



HAL
open science

Superposition for Full Higher-order Logic

Alexander Bentkamp, Jasmin Blanchette, Sophie Touret, Petar Vukmirović

► **To cite this version:**

Alexander Bentkamp, Jasmin Blanchette, Sophie Touret, Petar Vukmirović. Superposition for Full Higher-order Logic. CADE 2021 - 28th International Conference on Automated Deduction, Jul 2021, Pittsburgh, PA / online, United States. pp.396-412, 10.1007/978-3-030-79876-5_23 . hal-03364032

HAL Id: hal-03364032





<https://inria.hal.science/hal-03364032>

Submitted on 4 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Superposition for Full Higher-Order Logic

Alexander Bentkamp¹ , Jasmin Blanchette^{1,2,3} ,
Sophie Tourret^{2,3} , and Petar Vukmirović¹ 

¹ Vrije Universiteit Amsterdam, Amsterdam, the Netherlands

{a.bentkamp, j.c.blanchette, p.vukmirovic}@vu.nl

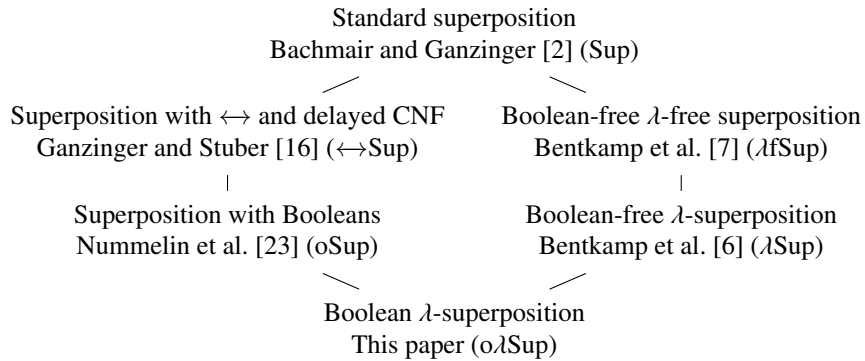
² Université de Lorraine, CNRS, Inria, LORIA, Nancy, France
sophie.tourret@inria.fr

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

Abstract. We recently designed two calculi as stepping stones towards superposition for full higher-order logic: Boolean-free λ -superposition and superposition for first-order logic with interpreted Booleans. Stepping on these stones, we finally reach a sound and refutationally complete calculus for higher-order logic with polymorphism, extensionality, Hilbert choice, and Henkin semantics. In addition to the complexity of combining the calculus's two predecessors, new challenges arise from the interplay between λ -terms and Booleans. Our implementation in Zipperposition outperforms all other higher-order theorem provers and is on a par with an earlier, pragmatic prototype of Booleans in Zipperposition.

1 Introduction

Superposition is a leading calculus for first-order logic with equality. We have been wondering for some years whether it would be possible to gracefully generalize it to extensional higher-order logic and use it as the basis of a strong higher-order automatic theorem prover. Towards this goal, we have, together with colleagues, designed superposition-like calculi for three intermediate logics between first-order and higher-order logic. Now we are finally ready to assemble a superposition calculus for full higher-order logic. The filiation of our new calculus from Bachmair and Ganzinger's standard first-order superposition is as follows:



Our goal was to devise an efficient calculus for higher-order logic. To achieve it, we pursued two objectives. First, the calculus should be refutationally complete. Second, the calculus should coincide as much as possible with its predecessors oSup and λSup on the respective fragments of higher-order logic (which in turn essentially coincide with Sup on first-order logic). Achieving these objectives is the main contribution of this paper. We made an effort to keep the calculus simple, but often the refutational completeness proof forced our hand to add conditions or special cases.

Like oSup , our calculus $\text{o}\lambda\text{Sup}$ operates on clauses that can contain Boolean subterms, and it interleaves clausification with other inferences. Like λSup , $\text{o}\lambda\text{Sup}$ eagerly $\beta\eta$ -normalizes terms, employs full higher-order unification, and relies on a fluid subterm superposition rule (FLUIDSUP) to simulate superposition inferences below applied variables—i.e., terms of the form $y t_1 \dots t_n$ for $n \geq 1$.

Because oSup contains several superposition-like inference rules for Boolean subterms, our completeness proof requires dedicated *fluid Boolean subterm hoisting rules* (FLUIDBOOLHOIST , FLUIDLOOBHOIST), which simulate Boolean inferences below applied variables, in addition to FLUIDSUP , which simulates superposition inferences.

Due to restrictions related to the term order that parameterizes superposition, it is difficult to handle variables bound by unclassified quantifiers if these variables occur applied or in arguments of applied variables. We solve the issue by replacing such quantified terms $\forall y. t$ by equivalent terms $(\lambda y. t) \approx (\lambda y. \mathbf{T})$ in a preprocessing step.

We implemented our calculus in the Zipperposition prover and evaluated it on TPTP and Sledgehammer benchmarks. The new Zipperposition outperforms all other higher-order provers and is on a par with an ad hoc implementation of Booleans in the same prover by Vukmirović and Nummelin [30]. We refer to the technical report [8] for the completeness proof and a more detailed account of the calculus and its evaluation.

2 Logic

Our logic is higher-order logic (simple type theory) with rank-1 polymorphism, Hilbert choice, and functional and Boolean extensionality. Its syntax mostly follows Gordon and Melham [17]. We use the notation \bar{a}_n or \bar{a} to stand for the tuple (a_1, \dots, a_n) where $n \geq 0$. Deviating from Gordon and Melham, type arguments are explicit, written as $c \langle \bar{\tau}_m \rangle$ for a symbol $c : \Pi \bar{a}_m. \nu$ and types $\bar{\tau}_m$. In the type signature Σ_{ty} , we require the presence of a nullary Boolean type constructor \mathbf{o} and a binary function type constructor \rightarrow . In the term signature Σ , we require the presence of the logical symbols \mathbf{T} , \perp , \neg , \wedge , \vee , \rightarrow , \forall , \exists , \approx , and $\not\approx$. The logical symbols are shown in bold to distinguish them from the notation used for clauses below. Moreover, we require the presence of the Hilbert choice operator $\varepsilon \in \Sigma$. Although ε is interpreted in our semantics, we do not consider it a logical symbol. Our calculus will enforce the semantics of ε by an axiom, whereas the semantics of the logical symbols will be enforced by inference rules. We write \mathcal{V} for the set of (term) variables. We use Henkin semantics, in the style of Fitting [15], with respect to which we can prove our calculus refutationally complete. In summary, our logic essentially coincides with the TPTP TH1 format [20].

We generally view terms modulo $\alpha\beta\eta$ -equivalence. When defining operations that need to analyze the structure of terms, however, we use a custom normal form as the

default representative of a $\beta\eta$ -equivalence class: The $\beta\eta Q_\eta$ -normal form $t \downarrow_{\beta\eta Q_\eta}$ of a term t is obtained by bringing the term into η -short β -normal form and finally applying the rewrite rule $Q\langle\tau\rangle s \rightarrow_{Q_\eta} Q\langle\tau\rangle (\lambda x. s x)$ exhaustively whenever s is not a λ -expression. Here and elsewhere, Q stands for either \forall or \exists .

On top of the standard higher-order terms, we install a clausal structure that allows us to formulate calculus rules in the style of first-order superposition. A literal $s \approx t$ is an equation $s \approx t$ or disequation $s \not\approx t$ of terms s and t ; both equations and disequations are unordered pairs. A clause $L_1 \vee \dots \vee L_n$ is a finite multiset of literals L_j . The empty clause is written as \perp . This clausal structure does not restrict the logic, because an arbitrary term t of Boolean type can be written as the clause $t \approx \top$.

We considered excluding negative literals by encoding them as $(s \approx t) \approx \perp$, following \leftrightarrow Sup [16]. However, this approach would make the conclusion of the equality factoring rule (EFACT) too large for our purposes. Regardless, the simplification machinery will allow us to reduce negative literals $t \not\approx \perp$ and $t \not\approx \top$ to $t \approx \top$ and $t \approx \perp$, respectively, thereby eliminating redundant representations of nonequational literals.

