

Verifying Saturation Provers Modularly

Anonymous Author(s)

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, completeness of various prover architectures implementing the calculus. It also helps derive completeness of calculi obtained by lifting a ground (i.e., variable-free) calculus. 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 [17], the first-order provers E [25], SPASS [31], and Vampire [19], and the higher-order provers Leo-III [26] and Zipperposition [7]. A saturation prover starts with a problem given in conjunctive normal form and systematically draws inferences from these 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 ultimately 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 index i represents the prover’s state at a point in time: If N_0 is unsatisfiable (i.e., provable), 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, which operate on a single clause set, and realistic provers with multiple clause sets. The *given clause procedure* [20], 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’s celebrations in Romania: 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 σ ; for example, $p(x)$ subsumes $p(a)$ and $p(x) \vee q(y)$. Regrettably, much of the saturation provers literature ignores subsumption and dynamic completeness. In the refinement chain leading to a verified prover, such a simplification turns 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, 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 non-ground lifting, subsumption, and the given clause procedure. The argument is hard to follow on paper, and even harder to formalize, as Schlichtkrull et al. found out [23].

To remedy this situation, Waldmann et al. [30] recently designed a comprehensive framework that applies to a wide range of calculi, including ordered resolution [4], unifying completion [2], standard superposition [3], theory superposition [29], hierarchic superposition [5], λ -free superposition [8], λ -superposition [7], and combinatory superposition [10]. This *saturation framework* can be used to lift ground static completeness to nonground static completeness and then 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.

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>

dynamic completeness of minimalistic proving processes or given clause provers with support for clause subsumption.

The framework was initially developed in L^AT_EX, but we also managed to formalize it in Isabelle/HOL in time to mention it briefly elsewhere, as four paragraphs in a conference paper. The present paper aims at introducing this work to the interactive theorem proving community. In addition to 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. The new proofs replace previous monolithic arguments.

We used the declarative Isar [32] in the hope that practitioners who are not Isabelle experts can follow 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 related to 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. The Isabelle theory files are available in the *Archive of Formal Proofs* and as an anonymized archive; 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 build on an abstract notion of formulas, which generalizes that of 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 = {⊥}`, 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`

`⊢ :: 'f set ⇒ 'f set ⇒ bool`

assumes

`Bot ≠ ∅ and`

`B ∈ Bot ⇒ {B} ⊢ N1 and`

`N2 ⊆ N1 ⇒ N1 ⊢ N2 and`

`(∀ C ∈ N2. N1 ⊢ {C}) ⇒ N1 ⊢ N2 and`

`N1 ⊢ N2 ⇒ N2 ⊢ N3 ⇒ N1 ⊢ N3`

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 ι 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 systems, 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 ⇒ 'f inference set where`

`Inf_from N = {ι ∈ Inf. prems_of ι ⊆ N}`

definition `Inf_between :: 'f set ⇒ 'f set ⇒ 'f inference set where`

`Inf_between N M =`

`Inf_from (N ∪ M) - Inf_from (N - M)`

The auxiliary function `Inf_from N` returns the set of possible inferences from premises in N . The variant `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 system indicates which inferences must a priori be drawn. The criterion specifies inferences that can a posteriori be omitted; it also identifies deletable formulas. Formally:

locale `calculus =`

`inference_system Inf +`

221 consequence_relation Bot (\models) +
 222 **fixes**
 223 $\text{Red}_I :: 'f \text{ set} \Rightarrow 'f \text{ inference set and}$
 224 $\text{Red}_F :: 'f \text{ set} \Rightarrow 'f \text{ set}$
 225 **assumes**
 226 $\text{Red}_I N \subseteq \text{Inf and}$
 227 $B \in \text{Bot} \Rightarrow N \models \{B\} \Rightarrow N - \text{Red}_F N \models \{B\} \text{ and}$
 228 $N \subseteq N' \Rightarrow \text{Red}_F N \subseteq \text{Red}_F N' \text{ and}$
 229 $N \subseteq N' \Rightarrow \text{Red}_I N \subseteq \text{Red}_I N' \text{ and}$
 230 $N' \subseteq \text{Red}_F N \Rightarrow \text{Red}_F N \subseteq \text{Red}_F (N - N') \text{ and}$
 231 $N' \subseteq \text{Red}_F N \Rightarrow \text{Red}_I N \subseteq \text{Red}_I (N - N') \text{ and}$
 232 $\iota \in \text{Inf} \Rightarrow \text{concl_of } \iota \in N \Rightarrow \iota \in \text{Red}_I N$

233 The locale inherits Inf (and the auxiliaries Inf_from and
 234 Inf_between) from inference_system as well as Bot and \models
 235 from consequence_relation. An inference ι is redundant if
 236 $\iota \in \text{Red}_I N$; a formula C is redundant if $C \in \text{Red}_F N$. The
 237 assumptions are needed to establish dynamic completeness.
 238 For example, the second assumption ensures that deleting a
 239 redundant formula from an unsatisfiable set preserves the
 240 set's unsatisfiability (i.e., its provability).

2.4 Static and Dynamic Completeness

243 Although dynamic completeness is the ultimate goal, static
 244 completeness is interesting because it is easier to formulate
 245 and prove. It is expressed in terms of a single formula set N ,
 246 which must be saturated. Saturation means that all inferences
 247 from N are redundant:

248 **definition** saturated :: $'f \text{ set} \Rightarrow \text{bool where}$
 249 saturated $N \longleftrightarrow \text{Inf_from } N \subseteq \text{Red}_I N$

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

252 **locale** statically_complete_calculus = calculus +
 253 **assumes** $B \in \text{Bot} \Rightarrow \text{saturated } N \Rightarrow N \models \{B\} \Rightarrow$
 254 $\exists B' \in \text{Bot}. B' \in N$

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

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

264 **inductive** $\triangleright :: 'f \text{ set} \Rightarrow 'f \text{ set} \Rightarrow \text{bool where}$
 265 $M - N \subseteq \text{Red}_F N \Rightarrow M \triangleright N$

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

276 We write chain (\triangleright) Ns to express that the sequence Ns is
 277 a derivation. The chain predicate is as in Schlichtkrull et al.
 278 Dynamic completeness applies only to fair derivations:

279 **definition** fair :: $'f \text{ set llist} \Rightarrow \text{bool where}$
 280 fair $Ns \longleftrightarrow$
 281 $\text{Inf_from } (\text{Liminf } Ns) \subseteq \text{Sup } (\text{Imap } \text{Red}_I Ns)$

282 Above, Liminf Ns denotes the limit inferior $\bigcup_i \bigcap_{j \geq i} Ns ! j$,
 283 and Sup Ns denotes the supremum $\bigcup_i Ns ! i$. The infix op-
 284 erator $!$ returns the sequence element at the given index. In
 285 accordance with the automated reasoning literature, we will
 286 refer to Liminf Ns as the *limit* of Ns .

