

Can a Higher-Order and a First-Order Theorem Prover Cooperate?*

Christoph Benz Müller¹ Volker Sorge² Mateja Jamnik³ Manfred Kerber²

¹Fachbereich Informatik, Universität des Saarlandes
66041 Saarbrücken, Germany (www.ags.uni-sb.de/~chris)

²School of Computer Science, The University of Birmingham
Birmingham B15 2TT, England, UK (www.cs.bham.ac.uk/~vxs|~mmk)

³University of Cambridge Computer Laboratory
Cambridge CB3 0FD, England, UK (www.cl.cam.ac.uk/~mj201)

Abstract. State-of-the-art first-order automated theorem proving systems have reached considerable strength over recent years. However, in many areas of mathematics they are still a long way from reliably proving theorems that would be considered relatively simple by humans. For example, when reasoning about sets, relations, or functions, first-order systems still exhibit serious weaknesses. While it has been shown in the past that higher-order reasoning systems can solve problems of this kind automatically, the complexity inherent in their calculi and their inefficiency in dealing with large numbers of clauses prevent these systems from solving a whole range of problems.

We present a solution to this challenge by combining a higher-order and a first-order automated theorem prover, both based on the resolution principle, in a flexible and distributed environment. By this we can exploit concise problem formulations without forgoing efficient reasoning on first-order subproblems. We demonstrate the effectiveness of our approach on a set of problems still considered non-trivial for many first-order theorem provers.

1 Introduction

When dealing with problems containing higher-order concepts, such as sets, functions, or relations, today's state-of-the-art first-order automated theorem provers (ATPs) still exhibit weaknesses on problems considered relatively simple by humans (cf. [14]). One reason is that the problem formulations use an encoding in a first-order set theory, which makes it particularly challenging when trying to prove theorems from first principles, that is, basic axioms. Therefore, to aid ATPs in finding proofs, problems are often enriched by hand-picked additional lemmata, or axioms of the selected set theory are dropped leaving the theory

* This work was supported by EPSRC grant GR/M22031 and DFG-SFB 378 (first author), EU Marie-Curie-Fellowship HPMF-CT-2002-01701 (second author), and EPSRC Advanced Research Fellowship GR/R76783 (third author).

incomplete. This has recently motivated extensions of state-of-the-art first-order calculi and systems, as for example presented in [14] for the SATURATE system. The extended SATURATE system can solve some problems from the SET domain in the TPTP [24] which VAMPIRE [21] and E-SETHEO's [23] cannot solve.

While it has already been shown in [6, 2] that many problems of this nature can be easily proved from first principles using a concise higher-order representation and the higher-order resolution ATP LEO, the combinatorial explosion inherent in LEO's calculus prevents the prover from solving a whole range of possible problems with one universal strategy. Often higher-order problems require only relatively few but essential steps of higher-order reasoning, while the overwhelming part of the reasoning is first-order or even propositional level. This suggests that LEO's performance could be improved when combining it with a first-order ATP to search efficiently for a possible refutation in the subset of those clauses that are essentially first-order.

The advantages of such a combination — further discussed in Sec. 2 — are not only that many problems can still be efficiently shown from first principles in a general purpose approach, but also that problems can be expressed in a very concise way. For instance, we present 45 problems from the SET domain of the TPTP-v3.0.1, together with their entire formalisation in less than two pages in this paper, which is difficult to achieve within a framework that does not provide λ -abstraction. We use this problem set, which is an extension of the problems considered in [14], in Sec. 4 to show the effectiveness of our approach. While many of the considered problems can be proved by LEO alone with some strategy, the combination of LEO with the first-order ATP BLIKSEM [11] is not only able to show more problems, but also needs only a single strategy to solve them. Several of our problems are considered very challenging by the first-order community and five of them (of which LEO can solve four) have a TPTP rating of 1.00, saying that they cannot be solved by any TPTP prover to date.