We let $\text{CSU}(s, t)$ denote an arbitrary (preferably, minimal) complete set of unifiers for two terms s and t on the set of free variables of the clauses in which s and t occur. To compute such sets, Huet-style preunification [18] is not sufficient, and we must resort to a full unification procedure [19, 29]. To cope with the nontermination of such procedures, we use dovetailing as described by Vukmirović et al. [28, Sect. 5].

Some of the rules in our calculus introduce Skolem symbols, representing objects mandated by existential quantification. We assume that these symbols do not occur in the input problem. More formally, given a problem over a term signature Σ , our calculus operates on a Skolem-extended term signature Σ_{sk} that, in addition to all symbols from Σ , inductively contains symbols $\text{sk}_{\Pi\bar{\alpha}. \forall \bar{x}. \exists z. t z} : \Pi \bar{\alpha}. \bar{\tau} \rightarrow \nu$ for all types ν , variables $z : \nu$, and terms $t : \nu \rightarrow \text{o}$ over Σ_{sk} , where $\bar{\alpha}$ are the free type variables occurring in t and $\bar{x} : \bar{\tau}$ are the free term variables occurring in t , both in order of first occurrence.

3 The Calculus

The $\text{o}\lambda\text{Sup}$ calculus closely resembles λSup , augmented with rules for Boolean reasoning that are inspired by oSup . As in λSup , superposition-like inferences are restricted to certain first-order-like subterms, the *green subterms*, which we define inductively as follows: Every term t is a green subterm of t , and for all symbols $f \in \Sigma \setminus \{\forall, \exists\}$, if t is a green subterm of u_i for some i , then t is a green subterm of $f\langle\bar{\tau}\rangle \bar{u}$. For example, the green subterms of $f(g(\neg p))(\forall\langle\tau\rangle(\lambda x. q))(y a)(\lambda x. h b)$ are the term itself, $g(\neg p)$, $\neg p$, $\forall\langle\tau\rangle(\lambda x. q)$, $y a$, and $\lambda x. h b$. We write $s\langle t \rangle$ to denote a term s with a green subterm t and call the first-order-like context $s\langle \rangle$ a *green context*.

Following λSup , we call a term t *fluid* if (1) $t \downarrow_{\beta\eta Q_\eta}$ is of the form $y \bar{u}_n$ where $n \geq 1$, or (2) $t \downarrow_{\beta\eta Q_\eta}$ is a λ -expression and there exists a substitution σ such that $t\sigma \downarrow_{\beta\eta Q_\eta}$ is not a λ -expression (due to η -reduction). Intuitively, fluid terms are terms whose normal form can change radically as a result of instantiation.

We define deeply occurring variables as in λSup , but exclude λ -expressions directly below quantifiers: A variable *occurs deeply* in a clause C if it occurs inside an argument of an applied variable or inside a λ -expression that is not directly below a quantifier.

Preprocessing. Our completeness theorem requires that quantified variables do not appear in certain higher-order contexts. We use preprocessing to eliminate problematic occurrences of quantifiers. The rewrite rules \forall_{\approx} and \exists_{\approx} , which we collectively denote by Q_{\approx} , are defined as $\forall\langle\tau\rangle \rightarrow_{\forall_{\approx}} \lambda y. y \approx (\lambda x. \top)$ and $\exists\langle\tau\rangle \rightarrow_{\exists_{\approx}} \lambda y. y \not\approx (\lambda x. \perp)$ where the rewritten occurrence of $Q\langle\tau\rangle$ is unapplied or has an argument of the form $\lambda x. v$ such that x occurs as a nongreen subterm of v . If either of these rewrite rules can be applied to a given term, the term is Q_{\approx} -reducible; otherwise, it is Q_{\approx} -normal.

For example, the term $\lambda y. \exists\langle\iota \rightarrow \iota\rangle (\lambda x. gxy(z y)(f x))$ is Q_{\approx} -normal. A term may be Q_{\approx} -reducible because a quantifier appears unapplied (e.g., $g\exists\langle\iota\rangle$); a quantified variable occurs applied (e.g., $\exists\langle\iota \rightarrow \iota\rangle (\lambda x. x a)$); a quantified variable occurs inside a nested λ -expression (e.g., $\forall\langle\iota\rangle (\lambda x. f(\lambda y. x))$); or a quantified variable occurs in the argument of a variable, either a free variable (e.g., $\forall\langle\iota\rangle (\lambda x. z x)$) or a variable bound above the quantifier (e.g., $\lambda y. \exists\langle\iota\rangle (\lambda x. y x)$).

A preprocessor Q_{\approx} -normalizes the input problem. Although inferences may produce Q_{\approx} -reducible clauses, we do not Q_{\approx} -normalize during the derivation process itself. Instead, Q_{\approx} -reducible ground instances of clauses will be considered redundant by the redundancy criterion. Thus, clauses whose ground instances are all Q_{\approx} -reducible can be deleted. However, there are Q_{\approx} -reducible clauses, such as $x\forall\langle\iota\rangle \approx a$, that nevertheless have Q_{\approx} -normal ground instances. Such clauses must be kept because the completeness proof relies on their Q_{\approx} -normal ground instances.

In principle, we could omit the side condition of the Q_{\approx} -rewrite rules and eliminate all quantifiers. However, the calculus (especially, the redundancy criterion) performs better with quantifiers than with λ -expressions, which is why we restrict Q_{\approx} -normalization as much as the completeness proof allows. Extending the preprocessing to eliminate all Boolean terms as in Kotelnikov et al. [21] does not work for higher-order logic because Boolean terms can contain variables bound by enclosing λ -expressions.

Term Order. The calculus is parameterized by a well-founded strict total order \succ on ground terms satisfying these four criteria: (O1) compatibility with green contexts—i.e., $s' \succ s$ implies $t\langle s'\rangle \succ t\langle s\rangle$; (O2) green subterm property—i.e. $t\langle s\rangle \succeq s$ where \succeq is the reflexive closure of \succ ; (O3) $u \succ \perp \succ \top$ for all terms $u \notin \{\top, \perp\}$; (O4) $Q\langle\tau\rangle t \succ t u$ for all types τ , terms t , and terms u such that $Q\langle\tau\rangle t$ and u are Q_{\approx} -normal and the only Boolean green subterms of u are \top and \perp . The restriction of (O4) to Q_{\approx} -normal terms ensures that term orders fulfilling the requirements exist, but it forces us to preprocess the input problem. We extend \succ to literals and clauses via the multiset extensions in the standard way [2, Sect. 2.4].

For nonground terms, \succ is required to be a strict partial order such that $t \succ s$ implies $t\theta \succ s\theta$ for all grounding substitutions θ . As in λSup , we also introduce a nonstrict variant \succeq for which we require that $t\theta \succeq s\theta$ for all grounding substitutions θ whenever $t \succeq s$, and similarly for literals and clauses.

To construct a concrete order fulfilling these requirements, we define an encoding into untyped first-order terms, and compare these using a variant of the Knuth–Bendix order. In a first step, denoted O , the encoding translates fluid terms t as fresh variables z_i ; nonfluid λ -expressions $\lambda x:\tau. u$ as $\text{lam}(O(\tau), O(u))$; applied quantifiers $Q\langle\tau\rangle(\lambda x:\tau. u)$ as $Q_1(O(\tau), O(u))$; and other terms $f\langle\bar{\tau}\rangle \bar{u}_k$ as $f_k(O(\bar{\tau}), O(\bar{u}_k))$. Bound variables are encoded as constants db^i corresponding to De Bruijn indices. In a second step, denoted \mathcal{P} , the

encoding replaces Q_1 by Q'_1 and variables z by z' whenever they occur below lam. For example, $\mathbf{V}(\iota)(\lambda x. p y y (\lambda u. f y y (\mathbf{V}(\iota)(\lambda v. u))))$ is encoded as $\mathbf{V}_1(\iota, p_3(y, y, \text{lam}(o, f_3(y', y', \mathbf{V}'_1(\iota, \text{db}^1))))$). The first-order terms can then be compared using a transfinite Knuth–Bendix order \succ_{kb} [22]. Let the weight of \mathbf{V}_1 and $\mathbf{\Xi}_1$ be ω , the weight of \mathbf{T}_0 and \mathbf{L}_0 be 1, and the weights of all other symbols be less than ω . Let the precedence $>$ be total and $\mathbf{L}_0, \mathbf{T}_0$ be the symbols of lowest precedence, with $\mathbf{L}_0 > \mathbf{T}_0$. Then let $t \succ s$ if $O(\mathcal{P}(t)) \succ_{\text{kb}} O(\mathcal{P}(s))$ and $t \succeq s$ if $O(\mathcal{P}(t)) \succeq_{\text{kb}} O(\mathcal{P}(s))$.

