

# Extending a Brainiac Prover to Lambda-Free Higher-Order Logic

Petar Vukmirović<sup>1</sup>(✉), Jasmin Christian Blanchette<sup>1,2</sup>,  
Simon Cruanes<sup>3</sup>, and Stephan Schulz<sup>4</sup>

<sup>1</sup> Vrije Universiteit Amsterdam, Amsterdam, The Netherlands  
{p.vukmirovic,j.c.blanchette}@vu.nl

<sup>2</sup> Max-Planck-Institut für Informatik, Saarland Informatics Campus,  
Saarbrücken, Germany

<sup>3</sup> Aesthetic Integration, Austin, Texas, USA

<sup>4</sup> DHBW Stuttgart, Stuttgart, Germany

**Abstract.** Decades of work have gone into developing efficient proof calculi, data structures, algorithms, and heuristics for first-order automatic theorem proving. Higher-order provers lag behind in terms of efficiency. Instead of developing a new higher-order prover from the ground up, we propose to start with a state-of-the-art superposition-based prover, E, and gradually enrich it with higher-order features. We explain how to extend the prover’s data structures, algorithms, and heuristics to  $\lambda$ -free higher-order logic, a formalism that supports partial application and applied variables. Our extension outperforms the traditional encoding and appears promising as a stepping stone towards full higher-order logic.

## 1 Introduction

Superposition provers, such as E [27], SPASS [34], and Vampire [18], are among the most successful first-order reasoning systems today. They serve as backends in various frameworks, ranging from software verifiers (e.g., Why3 [15]), to automatic higher-order theorem provers (e.g., Leo-III [28], Satalax [13]), to “hammers” in proof assistants (e.g., HOL Light [17], Isabelle [22]). Decades of research have gone into refining the proof calculus, devising efficient data structures and algorithms, and conceiving heuristics to guide proof search. This work has mostly focused on first-order logic with equality, with or without arithmetic.

Research on higher-order automatic provers has resulted in systems such as LEO [8], LEO-II [9], and Leo-III [28], based on paramodulation, and Satalax [13], based on analytic tableaux and SAT solving. These provers feature a “cooperative” architecture, pioneered by LEO: They are full-fledged higher-order provers that regularly invoke an external first-order prover with a low time limit as a terminal procedure, in an attempt to finish the proof quickly using only first-order reasoning. However, the first-order backend will succeed only if all the necessary higher-order reasoning has been performed, meaning that much of the first-order reasoning is carried out by the slower higher-order prover. As a result, this architecture leads to suboptimal performance on first-order problems and on problems with a large first-order component. For example, at the 2017

edition of the CADE ATP System Competition (CASC) [31], Leo-III proved 652 out of 2000 first-order problems originating from Isabelle, compared with 1433 for Vampire and 1185 for E.

To improve the situation, we propose to start from an optimized first-order prover and extend it to full higher-order logic one feature at a time. Our goal is to extend the prover in a *graceful* manner, so that it behaves as before on first-order problems, performs mostly like a first-order prover on typical, mildly higher-order problems, and scales up to arbitrary higher-order problems, in keeping with the zero-overhead principle: *What you don't use, you don't pay for*.

As a stepping stone towards full higher-order logic, we initially restrict our focus to a higher-order logic without  $\lambda$ -expressions (Sect. 2). Compared with first-order logic, its distinguishing features are partial application and applied variables. This formalism is rich enough to express the recursive equations of higher-order combinators, such as the `map` operation on finite lists:

$$\text{map } f \text{ nil} \approx \text{nil} \qquad \text{map } f \text{ (cons } x \text{ xs)} \approx \text{cons } (f \ x) \text{ (map } f \ \text{xs)}$$

Our vehicle is E, an open source prover developed primarily by Schulz. It is written in C and offers respectable performance, with the emphasis on “brainiac” heuristics rather than raw speed. It regularly secures second places at CASC and serves as a backend to competitive higher-order provers. We refer to our extended version of E as Ehoh. The source code is publicly available.<sup>1</sup> An earlier version of Ehoh, mischievously named hoE, is described in Vukmirović’s MSc thesis [32].

The three main challenges we faced were generalizing the type and term representation (Sect. 3), extending the core algorithms (Sect. 4), and generalizing the indexing data structures (Sect. 5). Some effort also went into adapting the inference rules (Sect. 6) and the heuristics (Sect. 7). We explain the key ideas; further details, including correctness proofs, are given in our technical report [33].

An innovative aspect of our work is a set of techniques we call *prefix optimization*. Higher-order terms contain twice as many proper subterms as first-order terms; for example, the term `f (g a) b` contains not only the argument subterms `g a`, `a`, and `b` but also the “prefix” subterms `f`, `f (g a)`, and `g`. Many operations require traversing all subterms of a term. With the prefix optimization, the prover traverses the subterms recursively in a first-order fashion, considering all the prefixes of the current subterm together, at no additional cost. Our experiments (Sect. 8) show that Ehoh is virtually as fast as E on first-order problems and can also prove higher-order problems that do not require synthesizing  $\lambda$ -terms. In future work, we plan to integrate Ehoh into the official E and consider  $\lambda$ -terms.

An alternative to using Ehoh consists of applying the *applicative encoding*: Every  $n$ -ary function symbol is converted to a nullary symbol, and application is represented by a distinguished binary symbol `@`. For example, the higher-order term `f (x a) b` can be encoded as the first-order term `@(@(f, @(x, a)), b)`. However, this representation is not graceful; it clutters data structures and interferes with proof search in subtle ways, leading to poorer performance. For these and further reasons (Sect. 9), it is not an ideal basis for higher-order reasoning.

<sup>1</sup> <https://bitbucket.org/petarvukmirovic/ehoh>

## 2 Logic

Our logic corresponds to the intensional  $\lambda$ -free higher-order logic ( $\lambda$ fHOL) described by Bentkamp, Blanchette, Cruanes, and Waldmann [7, Sect. 2]. Another possible name for this logic would be “applicative first-order logic.” Extensionality can be obtained by adding suitable axioms [7, Sect. 3.1].

A type is either an atomic type  $\iota$  or a function type  $\tau \rightarrow v$ , where  $\tau$  and  $v$  are themselves types. Terms, ranged over by  $s, t, u, v$ , are either variables  $x, y, z, \dots$ , symbols  $a, b, c, d, f, g, \dots$ , or binary applications of the form  $s t$ . The function arrow associates to the right, whereas application associates to the left. The typing rules are as for the simply typed  $\lambda$ -calculus. A term’s *arity* is the number of extra arguments it can take; thus, if  $f$  has type  $\iota \rightarrow \iota \rightarrow \iota$  and  $a$  has type  $\iota$ , then  $f$  has an arity of 2,  $f a$  has an arity of 1, and  $f a a$  has an arity of 0.

Following the *flattened* notation, terms have a unique decomposition of the form  $\zeta s_1 \dots s_m$ , where  $\zeta$ , the *head*, is a variable or symbol and  $s_1, \dots, s_m$ , the arguments, are arbitrary terms. We abbreviate  $m$ -tuples  $(a_1, \dots, a_m)$  to  $\overline{a_m}$  or  $\bar{a}$ ; abusing notation, we write  $\zeta \overline{s_m}$  for the curried application  $\zeta s_1 \dots s_m$ .