Technically, the combination — described in more detail in Sec. 3 — has been realised in the concurrent reasoning system OANTS [22, 8] which enables the cooperation of hybrid reasoning systems to construct a common proof object. In our past experiments, OANTS has been successfully employed to check the validity of set equations using higher-order and first-order ATPs, model generation, and computer algebra [5]. While this already enabled a cooperation between LEO and a first-order ATP, the proposed solution could not be classified as a general purpose approach. A major shortcoming was that all communication of partial results had to be conducted via the common proof object, which was very inefficient for hard examples. Thus, the solved examples from set theory were considered too trivial, albeit they were often similar to those still considered challenging in the TPTP in the first-order context. In this paper we now present a novel approach to the cooperation between LEO and BLIKSEM inside OANTS by decentralising communication. This leads not only to a higher overall efficiency — Sec. 4 details our results — but also to a general purpose approach based on a single strategy in LEO.

2 Why Linking Higher-Order and First-Order?

Existing higher-order ATPs generally exhibit deficits in efficiently reasoning with first-order problems for several reasons. Unlike in the case of first-order provers, for which sophisticated calculi and strategies, as well as advanced implementation techniques, such as term indexing [19], have been developed, fully mechanisable higher-order calculi are still at a comparably early stage of development. Some problems are much harder in higher-order, for instance, unification is undecidable, strong constraining term- and literal-orderings are not available, extensionality reasoning and set variable instantiation has to be addressed. Nevertheless, for some mathematical problem domains, such as naive set theory, for instance, automated higher-order reasoning performs very well.

We motivate the need for linking higher-order and first-order ATPs with some examples from Table 1. It contains a range of challenging problems taken from the TPTP, against which we will evaluate our system in Sec. 4. The problems are given by the identifiers used in the SET domain of the TPTP, and are formalised in a variant of Church’s simply typed λ -calculus with prefix polymorphism. In classical type theory terms and all their sub-terms are typed. Polymorphism allows the introduction of type variables such that statements can be made for all types. For instance, in problem SET014+4 the universally quantified variable $X_{o\alpha}$ denotes a mapping from objects of type α to objects of type o . We use Church’s notation $o\alpha$, which stands for the functional type $\alpha \rightarrow o$. The reader is referred to [1] for a more detailed introduction. In the remainder, o will denote the type of truth values, and small Greek letters will denote arbitrary types. Thus, $X_{o\alpha}$ (resp. its η -longform $\lambda y_{\alpha}.Xy$) is actually a characteristic function denoting the set of elements of type α , for which the predicate associated with X holds. As further notational convention, we use capital letter variables to denote sets, functions, or relations, while lower case letters denote individuals. Types are usually only given in the first occurrence of a variable and omitted if inferable from the context.

The problems in Table 1 employ defined concepts that are specified in a knowledge base of hierarchical theories that LEO has access to. All concepts necessary for defining our problems in Table 1 are given in Table 2. Concepts are defined in terms of λ -expressions and they may contain other, already specified concepts. For presentation purposes, we use customary mathematical symbols \cup, \cap , etc., for some concepts like *union*, *intersection*, etc., and we also use infix notation. For instance, the definition of union on sets can be easily read in its more common mathematical representation $A \cup B := \{x | x \in A \vee x \in B\}$. Before proving a problem, LEO always expands — recursively, if necessary — all occurring concepts. This straightforward expansion to first principles is realised by an automated preprocess in our current approach.

SET171+3 We first discuss example SET171+3 to contrast our formalisation to a standard first-order one. After recursively expanding the input problem, that is, completely reducing it to first principles, LEO turns it into a negated unit clause. Since this initial clause is not in normal form, LEO first normalises it with explicit

