S_M$ , where  $M$  is the limit, the nonground calculus is lifted simultaneously from all selection functions of the form  $S_M$ , where  $M$  is an arbitrary clause set. This was achieved formally by instantiating the lifting\_intersection\_locale (Section 2.5) and showing that the grounding functions  $\mathcal{G}_{\text{F}}$  and  $\mathcal{G}_{\text{I}}$  satisfy basic properties.

We could then instantiate statically\_complete\_calculus (Section 2.4) to lift the ground calculus to the nonground level. The key lemma states that any inference  $\iota_G$  from the grounded clauses is approximated at the nonground level by an inference  $\iota_{\text{F}}$ :

**lemma**  $G\_Inf\_overapprox\_F\_Inf$ :  
 $\iota_G \in G\_Inf\_from\ M\ (\bigcup (\mathcal{G}_{\text{F}}\ 'M)) \implies$   
 $\exists \iota_{\text{F}} \in F\_Inf\_from\ M. \iota_G \in \mathcal{G}_{\text{I}}\ M\ \iota_{\text{F}}$

For the core of the proof, we could again reuse a result from Schlichtkrull et al. This illustrates again how the saturation framework takes care of the bureaucracy and allows the user to focus on the calculus-specific reasoning.

Step 3 is about deriving dynamic completeness of a given clause procedure. We needed to define a suitable equivalence relation  $\doteq$  and tiebreaker order  $\triangleright$  in formulas. We also defined the labels New, Processed, and Old (= active) and defined the order  $\sqsupset$  such that New  $\sqsupset$  Processed  $\sqsupset$  Old. With these definitions in place, we instantiated given\_clause and retrieved two main results: (1) the induced  $\rightsquigarrow_{\text{GC}}$  refines  $\triangleright$  and (2) the induced  $\rightsquigarrow_{\text{GC}}$  is dynamically complete.

Step 4 was essentially a refinement proof: RP's states, which are triples of clause sets, must be converted to GC's labeled clause sets:

**definition** lcls\_of :: 'a state  $\implies$  ('a clause  $\times$  label) set  
**where**

$$\begin{aligned} \text{lcls\_of } \mathcal{S} = & (\lambda C. (C, \text{New}))\ 'N\_of\ \mathcal{S} \\ & \cup (\lambda C. (C, \text{Processed}))\ 'P\_of\ \mathcal{S} \\ & \cup (\lambda C. (C, \text{Old}))\ 'O\_of\ \mathcal{S} \end{aligned}$$

The first eight transition rules of RP can be simulated by the *process* rule of GC, and the last rule of RP can be simulated by *infer*. To illustrate this, we will sketch the proof for the first RP rule, which deletes tautologies.

Consider the transition  $(N \cup \{C\}, P, O) \rightsquigarrow_{\text{RP}} (N, P, O)$ , where  $\{\text{Neg } A, \text{Pos } A\} \subseteq C$ . We need to prove that there exist labeled clause sets  $N', M, M'$  such that  $N' \cup M \rightsquigarrow_{\text{GC}} N' \cup M'$

by the *process* rule. This amounts to showing that (1)  $M$  is redundant w.r.t.  $N' \cup M'$  and (2) no clause from  $M'$  is labeled active. We take  $N' := \text{lclss\_of}(N, P, O)$ ,  $M := \{(C, \text{New})\}$ , and  $M' := \emptyset$ . Condition (1) is obvious since  $C$  is a tautology, and (2) is vacuously true.

After completing step 4, we had all the results we needed to re-prove the main theorem, *RP\_complete\_if\_fair*—with one exception. To lighten the presentation, we used our notations for basic concepts such as inferences, redundancy, and fairness. However, Schlichtkrull et al. used slightly different formulations, based on Bachmair and Ganzinger. Their notions are restricted to clauses, which is not an issue for RP, but the side premises of an inference are stored in a multiset and thus unordered. This affects all definitions based on inferences, from redundancy to saturation. We first proved completeness and related results using our concepts; then we restated the results and re-proved them using theirs, as a sanity check. And although we were mostly interested in completeness, we also proved the soundness of RP derivations, exploiting the lemma *unsat\_limit\_iff* (Section 2.6).

Is the new proof an improvement over the one by Schlichtkrull et al.? It certainly is in terms of lines of source text. Let us ignore soundness and the conversions between basic concepts, and all the definitions and lemmas that are shared between the two proofs, including the definition of RP itself. Then we arrive at around 700 lines for the new proof compared with 1200 lines for the original proof. Among the 700 lines, nearly half are concerned with the refinement from GC to RP. But regardless of the line counts, we are convinced that the new proof's modularity makes it more intelligible and easier to teach, and easier to mimic to formally verify other saturation provers.