An equation  $s \approx t$  corresponds to an unordered pair of terms. A literal  $L$  is either an equation or its negation. Clauses  $C, D$  are finite multisets of literals, interpreted disjunctively:  $L_1 \vee \dots \vee L_n$ .

## 3 Types and Terms

The term representation is a fundamental question when building a theorem prover. Delicate changes to E’s term representation were needed to support partial application and especially applied variables. In contrast, the introduction of a higher-order type system had a less dramatic impact on the prover’s code.

**Types.** For most of its history, E supported only untyped first-order logic. Cruanes implemented support for atomic types for E 2.0. Symbols  $f$  are declared with a type signature:  $f : \tau_1 \times \dots \times \tau_m \rightarrow \tau$ . Atomic types are represented by integers in memory, leading to efficient type comparisons.

In  $\lambda$ fHOL, a type signature consists of types  $\tau$ , in which the function type constructor  $\rightarrow$  can be nested arbitrarily—e.g.,  $(\iota \rightarrow \iota) \rightarrow \iota \rightarrow \iota$ . A natural way to represent such types is to mimic their recursive structures using tagged unions. However, this leads to memory fragmentation, and a simple operation such as querying the type of a function’s  $i$ th argument would require traversing  $i$  nodes in memory. We prefer a flattened representation, in which a type  $\tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow \iota$  is represented by a single node labeled with  $\rightarrow$  and pointing to the array  $(\tau_1, \dots, \tau_n, \iota)$ . Applying  $k \leq n$  arguments to a function of the above type yields a term of type  $\tau_{k+1} \rightarrow \dots \rightarrow \tau_n \rightarrow \iota$ . In memory, this corresponds to skipping the first  $k$  array elements.

To speed up type comparisons, Ehoh stores all types in a shared bank and implements perfect sharing, ensuring that types that are structurally the same are represented by the same object in memory. Type equality can then be implemented as a pointer comparison.

**Terms.** In E, terms are represented as perfectly shared directed acyclic graphs. Each node, or *cell*, contains 11 fields, including `f_code`, an integer that identifies the term’s head symbol (if  $\geq 0$ ) or variable (if  $< 0$ ); `arity`, an integer corresponding to the number of arguments passed to the head symbol; `args`, an array of size `arity` consisting of pointers to argument terms; and `binding`, which possibly stores a substitution for a variable used for unification and matching.

In higher-order logic, variables may have function type and be applied, and symbols can be applied to fewer arguments than specified by their type signatures. A natural representation of  $\lambda$ HOL terms as tagged unions would distinguish between variables  $x$ , symbols  $f$ , and binary applications  $s\ t$ . However, this scheme would lead to memory fragmentation and linear-time access to the  $i$ th argument of a function, affecting performance on purely or mostly first-order problems. Instead, we propose a flattened representation, as a generalization of E’s existing data structures: Allow arguments to variables, and for symbols let `arity` be the number of *actual* arguments.

A side effect of the flattened representation is that prefix subterms are not shared. For example, the terms  $f\ a$  and  $f\ a\ b$  correspond to the flattened cells  $f(a)$  and  $f(a, b)$ . The argument subterm  $a$  is shared, but not the prefix  $f\ a$ . Similarly,  $x$  and  $x\ b$  are represented by two distinct cells,  $x()$  and  $x(b)$ , and there is no connection between the two  $x$ ’s. In particular, despite perfect sharing, their `binding` fields are unconnected, leading to inconsistencies.

A potential solution would be to systematically traverse a clause and set the `binding` fields of all cells of the form  $x(\bar{s})$  whenever a variable  $x$  is bound, but this would be inefficient and inelegant. Instead, we implemented a hybrid approach: Variables are applied by an explicit application operator `@`, to ensure that they are always perfectly shared. Thus,  $x\ b\ c$  is represented by the cell  $@(x, b, c)$ , where  $x$  is a shared subcell. This is perfectly graceful, since variables never occur applied in first-order terms. The main drawback of this technique is that some normalization is necessary after substitution: Whenever a variable is instantiated by a term with a symbol head, the `@` symbol must be eliminated. Applying the substitution  $\{x \mapsto f\ a\}$  to the cell  $@(x, b, c)$  must produce the cell  $f(a, b, c)$  and not  $@(f(a), b, c)$ , for consistency with other occurrences of  $f\ a\ b\ c$ .

There is one more complication related to the `binding` field. In E, it is easy and useful to traverse a term as if a substitution has been applied, by following all set `binding` fields. In Ehoh, this is not enough, because cells must also be normalized. To avoid repeatedly creating the same normalized cells, we introduced a `binding_cache` field that connects a  $@(x, \bar{s})$  cell with its substitution. However, this cache can easily become stale when the `binding` pointer is updated. To detect this situation, we store  $x$ ’s `binding` value in the  $@(x, \bar{s})$  cell’s (otherwise unused) `binding` field. To find out whether the cache is valid, it suffices to check that the `binding` fields of  $x$  and  $@(x, \bar{s})$  are equal.

**Term Orders.** Superposition provers rely on a term order to prune the search space. To ensure completeness, the order must be a simplification order that is total (linear) for variable-free terms. The Knuth–Bendix order (KBO) and the lexicographic path order (LPO) meet this criterion. KBO is generally consid-

ered the most robust and efficient option for superposition. E implements both. In earlier work, Blanchette and colleagues have shown that only KBO can be generalized gracefully while preserving all the necessary properties for superposition [5, 11]. For this reason, we focus on KBO.

E implements the linear-time algorithm for KBO described by Löchner [19], which relies on the tupling method to store intermediate results, avoiding repeated computations. Once we understood the algorithm, it was easy to gracefully generalize it to implement the  $\lambda$ FHOL version of KBO [5]. The main difference is that when comparing two terms  $f \overline{s_m}$  and  $f \overline{t_n}$ , because of partial application we may now have  $m \neq n$ ; this required changing the implementation to perform a length-lexicographic comparison of the tuples  $\overline{s_m}$  and  $\overline{t_n}$ .

## 4 Unification and Matching

Syntactic unification of  $\lambda$ FHOL terms has a definite first-order flavor. It is decidable, and most general unifiers (MGUs) are unique up to variable renaming. For example, the unification constraint  $f(y \ a) \stackrel{?}{=} f(a)$  has the MGU  $\{y \mapsto f\}$ , whereas in full higher-order logic it would admit infinitely many independent solutions of the form  $\{y \mapsto \lambda x. f(f(\dots(f x)\dots))\}$ . Matching is a special case of unification where only the variables on the left-hand side can be instantiated.

An easy but inefficient way to implement unification and matching for  $\lambda$ FHOL is to apply the applicative encoding (Sect. 1), perform first-order unification or matching, and decode the result. Instead, we propose to generalize the first-order unification and matching procedures to operate directly on  $\lambda$ FHOL terms.