Selection Functions. The calculus is also parameterized by a literal selection function and a Boolean subterm selection function. We define an element x of a multiset M to be \triangleright -maximal for some relation \triangleright if for all $y \in M$ with $y \triangleright x$, we have $y = x$. It is *strictly* \triangleright -maximal if it is \triangleright -maximal and occurs only once in M .

The literal selection function *HLitSel* maps each clause to a subset of *selected literals*. A literal may not be selected if it is positive and neither side is \mathbf{L} . Moreover, a literal $L \langle y \rangle$ may not be selected if $y \bar{u}_n$, with $n \geq 1$, is a \succeq -maximal term of the clause.

The Boolean subterm selection function *HBoolSel* maps each clause C to a subset of *selected subterms* in C . Selected subterms must be green subterms of Boolean type. Moreover, a subterm s must not be selected if $s = \mathbf{T}$, if $s = \mathbf{L}$, if s is a variable-headed term, if s is at the topmost position on either side of a positive literal, or if s contains a variable y as a green subterm, and $y \bar{u}_n$, with $n \geq 1$, is a \succeq -maximal term of the clause.

Eligibility. A literal L is (*strictly*) *eligible* w.r.t. a substitution σ in C if it is selected in C or there are no selected literals and no selected Boolean subterms in C and $L\sigma$ is (*strictly*) \succeq -maximal in $C\sigma$.

The eligible subterms of a clause C w.r.t. a substitution σ are inductively defined as follows: Any selected subterm is eligible. If a literal $L = s \approx t$ with $s\sigma \not\prec t\sigma$ is either eligible and negative or strictly eligible and positive, then the subterm s is eligible. If a subterm t is eligible and the head of t is not \approx or $\not\approx$, all direct green subterms of t are eligible. If a subterm t is eligible and t is of the form $u \approx v$ or $u \not\approx v$, then u is eligible if $u\sigma \not\prec v\sigma$ and v is eligible if $u\sigma \not\prec v\sigma$.

The Core Inference Rules. The calculus consists of the following core inference rules. The first five rules stem from λSup , with minor adaptations concerning Booleans:

$$\frac{\overbrace{D' \vee t \approx t'}^D \quad C \langle u \rangle}{(D' \vee C \langle t' \rangle)\sigma} \text{SUP} \quad \frac{\overbrace{C' \vee u \not\approx u'}^C}{C'\sigma} \text{ERES} \quad \frac{\overbrace{C' \vee u' \approx v' \vee u \approx v}^C}{(C' \vee v \not\approx v' \vee u \approx v')\sigma} \text{EFACT}$$

$$\frac{\overbrace{D' \vee t \approx t'}^D \quad C \langle u \rangle}{(D' \vee C \langle z t' \rangle)\sigma} \text{FLUIDSUP} \quad \frac{\overbrace{C' \vee s \approx s'}^C}{C'\sigma \vee s\sigma \bar{x}_n \approx s'\sigma \bar{x}_n} \text{ARGCONG}$$

SUP 1. u is not fluid; 2. u is not a variable deeply occurring in C ; 3. if u is a variable y , there must exist a grounding substitution θ such that $t\sigma\theta \succ t'\sigma\theta$ and $C\sigma\theta \prec C''\sigma\theta$, where $C'' = C\{y \mapsto t'\}$; 4. $\sigma \in \text{CSU}(t, u)$; 5. $t\sigma \not\prec t'\sigma$; 6. u is eligible in C w.r.t. σ ; 7. $C\sigma \not\prec D\sigma$; 8. $t \approx t'$ is strictly eligible in D w.r.t. σ ; 9. $t\sigma$ is not a fully applied logical symbol; 10. if $t'\sigma = \mathbf{L}$, the subterm u is at the top level of a positive literal.

- ERES 1. $\sigma \in \text{CSU}(u, u')$; 2. $u \not\approx u'$ is eligible in C w.r.t. σ .
- EFACT 1. $\sigma \in \text{CSU}(u, u')$; 2. $u\sigma \not\approx v\sigma$; 3. $(u \approx v)\sigma$ is \approx -maximal in $C\sigma$; 4. $u\sigma \not\approx v\sigma$; 5. nothing is selected in C .
- FLUIDSUP 1. u is a variable deeply occurring in C or u is fluid; 2. z is a fresh variable; 3. $\sigma \in \text{CSU}(zt, u)$; 4. $(zt')\sigma \neq (zt)\sigma$; 5.–10. as for SUP.
- ARGCONG 1. $n > 0$; 2. σ is the most general type substitution that ensures well-typedness of the conclusion for a given n ; 3. \bar{x}_n is a tuple of distinct fresh variables; 4. the literal $s \approx s'$ is strictly eligible in C w.r.t. σ .

The following rules are concerned with Boolean reasoning and originate from oSup. They have been adapted to support polymorphism and applied variables.

$$\begin{array}{c}
\frac{C\langle u \rangle}{(C\langle \perp \rangle \vee u \approx \top)\sigma} \text{BOOLHOIST} \\
\frac{C\langle u \rangle}{(C\langle \perp \rangle \vee x \approx y)\sigma} \text{EQHOIST} \\
\frac{C\langle u \rangle}{(C\langle \top \rangle \vee x \approx y)\sigma} \text{NEQHOIST} \\
\frac{C\langle u \rangle}{(C\langle \perp \rangle \vee yx \approx \top)\sigma} \text{FORALLHOIST} \\
\frac{C\langle u \rangle}{(C\langle \top \rangle \vee yx \approx \perp)\sigma} \text{EXISTSHOIST}
\end{array}
\qquad
\begin{array}{c}
\frac{\overbrace{C' \vee s \approx s'}^C}{C'\sigma} \text{FALSEELIM} \\
\frac{C\langle u \rangle}{C\langle t' \rangle\sigma} \text{BOOLRW} \\
\frac{C\langle u \rangle}{C\langle y(\text{sk}_{\Pi \bar{\alpha}. \forall \bar{x}. \exists z. \neg y\sigma z}(\bar{\alpha}) \bar{x}) \rangle\sigma} \text{FORALLRW} \\
\frac{C\langle u \rangle}{C\langle y(\text{sk}_{\Pi \bar{\alpha}. \forall \bar{x}. \exists z. y\sigma z}(\bar{\alpha}) \bar{x}) \rangle\sigma} \text{EXISTSRW}
\end{array}$$

- BOOLHOIST 1. σ is a type unifier of the type of u with the Boolean type \circ (i.e., the identity if u is Boolean or $\{\alpha \mapsto \circ\}$ if u is of type α for some type variable α); 2. the head of u is neither a variable nor a logical symbol; 3. u is eligible in C ; 4. the occurrence of u is not at the top level of a positive literal.
- EQHOIST, NEQHOIST, FORALLHOIST, EXISTSHOIST 1. $\sigma \in \text{CSU}(u, x \approx y)$, $\sigma \in \text{CSU}(u, x \not\approx y)$, $\sigma \in \text{CSU}(u, \forall(\alpha) y)$, or $\sigma \in \text{CSU}(u, \exists(\alpha) y)$, respectively; 2. x , y , and α are fresh variables; 3. u is eligible in C w.r.t. σ ; 4. if the head of u is a variable, it must be applied and the affected literal must be of the form $u \approx \top$, $u \approx \perp$, or $u \approx v$ where v is a variable-headed term.
- FALSEELIM 1. $\sigma \in \text{CSU}(s \approx s', \perp \approx \top)$; 2. $s \approx s'$ is strictly eligible in C w.r.t. σ .
- BOOLRW 1. $\sigma \in \text{CSU}(t, u)$ and (t, t') is one of the following pairs, where y is a fresh variable: $(\neg \perp, \top)$, $(\neg \top, \perp)$, $(\perp \wedge \perp, \perp)$, $(\top \wedge \perp, \perp)$, $(\perp \wedge \top, \perp)$, $(\top \wedge \top, \top)$, $(\perp \vee \perp, \perp)$, $(\top \vee \perp, \top)$, $(\perp \vee \top, \top)$, $(\top \vee \top, \top)$, $(\perp \rightarrow \perp, \top)$, $(\top \rightarrow \perp, \perp)$, $(\perp \rightarrow \top, \top)$, $(\top \rightarrow \top, \top)$, $(y \approx y, \top)$, $(y \not\approx y, \perp)$; 2. u is not a variable; 3. u is eligible in C w.r.t. σ ; 4. if the head of u is a variable, it must be applied and the affected literal must be of the form $u \approx \top$, $u \approx \perp$, or $u \approx v$ where v is a variable-headed term.
- FORALLRW, EXISTSRW 1. $\sigma \in \text{CSU}(\forall(\beta) y, u)$ and $\sigma \in \text{CSU}(\exists(\beta) y, u)$, respectively, where β is a fresh type variable, y is a fresh term variable, $\bar{\alpha}$ are the free type variables and \bar{x} are the free term variables occurring in $y\sigma$ in order of first occurrence;

