

# Boolean Reasoning in a Higher-Order Superposition Prover

Petar Vukmirović<sup>a</sup>, Visa Nummelin<sup>a</sup>

<sup>a</sup>Vrije Universiteit Amsterdam, De Boelelaan 1105, 1081 HV Amsterdam, The Netherlands

## Abstract

We present a pragmatic approach to extending a Boolean-free higher-order superposition calculus to support Boolean reasoning. Our approach extends inference rules that have been used only in a first-order setting, uses some well-known rules previously implemented in higher-order provers, as well as new rules. We have implemented the approach in the Zipperposition theorem prover. The evaluation shows highly competitive performance of our approach and clear improvement over previous techniques.

## Keywords

higher-order logic, theorem proving, superposition

## 1. Introduction

In the last decades, automatic theorem provers have been used successfully as backends to “hammers” in proof assistants [1, 2] and to software verifiers [3]. Most advanced provers, such as CVC4 [4], E [5], and Vampire [6], are based on first-order logic, whereas most frontends that use them are based on versions of higher-order logic. Thus, there is a large gap in expressiveness between front- and backends. This gap is bridged using well-known translations from higher-order to first-order logic [7, 8]. However, translations are usually less efficient than native support [9, 10, 11]. The distinguishing features of higher-order logic used by proof assistants that the translation must eliminate include  $\lambda$ -binders, function extensionality – the property that functions are equal if they agree on every argument, described by the axiom  $\forall(x, y : \tau \rightarrow \nu). (\forall(z : \tau). xz \approx yz) \Rightarrow x \approx y$ , and formulas occurring as arguments of function symbols [8].

A group of authors including Vukmirović [10] recently designed a complete calculus for extensional Boolean-free higher-order logic. This calculus is an extension of superposition, the calculus used in most successful provers such as E or Vampire. The extension removes the need to translate the first two above mentioned features of higher-order logic. Kotelnikov et al. [12, 13] extended the language of first-order logic to support the third feature of higher-order logic that requires translation. They described two approaches: one based on calculus-level treatment of Booleans and the other, which requires no changes to the calculus, based on preprocessing.

---


PAAR 2020: Seventh Workshop on Practical Aspects of Automated Reasoning, June 29–30, 2020, Paris, France (virtual)

✉ p.vukmirovic@vu.nl (P. Vukmirović); visa.nummelin@vu.nl (V. Nummelin)

🌐 <https://petarvukmirovic.github.io/home/> (P. Vukmirović)

🆔 0000-0001-7049-6847 (P. Vukmirović); 0000-0003-0078-790X (V. Nummelin)

© 2020 Copyright for this paper by its authors.  
Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

To fully bridge the gap between higher-order and first-order tools, we combine the two approaches: we use the efficient higher-order superposition calculus and extend it with inference rules that reason with Boolean terms. In early work, Kotelnikov et al. [12] have described a *FOOL paramodulation* rule that, under some order requirements, removes the need for the axiom describing the Boolean domain –  $\forall(p : o). p \approx \top \vee p \approx \perp$ . In this approach, it is assumed that a problem with formulas occurring as arguments of symbols is translated to first-order logic.

The backbone of our approach is based on an extension of this rule to higher-order logic. Namely, we do not translate away any Boolean structure that is nested inside non-Boolean terms and allow our rule to hoist the nested Booleans to the literal level. Then, we clausify the resulting formula (i.e., a clause that contains formulas in literals) using a new rule.

An important feature that we inherit by building on top of Bentkamp et al. [10] is support for (function) extensionality. Moving to higher-order logic with Booleans also means that we need to consider *Boolean extensionality*:  $\forall(p : o)(q : o). (p \Leftrightarrow q) \Rightarrow p \approx q$ . We extend the rules of Bentkamp et al. that treat function extensionality to support Boolean extensionality as well.

Rules that extend the two orthogonal approaches form the basis of our support for Boolean reasoning (Section 3). In addition, we have implemented rules that are inspired by the ones used in the higher-order provers Leo-III [14] and Satallax [15], such as elimination of Leibniz equality, primitive instantiation and treatment of choice operator [16]. We have also designed new rules that use higher-order unification to resolve Boolean formulas that are hoisted to literal level, delay clausification of non-atomic literals, reason about formulas under  $\lambda$ -binders, and many others. Even though the rules we use are inspired by the ones of refutationally complete higher-order provers, we do not guarantee completeness of our extension of  $\lambda$ -superposition.

We compare our native approach with two alternatives based on preprocessing (Section 4). First, we compare it to an axiomatization of the theory of Booleans. Second, inspired by work of Kotelnikov et al. [13], we implemented the preprocessing approach that does not require introduction of Boolean axioms. We also discuss some examples, coming from TPTP [17], that illustrate advantages and disadvantages of our approach (Section 5).

Our approach is implemented in the Zipperposition theorem prover [18, 19]. Zipperposition is an easily extensible open source prover that Bentkamp et al. used to implement their higher-order superposition calculus. We further extend their implementation.

We performed an extensive evaluation of our approach (Section 6). In addition to evaluating different configurations of our new rules, we have compared them to full higher-order provers CVC4, Leo-III, Satallax and Vampire. The results suggest that it is beneficial to natively support Boolean reasoning – the approach outperforms preprocessing-based approaches. Furthermore, it is very competitive with state-of-the-art higher order provers. We discuss the differences between our approach and the approaches we base on, as well as related approaches (Section 7).

## 2. Background

We base our work on Bentkamp et al.’s [10] extensional polymorphic clausal higher-order logic. We extend the syntax of this logic by adding logical connectives to the language of terms. The semantic of the logic is extended by interpreting Boolean type  $o$  as a two-element domain. This amounts to extending Bentkamp et al.’s fragment of higher-order logic to full-higher order logic

(HOL). Our notation, definitions and the following text are largely based on Bentkamp et al.'s.

A signature is a quadruple  $(\Sigma_{\text{ty}}, \mathcal{V}_{\text{ty}}, \Sigma, \mathcal{V})$  where  $\Sigma_{\text{ty}}$  is a set of type constructors,  $\mathcal{V}_{\text{ty}}$  is a set of type variables and  $\Sigma$  and  $\mathcal{V}$  are sets of constants and term variables, respectively. We require nullary type constructors  $\iota$  and  $o$ , as well as binary constructor  $\rightarrow$  to be in  $\Sigma_{\text{ty}}$ . A type  $\tau, v$  is either a type variable  $\alpha \in \mathcal{V}_{\text{ty}}$  or of the form  $\kappa(\tau_1, \dots, \tau_n)$  where  $\kappa$  is an  $n$ -ary type constructor. We write  $\kappa$  for  $\kappa()$ ,  $\tau \rightarrow v$  for  $\rightarrow(\tau, v)$ , and we abbreviate tuples  $(a_1, \dots, a_n)$  as  $\bar{a}_n$  for  $n \geq 0$ . Similarly, we drop parentheses to shorten  $\tau_1 \rightarrow (\dots \rightarrow (\tau_{n-1} \rightarrow \tau_n) \dots)$  into  $\tau_1 \rightarrow \dots \rightarrow \tau_n$ . Each symbol in  $\Sigma$  is assigned a type declaration of the form  $\Pi \bar{\alpha}_n. \tau$  where all variables occurring in  $\tau$  are among  $\bar{\alpha}_n$ .

Function symbols  $a, b, f, g, \dots$  are elements of  $\Sigma$ ; their type declarations are written as  $f : \Pi \bar{\alpha}_n. \tau$ . Term variables from the set  $\mathcal{V}$  are written  $x, y, z, \dots$  and we denote their types as  $x : \tau$ . When the type is not important, we omit type declarations. We assume that symbols  $\top, \perp, \neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$  with their standard meanings and type declarations are elements of  $\Sigma$ . Furthermore, we assume that polymorphic symbols  $\forall$  and  $\exists$  with type declarations  $\Pi \alpha. (\alpha \rightarrow o) \rightarrow o$  and  $\approx : \Pi \alpha. \alpha \rightarrow \alpha \rightarrow o$  are in  $\Sigma$ , with their standard meanings. All these symbols are called *logical symbols*. We write binary logical symbols in infix notation.