Looking at the line counts, one might be led to believe that the new proof was easy to develop. Unfortunately, this was not the case, because we worked with a less mature version of the framework. The circularity issue noted above led to a one-year hiatus, and when we resumed, we still had to find our way across a web of locales. The main difficulty we faced is that several copies of the same locale instances emerged, identical up to unfolding of definitions but distinct as far as the locale machinery is concerned. Instantiating the *given\_clause* locale, which pulls in a wide range of concepts, initially produced around 40 subgoals instead of the 14 we now get. There were so many layers of definitions that Sledgehammer was helpless.

One reason locales became so complicated is that we rigorously followed Waldmann et al. and their proofs by reduction. For example, lifting with a tiebreaker order is reduced to lifting with  $\emptyset$  as the tiebreaker, which in turn amounts to a standard lifting. A more direct proof would consider only the most general concept and avoid the replication of concepts. Fortunately, we found a way to simplify the locales while keeping compatibility with Waldmann et al. Often, it sufficed to replace an opaque **definition** by a transparent

**abbreviation**, or to choose a canonical locale instance and refrain from using the other definitionally equal instances.

Locales are truly a double-edged sword. On the one hand, they conveniently keep track of parameters and dependencies; without them, we would need to clutter definitions with extra arguments and lemmas with extra assumptions. On the other hand, locales often surprised us and provide too few introspection methods. Two examples: First, when instantiating a locale, there is no easy way to figure where the  $n$  subgoals come from in the locale hierarchy. Second, after instantiating a locale, the command **print\_theorems**, which is designed to display the lemmas introduced by the previous command, prints nothing. So after closing the 14 subgoals of *given\_clause*, we still have to search for what we have proved. Clearly, locales are immensely useful, but they could be made easier to use and debug.

## 6 Concluding Remarks

We formalized in Isabelle/HOL a framework designed to support the verification of automatic provers based on saturation calculi. This work joins a long list of verifications of logical calculi and theorem provers; we refer to Blanchette [12, Section 5] for a recent overview of related work. But unlike virtually all the previous work, instead of focusing on a single calculus or prover, we mechanized a framework applicable to a wide range of provers.

Because the informal proof and its formalization took place partly in parallel, they could benefit from each other. The precise and reliable informal proof was mostly a joy to translate into Isar. The formalization did unveil an unpleasant circularity in the application of the framework to RP, which was repaired by introducing a new concept. But we also encountered the opposite situation, where it was the translation to Isabelle that contained flawed conditions. We uncovered this also thanks to the RP case study.

The pen-and-paper version of the framework, due to Waldmann et al., is already being used in six separate informal proofs in ongoing work by colleagues and ourselves. If we want our formalized framework to be useful in the same way, we first need to prove static ground completeness for various interesting saturation calculi. So far, only ordered resolution and, thanks to Peltier [23], a variant of standard superposition have been formalized in Isabelle.

Beyond our framework, the IsaFoR (Isabelle Formalization of Rewriting) library [28] and the RP formalization [24] offer many definitions and lemmas related to first-order terms and clauses. Starting from RP, Schlichtkrull et al. [25] performed three further refinement steps to derive an executable functional prover, RPx, in Standard ML. With no effort, RPx could be rebased to use our RP proof. In principle, it would be possible to refine RPx further to use optimized data structures and algorithms, following the lines of Fleury's verification of an imperative SAT solver using Isabelle [17].

Waldmann et al. present a wealth of examples, especially in their technical report. For future work, some of these could be formalized. Particularly useful would be their three prover loops (Otter, DISCOUNT, and Zipperposition) based on the given clause procedures GC and LGC. Choosing one of these loops instead of GC or LGC as the basis of a prover could reduce the refinement burden. Using Peltier's formalization of superposition, it should now be possible to verify a superposition prover in the style of RP<sub>x</sub> in a few thousand lines of Isabelle. Integrating imperative data structures and algorithms, however, would involve much more work, perhaps of the order of a PhD project.

The saturation framework is very flexible, but it does not capture the important technique of clause splitting as implemented in SPASS and Vampire [16, 30]. The notion of redundancy criterion is also too weak to justify pure literal and blocked clause elimination [19]. To lift these restrictions, more theoretical research is necessary. If the present work is an indication of anything, this research should be carried out at least in part using a proof assistant.

## Acknowledgments

We thank Alexander Bentkamp, Robert Lewis, Simon Robillard, Anders Schlichtkrull, Mark Summerfield, Dmitriy Traytel, and Uwe Waldmann for their comments and suggestions. Blanchette's research has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement No. 713999, Matryoshka). He also benefited from the Netherlands Organization for Scientific Research (NWO) Incidental Financial Support scheme and he has received funding from the NWO under the Vidi program (project No. 016.Vidi.189.037, Lean Forward).