We present our unification procedure as a transition system, generalizing Baader and Nipkow [3]. A unification problem consists of a finite set  $S$  of unification constraints  $s_i \stackrel{?}{=} t_i$ , where  $s_i$  and  $t_i$  are of the same type. A problem is in *solved form* if it has the form  $\{x_1 \stackrel{?}{=} t_1, \dots, x_n \stackrel{?}{=} t_n\}$ , where the  $x_i$ 's are distinct and do not occur in the  $t_j$ 's. The corresponding unifier is  $\{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$ . The transition rules attempt to bring the input constraints in solved form.

The first group of rules consists of operations that focus on a single constraint and replace it with a new (possibly empty) set of constraints:

$$\begin{array}{ll}
\text{Delete} & \{t \stackrel{?}{=} t\} \uplus S \Longrightarrow S \\
\text{Decompose} & \{f \overline{s_m} \stackrel{?}{=} f \overline{t_m}\} \uplus S \Longrightarrow S \cup \{s_1 \stackrel{?}{=} t_1, \dots, s_m \stackrel{?}{=} t_m\} \\
\text{DecomposeX} & \{x \overline{s_m} \stackrel{?}{=} u \overline{t_m}\} \uplus S \Longrightarrow S \cup \{x \stackrel{?}{=} u, s_1 \stackrel{?}{=} t_1, \dots, s_m \stackrel{?}{=} t_m\} \\
& \text{if } x \text{ and } u \text{ have the same type and } m > 0 \\
\text{Orient} & \{f \overline{s} \stackrel{?}{=} x \overline{t}\} \uplus S \Longrightarrow S \cup \{x \overline{t} \stackrel{?}{=} f \overline{s}\} \\
\text{OrientXY} & \{x \overline{s_m} \stackrel{?}{=} y \overline{t_n}\} \uplus S \Longrightarrow S \cup \{y \overline{t_n} \stackrel{?}{=} x \overline{s_m}\} \quad \text{if } m > n \\
\text{Eliminate} & \{x \stackrel{?}{=} t\} \uplus S \Longrightarrow \{x \stackrel{?}{=} t\} \cup \{x \mapsto t\}(S) \quad \text{if } x \in \mathcal{V}ar(S) \setminus \mathcal{V}ar(t)
\end{array}$$

The Delete, Decompose, and Eliminate rules are essentially as for first-order terms. The Orient rule is generalized to allow applied variables and complemented by a new OrientXY rule. DecomposeX, also a new rule, can be seen as a variant of Decompose that analyzes applied variables; the term  $u$  may be an application.

The rules belonging to the second group detect unsolvable constraints:

$$\begin{array}{ll}
\text{Clash} & \{f \bar{s} \stackrel{?}{=} g \bar{t}\} \uplus S \Longrightarrow \perp \quad \text{if } f \neq g \\
\text{ClashTypeX} & \{x \bar{s}_m \stackrel{?}{=} u \bar{t}_m\} \uplus S \Longrightarrow \perp \quad \text{if } x \text{ and } u \text{ have different types} \\
\text{ClashLenXF} & \{x \bar{s}_m \stackrel{?}{=} f \bar{t}_n\} \uplus S \Longrightarrow \perp \quad \text{if } m > n \\
\text{OccursCheck} & \{x \stackrel{?}{=} t\} \uplus S \Longrightarrow \perp \quad \text{if } x \in \mathcal{Var}(t) \text{ and } x \neq t
\end{array}$$

The derivations below demonstrate the computation of MGUs for the unification problems  $\{f(y \ a) \stackrel{?}{=} y(f \ a)\}$  and  $\{x(z \ b \ c) \stackrel{?}{=} g \ a \ (y \ c)\}$ :

$$\begin{array}{llll}
& \{f(y \ a) \stackrel{?}{=} y(f \ a)\} & & \{x(z \ b \ c) \stackrel{?}{=} g \ a \ (y \ c)\} \\
\Longrightarrow_{\text{Orient}} & \{y(f \ a) \stackrel{?}{=} f(y \ a)\} & \Longrightarrow_{\text{DecomposeX}} & \{x \stackrel{?}{=} g \ a, z \ b \ c \stackrel{?}{=} y \ c\} \\
\Longrightarrow_{\text{DecomposeX}} & \{y \stackrel{?}{=} f, f \ a \stackrel{?}{=} y \ a\} & \Longrightarrow_{\text{OrientXY}} & \{x \stackrel{?}{=} g \ a, y \ c \stackrel{?}{=} z \ b \ c\} \\
\Longrightarrow_{\text{Eliminate}} & \{y \stackrel{?}{=} f, f \ a \stackrel{?}{=} f \ a\} & \Longrightarrow_{\text{DecomposeX}} & \{x \stackrel{?}{=} g \ a, y \stackrel{?}{=} z \ b, c \stackrel{?}{=} c\} \\
\Longrightarrow_{\text{Delete}} & \{y \stackrel{?}{=} f\} & \Longrightarrow_{\text{Delete}} & \{x \stackrel{?}{=} g \ a, y \stackrel{?}{=} z \ b\}
\end{array}$$

The implementation stores open constraints on a stack. It focuses on the topmost constraint, detecting unsolvable cases early, following the above rules. New constraints that emerge are put on top of the stack. Correctness proofs as well as the pseudocode for unification and matching algorithms based on these ideas are included in our technical report [33].

During proof search, E repeatedly needs to test for unifiability of a term  $s$  not only with some other term  $t$  but also with  $t$ 's subterms. The prefix optimization speeds up this test: The subterms of  $t$  are traversed in a first-order fashion; for each such subterm  $\zeta \bar{t}_n$ , at most one prefix  $\zeta \bar{t}_k$ , with  $k \leq n$ , is possibly unifiable with  $s$ , by virtue of their having the same arity. Using this technique, Ehoh exhibits the same time complexity as E on first-order terms.

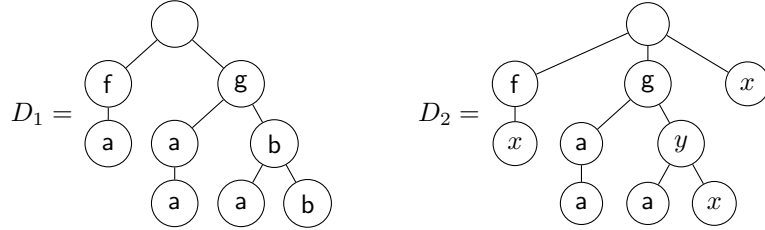
## 5 Indexing Data Structures

Superposition provers like E work by saturation. Their main loop heuristically selects a clause and searches for potential inference partners among a possibly large set of already derived clauses. Mechanisms such as simplification and subsumption also require locating terms in a large clause set. For example, when E derives a new equation  $s \approx t$ , if  $s$  is larger than  $t$  according to the term order, it will rewrite all instances  $\sigma(s)$  of  $s$  to  $\sigma(t)$  in existing clauses.

To avoid iterating over all terms (including subterms) in large clause sets, superposition provers store the potential inference partners in indexing data structures. A term index stores a set of terms  $S$ . Given a *query term*  $t$ , a query returns all terms  $s \in S$  that satisfy a given *retrieval condition*:  $\sigma(s) = \sigma(t)$  ( $s$  and  $t$  are unifiable),  $\sigma(s) = t$  ( $s$  generalizes  $t$ ), or  $s = \sigma(t)$  ( $s$  is an instance of  $t$ ), for some substitution  $\sigma$ . Perfect indices return the subset of terms satisfying the retrieval condition. In contrast, imperfect indices overapproximate the result; a second pass is needed to filter out the terms violating the retrieval condition.