Terms are defined inductively as follows. Variables  $x : \tau$  are terms of type  $\tau$ . If  $f : \Pi \bar{\alpha}_n. \tau$  is in  $\Sigma$  and  $\bar{v}_n$  is a tuple of types, called type arguments, then  $f\langle \bar{v}_n \rangle$  (written as  $f$  if  $n = 0$ , or if type arguments can be inferred from the context) is a term of type  $\tau\{\bar{\alpha}_n \mapsto \bar{v}_n\}$ , called constant. If  $x$  is a variable of type  $\tau$  and  $s$  is a term of type  $v$  then  $\lambda x. s$  is a term of type  $\tau \rightarrow v$ . If  $s$  and  $t$  are of type  $\tau \rightarrow v$  and  $\tau$ , respectively, then  $st$  is a term of type  $v$ . We call terms of Boolean type ( $o$ ) *formulas* and denote them by  $f, g, h, \dots$ ; we use  $p, q, r, \dots$  for variables whose result type is  $o$  and  $\mathfrak{p}, \mathfrak{q}, \mathfrak{r}$  for constants with the same result type. We shorten iterated lambda abstraction  $\lambda x_1. \dots \lambda x_n. s$  to  $\lambda \bar{x}_n. s$ , and iterated application  $(s t_1) \dots t_n$  to  $s \bar{t}_n$ . We assume the standard notion of free and bound variables, capture-avoiding substitutions  $\sigma, \rho, \theta, \dots$ , and  $\alpha$ -,  $\beta$ -,  $\eta$ -conversion. Unless stated otherwise, we view terms as  $\alpha\beta\eta$ -equivalence classes, with  $\eta$ -long  $\beta$ -reduced form as the representative. Each term  $s$  can be uniquely written as  $\lambda \bar{x}_m. a \bar{t}_n$  where  $a$  is either variable or constant and  $m, n \geq 0$ ; we call  $a$  the *head* of  $s$ . We say that a term  $a \bar{t}_n$  is written in *spine notation* [20]. Following our previous work [21], we define nonstandard notion of subterms and positions inductively as a graceful extension of the first-order counterparts: a term  $s$  is a subterm of itself at position  $\varepsilon$ . If  $s$  is a subterm of  $t_i$  at position  $p$  then  $s$  is a subterm of  $a \bar{t}_n$  at position  $i.p$ , where  $a$  is a head. If  $s$  is a subterm of  $t$  at position  $p$  then  $s$  is a subterm of  $\lambda x. t$  at position  $1.p$ . We use  $s|_p$  to denote subterm of  $s$  at position  $p$ .

Given a formula  $f$  we call its Boolean subterm  $f|_p$  a *top-level Boolean* if for all proper prefixes  $q$  of  $p$ , the head of  $f|_q$  is a logical constant. Otherwise, we call it a *nested Boolean*. For example, in the formula  $f = h a \approx g (\mathfrak{p} \Rightarrow \mathfrak{q}) \vee \neg \mathfrak{p}$ ,  $f|_1$  and  $f|_2$  are top-level Booleans, whereas  $f|_{1.2.1}$  is a nested Boolean (as well as its subterms). Only top-level Booleans are allowed in first-order logic, whereas nested Booleans are characteristic for higher-order logic. A formula is called an *atom* if it is of the form  $a \bar{t}_n$ , where  $a$  is a non-logical head, or of the form  $s \approx t$ , where if  $s$  or  $t$  are of type  $o$ , and one of them has a logical head, the other one must be  $\top$  or  $\perp$ . A *literal*  $L$  is an atom or its negation. A *clause*  $C$  is a multiset of literals, interpreted and written (abusing  $\vee$ ) disjunctively as  $L_1 \vee \dots \vee L_n$ . We write  $s \not\approx t$  for  $\neg(s \approx t)$ . We say a variable is *free* in a clause  $C$  if it is not bound inside any subterm of a literal in  $C$ .

### 3. The Native Approach

Some support for Booleans was already present in Zipperposition before we started extending the calculus of Bentkamp et al. In this section, we start by describing the internals of Zipperposition responsible for reasoning with Booleans. We continue by describing 15 rules that we have implemented. For ease of presentation we divide them in three categories. We assume some familiarity with the superposition calculus [22] and adopt the notation used by Schulz [23].

#### 3.1. Support for Booleans in Zipperposition

Zipperposition is an open source<sup>1</sup> prover written in OCaml. From its inception, it was designed as a prover that supports easy extension of its base superposition calculus to various theories, including arithmetic, induction and limited support for higher-order logic [18].

Zipperposition internally stores applications in flattened, spine notation. It also exploits associativity of  $\wedge$  and  $\vee$  to flatten nested applications of these symbols. Thus, the terms  $p \wedge (q \wedge r)$  and  $(p \wedge q) \wedge r$  are represented as  $\wedge p q r$ . The prover's support for  $\lambda$ -terms is used to represent quantified nested Booleans: formulas  $\forall x. f$  and  $\exists x. f$  are represented as  $\forall (\lambda x. f)$  and  $\exists (\lambda x. f)$ . After clausification of the input problem, no nested Booleans will be modified or renamed using fresh predicate symbols.

The version of Zipperposition preceding our modifications distinguished between equational and non-equational literals. Following E [5], we modified Zipperposition to represent all literals equationally: a non-equational literal  $f$  is stored as  $f \approx \top$ , whereas  $\neg f$  is stored as  $f \not\approx \top$ . Equations of the form  $f \approx \perp$  and  $f \not\approx \perp$  are transformed into  $f \not\approx \top$  and  $f \approx \top$ , respectively.

#### 3.2. Core Rules

Kotelnikov et al. [12], to the best of our knowledge, pioneered the approach of extending a first-order superposition prover to support nested Booleans. They call effects of including the axiom  $\forall(p : o). p \approx \top \vee p \approx \perp$  a “recipe for disaster”. To combat the explosive behavior of the axiom, they imposed the following two requirements to the simplification order  $>$  (which is a parameter to the superposition calculus):  $\top > \perp$  and  $\top$  and  $\perp$  are two smallest ground terms with respect to  $>$ . If these two requirements are met, there is no self-paramodulation of the clause and only paramodulation possible is from literal  $p \approx \top$  of the mentioned axiom into a Boolean subterm of another clause. Finally, Kotelnikov et al. replace the axiom with the inference rule *FOOL Paramodulation* (FP):

$$\frac{C[f]}{C[\top] \vee f \approx \perp} \text{FP}$$

where  $f$  is a nested non-variable Boolean subterm of clause  $C$ , different from  $\top$  and  $\perp$ . In addition, they translate the initial problem containing nested Booleans to first-order logic without interpreted Booleans; this translation introduces proxy symbols for  $\top$  and  $\perp$ , and proxy type for  $o$ .

<sup>1</sup><https://github.com/sneeuwballen/zipperposition>

We created two rules that are syntactically similar to FP but are adapted for higher-order logic with one key distinction – we do not perform any translation:

$$\frac{C[f]}{C[\perp] \vee f \approx \top} \text{CASES} \qquad \frac{C[f]}{C[\perp] \vee f \approx \top \quad C[\top] \vee f \not\approx \top} \text{CASESSIMP}$$

The double line in the definition of CASESSIMP denotes that the premise is replaced by conclusions; obviously, the prover that uses the rules should not include them both at the same time. In addition, since literals  $f \approx \perp$  are represented as negative equations  $f \not\approx \top$ , which cannot be used to paramodulate from, we change the first requirement on the order to  $\perp > \top$ .

These two rules hoist Boolean subterms  $f$  to the literal level; therefore, some results of CASES and CASESSIMP will have literals of the form  $f \approx \top$  (or  $f \not\approx \top$ ) where  $f$  is not an atom. This introduces the need for the rule called eager clausification (EC):

$$\frac{C}{D_1 \cdots D_m} \text{EC}$$

We say that a clause is *standard* if all of its literals are of the form  $s \approx t$ , where  $s$  and  $t$  are not Booleans or of the form  $f \approx \top$ , where the head of  $f$  is not a logical symbol and  $\approx$  denotes  $\approx$  or  $\not\approx$ . The rule EC is applicable if clause  $C = L_1 \vee \cdots \vee L_n$  is not standard. The resulting clauses  $\overline{D}_m$  represent the result of clausification of the formula  $\forall \bar{x}. L_1 \vee \cdots \vee L_n$  where  $\bar{x}$  are all free variables of  $C$ .

An advantage of leaving nested Booleans unmodified is that the prover will be able to prove some problems containing them without using the prolific rules described above. For example, given two clauses  $f(p x \Rightarrow p y) \approx a$  and  $f(p a \Rightarrow p b) \not\approx a$ , the empty clause can easily be derived without the above rules. A disadvantage of this approach is that the proving process will periodically be interrupted by expensive calls to the clausification algorithm.