## References

- [1] Andrea Asperti, Wilmer Ricciotti, Claudio Sacerdoti Coen, and Enrico Tassi. 2011. The Matita Interactive Theorem Prover. In *CADE-23 (LNCS, Vol. 6803)*, Nikolaj Bjørner and Viorica Sofronie-Stokkermans (Eds.). Springer, 64–69.
- [2] Leo Bachmair, Nachum Dershowitz, and David A. Plaisted. 1989. Completion without failure. In *Rewriting Techniques—Resolution of Equations in Algebraic Structures*, Hassan Aït-Kaci and Maurice Nivat (Eds.), Vol. 2. Academic Press, 1–30.
- [3] Leo Bachmair and Harald Ganzinger. 1994. Rewrite-based equational theorem proving with selection and simplification. *J. Log. Comput.* 4, 3 (1994), 217–247.
- [4] Leo Bachmair and Harald Ganzinger. 2001. Resolution theorem proving. In *Handbook of Automated Reasoning*, Alan Robinson and Andrei Voronkov (Eds.), Vol. I. Elsevier and MIT Press, 19–99.
- [5] Leo Bachmair, Harald Ganzinger, and Uwe Waldmann. 1994. Refutational theorem proving for hierarchic first-order theories. *Appl. Algebra Eng. Commun. Comput.* 5 (1994), 193–212.
- [6] Clemens Ballarín. 2014. Locales: A module system for mathematical theories. *J. Autom. Reason.* 52, 2 (2014), 123–153.
- [7] Alexander Bentkamp, Jasmin Blanchette, Sophie Tournet, Petar Vukmirović, and Uwe Waldmann. 2019. Superposition with lambdas. In *CADE-27 (LNCS, Vol. 11716)*, Pascal Fontaine (Ed.). Springer, 55–73.
- [8] Alexander Bentkamp, Jasmin Christian Blanchette, Simon Cruanes, and Uwe Waldmann. 2018. Superposition for lambda-free higher-order logic. In *IJCAR 2018 (LNCS, Vol. 10900)*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). Springer, 28–46.
- [9] Yves Bertot and Pierre Castéran. 2004. *Interactive Theorem Proving and Program Development—Coq'Art: The Calculus of Inductive Constructions*. Springer.
- [10] Ahmed Bhayat and Giles Reger. 2020. A combinator-based superposition calculus for higher-order logic. In *IJCAR 2020, Part I (LNCS, Vol. 12166)*, Nicolas Peltier and Viorica Sofronie-Stokkermans (Eds.). Springer, 278–296.
- [11] Jasmin Blanchette and Sophie Tournet. 2020. Extensions to the comprehensive framework for saturation theorem proving. *Archive of Formal Proofs 2020 (2020)*. [https://www.isa-afp.org/entries/Saturation\\_Framework\\_Extensions.shtml](https://www.isa-afp.org/entries/Saturation_Framework_Extensions.shtml)
- [12] Jasmin Christian Blanchette. 2019. Formalizing the metatheory of logical calculi and automatic provers in Isabelle/HOL (invited talk). In *CPP 2019*, Assia Mahboubi and Magnus O. Myreen (Eds.). ACM, 1–13.
- [13] Jasmin Christian Blanchette, Mathias Fleury, Peter Lammich, and Christoph Weidenbach. 2018. A Verified SAT Solver Framework with Learn, Forget, Restart, and Incrementality. *J. Autom. Reason.* 61, 3 (2018), 333–366.
- [14] Ana Bove, Peter Dybjer, and Ulf Norell. 2009. A Brief Overview of Agda—A Functional Language with Dependent Types. In *TPHOLS 2009 (LNCS, Vol. 5674)*, Stefan Berghofer, Tobias Nipkow, Christian Urban, and Makarius Wenzel (Eds.). Springer, 73–78.
- [15] Leonardo de Moura, Soonho Kong, Jeremy Avigad, Floris van Doorn, and Jakob von Raumer. 2015. The Lean Theorem Prover (System Description). In *CADE-25 (LNCS, Vol. 9195)*, Amy P. Felty and Aart Middeldorp (Eds.). Springer, 378–388.
- [16] Arnaud Fietzke and Christoph Weidenbach. 2009. Labelled splitting. *Ann. Math. Artif. Intell.* 55, 1–2 (2009), 3–34.
- [17] Mathias Fleury. 2019. Optimizing a Verified SAT Solver. In *NFM 2019 (LNCS, Vol. 11460)*, Julia M. Badger and Kristin Yvonne Rozier (Eds.). Springer, 148–165.
- [18] Thomas Hillenbrand and Bernd Löchner. 2002. The Next WALDMEISTER Loop. In *CADE-18 (LNCS, Vol. 2392)*, Andrei Voronkov (Ed.). Springer, 486–500.
- [19] Matti Jarvisalo, Armin Biere, and Marijn Heule. 2010. Blocked Clause Elimination. In *TACAS 2010 (LNCS, Vol. 6015)*, Javier Esparza and Rupak Majumdar (Eds.). Springer, 129–144.
- [20] Laura Kovács and Andrei Voronkov. 2013. First-order theorem proving and Vampire. In *CAV 2013 (LNCS, Vol. 8044)*, Natasha Sharygina and Helmut Veith (Eds.). Springer, 1–35.
- [21] William McCune and Larry Wos. 1997. Otter—the CADE-13 competition incarnations. *J. Autom. Reason.* 18, 2 (1997), 211–220.
- [22] Lawrence C. Paulson and Jasmin Christian Blanchette. 2012. Three Years of Experience with Sledgehammer, a Practical Link Between Automatic and Interactive Theorem Provers. In *IWIL-2010 (EPIc, Vol. 2)*, Geoff Sutcliffe, Stephan Schulz, and Eugenia Ternovska (Eds.). EasyChair, 1–11.
- [23] Nicolas Peltier. 2016. A Variant of the Superposition Calculus. *Archive of Formal Proofs 2016 (2016)*. <https://www.isa-afp.org/entries/SuperCalc.shtml>
- [24] Anders Schlichtkrull, Jasmin Blanchette, Dmitriy Traytel, and Uwe Waldmann. 2020. Formalizing Bachmair and Ganzinger's Ordered Resolution Prover. *J. Autom. Reasoning (2020)*. online first.
- [25] Anders Schlichtkrull, Jasmin Christian Blanchette, and Dmitriy Traytel. 2019. A verified prover based on ordered resolution. In *CPP 2019*, Assia Mahboubi and Magnus O. Myreen (Eds.). ACM, 152–165.
- [26] Stephan Schulz, Simon Cruanes, and Petar Vukmirović. 2019. Faster, higher, stronger: E 2.3. In *CADE-27 (LNCS, Vol. 11716)*, Pascal Fontaine (Ed.). Springer, 495–507.