E relies on two term indexing data structures, which needed to be generalized to  $\lambda$ fHOL: perfect discrimination trees [21] and fingerprint indices [25]. It also uses feature vector indices [26] for clauses, but these required no changes.

**Perfect Discrimination Trees.** Discrimination trees [21] are tries in which every node is labeled with a symbol or a variable. A path from the root to a leaf node corresponds to a “serialized term”—a term expressed without parentheses and commas. Consider the following discrimination trees:



Assuming  $a, b, x, y : \iota$ ,  $f : \iota \rightarrow \iota$ , and  $g : \iota^2 \rightarrow \iota$ , the trees  $D_1$  and  $D_2$  represent the term sets  $\{f(a), g(a, a), g(b, a), g(b, b)\}$  and  $\{f(x), g(a, a), g(y, a), g(y, x), x\}$ .

E uses perfect discrimination trees for finding generalizations of query terms. For example, if the query term is  $g(a, a)$ , it would follow the path  $g.a.a$  in the tree  $D_1$  and return  $\{g(a, a)\}$ . For  $D_2$ , it would also explore paths labeled with variables, binding them as it proceeds, and return  $\{g(a, a), g(y, a), g(y, x), x\}$ .

The data structure relies on the observation that serializing is unambiguous. Conveniently, this property also holds for  $\lambda$ HOL terms. Assume that two distinct  $\lambda$ HOL terms yield the same serialization. Clearly, they must disagree on parentheses; one will have the subterm  $s t u$  where the other has  $s (t u)$ . However, these two subterms cannot both be well typed.

When generalizing the data structure to  $\lambda$ HOL, we face a slight complication due to partial application. First-order terms can only be stored in leaf nodes, but in E<sub>hoh</sub> we must also be able to represent partially applied terms, such as  $f$ ,  $g$ , or  $g a$  (assuming, as above, that  $f$  is unary and  $g$  is binary). Conceptually, this can be solved by storing a Boolean on each node indicating whether it is an accepting state. In the implementation, the change is more subtle, because several parts of E’s code implicitly assume that only leaf nodes are accepting.

The main difficulty specific to  $\lambda$ HOL concerns applied variables. To enumerate all generalizing terms, E needs to backtrack from child to parent nodes. To achieve this, it relies on two stacks that store subterms of the query term: `term_stack` stores the terms that must be matched in turn against the current subtree; and `term_proc` stores, for each node from the root to the current subtree, the corresponding processed term, including any arguments yet to be matched.

The matching procedure starts at the root with an empty substitution  $\sigma$ . Initially, `term_stack` contains the query term, and `term_proc` is empty. The procedure advances by moving to a suitable child node:

- A. If the node is labeled with a symbol  $f$  and the top item  $t$  of `term_stack` is  $f(\overline{t_n})$ , replace  $t$  by  $n$  new items  $t_1, \dots, t_n$ , and push  $t$  onto `term_proc`.
- B. If the node is labeled with a variable  $x$ , there are two subcases. If  $x$  is already bound, check that  $\sigma(x) = t$ ; otherwise, extend  $\sigma$  so that  $\sigma(x) = t$ . Next, pop a term  $t$  from `term_stack` and push it onto `term_proc`.

The goal is to reach an accepting node. If the query term and all the terms

stored in the tree are first-order, `term_stack` will then be empty, and the entire query term will have been matched.

Backtracking works in reverse: Pop a term  $t$  from `term_proc`; if the current node is labeled with an  $n$ -ary symbol, discard `term_stack`'s topmost  $n$  items; finally, push  $t$  onto `term_stack`. Variable bindings must also be undone.

As an example, looking up  $g(b, a)$  in the tree  $D_1$  would result in the following succession of stack states, starting from the root  $\epsilon$  along the path  $g.b.a$ :

	$\epsilon$	$g$	$g.b$	$g.b.a$
<code>term_stack:</code>	[ $g(b, a)$ ]	[ $b, a$ ]	[ $a$ ]	[]
<code>term_proc:</code>	[]	[ $g(b, a)$ ]	[ $b, g(b, a)$ ]	[ $a, b, g(b, a)$ ]

(The notation  $[a_1, \dots, a_n]$  represents the  $n$ -item stack with  $a_1$  as the top item.) Backtracking amounts to moving leftwards: When backtracking from the node  $g$  to the root, we pop  $g(b, a)$  from `term_proc`, we discard two items from `term_stack`, and we push  $g(b, a)$  onto `term_stack`.

To adapt the procedure to  $\lambda$ HOL, the key idea is that an applied variable is not very different from an applied symbol. A node labeled with an  $n$ -ary symbol or variable  $\zeta$  matches a prefix  $t'$  of the  $k$ -ary term  $t$  popped from `term_stack` and leaves  $n - k$  arguments  $\bar{u}$  to be pushed back, with  $t = t' \bar{u}$ . If  $\zeta$  is a variable, it must be bound to the prefix  $t'$ . Backtracking works analogously: Given the arity  $n$  of the node label  $\zeta$  and the arity  $k$  of the term  $t$  popped from `term_proc`, we discard the topmost  $n - k$  items  $\bar{u}$  from `term_proc`.

To illustrate the procedure, we consider the tree  $D_2$  but change  $y$ 's type to  $\iota \rightarrow \iota$ . This tree represents the set  $\{f\ x, g\ a\ a, g\ (y\ a), g\ (y\ x), x\}$ . Let  $g\ (g\ a\ b)$  be the query term. We have the following sequence of substitutions and stacks:

	$\epsilon$	$g$	$g.y$	$g.y.x$
$\sigma:$	$\emptyset$	$\emptyset$	$\{y \mapsto g\ a\}$	$\{y \mapsto g\ a, x \mapsto b\}$
<code>term_stack:</code>	[ $g\ (g\ a\ b)$ ]	[ $g\ a\ b$ ]	[ $b$ ]	[]
<code>term_proc:</code>	[]	[ $g\ (g\ a\ b)$ ]	[ $g\ a\ b, g\ (g\ a\ b)$ ]	[ $b, g\ a\ b, g\ (g\ a\ b)$ ]

Finally, to avoid traversing twice as many subterms as in the first-order case, we can implement the prefix optimization: Given a query term  $\zeta\ \bar{t}_n$ , we can also match prefixes  $\zeta\ \bar{t}_k$  by allowing `term_stack` to be nonempty at the end.

**Fingerprint Indices.** Fingerprint indices [25] trade perfect indexing for a compact memory representation and more flexible retrieval conditions. The basic idea is to compare terms by looking only at a few predefined sample positions. If we know that term  $s$  has symbol  $f$  at the head of the subterm at 2.1 and term  $t$  has  $g$  at the same position, we can directly conclude that  $s$  and  $t$  are not unifiable.