Naive application of CASES and CASESSIMP rules can result in many redundant clauses. Consider a clause  $C = p(p(p(pa))) \approx \top$  where  $p : o \rightarrow o$ ,  $a : o$ . Then, the clause  $D = a \approx \top \vee p \perp \approx \top$  can be derived from  $C$  in eight ways using the rules, depending on which nested Boolean subterm was chosen for the inference. In general, if a clause has a subterm occurrence of the form  $p^n a$ , where both  $p$  and  $a$  have result type  $o$ , the clause  $a \approx \top \vee p \perp \approx \top$  can be derived in  $2^{n-1}$  ways. To combat these issues we implemented pragmatic restrictions of the rule: only the leftmost outermost (or innermost) eligible subterm will be considered. With this modification  $D$  can be derived in only one way. Furthermore, some intermediate conclusions of the rules will not be derived, pruning the search space.

The clausification algorithm by Nonnengart and Weidenbach [24] aggressively simplifies the input problem using well-known Boolean equivalences before clausifying it. For example, the formula  $p \wedge \top$  will be replaced by  $p$ . To simplify nested Booleans we implemented the rule

$$\frac{C[f\sigma]}{C[g\sigma]} \text{BOOLSIMP}$$

where  $f \rightarrow g \in E$  runs over fixed set of rewrite rules  $E$ , and  $\sigma$  is any substitution. In the current implementation of Zipperposition,  $E$  consists of the rules described by Nonnengart

and Weidenbach [24, Section 3]. This set contains the rules describing how each logical symbol behaves when either of its argument is  $\top$  or  $\perp$ : for example, it includes  $\top \Rightarrow p \longrightarrow p$  and  $p \Rightarrow \top \longrightarrow \top$ . Leo-III implements a similar rule, called *simp* [25, Section 4.2.1].

Our decision to represent negative atoms as negative equations was motivated by the need to alter Zipperposition’s earlier behavior as little as possible. Namely, negative atoms were not used as literals that can be used to paramodulate from, and as such added to the laziness of the superposition calculus. However, it might be useful to consider unit clauses of the form  $f \not\approx \top$  as  $f \approx \perp$  to strengthen demodulation. To that end, we have introduced the following rule:

$$\frac{f \not\approx \top \quad C[f\sigma]}{f \not\approx \top \quad C[\perp]} \text{BOOLDEMOD}$$

### 3.3. Higher-Order Considerations

To achieve refutational completeness of higher-order resolution and similar calculi it is necessary to instantiate variables with result type  $o$ , *predicate variables*, with arbitrary formulas [25, 16]. Fortunately, we can approximate the formulas using a complete set of logical symbols (e.g.,  $\neg$ ,  $\forall$ , and  $\wedge$ ). Since such an approximation is not only necessary for completeness of some calculi, but very useful in practice, we implemented the *primitive instantiation* (PI) rule:

$$\frac{C \vee \lambda \bar{x}_m. p \bar{s}_n \approx t}{(C \vee \lambda \bar{x}_m. p \bar{s}_n \approx t)\{p \mapsto f\}} \text{PI}$$

where  $p$  is a free variable of the type  $\tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow o$ . Choosing a different  $f$  that instantiates  $p$ , we can balance between explosiveness of approximating a complete set of logical symbols and incompleteness of pragmatic approaches. We borrow the notion of imitation from higher-order unification jargon [21]: we say that the term  $\lambda \bar{x}_m. f (y_1 \bar{x}_m) \dots (y_n \bar{x}_m)$  is an *imitation* of constant  $f : \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow \tau$  for some variable  $z$  of type  $\nu_1 \rightarrow \dots \rightarrow \nu_m \rightarrow \tau$ . Variables  $\bar{y}_n$  are fresh free variables, where each  $y_i$  has the type  $\nu_1 \rightarrow \dots \rightarrow \nu_m \rightarrow \tau_i$ ; variable  $x_i$  is of type  $\nu_i$ .

Rule PI was already implemented by Simon Cruanes in Zipperposition, before we started our modifications. The rule has different modes that generate sets of possible terms  $f$  for  $p : \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow o$ : *Full*, *Pragmatic*, and *Imit $\star$*  where  $\star$  is an element of a set of logical constants  $P = \{\wedge, \vee, \approx, \langle \alpha \rangle, \neg, \forall, \exists\}$ . Mode *Full* contains imitations (for  $p$ ) of all elements of  $P$ . Mode *Pragmatic* contains imitations of  $\neg$ ,  $\top$  and  $\perp$ ; if there exist indices  $i, j$  such that  $i \neq j$  and  $\tau_i = \tau_j$ , it contains  $\lambda \bar{x}_n. x_i \approx x_j$ ; if there exist indices  $i, j$  such that  $i \neq j$ , and  $\tau_i = \tau_j = o$ , then it contains  $\lambda \bar{x}_n. x_i \wedge x_j$  and  $\lambda \bar{x}_n. x_i \vee x_j$ ; if for some  $i, \tau_i = o$ , then it contains  $\lambda \bar{x}_n. x_i$ . Mode *Imit $\star$*  contains imitations of  $\top$ ,  $\perp$  and  $\star$  (except for *Imit $\forall\exists$*  which contains imitations of both  $\forall$  and  $\exists$ ).

While experimenting with our implementation we have noticed some proof patterns that led us to come up with the following modifications. First, it often suffices to perform PI only on initial clauses – which is why we allow the rule to be applied only to the clauses created using at most  $k$  generating inferences. Second, if the rule was used in the proof, its premise is usually only used as part of that inference – which is why we implemented a version of PI that removes

the clause after all possible PI inferences have been performed. We observed that the mode *Imit\** is useful in practice since often only a single approximation of a logical symbol is necessary.

Efficiently treating axiom of choice is notoriously difficult for higher-order provers. Andrews formulates this axiom as  $\forall(p : \alpha \rightarrow o). (\exists(x : \alpha). px) \Rightarrow p(\varepsilon p)$ , where  $\varepsilon : \Pi\alpha. (\alpha \rightarrow o) \rightarrow \alpha$  denotes the *choice operator* [16]. After clausification, this axiom becomes  $px \approx \top \vee p(\varepsilon p) \approx \top$ . Since term  $px$  matches any Boolean term in the proof state, this axiom is very explosive. Therefore, Leo-III [14] deals with the choice operator on the calculus level. Namely, whenever a clause  $C = px \approx \top \vee p(fp) \approx \top$  is chosen for processing,  $C$  is removed from the proof state and  $f$  is added to set of choice functions  $CF$  (which initially contains just  $\varepsilon$ ). Later, elements of  $CF$  will be used to heuristically instantiate the axiom of choice. We reused the method of recognizing choice functions, but generalized the rule for creating the instance of the axiom (assuming  $\xi \in CF$ ):

$$\frac{C[\xi t]}{x(ty) \approx \top \vee x(t(\xi(\lambda z. x(tz)))) \approx \top} \text{CHOICE}$$

Let  $D$  be the conclusion of CHOICE. The fresh variable  $x$  in  $D$  acts as arbitrary context around  $t$ , the chosen instantiation for  $p$  from axiom of choice; the variable  $x$  can later be replaced by imitation of logical symbols to create more complex instantiations of the choice axiom. To generate useful instances early, we create  $D\{x \mapsto \lambda z. z\}$  and  $D\{x \mapsto \lambda z. \neg z\}$ . Then, based on Zipperposition parameters,  $D$  will either be deleted or kept. Note that  $D$  will not subsume its instances, since the matching algorithm of Zipperposition is too weak for this.