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

292 Dynamic completeness is formulated in its own locale:

293 **locale** dynamically_complete_calculus = calculus +
 294 **assumes**
 295 $B \in \text{Bot} \Rightarrow \text{chain } (\triangleright) Ns \Rightarrow \text{fair } Ns \Rightarrow$
 296 $Ns ! 0 \models \{B\} \Rightarrow$
 297 $\exists i < \text{llength } Ns. \exists B' \in \text{Bot}. B' \in Ns ! i$

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

300 **sublocale** statically_complete_calculus
 301 \subseteq dynamically_complete_calculus
 302 **sublocale** dynamically_complete_calculus
 303 \subseteq statically_complete_calculus

304 The above definitions are all based on Waldmann et al.,
 305 which in turn largely follow Bachmair and Ganzinger while
 306 generalizing them. The literature is not entirely consistent
 307 on notions such as redundancy, saturation, and fairness.
 308 Waldmann et al. studied the variation and showed that com-
 309 pleteness and “reduced” completeness coincide. We also for-
 310 malized these proofs (in Calculus_Variations.thy). They
 311 make it possible to mix and match results. For example, sup-
 312 pose that a calculus is proved “reducedly” statically complete
 313 in the literature. Instead of redoing the proof for standard
 314 static completeness, we can reuse the reduced one as is and
 315 derive dynamic completeness using the framework.

2.5 Intersection of Redundancy Criteria

316 There are various situations where a prover may need to
 317 apply several redundancy criteria in parallel and identify an
 318 inference or formula as redundant only if all the criteria agree.
 319 This can be incorporated into the framework as a single
 320 intersection criterion. One scenario where this is useful is
 321 when redundancy is parameterized by a value $q :: 'q$ that
 322 is only known at the limit. The intersection amounts to a
 323 universal quantification “for all $q \in Q, \dots$ ” Once we know
 324 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 $Q :: 'q\ set$ and by a Q -indexed family of consequence relations $\text{entails}_q :: 'q \Rightarrow 'f\ set \Rightarrow 'f\ set \Rightarrow bool$ 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 consequence relation \models_Q as the intersection of the entails_q 's:

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

The result 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
 $\text{Red}_{Iq} :: 'q \Rightarrow 'f\ set \Rightarrow 'f\ inference\ set$ **and**
 $\text{Red}_{Fq} :: 'q \Rightarrow 'f\ set \Rightarrow 'f\ set$
assumes
 $Q \neq \emptyset$ **and**
 $\forall q \in Q. \text{calculus}\ \text{Bot}\ \text{Inf}\ (\text{entails}_q\ q)\ (\text{Red}_{Iq}\ q)\$
 $(\text{Red}_{Fq}\ q)$

The redundancy criterion is defined as a pair of intersections:

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

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

Intersection calculi qualify as calculi with the intersections Red_I and 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*:
 $(\forall N. \text{saturated}\ \text{Inf}\ \text{Red}_I\ N \wedge (\forall B \in \text{Bot}. B \notin N) \longrightarrow$
 $\exists B \in \text{Bot}. \exists q \in Q. \neg \text{entails}_q\ q\ N\ \{B\}) \implies$
 $\text{statically_complete_calculus}\ \text{Bot}\ \text{Inf}\ (\models_Q)\ \text{Red}_I\ \text{Red}_F$

2.6 Soundness

An inference system `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*:
 $\text{chain}\ (\triangleright_{\text{Red}})\ Ns \implies \text{chain}\ (\models)\ Ns \implies$
 $(\text{Liminf}\ Ns \models \text{Bot} \iff Ns\ !\ 0 \models \text{Bot})$

This result is stated in a locale requiring compactness, meaning that $N \models \text{Bot} \implies \exists M \subseteq N. \text{finite}\ M \wedge M \models \text{Bot}$. Informally, the proof 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. Thus, by transitivity of \models , the initial set $Ns\ !\ 0$ must be unsatisfiable.

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 `Inf` operating on formulas equipped with a well-order $<$ is counterexample-reducing if the following locale requirements can be met:

fixes $\text{l_of} :: 'f\ set \Rightarrow 'f\ set$
assumes
 $N \cap \text{Bot} = \emptyset \implies D \in N \implies \neg \text{l_of}\ N \models D \implies$
 $(\bigwedge C. C \in N \implies \neg \text{l_of}\ N \models C \implies D \leq C) \implies$
 $\exists \iota \in \text{Inf}. \text{prems_of}\ \iota \neq [] \wedge \text{main_prem_of}\ \iota = D$
 $\wedge \text{side_prems_of}\ \iota \subseteq N$
 $\wedge \text{l_of}\ N \models \text{side_prems_of}\ \iota$
 $\wedge \neg \text{l_of}\ N \models \text{concl_of}\ \iota \wedge \text{concl_of}\ \iota < D$

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. Given a formula set, the `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 bring them very far, however, because they could never connect that definition of `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 $\text{Red}_I :: 'f \text{ set} \Rightarrow 'f \text{ inference set}$ **where**

$$\begin{aligned} \text{Red}_I N = \{ & \iota \in \text{Inf}. \exists M \subseteq N. \\ & M \cup \text{side_prems_of } \iota \models \{\text{concl_of } \iota\} \\ & \wedge \forall D \in M. D < \text{main_prem_of } \iota \} \end{aligned}$$