| SET | Problem Formalisation |
|-------|---|
| 014+4 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, A_{\alpha\alpha} \bullet [X \subseteq A \wedge Y \subseteq A] \Rightarrow (X \cup Y) \subseteq A]$ |
| 017+1 | $\forall x_{\alpha}, y_{\alpha}, z_{\alpha} \bullet [UnOrderedPair(x, y) = UnOrderedPair(x, z) \Rightarrow y = z]$ |
| 066+1 | $\forall x_{\alpha}, y_{\alpha} \bullet [UnOrderedPair(x, y) = UnOrderedPair(y, x)]$ |
| 067+1 | $\forall x_{\alpha}, y_{\alpha} \bullet [UnOrderedPair(x, x) \subseteq UnOrderedPair(x, y)]$ |
| 076+1 | $\forall x_{\alpha}, y_{\alpha} \bullet \forall Z_{\alpha\alpha} \bullet x \in Z \wedge y \in Z \Rightarrow UnOrderedPair(x, y) \subseteq Z]$ |
| 086+1 | $\forall x_{\alpha} \bullet \exists y_{\alpha} \bullet [y \in Singleton(x)]$ |
| 096+1 | $\forall X_{\alpha\alpha}, y_{\alpha} \bullet [X \subseteq Singleton(y) \Rightarrow [X = \emptyset \vee X = Singleton(y)]]$ |
| 143+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [(X \cap Y) \cap Z = X \cap (Y \cap Z)]$ |
| 171+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [X \cup (Y \cap Z) = (X \cup Y) \cap (X \cup Z)]$ |
| 580+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, u_{\alpha} \bullet [u \in ExclUnion(X, Y) \Leftrightarrow [u \in X \Leftrightarrow u \notin Y]]$ |
| 601+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [(X \cap Y) \cup ((Y \cap Z) \cup (Z \cap X)) = (X \cup Y) \cap ((Y \cup Z) \cap (Z \cup X))]$ |
| 606+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [X \setminus (X \cap Y) = X \setminus Y]$ |
| 607+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [X \cup (Y \setminus X) = X \cup Y]$ |
| 609+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [X \setminus (Y \setminus Z) = (X \setminus Y) \cup (X \cap Z)]$ |
| 611+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [X \cap Y = \emptyset \Leftrightarrow X \setminus Y = X]$ |
| 612+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [X \setminus (Y \cup Z) = (X \setminus Y) \cap (X \setminus Z)]$ |
| 614+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [(X \setminus Y) \setminus Z = X \setminus (Y \cup Z)]$ |
| 615+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [(X \cup Y) \setminus Z = (X \setminus Z) \cup (Y \setminus Z)]$ |
| 623+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [ExclUnion(ExclUnion(X, Y), Z) = ExclUnion(X, ExclUnion(Y, Z))]$ |
| 624+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, Z_{\alpha\alpha} \bullet [Meets(X, (Y \cup Z)) \Leftrightarrow [Meets(X, Y) \vee Meets(X, Z)]]$ |
| 630+3 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [Misses(X \cap Y, ExclUnion(X, Y))]$ |
| 640+3 | $\forall R_{\alpha\beta\alpha}, Q_{\alpha\beta\alpha} \bullet [Subrel(R, Q) \Rightarrow Subrel(R, (\lambda u_{\alpha} \bullet \top) \times (\lambda v_{\beta} \bullet \top))]$ |
| 646+3 | $\forall x_{\alpha}, y_{\beta} \bullet [Subrel(Pair(x, y), (\lambda u_{\alpha} \bullet \top) \times (\lambda v_{\beta} \bullet \top))]$ |
| 647+3 | $\forall R_{\alpha\beta\alpha}, X_{\alpha\alpha} \bullet [(RDom(R) \subseteq X) \Rightarrow Subrel(R, X \times RCodom(R))]$ |
| 648+3 | $\forall R_{\alpha\beta\alpha}, Y_{\alpha\alpha} \bullet [(RCodom(R) \subseteq Y) \Rightarrow Subrel(R, RDom(R) \times Y)]$ |
| 649+3 | $\forall R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [RDom(R) \subseteq X \wedge RCodom(R) \subseteq Y] \Rightarrow Subrel(R, X \times Y)]$ |