Most provers natively support extensionality reasoning: Bhayat et al. [26] modify first-order unification to return unification constraints consisting of pairs of terms of functional type, whereas Steen relies on the unification rules of Leo-III's calculus [25, Section 4.3.3.] to deal with extensionality. Bentkamp et al [10] altered core generating inference rules of the superposition calculus to support extensionality. Instead of requiring that terms involved in the inference are unifiable, it is required that they can be decomposed into *disagreement pairs* such that at least one of the disagreement pairs is of functional type. Disagreement pairs of terms  $s$  and  $t$  of the same type are defined inductively using function  $\text{dp}$ :  $\text{dp}(s, t) = \emptyset$  if  $s$  and  $t$  are equal;  $\text{dp}(a\bar{s}_n, b\bar{t}_m) = \{(a\bar{s}_n, b\bar{t}_m)\}$  if  $a$  and  $b$  are different heads;  $\text{dp}(\lambda x. s, \lambda y. t) = \{(\lambda x. s, \lambda y. t)\}$ ;  $\text{dp}(a\bar{s}_n, a\bar{t}_n) = \bigcup_{i=1}^n \text{dp}(s_i, t_i)$ . Then the extensionality rules are stated as follows:

$$\frac{s \approx t \vee C \quad u[s'] \approx v \vee D}{(s_1 \approx s'_1 \vee \dots \vee s_n \approx s'_n \vee u[t] \approx v \vee C \vee D)\sigma} \text{EXTSUP}$$

$$\frac{s \approx s' \vee C}{(s_1 \approx s'_1 \vee \dots \vee s_n \approx s'_n \vee C)\sigma} \text{EXTER}$$

$$\frac{s \approx t \vee s' \approx u \vee C}{(s_1 \approx s'_1 \vee \dots \vee s_n \approx s'_n \vee t \approx u \vee s' \approx u \vee C)\sigma} \text{EXTEF}$$

Rules EXT<sub>SUP</sub>, EXT<sub>ER</sub>, and EXT<sub>EF</sub> are extensional versions of superposition, equality resolution and equality factoring [23]. The union of these three rules is denoted by EXT. In each rule,  $\sigma$  is a most general unifier of the types of  $s$  and  $s'$ , and  $\text{dp}(s\sigma, s'\sigma) = \{(s_1, s'_1), \dots, (s_n, s'_n)\}$ . All

side conditions for extensional rules are the same as for the standard rules, except that condition that  $s$  and  $s'$  are unifiable is replaced by the condition that at least one  $s_i$  is of functional type and that  $n > 0$ . This rule is easily extended to support Boolean extensionality by requiring that at least one  $s_i$  is of functional or type  $o$ , and adding the condition “ $\text{dp}(f, g) = \{(f, g)\}$  if  $f$  and  $g$  are different formulas” to the definition of  $\text{dp}$ .

Consider the clause  $f(\neg p \vee \neg q) \approx f(\neg(p \wedge q))$ . This problem is obviously unsatisfiable, since arguments of  $f$  on different sides of the disequation are extensionally equal; however, without EXT rules Zipperposition will rely on CASES(SIMP) and EC rules to derive the empty clause. Rule EXTER will generate  $C = \neg p \vee \neg q \approx \neg(p \wedge q)$ . Then,  $C$  will get classified using EC, effectively reducing the problem to  $\neg(\neg p \vee \neg q \Leftrightarrow \neg(p \wedge q))$ , which is first-order.

Zipperposition restricts EXTSUP by requiring that  $s$  and  $s'$  are not of function or Boolean types. If the terms are of function type, our experience is that better treatment of function extensionality is to apply fresh free variables (or Skolem terms, depending on the sign [10]) to both sides of a (dis)equation to reduce it to a first-order literal; Boolean extensionality is usually better supported by applying EC on the top-level Boolean term. Thus, for the following discussion we can assume  $s$  and  $s'$  are not  $\lambda$ -abstractions or formulas. Then, EXTSUP is applicable if  $s$  and  $s'$  have the same head, and a functional or Boolean subterm. To efficiently retrieve such terms, we added an index that maps symbols to positions in clauses where they appear as a head of a term that has a functional or Boolean subterm. This index will be empty for first-order problems, incurring no overhead if extensionality reasoning is not needed. Furthermore, we do not apply EXT rules if all disagreement pairs have at least one side whose head is a variable; those will be dealt with more efficiently using standard, non-extensional, versions of the rules. We also eagerly resolve literals  $s_i \approx s'_i$  using at most one unifier returned by terminating, pragmatic variant of unification algorithm by Vukmirović et al. [21].

Expressiveness of higher-order logic allows users to define equality using a single axiom, called Leibniz equality [16]:  $\forall(x : \alpha)(y : \alpha). (\forall(p : \alpha \rightarrow o). px \Rightarrow py) \Rightarrow x \approx y$ . Leibniz equality often appears in TPTP problems. Since modern provers have the native support for equality, it is usually beneficial to recognize and replace occurrences of Leibniz equality.

Before we began our modifications, Zipperposition had a powerful rule that recognizes clauses that contain variations of Leibniz equality and instantiates them with native equality. This rule was designed by Simon Cruanes, and to the best of our knowledge, it has not been documented so far. With his permission we describe this rule as follows:

$$\frac{p \bar{s}_n^1 \approx \top \vee \dots \vee p \bar{s}_n^i \approx \top \vee p \bar{t}_n^1 \approx \top \vee \dots \vee p \bar{t}_n^j \approx \top \vee C}{(p \bar{s}_n^1 \approx \top \vee \dots \vee p \bar{s}_n^i \approx \top \vee C)\sigma} \text{ELIMPREDVAR}$$

where  $p$  is a free variable,  $p$  does not occur in any  $s_k^l$  or  $t_k^l$ , or in  $C$ , and  $\sigma$  is defined as  $\{p \mapsto \lambda \bar{x}_n. \bigvee_{k=1}^j (\bigwedge_{l=1}^n x_l \approx t_l^k)\}$ .

To better understand how this rule removes variable-headed negative literals, consider the clause  $C = p a_1 a_2 \approx \top \vee p b_1 b_2 \approx \top \vee p c_1 c_2 \approx \top$ . Since all side conditions are fulfilled, the rule ELIMPREDVAR will generate  $\sigma = \{p \mapsto \lambda xy. (x \approx b_1 \wedge y \approx b_2) \vee (x \approx c_1 \wedge y \approx c_2)\}$ . After applying  $\sigma$  to  $C$  and subsequent  $\beta$ -reduction, negative literal  $p b_1 b_2 \approx \top$  will reduce to  $(b_1 \approx b_1 \wedge b_2 \approx b_2) \vee (b_1 \approx c_1 \wedge b_2 \approx c_2) \approx \top$ , which is equivalent to  $\perp$ . Thus, we can remove this literal and all negative literals of the form  $p \bar{t}_n \approx \top$  from  $C$  and apply  $\sigma$  to the remaining ones.



The previous rule removes all variables occurring in disequations in one attempt. We implemented two rules that behave more lazily, inspired by the ones present in Leo-III and Satallax:

$$\frac{p \bar{s}_n \approx \top \vee p \bar{t}_n \not\approx \top \vee C}{(s_i \approx t_i \vee C)\sigma} \text{ELIMLEIBNIZ+} \quad \frac{p \bar{s}_n \not\approx \top \vee p \bar{t}_n \approx \top \vee C}{(s_i \approx t_i \vee C)\sigma'} \text{ELIMLEIBNIZ-}$$

where  $p$  is a free variable,  $p$  does not occur in  $t_i$ ,  $\sigma = \{p \mapsto \lambda \bar{x}_n. x_i \approx t_i\}$  and  $\sigma' = \{p \mapsto \lambda \bar{x}_n. \neg(x_i \approx t_i)\}$ . This rule differs from ELIMPREDVAR in three ways. First, it acts on occurrences of variables in both positive and negative literals. Second, due to its simplicity, it usually does not require EC as the following step. Third, it imposes much weaker conditions on  $p$ . However, removing all negative variables in one step might improve performance. Coming back to example of the clause  $C = p a_1 a_2 \approx \top \vee p b_1 b_2 \not\approx \top \vee p c_1 c_2 \not\approx \top$ , we can apply ELIMLEIBNIZ+ using the substitution  $\sigma = \{\lambda xy. x \approx b_1\}$  to obtain the clause  $C' = a_1 \approx b_1 \vee a_1 \not\approx c_1$ .

### 3.4. Additional Rules

Zipperposition's unification algorithm [21] uses flattened representation of terms with logical operators  $\wedge$  and  $\vee$  for heads to unify terms that are not unifiable modulo  $\alpha\beta\eta$ -equivalence, but are unifiable modulo associativity and commutativity of  $\wedge$  and  $\vee$ . Let  $\diamond$  denote either  $\wedge$  or  $\vee$ . When the unification algorithm is given two terms  $\diamond \bar{s}_n$  and  $\diamond \bar{t}_n$ , where neither of  $\bar{s}_n$  nor  $\bar{t}_n$  contain duplicates, it performs the following steps: First, it removes all terms that appear in both  $\bar{s}_n$  and  $\bar{t}_n$  from the two argument tuples. Next, the remaining terms are sorted first by their head term and then by their weight. Finally, the sorted lists are unified pairwise. As an example, consider the problem of unifying the pair  $(\wedge (p a) (q (f a)), \wedge (q (f a)) (r (f (f a))))$  where  $r$  is a free variable. If the arguments of  $\wedge$  are simply sorted as described above, we would unsuccessfully try to unify  $p a$  with  $q (f a)$ . However, by removing term  $q (f a)$  from the argument lists, we will be left with the problem  $(p a, r (f (f a)))$  which has a unifier.