Let **A** (“at a variable”), **B** (“below a variable”), and **N** (“nonexistent”) be distinguished symbols. Given a term  $t$  and a position  $p$ , the *fingerprint function*  $\mathcal{G}fpf$  is defined as

$$\mathcal{G}fpf(t, p) = \begin{cases} \mathbf{f} & \text{if } t|_p \text{ has a symbol head } f \\ \mathbf{A} & \text{if } t|_p \text{ is a variable} \\ \mathbf{B} & \text{if } t|_q \text{ is a variable for some proper prefix } q \text{ of } p \\ \mathbf{N} & \text{otherwise} \end{cases}$$



Based on a fixed tuple of sample positions  $\overline{p_n}$ , the *fingerprint* of a term  $t$  is defined as  $\mathcal{F}p(t) = (\mathcal{G}fpf(t, p_1), \dots, \mathcal{G}fpf(t, p_n))$ . To compare two terms  $s$  and  $t$ , it suffices to check that their fingerprints are componentwise compatible using the following unification and matching matrices:

	f <sub>1</sub>	f <sub>2</sub>	A	B	N
f <sub>1</sub>		<b>X</b>			<b>X</b>
A					<b>X</b>
B					
N	<b>X</b>	<b>X</b>	<b>X</b>		

	f <sub>1</sub>	f <sub>2</sub>	A	B	N
f <sub>1</sub>		<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>
A				<b>X</b>	<b>X</b>
B					
N	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	

The rows and columns correspond to  $s$  and  $t$ , respectively. The metavariables  $f_1$  and  $f_2$  represent arbitrary distinct symbols. Incompatibility is indicated by **X**.

As an example, let  $(\epsilon, 1, 2, 1.1, 1.2, 2.1, 2.2)$  be the sample positions, and let  $s = f(a, x)$  and  $t = f(g(x), g(a))$  be the terms to unify. Their fingerprints are

$$\mathcal{F}p(s) = (f, a, A, N, N, B, B) \quad \mathcal{F}p(t) = (f, g, g, A, N, a, N)$$

Using the left matrix, we compute the compatibility vector  $(-, \mathbf{X}, -, \mathbf{X}, -, -, -)$ . The mismatches at positions 1 and 1.1 indicate that  $s$  and  $t$  are not unifiable.

A fingerprint index is a trie that stores a term set  $T$  keyed by fingerprint. The term  $f(g(x), g(a))$  above would be stored in the node addressed by  $f.g.g.A.N.a.N$ , possibly together with other terms that share the same fingerprint. This organization makes it possible to unify or match a query term  $s$  against all the terms  $T$  in one traversal. Once a node storing the terms  $U \subseteq T$  has been reached, due to overapproximation we must apply unification or matching on  $s$  and each  $u \in U$ .

When adapting this data structure to  $\lambda$ fHOL, we must first choose a suitable notion of position in a term. Conventionally, higher-order positions are strings over  $\{1, 2\}$  indicating, for each binary application  $t_1 t_2$ , which term  $t_i$  to follow. Given that this is not graceful, it seems preferable to generalize the first-order notion to flattened  $\lambda$ fHOL terms—e.g.,  $x a b|_1 = a$  and  $x a b|_2 = b$ . However, this approach fails on applied variables. For example, although  $x b$  and  $f a b$  are unifiable (using  $\{x \mapsto f a\}$ ), sampling position 1 would yield a clash between  $b$  and  $a$ . To ensure that positions remain stable under substitution, we propose instead to number arguments in reverse, from right to left:  $t|^\epsilon = t$  and  $\zeta t_n \dots t_1|^{i.p} = t_i|^{i.p}$  if  $1 \leq i \leq n$ . The operation is undefined for out-of-bound indices.

Let  $t|^{i.p}$  denote the subterm  $t|^{i.p}$  such that  $i.p$  is the longest prefix of  $p$  for which  $t|^{i.p}$  is defined. The  $\lambda$ fHOL version of the fingerprint function is defined as follows:

$$\mathcal{G}fpf'(t, p) = \begin{cases} f & \text{if } t|^{i.p} \text{ has a symbol head } f \\ A & \text{if } t|^{i.p} \text{ has a variable head} \\ B & \text{if } t|^{i.p} \text{ is undefined but } t|^{i.p} \text{ has a variable head} \\ N & \text{otherwise} \end{cases}$$

Except for the reversed numbering scheme,  $\mathcal{G}fpf'$  coincides with  $\mathcal{G}fpf$  on first-order terms. The fingerprint  $\mathcal{F}p'(t)$  of a term  $t$  is defined analogously as before, and the same compatibility matrices can be used.

The most interesting new case is that of an applied variable. Given the sample positions  $(\epsilon, 2, 1)$ , the fingerprint of  $x$  is  $(A, B, B)$  as before, whereas the finger-

print of  $x$   $c$  is  $(A, B, c)$ . As another example, let  $(\epsilon, 2, 1, 2.2, 2.1, 1.2, 1.1)$  be the sample positions, and let  $s = x (f b c)$  and  $t = g a (y d)$ . Their fingerprints are

$$\mathcal{Fp}(s) = (A, B, f, B, B, b, c) \quad \mathcal{Fp}(t) = (g, a, A, N, N, B, d)$$

The terms are not unifiable due to the incompatibility at position 1.1 ( $c$  versus  $d$ ).

We can easily support the prefix optimization for both terms  $s$  and  $t$  being compared: We ensure that  $s$  and  $t$  are fully applied, by adding enough fresh variables as arguments, before computing their fingerprints.

## 6 Inference Rules

E's main loop applies inference rules between a heuristically selected *given clause* and already derived clauses. There are two kinds of rules: The *generating rules* are necessary for completeness, whereas the *simplification rules* replace existing clauses by simpler clauses as an optimization. The main loop is where the algorithms and data structures described in Sects. 4 and 5 come into play.

Ehoh implements essentially the same logical calculus as E, except that it is generalized to  $\lambda$ HOL terms. The standard inference rules and completeness proof of superposition can be reused verbatim; the only changes concern the basic definitions of terms and substitutions [7, Sect. 1].

**The Generating Rules.** The superposition calculus consists of the following four core generating rules, whose conclusions are added to the proof state:

$$\begin{array}{c} \frac{s \not\approx s' \vee C}{\sigma(C)} \text{ER} \\ \frac{s \approx t \vee C \quad u[s'] \not\approx v \vee D}{\sigma(u[t] \not\approx v \vee C \vee D)} \text{SN} \end{array} \quad \begin{array}{c} \frac{s \approx t \vee s' \approx u \vee C}{\sigma(t \not\approx u \vee s \approx u \vee C)} \text{EF} \\ \frac{s \approx t \vee C \quad u[s'] \approx v \vee D}{\sigma(u[t] \approx v \vee C \vee D)} \text{SP} \end{array}$$

In each rule,  $\sigma$  denotes the MGU of  $s$  and  $s'$ . Not shown are side conditions that restrict the rules' applicability.

Equality resolution and factoring (ER and EF) work on entire terms that occur on either side of a literal occurring in the given clause. To generalize them, it suffices to disable the prefix optimization for our unification algorithm. By contrast, superposition into negative and positive literals (SN and SP) are more complex. As two-premise rules, they require the prover to find a partner for the given clause. There are two cases to consider, depending on whether the given clause acts as the first or second premise in an inference.