2. u is not a variable; 3. u is eligible in C w.r.t. σ ; 4. if the head of u is a variable, it must be applied and the affected literal must be of the form $u \approx \mathbf{T}$, $u \approx \mathbf{\perp}$, or $u \approx v$ where v is a variable-headed term; 5. for FORALLRW, the indicated occurrence of u is not in a literal $u \approx \mathbf{T}$, and for EXISTSRW, the indicated occurrence of u is not in a literal $u \approx \mathbf{\perp}$.

Like SUP, also the Boolean rules must be simulated in fluid terms. The following rules are Boolean counterparts of FLUIDSUP:

$$\frac{C\langle u \rangle}{(C\langle z \mathbf{\perp} \rangle \vee x \approx \mathbf{T})\sigma} \text{FLUID-BOOLHOIST} \quad \frac{C\langle u \rangle}{(C\langle z \mathbf{T} \rangle \vee x \approx \mathbf{\perp})\sigma} \text{FLUID-LOOBHOIST}$$

FLUIDBOOLHOIST 1. u is fluid; 2. z and x are fresh variables; 3. $\sigma \in \text{CSU}(z, x, u)$; 4. $(z \mathbf{\perp})\sigma \neq (zx)\sigma$; 5. $x\sigma \neq \mathbf{T}$ and $x\sigma \neq \mathbf{\perp}$; 6. u is eligible in C w.r.t. σ .

FLUIDLOOBHOIST Like the above but with $\mathbf{\perp}$ replaced by \mathbf{T} in condition 4.

In addition to the inference rules, our calculus relies on two axioms, below. Axiom (EXT), from λSup , embodies functional extensionality; the expression $\text{diff}\langle \alpha, \beta \rangle$ abbreviates $\text{sk}_{\prod \alpha \beta. \forall z y. \exists x. z x \neq y x}(\alpha, \beta)$. Axiom (CHOICE) characterizes the Hilbert choice operator ε .

$$z(\text{diff}\langle \alpha, \beta \rangle z y) \not\approx y(\text{diff}\langle \alpha, \beta \rangle z y) \vee z \approx y \quad (\text{EXT})$$

$$y x \approx \mathbf{\perp} \vee y(\varepsilon\langle \alpha \rangle y) \approx \mathbf{T} \quad (\text{CHOICE})$$

Rationale for the Rules. Most of the calculus's rules are adapted from its precursors. SUP, ERES, and EFACT are already present in Sup, with slightly different side conditions. Notably, as in λfSup and λSup , SUP inferences are required only into green contexts. Other subterms are accessed indirectly via ARGCONG and (EXT).

The rules BOOLHOIST, EQHOIST, NEQHOIST, FORALLHOIST, EXISTSHOIST, FALSEELIM, BOOLRW, FORALLRW, and EXISTSRW, concerned with Boolean reasoning, stem from oSup , which was inspired by $\leftrightarrow\text{Sup}$. Except for BOOLHOIST and FALSEELIM, these rules have a condition stating that “if the head of u is a variable, it must be applied and the affected literal must be of the form $u \approx \mathbf{T}$, $u \approx \mathbf{\perp}$, or $u \approx v$ where v is a variable-headed term.” The inferences at variable-headed terms permitted by this condition are our form of primitive substitution [1, 18], a mechanism that blindly substitutes logical connectives and quantifiers for variables z with a Boolean result type.

Example 1. Our calculus can prove that Leibniz equality implies equality (i.e., if two values behave the same for all predicates, they are equal) as follows:

$$\frac{\frac{\frac{\frac{z a \approx \mathbf{\perp} \vee z b \approx \mathbf{T}}{(x a \approx y a) \approx \mathbf{\perp} \vee \mathbf{\perp} \approx \mathbf{T} \vee x b \approx y b}}{\mathbf{T} \approx \mathbf{\perp} \vee \mathbf{\perp} \approx \mathbf{T} \vee w a b b \approx w b a b}}{\mathbf{\perp} \approx \mathbf{T} \vee w a b b \approx w b a b}}{a \not\approx b} \text{FALSEELIM}}{\frac{w a b b \approx w b a b}{a \not\approx a} \text{SUP}} \text{ERES}$$

The EQHOIST inference, applied on $z\ b$, illustrates how our calculus introduces logical symbols without a dedicated primitive substitution rule. Although \approx does not appear in the premise, we still need to apply EQHOIST on $z\ b$ with $\text{CSU}(z\ b, x_0 \approx y_0) = \{\{z \mapsto \lambda v. x\ v \approx y\ v, x_0 \mapsto x\ b, y_0 \mapsto y\ b\}\}$. Other calculi [1, 9, 18, 26] would apply an explicit primitive substitution rule instead, yielding essentially $(x\ a \approx y\ a) \approx \perp \vee (x\ b \approx y\ b) \approx \top$. However, in our approach this clause is subsumed and could be discarded immediately. By hoisting the equality to the clausal level, we bypass the redundancy criterion.

Next, BOOLRW can be applied to $x\ a \approx y\ a$ with $\text{CSU}(x\ a \approx y\ a, y_0 \approx y_0) = \{\{x \mapsto \lambda v. w\ a \vee v, y \mapsto \lambda v. w\ v\ a, y_0 \mapsto w\ a\ a\}\}$. The two FALSEELIM steps remove the $\perp \approx \top$ literals. Then SUP is applicable with the unifier $\{w \mapsto \lambda x_1\ x_2\ x_3. x_2\} \in \text{CSU}(b, w\ a\ b)$, and ERES derives the contradiction.

Like in λSup , the FLUIDSUP rule is responsible for simulating superposition inferences below applied variables, other fluid terms, and deeply occurring variables. Complementarily, FLUIDBOOLHOIST and FLUIDLOOBHOIST simulate the various Boolean inference rules below fluid terms. Initially, we considered adding a fluid version of each rule that operates on Boolean subterms, but we discovered that FLUID-BOOLHOIST and FLUIDLOOBHOIST suffice to achieve refutational completeness.

Example 2. The clause set consisting of $h(y\ b) \approx h(g\ \perp) \vee h(y\ a) \approx h(g\ \top)$ and $a \approx b$ highlights the need for FLUIDBOOLHOIST and its companion. The set is unsatisfiable because the instantiation $\{y \mapsto \lambda x. g(x\ \approx a)\}$ produces the clause $h(g(b\ \approx a)) \approx h(g\ \perp) \vee h(g(a\ \approx a)) \approx h(g\ \top)$, which is unsatisfiable in conjunction with $a \approx b$.