| 651+3 | $\forall R_{\alpha\beta\alpha} \bullet [RDom(R) \subseteq A_{\alpha\alpha} \Rightarrow Subrel(R, A \times (\lambda v_{\beta} \bullet \top))]$ |
| 657+3 | $\forall R_{\alpha\beta\alpha} \bullet [Field(R) \subseteq ((\lambda u_{\alpha} \bullet \top) \cup (\lambda v_{\beta} \bullet \top))]$ |
| 669+3 | $\forall R_{\alpha\alpha\alpha} \bullet [Subrel(Id(\lambda u_{\alpha} \bullet \top), R) \Rightarrow [(\lambda u_{\alpha} \bullet \top) \subseteq RDom(R) \wedge (\lambda u_{\alpha} \bullet \top) = RCodom(R)]]$ |
| 670+3 | $\forall Z_{\alpha\alpha}, R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \Rightarrow IsRelOn(RestrictRDom(R, Z), Z, Y)]$ |
| 671+3 | $\forall Z_{\alpha\alpha}, R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \wedge X \subseteq Z] \Rightarrow RestrictRDom(R, Z) = R]$ |
| 672+3 | $\forall Z_{\alpha\alpha}, R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \Rightarrow IsRelOn(RestrictRCodom(R, Z), X, Z)]$ |
| 673+3 | $\forall Z_{\alpha\alpha}, R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \wedge Y \subseteq Z] \Rightarrow RestrictRCodom(R, Z) = R]$ |
| 680+3 | $\forall R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \Rightarrow$ $[\forall u_{\alpha} \bullet u \in X \Rightarrow [u \in RDom(R) \Leftrightarrow \exists v_{\beta} \bullet v \in Y \wedge R(u, v)]]]$ |
| 683+3 | $\forall R_{\alpha\beta\alpha}, X_{\alpha\alpha}, Y_{\alpha\alpha} \bullet [IsRelOn(R, X, Y) \Rightarrow$ $[\forall v_{\beta} \bullet v \in Y \Rightarrow [v \in RCodom(R) \Rightarrow \exists u_{\alpha} \bullet u \in X \wedge u \in RDom(R)]]]$ |
| 684+3 | $\forall P_{\alpha\beta\alpha}, R_{\alpha\gamma\beta}, x_{\alpha}, z_{\gamma} \bullet [RelComp(P, R)xz \Leftrightarrow \exists y_{\beta} \bullet Pxy \wedge Ryz]$ |
| 686+3 | $\forall Z_{\alpha\alpha}, R_{\alpha\gamma\beta}, x_{\alpha} \bullet [x \in InverseImageR(R, Z) \Leftrightarrow \exists y_{\alpha} \bullet Rxy \wedge x \in Z]$ |
| 716+4 | $\forall F_{\beta\alpha}, G_{\gamma\beta} \bullet [[Inj(F) \wedge Inj(G)] \Rightarrow Inj(G \circ F)]$ |
| 724+4 | $\forall F_{\beta\alpha}, G_{\gamma\beta}, H_{\gamma\beta} \bullet [[F \circ G = F \circ H \wedge Surj(F)] \Rightarrow G = H]$ |
| 741+4 | $\forall F_{\beta\alpha}, G_{\gamma\beta}, H_{\alpha\gamma} \bullet [[Inj((F \circ G) \circ H) \wedge Surj((G \circ H) \circ F) \wedge Surj((H \circ F) \circ G)] \Rightarrow Bij(H)]$ |
| 747+4 | $\forall F_{\beta\alpha}, G_{\gamma\beta}, \triangleleft_{\alpha\alpha}^1, \triangleleft_{\alpha\alpha}^2, \triangleleft_{\alpha\alpha}^3, \triangleleft_{\gamma\gamma}^1, \triangleleft_{\gamma\gamma}^2, \triangleleft_{\gamma\gamma}^3 \bullet [[IncreasingF(F, \triangleleft^1, \triangleleft^2) \wedge DecreasingF(G, \triangleleft^2, \triangleleft^3)] \Rightarrow$ $DecreasingF(F \circ G, \triangleleft^1, \triangleleft^3)]$ |
| 752+4 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, F_{\beta\alpha} \bullet [ImageF(F, X \cup Y) = ImageF(F, X) \cup ImageF(F, Y)]$ |
| 753+4 | $\forall X_{\alpha\alpha}, Y_{\alpha\alpha}, F_{\beta\alpha} \bullet [ImageF(F, X \cap Y) \subseteq ImageF(F, X) \cap ImageF(F, Y)]$ |
| 764+4 | $\forall F_{\beta\alpha} \bullet [InverseImageF(F, \emptyset) = \emptyset]$ |
| 770+4 | $\forall R_{\alpha\beta\alpha}, Q_{\alpha\beta\alpha} \bullet [[EquivRel(R) \wedge EquivRel(Q)] \Rightarrow$ $[EquivClasses(R) = EquivClasses(Q) \vee Disjoint(EquivClasses(R), EquivClasses(Q))]]]$ |