The winner of THF division of CASC-27 [27], Satallax [15], has one crucial advantage over Zipperposition: it is based on higher-order tableaux, and as such it does not require formulas to be converted to clauses. The advantage of tableaux is that once it instantiates a variable with a term, this instantiation naturally propagates through the whole formula. In Zipperposition, which is based on higher-order superposition, the original formula is clausified and instantiating a variable in a clause  $C$  does not automatically instantiate it in all clauses that are results of clausification of the same formula as  $C$ . To mitigate this issue, we have created extensions of equality resolution and equality factoring that take Boolean extensionality into account:

$$\frac{s \approx s' \vee C}{C\sigma} \text{BOOLER} \quad \frac{p \bar{s}_n \approx \top \vee s' \not\approx \top \vee C}{(p \bar{s}_n \approx \neg s' \vee C)\sigma} \text{BOOLEF+-}$$

$$\frac{p \bar{s}_n \not\approx \top \vee s' \approx \top \vee C}{(p \bar{s}_n \approx \neg s' \vee C)\sigma} \text{BOOLEF-+} \quad \frac{p \bar{s}_n \not\approx \top \vee s' \not\approx \top \vee C}{(p \bar{s}_n \approx s' \vee C)\sigma} \text{BOOLEF--}$$

All side conditions except for the ones concerning the unifiability of terms are as in the original equality resolution and equality factoring rules. In rule BOOLER,  $\sigma$  unifies  $s$  and  $\neg s'$ . In the

$+ -$  and  $- +$  versions of **BOOLEF**,  $\sigma$  unifies  $p\bar{s}_n$  and  $\neg s'$ , and in the remaining version it unifies  $p\bar{s}_n$  and  $s'$ . Intuitively, these rules bring Boolean (dis)equations in the appropriate form for application of the corresponding base rules. It suffices to consider literals of the form  $s \approx s'$  for **BOOLEF** since Zipperposition rewrites  $s \Leftrightarrow t \approx \top$  and  $\neg(s \Leftrightarrow t) \not\approx \top$  to  $s \approx t$  (and does analogous rewriting into  $s \not\approx t$ ).

Another approach to mitigate harmful effects of eager clausification is to delay it as long as possible. Following the approach by Ganzinger and Stuber [28], we represent every input formula  $f$  as a unit clause  $f \approx \top$  and use the following lazy clausification (LC) rules:

$$\begin{array}{c} \frac{(f \wedge g) \approx \top \vee C}{f \approx \top \vee C \quad g \approx \top \vee C} \text{LC}_\wedge \quad \frac{(f \vee g) \approx \top \vee C}{f \approx \top \vee g \approx \top \vee C} \text{LC}_\vee \quad \frac{(f \Rightarrow g) \approx \top \vee C}{f \not\approx \top \vee g \approx \top \vee C} \text{LC}_\Rightarrow \\ \\ \frac{(\neg f) \approx \top \vee C}{f \not\approx \top \vee C} \text{LC}_\neg \quad \frac{(\forall x. f) \approx \top \vee C}{f\{x \mapsto y\} \vee C} \text{LC}_\forall \quad \frac{(\exists x. f) \approx \top \vee C}{f\{x \mapsto \text{sk}\langle \bar{\alpha} \rangle \bar{y}_n\} \vee C} \text{LC}_\exists \\ \\ \frac{f \approx g \vee C}{f \not\approx \top \vee g \approx \top \vee C \quad f \approx \top \vee g \not\approx \top \vee C} \text{LC}_\approx \end{array}$$

The rules described above are as given by Ganzinger and Stuber (adapted to our setting), with the omission of rules for negative literals ( $f \not\approx \top$ ), which are easy to derive and which can be found in their work [28]. In  $\text{LC}_\approx$  we require both  $f$  and  $g$  to be formulas and at least one of them not to be  $\top$ . In  $\text{LC}_\forall$ ,  $y$  is a fresh variable, and in  $\text{LC}_\exists$ ,  $\text{sk}$  is a fresh symbol and  $\bar{\alpha}$  and  $\bar{y}_n$  are all the type and term variables occurring freely in  $\exists x. f$ .

Naive application of the LC rules can result in exponential blowup in problem size. To avoid this, we rename formulas that have repeated occurrences. We keep the count of all non-atomic formulas occurring as either side of a literal. Before applying the LC rule on a clause  $f \approx \top \vee C$ , we check whether the number of  $f$ 's occurrences exceeds the threshold  $k$ . If it does, based on the polarity of the literal  $f \approx \top$ , we add the clause  $p\bar{y}_n \not\approx \top \vee f \approx \top$  (if the literal is positive) or  $p\bar{y}_n \approx \top \vee f \not\approx \top$  (if the literal is negative), where  $\bar{y}_n$  are all free variables of  $f$  and  $p$  is a fresh symbol. Then, we replace the clause  $f \approx \top \vee C$  by  $p\bar{y}_n \approx \top \vee C$ .

Before the number of occurrences of  $f$  is checked, we first check (using a fast, incomplete matching algorithm) if there is a formula  $g$ , for which definition was already introduced, such that  $g\sigma = f$ , for some substitution  $\sigma$ . This check can have three outcomes. First, if the definition  $q\bar{x}_n$  is already introduced for  $g$  with the polarity matching that of  $f \approx \top$ , then  $f$  is replaced by  $(q\bar{x}_n)\sigma$ . Second, if the definition was introduced, but with different polarity, we create the clause defining  $g$  with the missing polarity, and replace  $f$  with  $(q\bar{x}_n)\sigma$ . Last, if there is no renamed formula  $g$  generalizing  $f$ , then we perform the previously described check.

In addition to reusing names for formula definitions, we reuse the Skolem symbols introduced by the  $\text{LC}_\exists$  rule. When  $\text{LC}_\exists$  is applied to  $f = \exists x. f'$  we check if there is a Skolem  $\text{sk}\langle \bar{\alpha}_m \rangle \bar{y}_n$  introduced for a formula  $g = \exists x. g'$ , such that  $g\sigma = f$ . If so, the symbol  $\text{sk}$  is reused and  $\exists x. f'$  is replaced by  $f'\{x \mapsto (\text{sk}\langle \bar{\alpha}_m \rangle \bar{y}_n)\sigma\}$ . Renaming and name reusing techniques are inspired by the VCNF algorithm described by Reger et al. [29].

Rules **CASES** and **CASESIMP** deal with Boolean terms, but we need to rely on extensionality reasoning to deal with  $\lambda$ -abstractions whose body has type  $o$ . Using the observation that the

formula  $\forall \bar{x}_n. f$  implies that  $\lambda \bar{x}_n. f$  is extensionally equal to  $\lambda \bar{x}_n. \top$  (and similarly, if  $\forall \bar{x}_n. \neg f$ , then  $\lambda \bar{x}_n. f \approx \lambda \bar{x}_n. \perp$ ), we designed the following rule (where all free variables of  $f$  are  $\bar{x}_n$  and variables occurring freely in  $C$ ):

$$\frac{C[\lambda \bar{x}_n. f]}{(\forall \bar{x}_n. f) \approx \top \vee C[\lambda \bar{x}_n. \top] \quad (\forall \bar{x}_n. \neg f) \approx \top \vee C[\lambda \bar{x}_n. \perp]} \text{INTERPRET } \lambda$$

## 4. Alternative Approaches

An alternative to heavy modifications of the prover needed to support the rules described above is to treat Booleans as yet another theory. Since the theory of Booleans is finitely axiomatizable, simply stating those axioms instead of creating special rules might seem appealing. Another approach is to preprocess nested Booleans by hoisting them to the top level.

### 4.1. Axiomatization

A simple axiomatization of the theory of Booleans is given by Bentkamp et al. [10]. Following their approach, we introduce the proxy type *bool*, which corresponds to *o*, to the signature. We define proxy symbols *t*, *f*, *not*, *and*, *or*, *impl*, *equiv*, *forall*, *exists*, *choice*, and *eq* which correspond to the homologous logical constants from Section 2. In their type declarations *o* is replaced by *bool*.