The literal selection function can select either literal in the first clause. ERES is applicable in either case, but the unifiers $\{y \mapsto \lambda x. g\ \perp\}$ and $\{y \mapsto \lambda x. g\ \top\}$ do not lead to a contradiction. Instead, we need to apply FLUIDBOOLHOIST if the first literal is selected or FLUIDLOOBHOIST if the second literal is selected. In the first case, the derivation is as follows:

$$\begin{array}{c}
\frac{h(y\ b) \approx h(g\ \perp) \vee h(y\ a) \approx h(g\ \top)}{h(z'\ b\ \perp) \approx h(g\ \perp) \vee h(z'\ a(x'\ a)) \approx h(g\ \top) \vee x'\ b \approx \top} \text{FLUIDBOOLHOIST} \\
\frac{h(z'\ b\ \perp) \approx h(g\ \perp) \vee h(z'\ a(x'\ a)) \approx h(g\ \top) \vee x'\ b \approx \top}{h(g(x'\ a)) \approx h(g\ \top) \vee x'\ b \approx \top} \text{ERES} \\
\frac{a \approx b \quad h(g(x''\ a \approx x'''\ a)) \approx h(g\ \top) \vee \perp \approx \top \vee x''\ b \approx x'''\ b}{h(g(a \approx x'''\ a)) \approx h(g\ \top) \vee \perp \approx \top \vee a \approx x'''\ b} \text{EQHOIST} \\
\frac{h(g(a \approx x'''\ a)) \approx h(g\ \top) \vee \perp \approx \top \vee a \approx x'''\ b}{h(g\ \top) \approx h(g\ \top) \vee \perp \approx \top \vee a \approx a} \text{SUP} \\
\frac{h(g\ \top) \approx h(g\ \top) \vee \perp \approx \top \vee a \approx a}{\perp \approx \top \vee a \approx a} \text{BOOLRW} \\
\frac{\perp \approx \top \vee a \approx a}{\perp \approx \top} \text{ERES} \\
\frac{\perp \approx \top}{\perp} \text{FALSEELIM}
\end{array}$$

The FLUIDBOOLHOIST inference uses the unifier $\{y \mapsto \lambda u. z'\ u(x'\ u), z \mapsto \lambda u. z'\ b\ u, x \mapsto x'\ b\} \in \text{CSU}(z\ x, y\ b)$. We apply ERES to the first literal of the resulting clause, with unifier $\{z' \mapsto \lambda uv. g\ v\} \in \text{CSU}(h(z'\ b\ \perp), h(g\ \perp))$. Next, we apply EQHOIST with the unifier $\{x' \mapsto \lambda u. x''\ u \approx x'''\ u, w \mapsto x''\ b, w' \mapsto x'''\ b\} \in \text{CSU}(x'\ b, w \approx w')$ to the literal

created by FLUIDBOOLHOIST, effectively performing a primitive substitution. The resulting clause can superpose into $a \approx b$ with the unifier $\{x'' \mapsto \lambda u. u\} \in \text{CSU}(x'' b, b)$. The two sides of the interpreted equality in the first literal can then be unified, allowing us to apply BOOLRW with the unifier $\{y \mapsto a, x''' \mapsto \lambda u. a\} \in \text{CSU}(y \approx y, a \approx x''' b)$. Finally, applying ERES twice and FALSEELIM once yields the empty clause.

Remarkably, none of the provers that participated in the CASC-J10 competition can solve this two-clause problem within a minute. Satallax finds a proof after 72 s and LEO-II after over 7 minutes. Our new Zipperposition implementation solves it in 3 s.

The Redundancy Criterion. In first-order superposition, a clause is considered redundant if all its ground instances are entailed by \prec -smaller ground instances of other clauses. In essence, this will also be our definition, but we will use a different notion of ground instances and a different notion of entailment.

Given a clause C , let its *ground instances* $\mathcal{G}(C)$ be the set of all clauses of the form $C\theta$ for some substitution θ such that $C\theta$ is ground and Q_{\approx} -normal, and for all variables x occurring in C , the only Boolean green subterms of $x\theta$ are \mathbf{T} and \mathbf{F} . The rationale of this definition is to ensure that ground instances of the conclusion of FORALLHOIST, EXISTSHOIST, FORALLRW, and EXISTSRW inferences are smaller than the corresponding instances of their premise by property (O4).

The redundancy criterion's notion of entailment is defined via an encoding into a weaker logic, following λfSup and λSup . In this paper, the weaker logic is ground first-order logic with interpreted Booleans—the ground fragment of the logic of oSup . Its signature $(\Sigma_{\text{ty}}, \Sigma_{\text{GF}})$ is derived from our higher-order signature $(\Sigma_{\text{ty}}, \Sigma)$ as follows. The type constructors Σ_{ty} are the same in both signatures, but \rightarrow is an uninterpreted type constructor in first-order logic. For each ground instance $f(\bar{v}) : \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow \tau$ of a symbol $f \in \Sigma$, we introduce a first-order symbol $f_j^{\bar{v}} \in \Sigma_{\text{GF}}$ with argument types $\bar{\tau}_j$ and result type $\tau_{j+1} \rightarrow \dots \rightarrow \tau_n \rightarrow \tau$, for each j . Moreover, for each ground term $\lambda x. t$, we introduce a symbol $\text{lam}_{\lambda x. t} \in \Sigma_{\text{GF}}$ of the same type. The symbols $\mathbf{F}_0, \mathbf{T}_0, \neg_1, \mathbf{A}_2, \mathbf{V}_2, \rightarrow_2, \approx_2^{\tau}$, and \neq_2^{τ} are identified with the corresponding first-order logical symbols.

We define an encoding \mathcal{F} of Q_{\approx} -normal ground higher-order terms into this ground first-order logic recursively as follows: $\mathcal{F}(\forall(\tau)(\lambda x. t)) = \forall x. \mathcal{F}(t)$ and $\mathcal{F}(\exists(\tau)(\lambda x. t)) = \exists x. \mathcal{F}(t)$ for applied quantifiers; $\mathcal{F}(\lambda x. t) = \text{lam}_{\lambda x. t}$ for λ -expressions; and $\mathcal{F}(f(\bar{v}) \bar{s}_j) = f_j^{\bar{v}}(\mathcal{F}(\bar{s}_j))$ for other terms. For quantified variables, we define $\mathcal{F}(x) = x$. Here, Q_{\approx} -normality is crucial to ensure that bound variables do not occur applied or within λ -expressions. The definition of green subterms is devised such that green subterms correspond to first-order subterms via the encoding \mathcal{F} , with the exception of first-order subterms below quantifiers. The encoding \mathcal{F} is extended to clauses by mapping each literal and each side of a literal individually. From the entailment relation \models for the ground first-order logic, we derive an entailment relation $\models_{\mathcal{F}}$ on Q_{\approx} -normal ground higher-order clauses by defining $M \models_{\mathcal{F}} N$ if $\mathcal{F}(M) \models \mathcal{F}(N)$. This relation is weaker than standard higher-order entailment; for example, $\{f \approx g\} \not\models_{\mathcal{F}} \{f a \approx g a\}$ (because of the subscripts added by \mathcal{F}) and $\{\text{p}(\lambda x. \mathbf{T})\} \not\models_{\mathcal{F}} \{\text{p}(\lambda x. \neg \mathbf{F})\}$ (because of the lam symbols used by \mathcal{F}).

Using $\models_{\mathcal{F}}$, we define a clause C to be *redundant* w.r.t. a clause set N if for every $D \in \mathcal{G}(C)$, we have $\{E \in \mathcal{G}(N) \mid E \prec D\} \models_{\mathcal{F}} D$ or there exists a clause $C' \in N$ such that $C \sqsupset C'$ and $D \in \mathcal{G}(C')$. The tiebreaker \sqsupset can be an arbitrary well-founded partial order

on clauses; in practice, we use a well-founded restriction of the ill-founded strict subsumption relation [6, Sect. 3.4]. We denote the set of redundant clauses w.r.t. a clause set N by $Red_C(N)$. Note that $\models_{\mathcal{F}}$ is weak enough to ensure that the ARGCONG inference rule and axiom (EXT) are not immediately redundant and can fulfill their purpose.

For first-order superposition, an inference is considered redundant if for each of its ground instances, a premise is redundant or the conclusion is entailed by clauses smaller than the main premise. For most inference rules, our definition follows this idea, using $\models_{\mathcal{F}}$ for entailment; other rules need nonstandard notions of ground instances and redundancy. The definition of inference redundancy presented below is simpler than the more sophisticated notion in our technical report. Nonetheless, the redundant inferences below are a strict subset of the redundant inferences of our report and thus completeness also holds using the notion below. For the few prover optimizations based on inference redundancy that we know about (e.g., simultaneous superposition [4]), the following criterion suffices.