To cover the case where the given clause acts as the left premise, the prover relies on a fingerprint index to compute a set of clauses containing terms possibly unifiable with a side  $s$  of a positive literal of the given clause. Thanks to our generalization of fingerprints, in Ehoh this candidate set is guaranteed to over-approximate the set of all possible inference partners. The unification algorithm is then applied to filter out unsuitable candidates. Thanks to the prefix optimization, we can avoid gracelessly polluting the index with all prefix subterms.

For the case where the given clause is the right premise, the prover traverses its subterms  $s'$  looking for inference partners in another fingerprint index, which contains only entire left- and right-hand sides of equalities. Like E, Ehoh traverses subterms in a first-order fashion. If prefix unification succeeds, Ehoh determines the unified prefix and applies the appropriate inference instance.

**The Simplifying Rules.** Unlike the generating rules, the simplifying rules do not only add conclusions to the proof state—they also remove the premises. E has over a dozen simplifying rules, but it will suffice to consider a single example:

$$\frac{s \approx t \quad u[\sigma(s)] \approx u[\sigma(t)] \vee C}{s \approx t} \text{ES}$$

Given an equation  $s \approx t$ , equality subsumption (ES) removes a clause containing a literal whose two sides are equal except that an instance of  $s$  appears on one side where the corresponding instance of  $t$  appears on the other side.

E maintains a perfect discrimination tree that stores clauses of the form  $s \approx t$  indexed by  $s$  and  $t$ . When applying the ES rule, E considers each literal  $u \approx v$  of the given clause in turn. It starts by taking the left-hand side  $u$  as a query term. If an equation  $s \approx t$  (or  $t \approx s$ ) is found in the tree, with  $\sigma(s) = u$ , the prover checks whether  $\sigma'(t) = v$  for some extension  $\sigma'$  of  $\sigma$ . If so, ES is applicable. To consider nonempty contexts, the prover traverses the subterms  $u'$  and  $v'$  of  $u$  and  $v$  in lockstep, as long as they appear under identical contexts. Thanks to the prefix optimization, when Ehoh is given a subterm  $u'$ , it can find an equation  $s \approx t$  in the tree such that  $\sigma(s)$  is equal to some prefix of  $u'$ , with  $n$  arguments  $\bar{u}_n$  remaining as unmatched. Checking for equality subsumption then amounts to checking that  $v' = \sigma'(t) \bar{u}_n$ , for some extension  $\sigma'$  of  $\sigma$ .

For example, let  $f(g\ a\ b) \approx f(h\ g\ b)$  be the given clause, and suppose that  $x\ a \approx h\ x$  is indexed. Under context  $f[\ ]$ , Ehoh considers the subterms  $g\ a\ b$  and  $h\ x\ b$ . It finds the prefix  $g\ a$  of  $g\ a\ b$  in the tree, with  $\sigma = \{x \mapsto g\}$ . The prefix  $h\ g$  of  $h\ g\ b$  matches the indexed equation's right-hand side  $h\ x$  using the same substitution, and the remaining argument in both subterms,  $b$ , is identical.

## 7 Heuristics

E's heuristics are largely independent of the prover's logic and work unchanged for Ehoh. On first-order problems, Ehoh's behavior is virtually the same as E's. Yet, in preliminary experiments, we observed that some  $\lambda$ fHOL benchmarks were proved quickly by E in conjunction with the applicative encoding (Sect. 1) but timed out with Ehoh. Based on these observations, we extended the heuristics to exploit  $\lambda$ fHOL-specific features.

**Term Order Generation.** The inference rules and the redundancy criterion are parameterized by a term order—typically an instance of KBO or LPO (Sect. 3). E can generate a *symbol weight* function (for KBO) and a *symbol precedence* (for KBO and LPO) based on criteria such as the symbols' frequencies, their arities, and whether they appear in the conjecture.

In preliminary experiments, we discovered that the presence of an explicit application operator @ can be beneficial for some problems. With the applicative encoding, generation schemes can take the symbols  $@_{\tau,v}$  into account, effectively exploiting the type information carried by such symbols. To simulate this behavior, we introduced four generation schemes that extend E’s existing symbol-frequency-based schemes by partitioning the symbols by type. To each symbol, the new schemes assign a frequency corresponding to the sum of all symbol frequencies for its class. In addition, we designed four schemes that combine E’s type-agnostic and Ehoh’s type-aware approaches.

To generate symbol precedences, E can sort symbols by weight and use the symbol’s position in the sorted array as the basis for precedence. To account for the type information introduced by the applicative encoding, we implemented four type-aware precedence generation schemes.

**Literal Selection.** The side conditions of the superposition rules (SN and SP, Sect. 6) use a literal selection function to restrict which literals the rules need to consider, thereby pruning the search space. Given a clause  $L_1 \vee \dots \vee L_n$ , a literal selection function returns either  $\{L_i\}$  for some  $i$  or  $\emptyset$ . The most widely used function in E is probably `SelectMaxLComplexAvoidPosPred`, which we abbreviate to `SelectMLCAPP`. It implements a complex set of criteria that work well in practice. It only selects negative literals, and it strongly prefers literals that are maximal in the clause, a notion that is derived from the term order.

**Clause Selection.** E’s main loop repeatedly selects a clause—the given clause—and performs inferences involving it. The choice can be arbitrary, but to ensure completeness every candidate clause should be eventually selected.

E heuristically assigns *clause priorities* and *clause weights* to the candidates. E’s main loop visits, in round-robin fashion, a set of priority queues. From each queue, it selects a number of clauses with the highest priorities, breaking ties by preferring smaller weights.

E provides template weight functions that allow users to fine-tune parameters such as weights assigned to variables or function symbols. The most widely used template is `ConjectureRelativeSymbolWeight`. It computes term and clause weights according to eight parameters, notably *conj\_mul*, a multiplier applied to the weight of conjecture symbols. We implemented a new type-aware template function, called `ConjectureRelativeSymbolTypeWeight`, that applies the *conj\_mul* multiplier to all symbols whose type occurs in the conjecture.

**Configurations and Modes.** A combination of parameters—including term order, literal selection, and clause selection—is called a *configuration*. For years, E has provided an *auto* mode, which analyzes the input problem and chooses a configuration known to perform well on similar problems. More recently, E has been extended with an *autoschedule* mode, which applies a portfolio of configurations in sequence on the given problem. Configurations that perform well on a wide range of problems have emerged over time. One of them is the configuration that is most often chosen by E’s *auto* mode. We call it *boa* (“best of *auto*”).

## 8 Evaluation

How useful are Ehoh’s new heuristics? And how does Ehoh perform compared with E, used directly or in conjunction with the applicative encoding? To answer the first question, we evaluated each new parameter independently. From the empirical results, we derived a new configuration optimized for  $\lambda$ FHOL problems. To answer the second question, we compared Ehoh’s success rate on  $\lambda$ FHOL problems with E’s on their applicatively encoded counterparts. We also included first-order benchmarks to measure Ehoh’s overhead compared with E.