To make this paper self-contained we include the axioms from Bentkamp et al. [10]. Definitions of symbols are computational in nature: symbols are characterized by their behavior on *t* and *f*. This also reduces interferences between different axioms. Axioms are listed as follows:

$$\begin{array}{lll} t \approx f & \text{or } t x \approx t & \text{equiv } x y \approx \text{and } (\text{impl } x y) (\text{impl } y x) \\ x \approx t \vee x \approx f & \text{or } f x \approx x & \text{forall } \langle \alpha \rangle (\lambda x. t) \approx t \\ \text{not } t \approx f & \text{impl } t x \approx x & y \approx (\lambda x. t) \vee \text{forall } \langle \alpha \rangle y \approx f \\ \text{not } f \approx t & \text{impl } f x \approx t & \text{exists } \langle \alpha \rangle y \approx \\ \text{and } t x \approx x & x \approx y \vee \text{eq } \langle \alpha \rangle x y \approx t & \text{not } (\text{forall } \langle \alpha \rangle (\lambda x. \text{not } (y x))) \\ \text{and } f x \approx f & x \approx y \vee \text{eq } \langle \alpha \rangle x y \approx f & y x \approx f \vee y (\text{choice } \langle \alpha \rangle y) \approx t \end{array}$$

### 4.2. Preprocessing Booleans

Kotelnikov et al. extended VCNF, Vampire’s algorithm for clausification, to support nested Booleans [13]. Vukmirović et al. extended the clausification algorithm of Ehoh, the lambda-free higher-order version of E, to support nested Booleans inspired by VCNF extension [11, Section 8]. Zipperposition and Ehoh share the same clausification algorithm, enabling us to reuse the extension with one notable difference: unlike in Ehoh, not all nested Booleans different from variables,  $\top$  and  $\perp$  will be removed. Namely, Booleans that are below  $\lambda$ -abstraction and contain  $\lambda$ -bound variables will not be preprocessed. They cannot be easily hoisted to the level of an atom in which they appear, since this process might leak any variables bound in the context in which the nested Boolean appears. Similar preprocessing techniques are used in other higher-order provers [30].

## 5. Examples

The TPTP library contains thousands of higher-order benchmarks, many of them hand-crafted to point out subtle interferences of functional and Boolean properties of higher order logic. In this section we discuss some problems from the TPTP library that illustrate the advantages and disadvantages of our approach.

In the last five instances of the CASC theorem proving competition, the core calculus of the best performing higher-order prover was tableaux – a striking contrast from first-order part of the competition dominated by superposition-based provers. TPTP problem SET557<sup>^</sup>1 (a statement of Cantor’s theorem) might shed some light on why tableaux-based provers excel on higher-order problems. This problem conjectures that there is no surjection from a set to its power set:

$$\neg(\exists(x : \iota \rightarrow \iota \rightarrow o). \forall(y : \iota \rightarrow o). \exists(z : \iota). xz \approx y)$$

After negating the conjecture and clausification this problem becomes  $sk_1(sk_2 y) \approx y$  where  $sk_1$  and  $sk_2$  are Skolem symbols. Then, we can use ARGCONG rule [10] which applies fresh variable  $w$  to both sides of the equation, yielding clause  $C = sk_1(sk_2 y)w \approx yw$ . Most superposition- or paramodulation-based higher-order theorem provers (such as Leo-III, Vampire and Zipperposition) will split this clause into two clauses  $C_1 = sk_1(sk_2 y)w \not\approx \top \vee yw \approx \top$  and  $C_2 = sk_1(sk_2 y)w \approx \top \vee yw \not\approx \top$ . This clausification step makes the problem considerably harder. Namely, the clause  $C$  instantiated with the substitution  $\{y \mapsto \lambda x. \neg(sk_1 x x), w \mapsto sk_2(\lambda x. \neg(sk_1 x x))\}$  yields the empty clause. However, if the original clause is split into two as described above, Zipperposition will rely on PI rule to instantiate  $y$  with imitation of  $\neg$  and on equality factoring to further instantiate this approximation. These desired inferences need to be applied on both new clauses and represent only a fraction of inferences that can be done with  $C_1$  and  $C_2$ , reducing the chance of successful proof attempt. Rule BOOLER imitates the behavior of tableaux prover: it essentially rewrites the clause  $C$  into  $\neg(sk_1(sk_2 y)w) \not\approx yw$  which makes finding the necessary substitution easy and does not require a clausification step.

Combining rule (BOOL)ER with lazy clausification is very fruitful as the problem SYO033<sup>^</sup>1 illustrates. This problem also contains the single conjecture

$$\exists(x : (\iota \rightarrow o) \rightarrow o). \forall(y : \iota \rightarrow o). (xy \Leftrightarrow (\forall(z : \iota). yz))$$

The problem is easily solved if we instantiate variable  $x$  with the constant  $\forall$ . Moreover, the prover does not have to blindly guess this instantiation for  $x$ , but can obtain it by unifying  $xy$  with  $\forall y$  (which is the  $\eta$ -short form of  $\forall(z : \iota). yz$ ). However, when the problem is clausified, all quantifiers are removed. Then, Zipperposition only finds the proof if appropriate instantiation mode of PI is used, and if both clauses resulting from clausifying the negated conjecture are appropriately instantiated. In contrast, lazy clausification will derive the clause  $x(sk x) \not\approx \forall(sk x)$  from the negated conjecture in three steps. Then, equality resolution results in an empty clause, swiftly finishing the proof without any explosive inferences. This effect is even more pronounced on problems SYO287<sup>^</sup>5 and SYO288<sup>^</sup>5, in which critical proof step is instantiation of a variable with imitation of  $\vee$  and  $\wedge$ . In configurations that do not use lazy clausification and BOOLER, Zipperposition times out in any reasonable time limit; with those two options it solves mentioned problems in less than 100 ms.

In some cases, it is better to preprocess the problem. For example, TPTP problem SY0500^1.005 contains many nested Boolean terms:

$$f_0(f_1(f_1(f_1(f_2(f_3(f_3(f_3(f_4 a)))))))) \approx f_0(f_0(f_0(f_1(f_2(f_2(f_2(f_3(f_4(f_4(f_4 a))))))))))$$

In this problem, all functions  $f_i$  are of type  $o \rightarrow o$ , and constant  $a$  is of type  $o$ . FOOL unfolding of nested Boolean terms will result in exponential blowup in the problem size. However, superposition-based theorem provers are well-equipped for this issue: their CNF algorithms use smart simplifications and formula renaming to mitigate these effects. Moreover, when the problem is preprocessed, the prover is aware of the problem size before the proving process starts and can adjust its heuristics properly. E, Zipperposition and Vampire, instructed to perform FOOL unfolding, solve the problem swiftly, using their default modes. However, if problem is not preprocessed, Zipperposition struggles to prove it using CASES(SIMP) and due to the large number of (redundant) clauses it creates, succeeds only if specific heuristic choices are made.

## 6. Evaluation

We performed extensive evaluation to determine usefulness of our approach. As our benchmark set, we used all 2606 monomorphic theorems from the TPTP library, given in THF format. All of the experiments were performed on StarExec [31] servers with Intel Xeon E5-2609 0 CPUs clocked at 2.40 GHz. The evaluation is separated in two parts that answer different questions: How useful are the new rules? How does our approach compare with state-of-the-art higher-order provers?

### 6.1. Evaluation of the Rules

For this part of the evaluation, we fixed a single well-performing Zipperposition configuration called *base* ( $b$ ). Since we are testing a single configuration, we used the CPU time limit of 15 s – roughly the time a single configuration is given in a portfolio mode. Configuration  $b$  uses the pragmatic variant  $pv_{1121}^2$  of the unification algorithm given by Vukmirović et al. [21]. It enables BOOLSIMP rule, EC rule, PI rule in *Pragmatic* mode with  $k = 2$ , ELIMLEIBNIZ and ELIMPREDVAR rules, BOOLER rule, and BOOLEF rules. To evaluate the usefulness of all rules we described above, we enable, disable or change the parameters of a single rule, while keeping all other parameters of  $b$  intact. In figures that contain sufficiently different configurations, cells are of the form  $n(m)$  where  $n$  is the total number of proved problems by a particular configuration and  $m$  is the number of unique problems that a given configuration solved, compared to the other configurations in the same figure. Intersections of rows and columns denote corresponding combination of parameters. Result for the base configuration is written in *curly*; the best result is written in **bold**.