Table 1. Problems from TPTP for the evaluation of OANTS.

clause normalisation rules to reach some proper initial clauses. In our concrete case, this normalisation process leads to the following unit clause consisting of a (syntactically not solvable) unification constraint (here $B_{\alpha\alpha}, C_{\alpha\alpha}, D_{\alpha\alpha}$ are Skolem constants and Bx is obtained from expansion of $x \in B$):

$$[(\lambda x_{\alpha} \bullet Bx \vee (Cx \wedge Dx)) =^? (\lambda x_{\alpha} \bullet (Bx \vee Cx) \wedge (Bx \vee Dx))]$$

Note that negated primitive equations are generally automatically converted by LEO into unification constraints. This is why $[(\lambda x_{\alpha} \bullet Bx \vee (Cx \wedge Dx)) =^?$

| Defined Notions in Theory Typed Set | |
|-------------------------------------|---|
| $- \in -$ | $:= \lambda x_\alpha, A_{o\alpha} \cdot [A x]$ |
| \emptyset | $:= [\lambda x_\alpha \cdot \perp]$ |
| $- \subseteq -$ | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\forall x_\alpha \cdot x \in A \Rightarrow x \in B]$ |
| $- \cup -$ | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\lambda x_\alpha \cdot x \in A \vee x \in B]$ |
| $- \cap -$ | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\lambda x_\alpha \cdot x \in A \wedge x \in B]$ |
| $-$ | $:= \lambda A_{o\alpha} \cdot [\lambda x_\alpha \cdot x \notin A]$ |
| $- \setminus -$ | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\lambda x_\alpha \cdot x \in A \wedge x \notin B]$ |
| <i>ExclUnion</i> (-, -) | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [(A \setminus B) \cup (B \setminus A)]$ |
| <i>Disjoint</i> (-, -) | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [A \cap B = \emptyset]$ |
| <i>Meets</i> (-, -) | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\exists x_\alpha \cdot x \in A \wedge x \in B]$ |
| <i>Misses</i> (-, -) | $:= \lambda A_{o\alpha}, B_{o\alpha} \cdot [\neg \exists x_\alpha \cdot x \in A \wedge x \in B]$ |
| Defined Notions in Theory Relation | |
| <i>UnOrderedPair</i> (-, -) | $:= \lambda x_\alpha, y_\alpha \cdot [\lambda u_\alpha \cdot u = x \vee u = y]$ |
| <i>Singleton</i> (-) | $:= \lambda x_\alpha \cdot [\lambda u_\alpha \cdot u = x]$ |
| <i>Pair</i> (-, -) | $:= \lambda x_\alpha, y_\beta \cdot [\lambda u_\alpha, v_\beta \cdot u = x \wedge v = y]$ |
| $- \times -$ | $:= \lambda A_{o\alpha}, B_{o\beta} \cdot [\lambda u_\alpha, v_\beta \cdot u \in A \wedge v \in B]$ |
| <i>RDom</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\lambda x_\alpha \cdot \exists y_\beta \cdot R x y]$ |
| <i>RCodom</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\lambda y_\beta \cdot \exists x_\alpha \cdot R x y]$ |
| <i>Subrel</i> (-, -) | $:= \lambda R_{o\beta\alpha}, Q_{o\beta\alpha} \cdot [\forall x_\alpha \cdot \forall y_\alpha \cdot R x y \Rightarrow Q x y]$ |