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

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

For inference systems 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 as a finite multisets of literals introduced by Schlichtkrull et al. [23, Section 2.3]. The empty clause \perp is represented by \emptyset , but we will continue writing \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 $\text{J_of} :: 'a \text{ clause set} \Rightarrow 'a \text{ set}$

assumes

$$\begin{aligned} \perp \notin N & \Rightarrow D \in N \Rightarrow \neg \text{J_of } N \models D \Rightarrow \\ (\wedge C. C \in N & \Rightarrow \neg \text{J_of } N \models C \Rightarrow D \leq C) \Rightarrow \\ \exists \iota \in \text{Inf}. & \text{prems_of } \iota \neq [] \wedge \text{main_prem_of } \iota = D \\ & \wedge \text{side_prems_of } \iota \subseteq N \\ & \wedge \text{J_of } N \models \text{side_prems_of } \iota \\ & \wedge \neg \text{J_of } N \models \text{concl_of } \iota \wedge \text{concl_of } \iota < D \end{aligned}$$

For this locale, we could represent a candidate model more standardly as a set of true atoms, and we could use the modelhood relation \models instead of entailment \models of clauses.

We connected the two formulations of counterexample reduction so that end users do not have to. 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: $\text{saturated } N \Rightarrow \perp \notin N \Rightarrow \text{J_of } N \models N$.

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}_I 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 =`

$$\begin{aligned} & \text{inference_system } \text{Inf}_F + \\ & \text{calculus } \text{Bot}_G \text{ Inf}_G (\models) \text{Red}_{IG} \text{Red}_{FG} + \end{aligned}$$

fixes

$\mathcal{G}_F :: 'f \Rightarrow 'g \text{ set}$ **and**

$\mathcal{G}_I :: 'f \text{ inference} \Rightarrow 'g \text{ inference set option}$

assumes

$\text{Bot}_F \neq \emptyset$ **and**

$B \in \text{Bot}_F \Rightarrow \mathcal{G}_F B \neq \emptyset$ **and**

$B \in \text{Bot}_F \Rightarrow \mathcal{G}_F B \subseteq \text{Bot}_G$ **and**

$\mathcal{G}_F C \cap \text{Bot}_G \neq \emptyset \rightarrow C \in \text{Bot}_F$ **and**

$\iota \in \text{Inf}_F \Rightarrow \mathcal{G}_I \iota \neq \text{None} \Rightarrow$

$\text{the } (\mathcal{G}_I \iota) \subseteq \text{Red}_{IG} (\mathcal{G}_F (\text{concl_of } \iota))$

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 $\mathcal{G}_{Fset} :: 'fset \Rightarrow 'gset$ 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 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**

$$\text{Red}_I N = \{\iota \in \text{Inf}_F.$$

$$(\mathcal{G}_I \iota \neq \text{None} \wedge \text{the } (\mathcal{G}_I \iota) \subseteq \text{Red}_{IG} (\mathcal{G}_{Fset} N))$$

$$\vee (\mathcal{G}_I \iota = \text{None}$$

$$\wedge \mathcal{G}_F (\text{concl_of } \iota) \subseteq \mathcal{G}_{Fset} N \cup \text{Red}_{FG} (\mathcal{G}_{Fset} N)\}$$

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 BotF InfF (\models_G) RedI RedF`

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

`statically_complete_calculus BotG InfG (\models)`

`RedIG RedFG` **and**

`$\wedge N. \text{saturated } N \Rightarrow \text{ground_Inf_redundant } N$`

shows

`statically_complete_calculus BotF InfF (\models_G)`

`RedI RedF`

3.2 Adding Tiebreaker Orders

The completeness proof for standard superposition [3] and for 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. The idea is to use a well-founded order, or more generally a family of orders, which they call a tiebreaking order. The tiebreaker can be instantiated with subsumption. It can also serve other purposes; for example, in a given clause procedure, it can 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 Prec_{FG} of well-founded orders:

```

locale tiebreaker_lifting =
  std : standard_lifting BotF InfF BotG InfG (≡)
  RedIG RedFG  $\mathcal{G}_F$   $\mathcal{G}_I$  +
fixes PrecFG :: 'g ⇒ 'f ⇒ 'f ⇒ bool
assumes minimal_element (PrecFG g) UNIV

```

Redundancy of inferences Red_I is as for the standard lifting, but redundancy of formulas can use the tiebreaker:

```

definition RedF :: 'f set ⇒ 'f set where
  RedF N = {C. ∀D ∈  $\mathcal{G}_F$  C. D ∈ RedFG ( $\mathcal{G}_{\text{Fset}}$  N)
    ∨ ∃E ∈ N. PrecFG D E C ∧ D ∈  $\mathcal{G}_F$  E}

```

Based on this definition, we proved that `tiebreaker_lifting` produces a nonground calculus.

In Waldmann et al's paper, the proof that `standard_lifting` produces a nonground calculus is only mentioned as a simplification of the same proof for `tiebreaker_lifting`. To follow this structure, we added to `standard_lifting` an **interpretation** of `tiebreaker_lifting` where `PrecFG` is replaced by " $\lambda g C C'$. False", the empty order family. Relying on an **interpretation** instead of a **sublocale** conveniently avoids the issue of forbidden cyclic dependencies between locales, and is sufficient to derive the desired nonground calculus result in `standard_lifting` as a one-liner, avoiding a tedious copy-paste-simplify work on a 200 lines proof. Thanks to this construction, it was a breeze to show that the tiebreaker order can be ignored when proving static completeness and recovered for dynamic completeness, as Waldmann et al. do.

```

theorem static_to_dynamic:
  statically_complete_calculus BotF InfF (≡G)
  RedI std.RedF ⟷
  dynamically_complete_calculus BotF InfF (≡G)
  RedI RedF