First, we tested different parameters of CASES and CASESIMP rules. In Figure 1 we report the results. The columns correspond to three possible options to choose subterm on which the inference is performed: **a** stands for any eligible subterm, **lo** and **li** stands for leftmost outermost and leftmost innermost subterms, respectively. The rows correspond to two different rules: **b** is the base configuration, which uses CASESIMP, and **b<sub>c</sub>** swaps this rule for CASES. Although the margin is slim, the results show it is usually preferable to select leftmost-outermost subterm.

	a	lo	li
b	1646	<b>1648</b>	1640
b <sub>c</sub>	1644	1645	1644

**Figure 1:** Effect of the CASES(SIMP) rule on success rate

	−PI	b <sub>p</sub>	b <sub>f</sub>	b <sub>∧</sub>	b <sub>∨</sub>	b <sub>≈</sub>	b <sub>¬</sub>	b <sub>∀∃</sub>
k = 1		<b>1648</b>	1628	1637	1634	1630	1641	1637
k = 2	1636	1646	1629	1636	1631	1627	1638	1634
k = 8		1643	1625	1633	1631	1623	1637	1635

**Figure 2:** Effect of PI rule on success rate

	−EL	+EL	−BEF	+BEF
−EPV	1584 (0)	1644 (0)	−BER	1644 (2)
+EPV	1612 (0)	<b>1646 (0)</b>	+BER	1645 (0)

**Figure 3:** Effect of Leibniz equality elimination    **Figure 4:** Effect of BOOLER and BOOLEF

Second, we evaluated all the modes of PI rule with 3 values for parameter  $k$ : 1, 2, and 8 (Figure 2). The columns denote, from left to right: disabling the PI rule, *Pragmatic* mode, *Full* mode, and *Imit<sub>∗</sub>* modes with appropriate logical symbols. The rows denote different values of  $k$ . The results show that different values for  $k$  have a modest effect on success rate. The raw data reveal that when we focus our attention to configurations with  $k = 2$ , mode *Full* can solve 10 problems no other mode (including disabling PI rule) can. Modes *Imit<sub>∧</sub>* and *Pragmatic* solve 2 problems, whereas *Imit<sub>∨</sub>* solves one problem uniquely. This result suggests that, even though this is not evident from Figure 2, sets of problems solved by different modes somewhat differ.

Figure 3 gives results of evaluating rules that treat Leibniz equality on the calculus level: EL stands for ELIMLEIBNIZ, whereas EPV denotes ELIMPREDVAR; signs  $−$  and  $+$  denote that rule is removed from or added to configuration  $b$ , respectively. Disabling both rules severely lowers the success rate. The results suggest that including ELIMLEIBNIZ is beneficial to performance.

Similarly, Figure 4 discusses merits of including (+) or excluding (−) BOOLER (BER) and BOOLEF (BEF) rules. Our expectations were that inclusion of those two rules would make bigger impact on success rate. It turned out that, in practice, most of the effects of these rules could be achieved using a combination of the PI rule and basic superposition calculus rules.

Combining these two rules with lazy clausification is more useful: when the rule EC is replaced by the rule LC, the success rate increases (compared to 1646 problems solved by  $b$ ) to 1660 problems. We also discovered that reasoning with choice is useful: when rule CHOICE is enabled, the success rate increases to 1653. We determined that including or excluding the conclusion  $D$  of CHOICE, after it is simplified, makes no difference. Counterintuitively, disabling BOOLSIMP rule results in 1640 problems, which is only 6 problems short of configuration  $b$ . Disabling EXT and INTERPRET- $\lambda$  rules results in solving 25 and 31 problems less, respectively. Raw data show

	CVC4	Leo-III	Satallax	Vampire	Zipperposition
pure	1806 (5)	1627 (0)	2067 (0)	1924 (7)	1980 ( 0)
coop	–	2085 (3)	<b>2214 (9)</b>	–	2190 (17)

**Figure 5:** Comparison with other higher-order provers

that in total, using configurations from Figure 1 to Figure 4, 1682 problems can be solved.

Last, we compare our approach to alternatives. Axiomatizing Booleans brings Zipperposition down to a grinding halt: only 1106 problems can be solved using this mode. On the other hand, preprocessing is fairly competitive: it solves only 8 problems less than the *b* configuration.

## 6.2. Comparison with Other Higher-Order Provers

We compared Zipperposition with all higher-order theorem provers that took part in THF division of CASC-27[27]: CVC4 1.8 prerelease [4], Leo-III 1.4 [14], Satallax 3.4 [15], and Vampire-THF 4.4 [6]. In this part of the evaluation, Zipperposition is ran in portfolio mode that runs configurations in different time slices. We set the CPU time limit to 180 s, the time allotted to each prover at CASC-27.

Leo-III and Satallax are cooperative theorem provers – they periodically invoke first-order provers to finish the proof attempt. Leo-III uses CVC4, E and iProver [32] as backends, while Satallax uses Ehoh [11] as backend. Zipperposition can use Ehoh as backend as well. To test effectiveness of each calculus, we run the cooperative provers in two versions: *pure*, which disables backends, and *coop* which uses all supported backends.

In both pure and cooperative mode, Satallax comes out as the winner. Zipperposition comes in close second, showing that our approach is a promising basis for further extensions. Leo-III uses SMT solver CVC4, which features native support for Booleans, as a backend. It is possible that the use of CVC4 is one of the reasons for massive improvement in success rate of cooperative configuration of Leo-III, compared with the pure version. Therefore, we conjecture that including support for SMT backends in Zipperposition might be beneficial.

## 7. Discussion and Related Work

Our work is primarily motivated by the goal of closing the gap between higher-order “hammer” or software verifier frontends and first-order backends. Considerable amount of research effort has gone into making the translations of higher-order logic as efficient as possible. Descriptions of hammers like HOLyHammer [1] and Sledgehammer [2] for Isabelle contain details of these translations. Software verifiers Boogie [33] and Why3 [34] use similar translations.

Established higher-order provers like Leo-III and Satallax perform very well on TPTP benchmarks; however, recent evaluations show that on Sledgehammer problems they are outperformed by translations to first-order logic [10, 11, 9]. Those two provers are built from the ground up as higher-order provers – treatment of exclusively higher-order issues such as extensionality or choice is built into them usually using explosive rules. Those explosive rules might contribute

to their suboptimal performance on mostly first-order Sledgehammer problems.

In contrast, our approach is to start with a first-order prover and gradually extend it with higher-order features. The work performed in the context of Matryoshka project [35], in which both authors of this paper participate, resulted in adding support for  $\lambda$ -free higher-order logic with Booleans to E [11] and veriT [9], and adding support for Boolean-free higher-order logic to Zipperposition. Authors of many state-of-the-art first-order provers have implemented some form of support for higher-order reasoning. This is true both for SMT solvers, witnessed by the recent extension of CVC4 and veriT [9], and for superposition provers, witnessed by the extension of Vampire [36]. All of those approaches were arguably more focused on functional aspects of higher-order logic, such as  $\lambda$ -binders and function extensionality, than on Boolean aspects such as Boolean subterms and Boolean extensionality. A notable exception is work by Kotelnikov et al. that introduced support for Boolean subterms to first-order Vampire [12, 13].

The main merit of our approach is that it combines two successful complementary approaches to support features of higher-order logic that have not been combined before in a modular way. It is based on a higher-order superposition calculus that incurs around 1% of overhead on first-order problems compared with classic superposition [10]. We conjecture that it is this efficient reasoning base on which the approach is based that contributes to its competitive performance.

## 8. Conclusion

We presented a pragmatic approach to support Booleans in a modern automatic prover for clausal higher-order logic. Our approach combines previous research efforts that extended first-order provers with complementary features of higher-order logic. It also proposes some solutions for the issues that emerge with this combination. The implementation shows clear improvement over previous techniques and competitive performance.

What our work misses is an overview of heuristics that can be used to curb the explosion incurred by some of the rules described in this paper. In future work, we plan to address this issue. Similarly, unlike Bentkamp et al. [10], we do not give any completeness guarantees for our extension. We plan to develop a refutationally complete calculus that supports Booleans around core rules such as CASES and LC in future work.