For SUP, ERES, EFACT, BOOLHOIST, FALSEELIM, EQHOIST, NEQHOIST, and BOOLRW, we define ground instances as usual: *Ground instances* are all inferences obtained by applying a grounding substitution to premises and conclusion such that the result adheres to the conditions of the given rule w.r.t. selection functions that select literals and subterms as in the original premise. For FLUIDSUP and FLUIDBOOLHOIST, we define ground instances in the same way except that we require that ground instances adhere to the conditions of SUP or BOOLHOIST, respectively. For FORALLRW, EXISTSRW, FORALLHOIST, EXISTSHOIST, which do not have ground instances in the sense above, we define a *ground instance* as any inference that is obtained by applying the unifier σ to the premise and then applying a grounding substitution to premise and conclusion, regardless of whether the resulting inference is an inference of our calculus.

For all rules except FLUIDLOOBHOIST and ARGCONG, we define an inference to be *redundant* w.r.t. a clause set N if for each ground instance ι , a premise of ι is redundant w.r.t. $\mathcal{G}(N)$ or the conclusion of ι is entailed w.r.t. $\models_{\mathcal{F}}$ by clauses from $\mathcal{G}(N)$ that are smaller than the main (i.e., rightmost) premise of ι . For the rules FLUIDLOOBHOIST and ARGCONG, as well as axioms (EXT) and (CHOICE)—viewed as premiseless inferences—we define an inference to be *redundant* w.r.t. a clause set N if all ground instances of its conclusion are contained in $\mathcal{G}(N)$ or redundant w.r.t. $\mathcal{G}(N)$. We denote the set of redundant inferences w.r.t. N by $Red_I(N)$.

Simplification Rules. Our redundancy criterion is strong enough to support counterparts of most simplification rules implemented in Schulz’s first-order E [25, Sect. 2.3.1 and 2.3.2]. Deletion of duplicated literals, deletion of resolved literals, syntactic tautology deletion, negative simplify-reflect, and clause subsumption adhere to our redundancy criterion. Positive simplify-reflect, equality subsumption, and rewriting (de-modulation) of positive and negative literals are supported if they are applied on green subterms or on other subterms that are encoded into first-order subterms by \mathcal{G} and \mathcal{F} . Semantic tautology deletion can be applied as well, using $\models_{\mathcal{F}}$; moreover, for positive literals, the rewriting clause must be smaller than the rewritten clause.

Under some circumstances, inference rules can be applied as simplifications. The FALSEELIM and BOOLRW rules can be applied as a simplification if σ is the identity. If the head of u is \forall , FORALLHOIST and FORALLRW can both be applied and, together,

serve as one simplification rule. The same holds for EXISTS_{HOIST} and EXISTS_{RW} if the head of u is \exists . For all of these rules, the eligibility conditions can be ignored.

Clausification. Like oSup, our calculus does not require the input problem to be clausified during the preprocessing, and it supports higher-order analogues of the three inprocessing clausification methods introduced by Nummelin et al. *Inner delayed clausification* relies on our core calculus rules to destruct logical symbols. *Outer delayed clausification* adds the following clausification rules to the calculus:

$$\begin{array}{c}
 \frac{s \approx \mathbf{T} \vee C}{oc(s, C)} \text{ POSOUTERCLAUS} \qquad \frac{s \approx \mathbf{\perp} \vee C}{oc(\neg s, C)} \text{ NEGOUTERCLAUS} \\
 \\
 \frac{s \approx t \vee C}{s \approx \mathbf{\perp} \vee t \approx \mathbf{T} \vee C \quad s \approx \mathbf{T} \vee t \approx \mathbf{\perp} \vee C} \text{ EQOUTERCLAUS} \\
 \\
 \frac{s \not\approx t \vee C}{s \approx \mathbf{\perp} \vee t \approx \mathbf{\perp} \vee C \quad s \approx \mathbf{T} \vee t \approx \mathbf{T} \vee C} \text{ NEQOUTERCLAUS}
 \end{array}$$

The double bars identify simplification rules (i.e., the conclusions make the premise redundant and can replace it). The first two rules require that s has a logical symbol as its head, whereas the last two require that s and t are Boolean terms other than \mathbf{T} and $\mathbf{\perp}$. The function oc distributes the logical symbols over the clause—e.g., $oc(s \rightarrow t, C) = \{s \approx \mathbf{\perp} \vee t \approx \mathbf{T} \vee C\}$, and $oc(\neg(s \mathbf{V} t), C) = \{s \approx \mathbf{\perp} \vee C, t \approx \mathbf{\perp} \vee C\}$. It is easy to check that our redundancy criterion allows us to replace the premise of the OUTERCLAUS rules with their conclusion. Nonetheless, we apply EQOUTERCLAUS and NEQOUTERCLAUS as inferences because the premises might be useful in their original form.

Besides the two delayed clausification methods, a third inprocessing clausification method is *immediate* clausification. This clausifies the input problem’s outer Boolean structure in one swoop, resulting in a set of higher-order clauses. If unclausified Boolean terms rise to the top during saturation, the same algorithm is run to clausify them.

Unlike delayed clausification, immediate clausification is a black box and is unaware of the proof state other than the Boolean term it is applied to. Delayed clausification, on the other hand, clausifies the term step by step, allowing us to interleave clausification with the strong simplification machinery of superposition provers. It is especially powerful in higher-order contexts: Examples such as $y p q \not\approx (p \mathbf{V} q)$ can be refuted directly by equality resolution, rather than via more explosive rules on the clausified form.

4 Refutational Completeness

Our calculus is dynamically refutationally complete for problems in \mathbf{Q}_{\approx} -normal form. The full proof can be found in our technical report [8].

Theorem 3 (Dynamic refutational completeness). *Let $(N_i)_i$ be a derivation—i.e., $N_i \setminus N_{i+1} \subseteq \text{Red}_C(N_{i+1})$ for all i . Let N_0 be \mathbf{Q}_{\approx} -normal and such that $N_0 \models \perp$. Moreover, assume that $(N_i)_i$ is fair—i.e., all inferences from clauses in the limit inferior $\bigcup_i \bigcap_{j \geq i} N_j$ are contained in $\bigcup_i \text{Red}_1(N_i)$. Then we have $\perp \in N_i$ for some i .*

Following the completeness proof of λSup , our proof is structured in three levels of logics. For each, we define a calculus and show that it is refutationally complete: ground monomorphic first-order logic with an interpreted Boolean type (GF); the \mathcal{Q}_{\approx} -normal ground fragment of higher-order logic (GH); and higher-order logic (H).

The logic of the GF level is the ground fragment of oSup 's logic. The GF calculus is a ground version of oSup , which Nummelin et al. showed refutationally complete. It consists of ground first-order equivalents of our rules, excluding ARGCONG , FLUID-BOOLHOIST , and FLUIDLOOBHOIST , which are specific to higher-order logic. The counterparts to FORALLHOIST and EXISTSHOIST enumerate ground terms instead of producing free variables, to stay within the ground fragment. For compatibility with the nonground level, the conclusions of FORALLRW and EXISTSRW cannot contain concrete Skolem functions. Instead, the GF calculus is parameterized by a witness function that can assign an arbitrary term to each occurrence of a quantifier in a clause. This witness function is used to retrieve the Skolem terms in the GF equivalents of FORALLRW and EXISTSRW .

On the next level, the GH calculus includes inference rules isomorphic to the GF rules, transferred to higher-order logic via \mathcal{F}^{-1} . Moreover, it contains an ARGCONG variant that enumerates ground terms instead of introducing fresh variables, as well as rules enumerating ground instances of axioms (EXT) and (CHOICE). We prove refutational completeness of the GH calculus by constructing a higher-order interpretation based on the model constructed for the completeness proof of the GF level. This proof step is analogous to the corresponding step in λSup 's proof, but we must also consider \mathcal{Q}_{\approx} -normality and the logical symbols.

To lift completeness to the H level, we use the saturation framework of Waldmann et al. [31]. The main proof obligation it leaves us to show is that nonredundant GH inferences can be lifted to corresponding nonground H inferences. For this lifting, we must choose a suitable GH witness function and appropriate GH selection functions for literals and Boolean subterms, given a saturated clause set at the H level and the H selection functions. Then the saturation framework guarantees static refutational completeness w.r.t. Herbrand entailment, which is the entailment relation induced by the grounding function \mathcal{G} . We then show that this implies dynamic refutational completeness w.r.t. \models for \mathcal{Q}_{\approx} -normal initial clause sets.