We set a CPU time limit of 60 s per problem. The experiments were performed on StarExec [29] nodes equipped with Intel Xeon E5-2609 0 CPUs clocked at 2.40 GHz and with 8192 MB of memory. Our raw data is publicly available.<sup>2</sup>

We used the *boa* configuration as the basis to evaluate the new heuristic schemes. For each heuristic parameter we tuned, we changed only its value while keeping the other parameters the same as for *boa*.

All heuristic parameters were tested on a suite of 5012 problems generated using Sledgehammer, consisting of four versions of the Judgment Day [12] suite. The problems were given in native  $\lambda$ FHOL syntax.

Our main findings are as follows:

1. The combination of the weight generation scheme `invtypefreqrank` and the precedence generation scheme `invtypefreq` performs best.
2. The clause selection function `ConjectureRelativeSymbolTypeWeight` with `ConstPrio` priority and an `appv_mul` factor of 1.41 performs best.
3. The literal selection heuristics `SelectMLCAPP`, `SelectMLCAPPPreferAppVar`, and `SelectMLCAPPAvoidAppVar` give virtually the same results.

We derived *hoboa* from *boa* by enabling the features identified in points 1 and 2. Next, we present a more detailed evaluation of *hoboa*, along with other configurations, on a larger benchmark suite.

The benchmark problems are partitioned as follows:

- 1147 first-order TPTP [30] problems belonging to the FOF (untyped) and TF0 (monomorphic) categories, excluding arithmetic;
- 5012 Sledgehammer-generated problems from the Judgment Day [12] suite, targeting the monomorphic first-order logic embodied by TPTP TF0;
- all 355 monomorphic higher-order problems from the TH0 category of the TPTP library belonging to the  $\lambda$ FHOL fragment;
- 5012 Sledgehammer-generated problems from the Judgment Day suite, targeting the  $\lambda$ FHOL fragment of TPTP TH0.

For the first group of benchmarks, we randomly chose 1000 FOF problems (from about 8000) and all monomorphic TFF problems that are parsable by E within 60 s (amounting to 147 out of 231 monomorphic TFF problems). Both groups

<sup>2</sup> [http://matryoshka.gforge.inria.fr/pubs/ehoh\\_data/](http://matryoshka.gforge.inria.fr/pubs/ehoh_data/)

	First-order			Higher-order		
	TPTP	SH 32	SH 512	TPTP	SH 32	SH 512
E a	671	1095	1215			
E as	<b>717</b>	<b>1187</b>	<b>1310</b>			
E b	601	1117	1244			
@+E a	578	1095	1107	338	1120	1114
@+E as	623	1165	1153	<b>341</b>	1195	1155
@+E b	589	1095	1229	339	1121	1274
Ehoh a	668	1110	1213	339	1138	1229
Ehoh as	716	<b>1187</b>	1304	340	<b>1212</b>	<b>1311</b>
Ehoh b	599	1117	1243	339	1137	1243
Ehoh hb	552	1106	1225	338	1141	1256

**Fig. 1.** Number of solved problems

of Sledgehammer problems include two subgroups of 2506 problems, generated to include 32 or 512 Isabelle lemmas (SH 32 and SH 512), to represent both smaller and larger problems arising in interactive verification. Each subgroup itself consists of two sub-subgroups of 1253 problems, generated by using either  $\lambda$ -lifting or SK-style combinators to encode  $\lambda$ -expressions.

We evaluated Ehoh against a version of E, which we call @+E, that first performs the applicative encoding. We also evaluated E, @+E, and Ehoh on first-order problems. The number of problems each prover solved is given in Figure 1. We considered the modes *auto* (a) and *autoschedule* (as) and the configurations *boa* (b) and *hoboa* (hb). We observe the following:

- By comparing the Ehoh rows with the corresponding E rows, we see that Ehoh’s overhead is small. An inspection of the raw evaluation data reveals that most of these are problems solved by E shortly before the time limit, which Ehoh would presumably have solved if given a few more seconds.
- Ehoh generally outperforms the applicative encoding, on both first-order and higher-order problems. On Sledgehammer benchmarks, the best Ehoh mode (*autoschedule*) clearly outperforms all @+E modes and configurations.
- The *hoboa* configuration outperforms *boa* on higher-order problems, suggesting that it could be worthwhile to update *auto* and *autoschedule*.

For future work, we plan to design further heuristics and re-train *auto* and *autoschedule* based on  $\lambda$ HOL benchmarks.

## 9 Discussion and Related Work

Most higher-order provers were developed from the ground up. Notable exceptions are Otter- $\lambda$  by Beeson [6] and Zipperposition by Cruanes [14]. Otter- $\lambda$  adds  $\lambda$ -expressions and second-order unification to the superposition-based Otter [20]. The approach is resolutely pragmatic, with little emphasis on completeness or gracefulness. Zipperposition is a superposition-based prover written in OCaml. It

was initially designed for first-order logic but subsequently extended to higher-order logic. Its performance is a far cry from E’s, but it is easier to modify. It is used by Bentkamp et al. [7] for prototyping. Finally, there is noteworthy preliminary work by the developers of Vampire [10] and of CVC4 and veriT [4].

Native higher-order reasoning was pioneered by Robinson [23], Andrews [1], and Huet [16]. TPS, by Andrews et al. [2], was based on expansion proofs and let the user specify proof outlines. The Leo family of systems, developed by Benzmüller and his colleagues, is based on resolution and paramodulation. LEO [8] introduced the cooperative paradigm to integrate first-order provers. Leo-III [28] expands the cooperation with SMT solvers and introduces term orders. Brown’s Satallax [13] is based on a complete higher-order tableau, guided by a SAT solver; recent versions also cooperate with first-order provers.

An alternative to all of the above is to reduce higher-order logic to first-order logic by means of a translation. Robinson [24] outlined this approach decades before tools such as Sledgehammer [22] and HOLyHammer [17] popularized it in proof assistants. In addition to performing an applicative encoding, such translations must eliminate the  $\lambda$ -expressions and encode the type information.

By removing the need for the applicative encoding, our current work reduces the translation gap. The encoding buries the  $\lambda$ fHOL terms’ heads under layers of @ symbols, possibly influencing the heuristics. Terms double in size, cluttering the data structures, and twice as many subterm positions must be considered for inferences. In addition, the encoding must be undone in the proofs. But our main objection to the applicative encoding is that it is incompatible with interpreted symbols, notably for arithmetic. The traditional solution is to introduce proxies to connect an uninterpreted nullary symbol with its interpreted counterpart, but this is clumsy. Moreover, in a monomorphic logic, @ is actually a type-indexed family of symbols  $@_{\tau,v}$ , which must be correctly introduced and recognized. While it should be possible to base a higher-order prover on such an encoding, the prospect is aesthetically and technically unappealing.

## 10 Conclusion

Despite considerable progress since the 1970s, higher-order automated reasoning has not yet assimilated some of the most successful methods for first-order logic with equality, such as superposition. We presented a graceful extension of a state-of-the-art first-order theorem prover to a  $\lambda$ -free fragment of higher-order logic. Our work covers both theoretical and practical aspects. Experiments show promising results on higher-order problems and almost no overhead for first-order problems, as we would expect from a graceful generalization.