| <i>Id</i> (-) | $:= \lambda A_{o\alpha} \cdot [\lambda x_\alpha, y_\alpha \cdot x \in A \wedge x = y]$ |
| <i>Field</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [RDom(B) \cup RCodom(R)]$ |
| <i>IsRelOn</i> (-, -) | $:= \lambda R_{o\beta\alpha}, A_{o\alpha} \cdot \lambda B_{o\beta} \cdot [\forall x_\alpha, y_\beta \cdot R x y \Rightarrow (x \in A \wedge y \in B)]$ |
| <i>RestrictRCodom</i> (-, -) | $:= \lambda R_{o\beta\alpha}, A_{o\alpha} \cdot [\lambda x_\alpha, y_\beta \cdot x \in A \wedge R x y]$ |
| <i>RelComp</i> (-, -) | $:= \lambda R_{o\beta\alpha}, Q_{o\gamma\beta} \cdot [\lambda x_\alpha, z_\gamma \cdot \exists y_\beta \cdot R x y \wedge R y z]$ |
| <i>InverseImageR</i> (-, -) | $:= \lambda R_{o\beta\alpha}, B_{o\beta} \cdot [\lambda x_\alpha \cdot \exists y_\beta \cdot y \in B \wedge R x y]$ |
| <i>Reflexive</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\forall x_\alpha \cdot R x x]$ |
| <i>Symmetric</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\forall x_\alpha \cdot \forall y_\alpha \cdot R x y \Rightarrow R y x]$ |
| <i>Transitive</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\forall x_\alpha \cdot \forall y_\alpha \cdot \forall z_\alpha \cdot R x y \wedge R y z \Rightarrow R x z]$ |
| <i>EquivRel</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [Reflexive(R) \wedge Symmetric(R) \wedge Transitive(R)]$ |
| <i>EquivClasses</i> (-) | $:= \lambda R_{o\beta\alpha} \cdot [\lambda A_{o\alpha} \cdot \exists u_\alpha \cdot u \in A \wedge \forall v_\alpha \cdot v \in A \Leftrightarrow R u v]$ |
| Defined Notions in Theory Function | |
| <i>Inj</i> (-) | $:= \lambda F_{\beta\alpha} \cdot [\forall x_\alpha, y_\beta \cdot F(x) = F(y) \Rightarrow x = y]$ |
| <i>Surj</i> (-) | $:= \lambda F_{\beta\alpha} \cdot [\forall y_\beta \cdot \exists x_\alpha \cdot y = F(x)]$ |
| <i>Bij</i> (-) | $:= \lambda F_{\beta\alpha} \cdot [Surj(F) \wedge Inj(F)]$ |
| <i>ImageF</i> (-, -) | $:= \lambda F_{\beta\alpha}, A_{o\alpha} \cdot [\lambda y_\beta \cdot \exists x_\alpha \cdot x \in A \wedge y = F(x)]$ |
| <i>InverseImageF</i> (-, -) | $:= \lambda F_{\beta\alpha}, B_{o\beta} \cdot [\lambda x_\alpha \cdot \exists y_\beta \cdot y \in B \wedge y = F(x)]$ |
| $- \circ -$ | $:= \lambda F_{\beta\alpha}, G_{\gamma\beta} \cdot [\lambda x_\alpha \cdot G(F(x))]$ |
| <i>IncreasingF</i> (-, -) | $:= \lambda F_{\beta\alpha}, \triangleleft_{o\alpha\alpha}^1, \triangleleft_{o\beta\beta}^2 \cdot [\forall x_\alpha, y_\alpha \cdot x \triangleleft^1 y \Rightarrow F(x) \triangleleft^2 F(y)]$ |
| <i>DecreasingF</i> (-, -) | $:= \lambda F_{\beta\alpha}, \triangleleft_{o\alpha\alpha}^1, \triangleleft_{o\beta\beta}^2 \cdot [\forall x_\alpha, y_\alpha \cdot x \triangleleft^1 y \Rightarrow F(y) \triangleleft^2 F(x)]$ |

Table 2. Defined concepts occurring in problems from Table 1.