**Acknowledgment** We are grateful to the maintainers of StarExec for letting us use their service. We thank Alexander Bentkamp, Jasmin Blanchette and Simon Cruanes for many stimulating discussions and all the useful advice. Alexander Bentkamp suggested the rule `BOOLER`. We are also thankful to Ahmed Bhayat, Predrag Janičić and the anonymous reviewers for suggesting many improvements to this paper. We thank Evgeny Kotelnikov for patiently explaining the details of FOOL, Alexander Steen for clarifying Leo-III’s treatment of Booleans, and Giles Reger and Martin Suda for explaining the renaming mechanism implemented in VCNF. 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). Nummelin 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] C. Kaliszyk, J. Urban, Hol(y)hammer: Online ATP service for HOL light, *Mathematics in Computer Science* 9 (2015) 5–22.
- [2] L. C. Paulson, J. C. Blanchette, Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers, in: G. Sutcliffe, S. Schulz, E. Ternovska (Eds.), *IWIL-2010*, volume 2 of *EPiC*, EasyChair, 2012, pp. 1–11.
- [3] J. Filiâtre, A. Paskevich, Why3 - where programs meet provers, in: M. Felleisen, P. Gardner (Eds.), *Programming Languages and Systems - 22nd European Symposium on Programming, ESOP 2013, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2013, Rome, Italy, March 16-24, 2013. Proceedings*, volume 7792 of *LNCS*, Springer, 2013, pp. 125–128.
- [4] C. W. Barrett, C. L. Conway, M. Deters, L. Hadarean, D. Jovanovic, T. King, A. Reynolds, C. Tinelli, CVC4, in: G. Gopalakrishnan, S. Qadeer (Eds.), *CAV 2011*, volume 6806 of *LNCS*, Springer, 2011, pp. 171–177.
- [5] S. Schulz, S. Cruanes, P. Vukmirović, Faster, higher, stronger: E 2.3, in: P. Fontaine (Ed.), *CADE 2019*, volume 11716 of *LNCS*, Springer, 2019, pp. 495–507.
- [6] L. Kovács, A. Voronkov, First-order theorem proving and vampire, in: N. Sharygina, H. Veith (Eds.), *CAV 2013*, volume 8044 of *LNCS*, Springer, 2013, pp. 1–35.
- [7] J. Robinson, A note on mechanizing higher order logic, in: B. Meltzer, D. Michie (Eds.), *Machine Intelligence*, volume 5, Edinburgh University Press, 1970, pp. 121–135.
- [8] J. Meng, L. C. Paulson, Translating higher-order clauses to first-order clauses, *J. Autom. Reasoning* 40 (2008) 35–60.
- [9] H. Barbosa, A. Reynolds, D. E. Ouraoui, C. Tinelli, C. W. Barrett, Extending SMT solvers to higher-order logic, in: P. Fontaine (Ed.), *CADE 2019*, volume 11716 of *LNCS*, Springer, 2019, pp. 35–54.
- [10] A. Bentkamp, J. Blanchette, S. Tournet, P. Vukmirović, U. Waldmann, Superposition with lambdas, 2020. Submitted to a journal, [http://matryoshka.gforge.inria.fr/pubs/lamsup\\_report.pdf](http://matryoshka.gforge.inria.fr/pubs/lamsup_report.pdf).
- [11] P. Vukmirović, J. C. Blanchette, S. Cruanes, S. Schulz, Extending a brainiac prover to lambda-free higher-order logic, 2020. Submitted to a journal, [http://matryoshka.gforge.inria.fr/pubs/ehoh\\_article.pdf](http://matryoshka.gforge.inria.fr/pubs/ehoh_article.pdf).
- [12] E. Kotelnikov, L. Kovács, A. Voronkov, A first class boolean sort in first-order theorem proving and TPTP, *CoRR* abs/1505.01682 (2015). [arXiv:1505.01682](https://arxiv.org/abs/1505.01682).
- [13] E. Kotelnikov, L. Kovács, M. Suda, A. Voronkov, A clausal normal form translation for FOOL, in: C. Benzmüller, G. Sutcliffe, R. Rojas (Eds.), *GCAI 2016*, volume 41 of *EPiC Series in Computing*, EasyChair, 2016, pp. 53–71.
- [14] A. Steen, C. Benzmüller, The higher-order prover Leo-III, in: D. Galmiche, S. Schulz, R. Sebastiani (Eds.), *IJCAR 2018*, volume 10900 of *LNCS*, Springer, 2018, pp. 108–116.
- [15] C. E. Brown, Satallax: An automatic higher-order prover, in: B. Gramlich, D. Miller, U. Sattler (Eds.), *IJCAR 2012*, volume 7364 of *LNCS*, Springer, 2012, pp. 111–117.
- [16] P. B. Andrews, Classical type theory, in: J. A. Robinson, A. Voronkov (Eds.), *Handbook of Automated Reasoning*, volume II, Elsevier and MIT Press, 2001, pp. 965–1007.
- [17] G. Sutcliffe, The TPTP problem library and associated infrastructure - from CNF to TH0,

- TPTP v6.4.0, *J. Autom. Reasoning* 59 (2017) 483–502.
- [18] S. Cruanes, *Extending Superposition with Integer Arithmetic, Structural Induction, and Beyond*, Ph.D. thesis, École polytechnique, 2015.
- [19] S. Cruanes, *Superposition with structural induction*, in: C. Dixon, M. Finger (Eds.), *FroCoS 2017*, volume 10483 of *LNCS*, Springer, 2017, pp. 172–188.
- [20] I. Cervesato, F. Pfenning, *A linear spine calculus*, *J. Log. Comput.* 13 (2003) 639–688.
- [21] P. Vukmirović, A. Bentkamp, V. Nummelin, *Efficient full higher-order unification*, in: Z. M. Ariola (Ed.), *FSCD 2020*, volume 167 of *LIPICs*, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020, pp. 5:1–5:17.
- [22] L. Bachmair, H. Ganzinger, *Rewrite-based equational theorem proving with selection and simplification*, *J. Log. Comput.* 4 (1994) 217–247.
- [23] S. Schulz, *E - a brainiac theorem prover*, *AI Commun.* 15 (2002) 111–126.
- [24] A. Nonnengart, C. Weidenbach, *Computing small clause normal forms*, in: J. A. Robinson, A. Voronkov (Eds.), *Handbook of Automated Reasoning* (in 2 volumes), Elsevier and MIT Press, 2001, pp. 335–367.
- [25] A. Steen, *Extensional paramodulation for higher-order logic and its effective implementation Leo-III*, Ph.D. thesis, Free University of Berlin, Dahlem, Germany, 2018.
- [26] A. Bhayat, G. Reger, *A combinator-based superposition calculus for higher-order logic*, in: N. Peltier, V. Sofronie-Stokkermans (Eds.), *IJCAR 2020*, volume 12166 of *Lecture Notes in Computer Science*, Springer, 2020, pp. 278–296.
- [27] G. Sutcliffe, *The CADE-27 automated theorem proving system competition - CASC-27*, *AI Commun.* 32 (2019) 373–389.
- [28] H. Ganzinger, J. Stuber, *Superposition with equivalence reasoning and delayed clause normal form transformation*, *Inf. Comput.* 199 (2005) 3–23.
- [29] G. Reger, M. Suda, A. Voronkov, *New techniques in clausal form generation*, in: C. Benzmüller, G. Sutcliffe, R. Rojas (Eds.), *GCAI 2016*, volume 41 of *EPiC Series in Computing*, EasyChair, 2016, pp. 11–23.
- [30] M. Wisniewski, A. Steen, K. Kern, C. Benzmüller, *Effective normalization techniques for HOL*, in: N. Olivetti, A. Tiwari (Eds.), *IJCAR 2016*, volume 9706 of *LNCS*, Springer, 2016, pp. 362–370.
- [31] A. Stump, G. Sutcliffe, C. Tinelli, *StarExec: A cross-community infrastructure for logic solving*, in: S. Demri, D. Kapur, C. Weidenbach (Eds.), *IJCAR 2014*, volume 8562 of *LNCS*, Springer, 2014, pp. 367–373.
- [32] K. Korovin, *iprover - an instantiation-based theorem prover for first-order logic (system description)*, in: A. Armando, P. Baumgartner, G. Dowek (Eds.), *IJCAR 2008*, volume 5195 of *LNCS*, Springer, 2008, pp. 292–298.
- [33] K. R. M. Leino, P. Rümmer, *A polymorphic intermediate verification language: Design and logical encoding*, in: J. Esparza, R. Majumdar (Eds.), *TACAS 2010*, volume 6015 of *LNCS*, Springer, 2010, pp. 312–327.
- [34] F. Bobot, J.-C. Filliâtre, C. Marché, A. Paskevich, *Why3: Shepherd Your Herd of Provers*, in: *Boogie 2011: First International Workshop on Intermediate Verification Languages*, Wrocław, Poland, 2011, pp. 53–64.
- [35] J. Blanchette, P. Fontaine, S. Schulz, S. Tourret, U. Waldmann, *Stronger higher-order automation: A report on the ongoing matryoshka project*, in: M. Suda, S. Winkler (Eds.),

- EPTCS 311: Proceedings of the Second International Workshop on Automated Reasoning: Challenges, Applications, Directions, Exemplary Achievements - Natal, Brazil, August 26, 2019, Electronic Proceedings in Theoretical Computer Science, EPTCS, EPTCS, 2019, pp. 11–18.
- [36] A. Bhayat, G. Reger, Restricted combinatory unification, in: P. Fontaine (Ed.), CADE 2019, volume 11716 of *LNCS*, Springer, 2019, pp. 74–93.