5 Implementation

We implemented our calculus in the Zipperposition prover [14], whose OCaml source code makes it convenient to prototype calculus extensions. Except for the presence of axioms (EXT) and (CHOICE), the new code gracefully extends Zipperposition's implementation of oSup in the sense that $\text{o}\lambda\text{Sup}$ coincides with oSup on first-order problems. The same cannot be said w.r.t. λSup on Boolean-free problems because of the FLUIDBOOLHOIST and FLUIDLOOBHOIST rules, which are triggered by any applied variable. From the implementation of λSup , we inherit the given clause procedure, which supports infinitely branching inferences, as well as calculus extensions and heuristics [28]. From the implementation of oSup , we inherit the simplification rule BOOLSIMP , a mainstay of our Boolean simplification machinery.

As in the implementation of λ Sup, we approximate fluid terms as terms that are either nonground λ -expressions or terms of the form $x\bar{s}_n$ with $n > 0$. Two slight, accidental discrepancies are that we also count variable occurrences below quantifiers as deep and perform EFACT inferences even if the maximal literal is selected. Since we expect FLUIDBOOLHOIST and FLUIDLOOBHOIST to be highly explosive, we penalize them and all of their offspring. In addition to various λ Sup extensions [6, Sect. 5], we also use all the rules for Boolean reasoning described by Vukmirović and Nummelin [30] except for the BOOLEF rules.

6 Evaluation

We evaluate the calculus implementation in Zipperposition and compare it with other higher-order provers. Our experiments were performed on StarExec Miami servers equipped with Intel Xeon E5-2620 v4 CPUs clocked at 2.10 GHz. We used all 2606 TH0 theorems from the TPTP 7.3.0 library [27] and 1253 “Judgment Day” problems [12] generated using Sledgehammer (SH) [24] as our benchmark set. An archive containing the benchmarks and the raw evaluation results is publicly available [5].

Calculus Evaluation. In this first part, we evaluate selected parameters of Zipperposition by varying only the studied parameter in a fixed well-performing configuration. This base configuration disables axioms (CHOICE) and (EXT) and the FLUID- rules. It uses the unification procedure of Vukmirović et al. [29] in its complete variant—i.e., the variant that produces a complete set of unifiers. It uses none of the early Boolean rules described by Vukmirović and Nummelin [30]. The preprocessor Q_{∞} is disabled as well. All of the completeness-preserving simplification rules listed in Sect. 3 are enabled. The configuration uses immediate clausification. We set the CPU time limit to 30 s in all three experiments.

In the first experiment, we assess the overhead incurred by the FLUID- rules. These rules unify with a term whose head is a fresh variable. Thus, we expected that they needed to be tightly controlled to achieve good performance. To test our hypothesis, we simultaneously modified the parameters of these three rules. In Figure 1, the *off* mode simply disables the rules, the *pragmatic* mode uses a terminating incomplete unification algorithm (the pragmatic variant of Vukmirović et al. [29]), and the *complete* mode uses a complete unification algorithm. The results show that disabling FLUID- rules altogether achieves the best performance. However, on TPTP problems, *complete* finds 35 proofs not found by *off*, and *pragmatic* finds 22 proofs not found by *off*. On Sledgehammer benchmarks, this effect is much weaker, likely because the Sledgehammer benchmarks require less higher-order reasoning: *complete* finds only one new proof over *off*, and *pragmatic* finds only four.

In the second experiment, we explore the clausification methods introduced at the end of Sect. 3: *inner* delayed clausification, *outer* delayed clausification, and *immediate* clausification. The modes *inner* and *outer* employ oSup’s RENAME rule, which renames Boolean terms headed by logical symbols using a Tseitin-like transformation if they occur at least four times in the proof state. Vukmirović and Nummelin [30] observed that *outer* clausification can greatly help prove higher-order problems, and we expected

	<i>off</i>	<i>pragmatic</i>	<i>complete</i>
TPTP	1642	1591	1619
SH	467	431	437

Fig. 1. Evaluation of FLUID- rules

	<i>inner</i>	<i>outer</i>	<i>immediate</i>
TPTP	1323	1670	1642
SH	406	470	467

Fig. 2. Evaluation of classification method

	<i>off</i>	$p = 64$	$p = 16$	$p = 4$	$p = 1$
TPTP	1642	1617	1613	1615	1594
SH	467	458	458	459	445

Fig. 3. Evaluation of axiom (CHOICE)

	TPTP	ofSH	SH
CVC4 1.8	1796	680	619
Leo-III 1.5.2	2104	681	621
Vampire 4.5	2131	692	681
Satallax 3.5	2162	573	587
Zip (CASC-J10)	2301	734	736
New Zip	2320	724	720

Fig. 4. Evaluation of all competitive higher-order provers

it to perform well for our calculus, too. The results are shown in Figure 2. The results confirm our hypothesis: The *outer* mode outperforms *immediate* on both TPTP and Sledgehammer benchmarks. The *inner* mode performs worst, but on Sledgehammer benchmarks, it proves 17 problems beyond the reach of the other two. Interestingly, several of these problems contain axioms of the form $\phi \rightarrow \psi$, and applying superposition and demodulation to these axioms is preferable to classifying them.

In the third experiment, we investigate the effect of axiom (CHOICE), which is necessary to achieve refutational completeness. To evaluate (CHOICE), we either disabled it in a configuration labeled *off* or set the axiom’s penalty p to different values. In Zipperposition, penalties are propagated through inference and simplification rules and are used to increase the heuristic weight of clauses, postponing the selection of penalized clauses. The results are shown in Figure 3. As expected, disabling (CHOICE), or at least penalizing it heavily, improves performance. Yet enabling (CHOICE) can be crucial: For 19 TPTP problems, the proofs are found when (CHOICE) is enabled and $p = 4$, but not when the rule is disabled. On Sledgehammer problems, this effect is weaker, with only two new problems proved for $p = 4$.

Prover Comparison. In this second part, we compare Zipperposition’s performance with other higher-order provers. Like at CASC-J10, the wall-clock time limit was 120 s, the CPU time limit was 960 s, and the provers were run on StarExec Miami. We used the following versions of all systems that took part in the THF division: CVC4 1.8 [3], Leo-III 1.5.2 [26], Satallax 3.5 [13], and Vampire 4.5 [11]. The developers of Vampire have informed us that its higher-order schedule is optimized for running on a single core. As a result, the prover suffers some degradation of performance when running on multiple cores. We evaluate both the version of Zipperposition that took part in CASC-J10 (*Zip*) and the updated version of Zipperposition that supports our new calculus (*New Zip*). Zip’s portfolio of prover configurations is based on λ Sup and techniques described by Vukmirović and Nummelin [30]. New Zip’s portfolio is specially designed for our

new calculus and optimized for TPTP problems. To assess the performance of Boolean reasoning, we used Sledgehammer benchmarks generated both with native Booleans (SH) and with an encoding into Boolean-free higher-order logic (ofSH). For technical reasons, the encoding also performs λ -lifting, but this minor transformation should have little impact on results [6, Sect. 7].

The results are shown in Figure 4. The two versions of Zipperposition are ahead of all other provers on both benchmark sets. This shows that, with thorough parameter tuning, higher-order superposition outperforms tableaux, which had been the state of the art in higher-order reasoning for a decade. The updated version of New Zip beats Zip on TPTP problems but lags behind Zip on Sledgehammer benchmarks as we have yet to further explore more general heuristics that work well with our new calculus. The Sledgehammer benchmarks fail to demonstrate the superiority of native Booleans reasoning compared with an encoding, and in fact CVC4 and Leo-III perform dramatically better on the encoded Boolean problems, suggesting that there is room for tuning.

7 Conclusion

We have created a superposition calculus for higher-order logic that is refutationally complete. Most of the key ideas have been developed in previous work by us and colleagues, but combining them in the right way has been challenging. A key idea was to Q_{∞} -normalize away inconvenient terms.

Unlike earlier refutationally complete calculi for full higher-order logic based on resolution or paramodulation, our calculus employs a term order, which restricts the proof search, and a redundancy criterion, which can be used to add various simplification rules while keeping refutational completeness. These two mechanisms are undoubtedly major factors in the success of first-order superposition, and it is very fortunate that we could incorporate both in a higher-order calculus. An alternative calculus with the same two mechanisms could be achieved by combining oSup with Bhayat and Reger’s combinatory superposition [10]. The article on λ Sup [6, Sect. 8] discusses related work in more detail.