```

The theorem holds because the definitions of saturation, fairness, and completeness depend on Red_I but not on 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 following the same general approach. The locale declaration is as follows:

```

locale lifting_intersection =
  inference_system InfF +
  inference_system_family Q InfGq +
  consequence_relation_family BotG Q entailsq +
fixes
  BotF :: 'f set and
   $\mathcal{G}_{\text{Fq}}$  :: 'q ⇒ 'f ⇒ 'g set and
   $\mathcal{G}_{\text{Iq}}$  :: 'q ⇒ 'f inference ⇒ 'g inference set option and
  PrecFG :: 'g ⇒ 'f ⇒ 'f ⇒ bool
assumes
  ∀q ∈ Q. tiebreaker_lifting BotF InfF BotG
    (entailsq q) (InfGq q) (RedIq q) (RedFq q)
    ( $\mathcal{G}_{\text{Fq}}$  q) ( $\mathcal{G}_{\text{Iq}}$  q) PrecFG

```

For each index, we defined redundancy with and without tiebreaking order ($\text{Red}_{I\mathcal{G}_q}$, and resp. $\text{Red}_{F\mathcal{G}_q}$ and $\text{Red}_{F\mathcal{G}_q\emptyset}$) and entailment ($\text{entails}_{\mathcal{G}_q}$) such that the intersection of the corresponding nonground calculi are also calculi:

```

sublocale intersection_calculus BotF InfF Q entailsq
  RedIq RedFq
sublocale empty_ord: intersection_calculus BotF InfF Q
  entailsq RedIq RedFq RedFq

```

These commands introduce intersected notions of entailment ($\text{≡}_{\cap\mathcal{G}}$) and redundancies (Red_{IG} , and resp. Red_{FG} and $\text{empty_ord.Red}_{FG}$).

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`:

```

theorem stat_eq_dyn_ref_comp_fam_inter:
  statically_complete_calculus BotF InfF (≡∩G)
  RedI empty_ord.RedFG ⟷
  dynamically_complete_calculus BotF InfF (≡∩G)
  RedI RedF

```

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 Inf_{FL} such that whenever an Inf_F inference is possible, an Inf_{FL} is also possible, and vice versa:

```

locale labeled_tiebreaking_lifting =
  no_labels: tiebreaker_lifting BotF InfF BotG (≡)

```

661 $\text{Inf}_G \text{Red}_{IG} \text{Red}_{FG} \mathcal{G}_F \mathcal{G}_I \text{Prec}_{FG} +$
 662 **fixes** $\text{Inf}_{FL} :: ('f \times 'l) \text{ inference set}$
 663 **assumes**
 664 $\iota_F \in \text{Inf}_F \implies \text{length } Ls = \text{length } (\text{prems_of } \iota_F) \implies$
 665 $\exists L_0. \text{Infer } (\text{zip } (\text{prems_of } \iota_F) Ls)$
 666 $(\text{concl_of } \iota_F, L_0) \in \text{Inf}_{FL} \text{ and}$
 667 $\iota_{FL} \in \text{Inf}_{FL} \implies \text{Infer } (\text{map } \text{fst } (\text{prems_of } \iota_{FL}))$
 668 $(\text{fst } (\text{concl_of } \iota_{FL})) \in \text{Inf}_F$

669 $\text{Bot}_{FL} :: ('f \times 'l) \text{ set}$ is defined as $\text{Bot}_F \times \text{UNIV}$. The ground-
 670 ing functions (\mathcal{G}_{IL} and \mathcal{G}_{FL}) and the entailment relation ($\models_{\mathcal{G}_L}$)
 671 simply ignore labels. The main result lets us derive the static
 672 completeness of the labeled calculus from the static com-
 673 pleteness of the base calculus:
 674

675 **lemma** *stat_ref_comp_to_labeled_sta_ref_comp*:
 676 $\text{statically_complete_calculus } \text{Bot}_F \text{Inf}_F (\models_{\mathcal{G}})$
 677 $\text{no_labels.Red}_I \text{no_labels.Red}_F \implies$
 678 $\text{statically_complete_calculus } \text{Bot}_{FL} \text{Inf}_{FL} (\models_{\mathcal{G}_L})$
 679 $\text{Red}_I \text{Red}_F$

680 Labels can also be added to calculi that work with fami-
 681 lies of redundancy criteria. Since completeness results are
 682 independent of the tiebreaker order, we only considered the
 683 empty tiebreaker order:
 684

685 **locale** *labeled_lifting_intersection* =
 686 $\text{no_labels: lifting_intersection } \text{Inf}_F \text{Bot}_G \text{Q } \text{Inf}_{Gq}$
 687 $\text{entails}_q \text{Red}_{Iq} \text{Red}_{Fq} \text{Bot}_F \mathcal{G}_{Fq} \mathcal{G}_{Iq}$
 688 $(\lambda g \text{ CL } CL'. \text{False}) +$
 689 **fixes** $\text{Inf}_{FL} :: ('f \times 'l) \text{ inference set}$

690 We proved that static completeness with labels can be re-
 691 duced to static completeness without labels.
 692

693 4 Prover Architectures

694 Dynamic completeness is formulated in terms of an abstract
 695 single-formula-set prover \triangleright (Section 2.4). In practice, satura-
 696 tion provers implement the given clause procedure described
 697 in the introduction. This simplifies bookkeeping: Instead of
 698 having to keep track of all possible n -ary inferences, actual
 699 provers need only to track active formulas.
 700

701 Waldmann et al. describe two versions of the given
 702 clause procedure: an eager version called GC and a lazy
 703 version called LGC. LGC strictly generalizes GC, but be-
 704 cause GC is sufficient for most applications and simpler
 705 to use, we formalized both separately (in `Given-Clause-
 706 Architectures.thy` and `Given-Clause-Architectures-
 707 Revisited.thy`).
 708

709 4.1 Basic Setup

710 Much of the setup between GC and LGC can be shared,
 711 in a locale called `given_clause_basis` that extends `labeled_
 712 lifting_intersection`. The locale assumes that the unlabeled
 713 version of the intersection calculus is complete:
 714