The resulting Ehoh prover will form the basis of our work towards strong higher-order automation. Our aim is to turn it into a prover that excels on proof obligations emerging from interactive verification. In our experience, these tend to be large but only mildly higher-order. Most of our techniques, including the prefix optimization, are applicable to other superposition-based provers, such as SPASS and Vampire, and to the paramodulation-based higher-order Leo-III.

**Acknowledgment.** We are grateful to the maintainers of StarExec for letting us use their service. We thank Ahmed Bhayat, Alexander Bentkamp, Daniel El Ouraoui, Michael Färber, Pascal Fontaine, Predrag Janičić, Tomer Libal, Giles Reger, Hans-Jörg Schurr, and Alexander Steen for suggesting many textual improvements to this paper and to Vukmirović’s MSc thesis. We also want to thank the other members of the Matryoshka team, including Sophie Tourret and Uwe Waldmann, as well as Christoph Benzmüller, Andrei Voronkov, and Daniel Wand, for many stimulating discussions. Vukmirović and Blanchette’s research has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 713999, Matryoshka).

## References

- [1] Andrews, P.B.: Resolution in type theory. *J. Symb. Log.* 36(3), 414–432 (1971)
- [2] Andrews, P.B., Bishop, M., Issar, S., Nesmith, D., Pfenning, F., Xi, H.: TPS: A theorem-proving system for classical type theory. *J. Autom. Reason.* 16(3), 321–353 (1996)
- [3] Baader, F., Nipkow, T.: *Term Rewriting and All That*. Cambridge University Press (1998)
- [4] Barbosa, H., Reynolds, A., Fontaine, P., Ouraoui, D.E., Tinelli, C.: Higher-order SMT solving. In: Dimitrova, R., D’Silva, V. (eds.) *SMT 2018* (2018)
- [5] Becker, H., Blanchette, J.C., Waldmann, U., Wand, D.: A transfinite Knuth–Bendix order for lambda-free higher-order terms. In: de Moura, L. (ed.) *CADE-26*. LNCS, vol. 10395, pp. 432–453. Springer (2017)
- [6] Beeson, M.: Lambda logic. In: Basin, D.A., Rusinowitch, M. (eds.) *IJCAR 2004*. LNCS, vol. 3097, pp. 460–474. Springer (2004)
- [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] Benzmüller, C., Kohlhase, M.: System description: LEO—a higher-order theorem prover. In: Kirchner, C., Kirchner, H. (eds.) *CADE-15*. LNCS, vol. 1421, pp. 139–144. Springer (1998)
- [9] Benzmüller, C., Sultana, N., Paulson, L.C., Theiss, F.: The higher-order prover LEO-II. *J. Autom. Reason.* 55(4), 389–404 (2015)
- [10] Bhayat, A., Reger, 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] Blanchette, J.C., Waldmann, U., Wand, D.: A lambda-free higher-order recursive path order. In: Esparza, J., Murawski, A.S. (eds.) *FoSSaCS 2017*. LNCS, vol. 10203, pp. 461–479. Springer (2017)
- [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.: Superposition with structural induction. In: Dixon, C., Finger, M. (eds.) *FroCoS 2017*. LNCS, vol. 10483, pp. 172–188. Springer (2017)
- [15] Filliâtre, J.-C., Paskevich, A.: Why3—where programs meet provers. In: Felleisen, M., Gardner, P. (eds.) *ESOP 2013*. LNCS, vol. 7792, pp. 125–128. Springer (2013)
- [16] Huet, G.P.: A mechanization of type theory. In: Nilsson, N.J. (ed.) *IJCAI-73*. pp. 139–146. William Kaufmann (1973)



- [17] Kaliszyk, C., Urban, J.: HOL(y)Hammer: Online ATP service for HOL Light. *Math. Comput. Sci.* 9(1), 5–22 (2015)
- [18] Kovács, L., Voronkov, A.: First-order theorem proving and Vampire. In: Sharygina, N., Veith, H. (eds.) CAV 2013. LNCS, vol. 8044, pp. 1–35. Springer (2013)
- [19] Löchner, B.: Things to know when implementing KBO. *J. Autom. Reason.* 36(4), 289–310 (2006)
- [20] McCune, W.: OTTER 2.0. In: Stickel, M.E. (ed.) CADE-10. LNCS, vol. 449, pp. 663–664. Springer (1990)
- [21] McCune, W.: Experiments with discrimination-tree indexing and path indexing for term retrieval. *J. Autom. Reason.* 9(2), 147–167 (1992)
- [22] 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)
- [23] Robinson, J.: Mechanizing higher order logic. In: Meltzer, B., Michie, D. (eds.) *Machine Intelligence*, vol. 4, pp. 151–170. Edinburgh University Press (1969)
- [24] Robinson, J.: A note on mechanizing higher order logic. In: Meltzer, B., Michie, D. (eds.) *Machine Intelligence*, vol. 5, pp. 121–135. Edinburgh University Press (1970)
- [25] Schulz, S.: Fingerprint indexing for paramodulation and rewriting. In: Gramlich, B., Miller, D., Sattler, U. (eds.) IJCAR 2012. LNCS, vol. 7364, pp. 477–483. Springer (2012)
- [26] Schulz, S.: Simple and efficient clause subsumption with feature vector indexing. In: Bonacina, M.P., Stickel, M.E. (eds.) *Automated Reasoning and Mathematics—Essays in Memory of William W. McCune*. LNCS, vol. 7788, pp. 45–67. Springer (2013)
- [27] Schulz, S.: System description: E 1.8. In: McMillan, K.L., Middeldorp, A., Voronkov, A. (eds.) LPAR-19. LNCS, vol. 8312, pp. 735–743. Springer (2013)
- [28] Steen, A., Benzmüller, 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)
- [29] Stump, A., Sutcliffe, G., Tinelli, C.: StarExec: A cross-community infrastructure for logic solving. In: Demri, S., Kapur, D., Weidenbach, C. (eds.) IJCAR 2014. LNCS, vol. 8562, pp. 367–373. Springer (2014)
- [30] 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)
- [31] Sutcliffe, G.: The CADE-26 automated theorem proving system competition—CASC-26. *AI Commun.* 30(6), 419–432 (2017)
- [32] Vukmirović, P.: Implementation of Lambda-Free Higher-Order Superposition. MSc thesis, Vrije Universiteit Amsterdam (2018), [http://matryoshka.gforge.inria.fr/pubs/vukmirovic\\_msc\\_thesis.pdf](http://matryoshka.gforge.inria.fr/pubs/vukmirovic_msc_thesis.pdf)
- [33] Vukmirović, P., Blanchette, J.C., Cruanes, S., Schulz, S.: Extending a brainiac prover to lambda-free higher-order logic (technical report). Technical report (2018), [http://matryoshka.gforge.inria.fr/pubs/ehoh\\_report.pdf](http://matryoshka.gforge.inria.fr/pubs/ehoh_report.pdf)
- [34] Weidenbach, C., Dimova, D., Fietzke, A., Kumar, R., Suda, M., Wischniewski, P.: SPASS version 3.5. In: Schmidt, R.A. (ed.) CADE-22. LNCS, vol. 5663, pp. 140–145. Springer (2009)