- [27] Alexander Steen and Christoph Benzmüller. 2018. The higher-order prover Leo-III. In *IJCAR 2018 (LNCS, Vol. 10900)*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). Springer, 108–116.
- [28] René Thiemann and Christian Sternagel. 2009. Certification of Termination Proofs Using CeTA. In *TPHOLS 2009 (LNCS, Vol. 5674)*, Stefan Berghofer, Tobias Nipkow, Christian Urban, and Makarius Wenzel (Eds.). Springer, 452–468.
- [29] Sophie Tourret. 2020. A comprehensive framework for saturation theorem proving. *Archive of Formal Proofs 2020* (2020). [https://www.isa-afp.org/entries/Saturation\\_Framework.shtml](https://www.isa-afp.org/entries/Saturation_Framework.shtml)
- [30] Andrei Voronkov. 2014. AVATAR: The Architecture for First-Order Theorem Provers. In *CAV 2014 (LNCS, Vol. 8559)*, Armin Biere and Roderick Bloem (Eds.). Springer, 696–710.
- [31] Uwe Waldmann. 2002. Cancellative abelian monoids and related structures in refutational theorem proving (part I). *J. Symb. Comput.* 33, 6 (2002), 777–829.
- [32] Uwe Waldmann, Sophie Tourret, Simon Robillard, and Jasmin Blanchette. 2020. A comprehensive framework for saturation theorem proving. In *IJCAR 2020, Part I (LNCS, Vol. 12166)*, Nicolas Peltier and Viorica Sofronie-Stokkermans (Eds.). Springer, 316–334.
- [33] Christoph Weidenbach, Dilyana Dimova, Arnaud Fietzke, Rohit Kumar, Martin Suda, and Patrick Wischniewski. 2009. SPASS version 3.5. In *CADE-22 (LNCS, Vol. 5663)*, Renate A. Schmidt (Ed.). Springer, 140–145.
- [34] Makarius Wenzel. 2007. Isabelle/Isar—a generic framework for human-readable proof documents. In *From Insight to Proof: Festschrift in Honour of Andrzej Trybulec*, Roman Matuszewski and Anna Zalewska (Eds.). Studies in Logic, Grammar, and Rhetoric, Vol. 10(23). University of Białystok.