715

assumes
 716 $\text{statically_complete_calculus } \text{Bot}_F \text{Inf}_F (\models_{FG})$
 717 $\text{no_labels.Red}_I \text{no_labels.Red}_F$
 718

719 It introduces an equivalence relation \doteq and an order \prec on
 720 unlabeled formulas that must be compatible with each other
 721 and with the grounding function:
 722

assumes
 723 $C_1 \doteq D_1 \implies C_2 \doteq D_2 \implies C_1 \prec C_2 \implies D_1 \prec D_2 \text{ and}$
 724 $q \in \text{Q} \implies C_1 \doteq C_2 \implies \mathcal{G}_{Fq} q C_1 \implies \mathcal{G}_{Fq} q C_2 \text{ and}$
 725 $q \in \text{Q} \implies C_2 \prec C_1 \implies \mathcal{G}_{Fq} q C_1 \implies \mathcal{G}_{Fq} q C_2$
 726

727 It also introduces an order \sqsubset_L on labels and a distinguished
 728 label active that must be minimal w.r.t. \sqsubset_L . Both orders are
 729 assumed to be well founded. Inside the locale, we used these
 730 relations to define a tiebreaker order \sqsubset on labeled formulas
 731 and showed it well founded:
 732

definition $\sqsubset :: ('f \times 'l) \implies ('f \times 'l) \implies \text{bool}$ **where**
 733 $CL_1 \sqsubset CL_2 \iff \text{fst } CL_1 \prec \text{fst } CL_2$
 734 $\vee (\text{fst } CL_1 \doteq \text{fst } CL_2 \wedge \text{snd } CL_1 \sqsubset_L \text{snd } CL_2)$
 735

lemma *wf_prec_FL*: $\text{minimal_element } (\sqsubset) \text{UNIV}$
 736

737 This order can be used to strengthen the redundancy cri-
 738 terion of the labeled intersection calculus. Below, Red_I and
 739 Red_F refer to the strengthened criterion.
 740

741 Finally, we present two lemmas that make explicit the rela-
 742 tion between the labeled and unlabeled redundancy criteria:
 743

lemma *labeled_red_inf_eq_red_inf*:
 744 $\iota \in \text{Inf}_{FL} \implies \iota \in \text{Red}_I N \iff$
 745 $\text{to}_F \iota \in \text{no_labels.Red}_I (\text{fst } \iota N)$

lemma *red_labeled_clauses*:

746 $C \in \text{no_labels.Red}_F (\text{fst } \iota N) \vee (\exists C' \in \text{fst } \iota N. C' \prec C)$
 747 $\vee (\exists (C', L') \in N. L' \sqsubset_L L \wedge C' \prec C) \implies$
 748 $(C, L) \in \text{Red}_F N$
 749

750 Above, to_F denotes a function that erases all labels, and the
 751 infix operator ι denotes the image of a set under a function.
 752

753 Although the first lemma is obvious when we compare the
 754 definitions of Red_I and no_labels.Red_I , the corresponding
 755 Isabelle proof is over 100 lines long. This is an unfortunate
 756 consequence of the deep nesting of locales. The identity of
 757 notions is hidden under multiple layers of names, and Isabelle
 758 provides no automation that can unfold these definitions au-
 759 tomatically so as to uncover such connections. We conjecture
 760 that proof assistants based on computational type theories
 761 such as Agda [13], Coq [9], Lean [14], and Matita [1] would
 762 cope more gracefully with these definitional equalities.
 763

764 The `given_clause_basis` locale also ensures that inferences
 765 never produce active formulas. Only the procedure itself is
 766 allowed to make a formula active.
 767

768 When we developed this locale and the two locales based
 769 on it, we did not immediately try to instantiate them. This
 770 led to an amusing incident when we attempted to verify RP.
 771 Of course, the `given_clause_basis` locale axioms were flawed.
 772 First, we had stated that \succ must be well founded when we
 773

771 meant \prec . Strangely enough, we had successfully proved \sqsubseteq
 772 well founded, even though it is based on \prec .

773 How could that be? The answer is that we had relied on
 774 very powerful proof automation: Sledgehammer [21]. That
 775 tool had discovered *another* flawed assumption and exploited
 776 it. This time, we had misspelled a constant's name. In Isabelle,
 777 this creates an implicitly universal variable. This convention,
 778 which certainly saves a lot of typing and is considered by
 779 some as an advantage of Isabelle over its rivals, had led to
 780 an inconsistent assumption.

781 4.2 The Eager Given Clause Procedure

782 The eager given clause procedure GC is described as a tran-
 783 sition system that operates on labeled formula sets. This
 784 corresponds naturally to an inductive predicate \rightsquigarrow_{GC} :
 785

786 **inductive** $\rightsquigarrow_{GC} :: ('f \times 'l) \text{ set} \Rightarrow ('f \times 'l) \text{ set} \Rightarrow \text{bool}$
 787 **where**
 788 $\text{process: } N_1 = N \cup M \Rightarrow N_2 = N \cup M' \Rightarrow$
 789 $M \subseteq \text{Red}_F(N \cup M') \Rightarrow \text{active_subset } M' = \emptyset \Rightarrow$
 790 $N_1 \rightsquigarrow_{GC} N_2$
 791 $| \text{infer: } N_1 = N \cup \{(C, L)\} \Rightarrow$
 792 $N_2 = N \cup \{(C, \text{active})\} \cup M \Rightarrow L \neq \text{active} \Rightarrow$
 793 $\text{active_subset } M = \emptyset \Rightarrow$
 794 $\text{no_labels.Inf_between } (\text{fst } ' \text{ active_subset } N) \{C\}$
 795 $\subseteq \text{no_labels.Red}_1 (\text{fst } '(N \cup (C, \text{active}) \cup M)) \Rightarrow$
 796 $N_1 \rightsquigarrow_{GC} N_2$

797 The *process* rule can be used to add arbitrary passive for-
 798 mulas M' or to remove redundant formulas M . For example,
 799 the simplification of $p(2 + 2)$ to $p(4)$ would be modeled by
 800 taking $M := \{p(2 + 2)\}$ and $M' := \{p(4)\}$. The *infer* rule se-
 801 lects a passive formula C —called the given formula or given
 802 clause—and makes it active. In addition, it draws all infer-
 803 ences between C and all the active formulas; more precisely,
 804 the rule requires the addition of enough passive formulas M
 805 to make these inferences redundant. The locale assumes that
 806 the inference system Inf_F contains no nullary inferences.