The evaluation results show that our calculus is an excellent basis for higher-order theorem proving. In future work, we want to experiment further with the different parameters of the calculus (for example, with Boolean subterm selection heuristics) and implement it in a state-of-the-art prover such as E.

Acknowledgment. Uwe Waldmann provided advice and carefully checked the completeness proof. Visa Nummelin led the design of the oSup calculus. Simon Cruanes helped us with the implementation. Martin Desharnais generated the Sledgehammer benchmarks. Christoph Benzmüller, Ahmed Bhayat, Mathias Fleury, Herman Geuvers, Giles Reger, Alexander Steen, Mark Summerfield, Geoff Sutcliffe, and the anonymous reviewers helped us in various ways. We thank them all.

Bentkamp, Blanchette, and Vukmirović’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). Blanchette’s research has received funding from the Netherlands Organization for Scientific Research (NWO) under the Vidi program (project No. 016.Vidi.189.037, Lean Forward).

References

- [1] Andrews, P.B.: On connections and higher-order logic. *J. Autom. Reason.* **5**(3), 257–291 (1989)
- [2] Bachmair, L., Ganzinger, H.: Rewrite-based equational theorem proving with selection and simplification. *J. Log. Comput.* **4**(3), 217–247 (1994)
- [3] Barrett, C.W., Conway, C.L., Deters, M., Hadarean, L., Jovanovic, D., King, T., Reynolds, A., Tinelli, C.: CVC4. In: *CAV. LNCS*, vol. 6806, pp. 171–177. Springer (2011)
- [4] Benanav, D.: Simultaneous paramodulation. In: Stickel, M.E. (ed.) *CADE-10. LNCS*, vol. 449, pp. 442–455. Springer (1990)
- [5] Bentkamp, A., Blanchette, J., Tourret, S., Vukmirović, P.: Superposition for full higher-order logic (supplementary material), <https://doi.org/10.5281/zenodo.4534759>
- [6] Bentkamp, A., Blanchette, J., Tourret, S., Vukmirović, P., Waldmann, U.: Superposition with lambdas, accepted in *J. Autom. Reason.* Preprint at <https://arxiv.org/abs/2102.00453v1> (2021)
- [7] Bentkamp, A., Blanchette, J.C., Cruanes, S., Waldmann, U.: Superposition for lambda-free higher-order logic. In: Galmiche, D., Schulz, S., Sebastiani, R. (eds.) *IJCAR 2018. LNCS*, vol. 10900, pp. 28–46. Springer (2018)
- [8] Bentkamp, A., Blanchette, J.C., Tourret, S., Vukmirović, P.: Superposition for full higher-order logic (technical report). Technical report (2021), https://matryoshka-project.github.io/pubs/hosup_report.pdf
- [9] Benz Müller, C., Paulson, L.C., Theiss, F., Fietzke, A.: LEO-II—A cooperative automatic theorem prover for higher-order logic. In: Armando, A., Baumgartner, P., Dowek, G. (eds.) *IJCAR 2008. LNCS*, vol. 5195, pp. 162–170. Springer (2008)
- [10] Bhayat, A., Rege, G.: Set of support for higher-order reasoning. In: Konev, B., Urban, J., Rümmer, P. (eds.) *PAAR-2018. CEUR Workshop Proceedings*, vol. 2162, pp. 2–16. CEUR-WS.org (2018)
- [11] Bhayat, A., Rege, G.: A combinator-based superposition calculus for higher-order logic. In: Peltier, N., Sofronie-Stokkermans, V. (eds.) *IJCAR 2020, Part I. LNCS*, vol. 12166, pp. 278–296. Springer (2020)
- [12] Böhme, S., Nipkow, T.: Sledgehammer: Judgement Day. In: Giesl, J., Hähnle, R. (eds.) *IJCAR 2010. LNCS*, vol. 6173, pp. 107–121. Springer (2010)
- [13] Brown, C.E.: Satallax: An automatic higher-order prover. In: Gramlich, B., Miller, D., Sattler, U. (eds.) *IJCAR 2012. LNCS*, vol. 7364, pp. 111–117. Springer (2012)
- [14] Cruanes, S.: Extending Superposition with Integer Arithmetic, Structural Induction, and Beyond. Ph.D. thesis, École polytechnique (2015)
- [15] Fitting, M.: *Types, Tableaus, and Gödel’s God*. Kluwer (2002)
- [16] Ganzinger, H., Stuber, J.: Superposition with equivalence reasoning and delayed clause normal form transformation. *Information and Computation* **199**(1–2), 3–23 (2005)
- [17] Gordon, M.J.C., Melham, T.F. (eds.): *Introduction to HOL: A Theorem Proving Environment for Higher Order Logic*. Cambridge University Press (1993)
- [18] Huet, G.P.: A mechanization of type theory. In: Nilsson, N.J. (ed.) *IJCAI-73*. pp. 139–146. William Kaufmann (1973)
- [19] Jensen, D.C., Pietrzykowski, T.: Mechanizing ω -order type theory through unification. *Theor. Comput. Sci.* **3**(2), 123–171 (1976)
- [20] Kaliszyk, C., Sutcliffe, G., Rabe, F.: TH1: The TPTP typed higher-order form with rank-1 polymorphism. In: Fontaine, P., Schulz, S., Urban, J. (eds.) *PAAR-2016. CEUR Workshop Proceedings*, vol. 1635, pp. 41–55. CEUR-WS.org (2016)
- [21] Kotelnikov, E., Kovács, L., Suda, M., Voronkov, A.: A clausal normal form translation for FOOL. In: Benz Müller, C., Sutcliffe, G., Rojas, R. (eds.) *GCAI 2016. EPiC*, vol. 41, pp. 53–71. EasyChair (2016)

- [22] Ludwig, M., Waldmann, U.: An extension of the Knuth-Bendix ordering with LPO-like properties. In: Dershowitz, N., Voronkov, A. (eds.) LPAR-14. LNCS, vol. 4790, pp. 348–362. Springer (2007)
- [23] Nummelin, V., Bentkamp, A., Tourret, S., Vukmirović, P.: Superposition with first-class Booleans and inprocessing clausification. In: Platzer, A., Sutcliffe, G. (eds.) CADE-28. LNCS, Springer (2021)
- [24] Paulson, L.C., Blanchette, J.C.: Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In: Sutcliffe, G., Schulz, S., Ternovska, E. (eds.) IWIL-2010. EPIc, vol. 2, pp. 1–11. EasyChair (2012)
- [25] Schulz, S.: E - a brainiac theorem prover. *AI Commun.* **15**(2-3), 111–126 (2002)
- [26] Steen, A., Benz Müller, C.: The higher-order prover Leo-III. In: Galmiche, D., Schulz, S., Sebastiani, R. (eds.) IJCAR 2018. LNCS, vol. 10900, pp. 108–116. Springer (2018)
- [27] Sutcliffe, G.: The TPTP problem library and associated infrastructure—from CNF to TH0, TPTP v6.4.0. *J. Autom. Reason.* **59**(4), 483–502 (2017)
- [28] Vukmirović, P., Bentkamp, A., Blanchette, J., Cruanes, S., Nummelin, V., Tourret, S.: Making higher-order superposition work. In: Platzer, A., Sutcliffe, G. (eds.) CADE-28. LNCS, Springer (2021)
- [29] Vukmirović, P., Bentkamp, A., Nummelin, V.: Efficient full higher-order unification. In: Ariola, Z.M. (ed.) FSCD 2020. LIPIcs, vol. 167, pp. 5:1–5:17. Schloss Dagstuhl—Leibniz-Zentrum für Informatik (2020)
- [30] Vukmirović, P., Nummelin, V.: Boolean reasoning in a higher-order superposition prover. In: PAAR-2020. CEUR Workshop Proceedings, vol. 2752, pp. 148–166. CEUR-WS.org (2020)
- [31] Waldmann, U., Tourret, S., Robillard, S., Blanchette, J.: A comprehensive framework for saturation theorem proving. In: Peltier, N., Sofronie-Stokkermans, V. (eds.) IJCAR 2020, Part I. LNCS, vol. 12166, pp. 316–334. Springer (2020)