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

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

816 **lemma** *gc_fair*:
 817 **assumes**
 818 $\text{chain } (\rightsquigarrow_{GC}) \text{ } Ns$ **and**
 819 $\text{active_subset } (Ns ! 0) = \emptyset$ **and**
 820 $\text{passive_subset } (\text{Liminf } Ns) = \emptyset$
 821 **shows** $\text{fair } Ns$

822 The proof by Waldmann is 13 lines long, but we needed
 823 197 lines of Isabelle. The explanation for this poor De Bruijn
 824

825 factor is that the proof considers the maximal element of a
 826 set of integers, each representing a step in the derivation
 827 that is the first to satisfy a complex property for one of the
 828 premises of a given inference. The previous sentence already
 829 hints at the number of hidden reasoning steps necessary
 830 to obtain this particular element, although it is fairly easy
 831 to write the notion mathematically in a way that enables
 832 readers to build an adequate mental representation of the
 833 object. In Isabelle, however, every intermediate object in this
 834 construction must be showed to exist.
 835

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

845 We were not satisfied with this proof. A proof that is te-
 846 dious to write is also tedious to read, and hides its insights
 847 from the reader. The situation reminded us of a 700-line Isa-
 848 belle formalization by Fleury of a 8-line pen-and-paper proof
 849 by Weidenbach [12, Section 4.3]. The solution in that case
 850 was to reason by invariance.

851 Can invariance help prove *gc_fair* in a more perspicuous
 852 way? It is widely known that the given clause procedure
 853 satisfies the following key invariant:

854 All inferences from active formulas are redun-
 855 dant w.r.t. the current set of formulas.

856 Initially, all formulas are passive, so the invariant holds vac-
 857 uously. Whenever a given formula is made active, all in-
 858 ferences with active partners are made redundant, thereby
 859 maintaining the invariant. And if all formulas are active at
 860 the limit, then this means that all inferences from the en-
 861 tire limit N are redundant. Hence the derivation is fair and
 862 consequently N is saturated.

863 Paradoxically, although the invariant above is widely rec-
 864 ognized as an important property of the given clause pro-
 865 cedure, which is itself central in automated reasoning, we
 866 are not aware of any proofs in the literature that rely on it.
 867 Instead, the invariant is presented as a guide to intuition, if
 868 at all. For example, Bachmair and Ganzinger do not mention
 869 the invariant in their *Handbook* chapter.
 870

871 In Isabelle, we defined the invariant as follows:

872 **definition** *gc_invar* :: $('f \times 'l) \text{ set } llist \Rightarrow \text{enat} \Rightarrow \text{bool}$
 873 **where**
 874 $\text{gc_invar } Ns \ i \longleftrightarrow$
 875 $\text{Inf_from } (\text{active_subset } (\text{Liminf_upto } Ns \ i))$
 876 $\subseteq \text{Sup_upto } (\text{Imap } \text{Red}_{IG} \ Ns) \ i$
 877

878 The predicate *gc_invar* $Ns \ i$ determines whether the invari-
 879 ant holds at index i (which can be ∞). If k is within bounds,
 880

881 $\text{Liminf_upto } Ns \ k$ is defined as the bounded limit inferior
 882 $\bigcup_i \bigcap_{j=i}^k Ns \ ! \ j$, which is equal to the limit if $k = \infty$ and to
 883 $Ns \ ! \ k$ otherwise. Similarly, $\text{Sup_upto } Ns \ k$ is the bounded
 884 supremum $\bigcup_i^k Ns \ ! \ i$.

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

887 **lemma** *gc_invar_gc*:

888 $\text{chain } (\rightsquigarrow_{GC}) \ Ns \ \Longrightarrow \ \text{active_subset } (Ns \ ! \ 0) = \emptyset \ \Longrightarrow$
 889 $i < \text{llength } Ns \ \Longrightarrow \ \text{gc_invar } Ns \ i$

890 **lemma** *gc_invar_infinity*:

891 $Ns \neq [] \ \Longrightarrow \ (\forall i < \text{llength } Ns. \ \text{gc_invar } Ns \ i) \ \Longrightarrow$
 892 $\text{gc_invar } Ns \ \infty$

893 Assuming that all formulas eventually become active, the
 894 invariant at the limit simplifies to $\text{Inf_from } (\text{Liminf } Ns) \subseteq$
 895 $\text{Sup } (\text{lmap } \text{Red}_{LGC} \ Ns)$. This corresponds exactly to the defini-
 896 tion of fairness, which is what we want.

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

905 4.3 The Given Clause Procedure with Delayed 906 Inferences

907 The lazy given clause procedure LGC decouples inference
 908 scheduling and computation. Each inference ι is first added
 909 to a “to do” set T of scheduled inferences. At some later
 910 point, ι is computed, meaning that it is removed from T
 911 and its conclusion is added to the set N of labeled formulas.
 912 Alternatively, ι might be identified as redundant and deleted,
 913 a technique with the unfortunate name of “orphan deletion.”
 914 The formal definition of LGC follows:

915 **inductive**

916 $\rightsquigarrow_{LGC} :: ('f \ \text{inference set} \times ('f \times 'l) \ \text{set}) \ \text{llist} \Rightarrow \ \text{bool}$

917 **where**

918 $\text{process}: N_1 = N \cup M \ \Longrightarrow \ N_2 = N \cup M' \ \Longrightarrow$

919 $M \subseteq \text{Red}_F (N \cup M') \ \Longrightarrow \ \text{active_subset } M' = \emptyset \ \Longrightarrow$
 920 $(T, N_1) \rightsquigarrow_{LGC} (T, N_2)$

921 | $\text{schedule_infer}: T_2 = T_1 \cup T' \ \Longrightarrow \ N_1 = N \cup \{(C, L)\} \ \Longrightarrow$

922 $N_2 = N \cup \{(C, \text{active})\} \ \Longrightarrow \ L \neq \text{active} \ \Longrightarrow$

923 $T' = \text{no_labels. Inf_between}$

924 $(\text{fst } \text{active_subset } N) \ \{C\} \ \Longrightarrow$

925 $(T_1, N_1) \rightsquigarrow_{LGC} (T_2, N_2)$

926 | $\text{compute_infer}: T_1 = T_2 \cup \{\iota\} \ \Longrightarrow \ N_2 = N_1 \cup M \ \Longrightarrow$

927 $\text{active_subset } M = \emptyset \ \Longrightarrow$

928 $\iota \in \text{no_labels.Red}_1 (\text{fst } \text{active_subset } N) \ \Longrightarrow$

929 $(T_1, N_1) \rightsquigarrow_{LGC} (T_2, N_2)$

930 | $\text{delete_orphans}: T_1 = T_2 \cup T' \ \Longrightarrow \ T' \cap$

931 $\text{no_labels. Inf_from } (\text{fst } \text{active_subset } N) = \emptyset \ \Longrightarrow$

932 $(T_1, N) \rightsquigarrow_{LGC} (T_2, N)$

933 LGC has several benefits beyond orphan deletion. The “to
 934 do” set T can be used to support nullary inferences, or axioms.
 935 It can be used to represent infinite streams of inferences,
 936 which might for example result from higher-order unification
 937 (as in Zipperposition). Finally, inferences can often be stored
 938 more compactly than their conclusions, so LGC can be used
 939 to enable a low-level optimization (as in Waldmeister).
 940 The key result is to show that \rightsquigarrow_{LGC} -derivations are fair
 941 under some reasonable assumptions:

942 **lemma** *lgc_fair*:

943 **assumes**

944 $\text{chain } (\rightsquigarrow_{LGC}) \ TNs \ \text{and}$

945 $\text{active_subset } (\text{snd } (TNs \ ! \ 0)) = \emptyset \ \text{and}$

946 $\text{passive_subset } (\text{Liminf } (\text{lmap } \text{snd } TNs)) = \emptyset \ \text{and}$

947 $\forall \iota \in \text{Inf}_F. \ \text{prems_of } \iota = [] \ \longrightarrow \ \iota \in \text{fst } (TNs \ ! \ 0) \ \text{and}$

948 $\text{Liminf } (\text{lmap } \text{fst } TNs) = \emptyset$

949 **shows** $\text{fair } (\text{lmap } \text{snd } TNs)$

950 Compared with \rightsquigarrow_{GC} , an extra condition is added to ensure
 951 that all premise-free inferences are scheduled up front. The
 952 proof has the same structure as that of *gc_fair*, but reason-
 953 ing steps are more complex due to the extra rules and the
 954 presence of scheduled inferences. We needed 326 lines.

955 Can invariance help re-prove *lgc_fair* in the same way
 956 that it helped us with *gc_fair*? The lazy given clause proce-
 957 dure is hardly discussed in the literature and was specified
 958 rigorously only by Waldmann et al., but it is easy to think
 959 up a suitable generalization of the previous invariant:

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

963 Initially, all formulas are passive, and all nullary inferences
 964 are in T , so the invariant holds. Whenever a given formula is
 965 made active, all inferences with active partners are scheduled,
 966 thereby maintaining the invariant. Removing an inference
 967 from T and computing it makes it redundant, preserving the
 968 invariant. Orphan deletion also preserves the invariant, be-
 969 cause by definition an orphan is not an inference from active
 970 formulas. And if T is empty and all formulas are active at
 971 the limit, then all inferences from the limit N are redundant,
 972 as required by fairness.

973 We defined the invariant as follows in Isabelle:

974 **definition** *lgc_invar* ::

975 $('f \ \text{inference set} \times ('f \times 'l) \ \text{set}) \ \text{llist} \Rightarrow \ \text{enat} \Rightarrow \ \text{bool}$

976 **where**

977 $\text{lgc_invar } TNs \ i \ \longleftrightarrow$

978 $\text{Inf_from } (\text{active_subset}$

979 $(\text{Liminf_upto } (\text{lmap } \text{snd } TNs) \ i))$

980 $\subseteq \text{from_F } \text{Liminf_upto } (\text{lmap } \text{fst } TNs) \ i$

981 $\cup \text{Sup_upto } (\text{lmap } (\text{Red}_{LGC} \circ \text{snd}) \ TNs) \ i$

982 The formal proof by invariance resembles that for GC but
 983 with more cases to consider. Compared with the monolithic
 984 proof, we have a small improvement in number of lines (301
 985

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. [30, 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 its formalization by Schlichtkrull et al. We lifted ground static completeness explicitly and re-proved 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

$$\frac{(C_i \vee A_{i1} \vee \dots \vee A_{ik_i})_{i=1}^n \quad \neg A_1 \vee \dots \vee \neg A_n \vee D}{(C_1 \vee \dots \vee C_n \vee D)\sigma}$$

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

RP is defined as a transition system $\rightsquigarrow_{\text{RP}}$ 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 $\rightsquigarrow_{\text{RP}} :: 'a \text{ state} \Rightarrow 'a \text{ state} \Rightarrow \text{bool}$ **where**
 $\{\text{Neg } A, \text{Pos } A\} \subseteq C \Rightarrow (N \cup \{C\}, P, O) \rightsquigarrow_{\text{RP}} (N, P, O)$
 $| D \in P \cup O \Rightarrow \text{subsumes } D C \Rightarrow$
 $(N \cup \{C\}, P, O) \rightsquigarrow_{\text{RP}} (N, P, O)$
 $| D \in N \Rightarrow \text{strict_subsumes } D C \Rightarrow$
 $(N, P \cup \{C\}, O) \rightsquigarrow_{\text{RP}} (N, P, O)$
 $| D \in N \Rightarrow \text{strict_subsumes } D C \Rightarrow$
 $(N, P, O \cup \{C\}) \rightsquigarrow_{\text{RP}} (N, P, O)$

$| D \in P \cup O \Rightarrow \text{reduces } D C L \Rightarrow$ 1046
 $(N \cup \{C \uplus \{L\}\}, P, O) \rightsquigarrow_{\text{RP}} (N \cup \{C\}, P, O)$ 1047
 $| D \in N \Rightarrow \text{reduces } D C L \Rightarrow$ 1048
 $(N, P \cup \{C \uplus \{L\}\}, O) \rightsquigarrow_{\text{RP}} (N, P \cup \{C\}, O)$ 1049
 $| D \in N \Rightarrow \text{reduces } D C L \Rightarrow$ 1050
 $(N, P, O \cup \{C \uplus \{L\}\}) \rightsquigarrow_{\text{RP}} (N, P \cup \{C\}, O)$ 1051
 $| (N \cup \{C\}, P, O) \rightsquigarrow_{\text{RP}} (N, P \cup \{C\}, O)$ 1052
 $| (\emptyset, P \cup \{C\}, O) \rightsquigarrow_{\text{RP}}$ 1053
 $(\text{concl_of } ' \text{Inf_between } O \{C\}, P, O \cup \{C\})$ 1054

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 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 $N_{\text{of}}(\text{Liminf } \mathcal{S}_s) = P_{\text{of}}(\text{Liminf } \mathcal{S}_s) = \emptyset$. The completeness theorem for RP can be stated as follows:

theorem *RP_complete_if_fair*:

assumes

chain $(\rightsquigarrow_{\text{RP}}) \mathcal{S}_s$ **and**
 $O_{\text{of}}(\mathcal{S}_s ! 0) = \emptyset$ **and**
 fair \mathcal{S}_s **and**
 $\forall I. \neg I \models \mathcal{G}_S(\mathcal{S}_s ! 0)$

shows $\perp \in O_{\text{of}}(\text{Liminf } \mathcal{S}_s)$

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 $\mathcal{G}_S \mathcal{S}$ to $\mathcal{G}_S \mathcal{S}'$, where \mathcal{G}_S returns, for a given state, the set consisting of the \mathcal{G}_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}_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 [23, Appendix A], the proof is hard

to penetrate and commit to memory. 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 α -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 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}_F on clauses and \mathcal{G}_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 concerns 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 \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 function 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}_F and \mathcal{G}_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 ι_G from the grounded clauses is approximated at the nonground level by an inference ι :

lemma $G_Inf_overapprox_F_Inf$:
 $\iota_G \in G_Inf_from\ M\ (\cup\ (\mathcal{G}_F\ 'M)) \implies$
 $\exists\ \iota_F \in F_Inf_from\ M.\ \iota_G \in \mathcal{G}_I\ M\ \iota_F$

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_{GC} refines \triangleright and (2) the induced \rightsquigarrow_{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 `lclss_of` :: *'a state* \implies (*'a clause* \times *label*) *set*
where
`lclss_of` $\mathcal{S} = (\lambda C.\ (C,\ \text{New}))\ 'N_of\ \mathcal{S}$
 $\cup\ (\lambda C.\ (C,\ \text{Processed}))\ 'P_of\ \mathcal{S}$
 $\cup\ (\lambda C.\ (C,\ \text{Old}))\ 'O_of\ \mathcal{S}$

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 $(N \cup \{C\}, P, O) \rightsquigarrow_{RP} (N, P, O)$, where $\{\text{Neg } A, \text{Pos } A\} \subseteq C$. We need to prove that there exist labeled clause sets N', M, M' such that $N' \cup M \rightsquigarrow_{GC} N' \cup M'$ 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 condition (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 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 750 lines for the new proof compared with 1200 lines for the original proof. Among the 750 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 even 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 to find 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 a recent CPP invited talk [11, Section 5] for an 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 highly 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 [22], a variant of standard superposition have been formalized in Isabelle.

Beyond our framework, the IsaFoR (Isabelle Formalization of Rewriting) library [27] and the RP formalization [23] offer many definitions and lemmas related to first-order terms and clauses. Starting from RP, Schlichtkrull et al. [24] 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 [16].

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 RPx 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 [15, 28]. The notion of redundancy criterion is also too weak to justify pure literal and blocked clause elimination [18]. 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.

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 Ait-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 Tourret, 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 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.
- [12] 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.
- [13] 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.
- [14] 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.
- [15] Arnaud Fietzke and Christoph Weidenbach. 2009. Labelled splitting. *Ann. Math. Artif. Intell.* 55, 1–2 (2009), 3–34.
- [16] 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.
- [17] Thomas Hillenbrand and Bernd Löchner. 2002. The Next WALDMEISTER Loop. In *CADE-18 (LNCS, Vol. 2392)*, Andrei Voronkov (Ed.). Springer, 486–500.
- [18] 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.
- [19] 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.
- [20] William McCune and Larry Wos. 1997. Otter—the CADE-13 competition incarnations. *J. Autom. Reason.* 18, 2 (1997), 211–220.
- [21] 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.
- [22] Nicolas Peltier. 2016. A Variant of the Superposition Calculus. *Archive of Formal Proofs* 2016 (2016). <https://www.isa-afp.org/entries/SuperCalc.shtml>
- [23] Anders Schlichtkrull, Jasmin Blanchette, Dmitriy Traytel, and Uwe Waldmann. 2020. Formalizing Bachmair and Ganzinger’s Ordered Resolution Prover. *J. Autom. Reasoning* (2020). online first.
- [24] 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.
- [25] 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.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] Uwe Waldmann. 2002. Cancellative abelian monoids and related structures in refutational theorem proving (part I). *J. Symb. Comput.* 33, 6 (2002), 777–829.
- [30] 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.
- [31] 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.
- [32] 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.