$(\lambda x_\alpha \cdot (B x \vee C x) \wedge (B x \vee D x))$ is generated, and not $[(\lambda x_\alpha \cdot B x \vee (C x \wedge D x)) = (\lambda x_\alpha \cdot (B x \vee C x) \wedge (B x \vee D x))]^F$. Observe, that we write $[.]^T$ and $[.]^F$ for positive and negative literals, respectively. LEO then applies its goal directed functional and Boolean extensionality rules which replace this unification constraint by the negative literal (where x is a Skolem constant):

$$[(B x \vee (C x \wedge D x)) \Leftrightarrow ((B x \vee C x) \wedge (B x \vee D x))]^F$$

This unit clause is again not normal; normalisation, factorisation and subsumption yield the following set of clauses:

$$[B x]^F \quad [B x]^T \vee [C x]^T \quad [B x]^T \vee [D x]^T \quad [C x]^F \vee [D x]^F$$

This set is essentially of propositional logic character and trivially refutable. LEO needs 0.56 seconds for solving the problem and generates a total of 36 clauses.

| | |
|---|-----|
| Assumptions: $\forall B, C, x. [x \in (B \cup C) \Leftrightarrow x \in B \vee x \in C]$ | (1) |
| $\forall B, C, x. [x \in (B \cap C) \Leftrightarrow x \in B \wedge x \in C]$ | (2) |
| $\forall B, C. [B = C \Leftrightarrow B \subseteq C \wedge C \subseteq B]$ | (3) |
| $\forall B, C. [B \cup C = C \cup B]$ | (4) |
| $\forall B, C. [B \cap C = C \cap B]$ | (5) |
| $\forall B, C. [B \subseteq C \Leftrightarrow \forall x. x \in B \Rightarrow x \in C]$ | (6) |
| $\forall B, C. [B = C \Leftrightarrow \forall x. x \in B \Leftrightarrow x \in C]$ | (7) |
| Proof Goal: $\forall B, C, D. [B \cup (C \cap D) = (B \cup C) \cap (B \cup D)]$ | (8) |

Table 3. TPTP problem SET171+3 — distributivity of \cup over \cap .

Let us consider now this same example SET171+3 in its first-order formulation from the TPTP (see Table 3). We can observe that the assumptions provide only a partial axiomatisation of naive set theory. On the other hand, the specification introduces lemmata that are useful for solving the problem. In particular, assumption (7) is trivially derivable from (3) with (6). Obviously, clausal normalisation of this first-order problem description yields a much larger and more difficult set of clauses. Furthermore, definitions of concepts are not directly expanded as in LEO. It is therefore not surprising that most first-order ATPs still fail to prove this problem. In fact, very few TPTP provers were successful in proving SET171+3. Amongst them are MUSCADET 2.4. [20], VAMPIRE 7.0, and SATURATE. The natural deduction system MUSCADET uses special inference rules for sets and needs 0.2 seconds to prove this problem. VAMPIRE needs 108 seconds. The SATURATE system [14] (which extends VAMPIRE with Boolean extensionality rules that are a one-to-one correspondence to LEO’s rules for Extensional Higher-Order Paramodulation [3]) can solve the problem in 2.9 seconds while generating 159 clauses. The significance of such comparisons is clearly limited since different systems are optimised to a different degree. One noted difference between the experiments with first-order provers listed above, and the experiments with LEO and LEO-BLIKSEM is that first-order systems often use a case tailored problem representation (e.g., by avoiding some base axioms of the addressed theory), while LEO and LEO-BLIKSEM have a harder task of dealing with a general (not specifically tailored) representation.

For the experiments with LEO and the cooperation of LEO with the first-order theorem prover BLIKSEM, λ -abstraction as well as the extensionality treatment inherent in LEO’s calculus [4] is used. This enables a theoretically¹ Henkin-complete proof system for set theory. In the above example SET171+3, LEO generally uses the application of functional extensionality to push extensional unification constraints down to base type level, and then eventually applies Boolean extensionality to generate clauses from them. These are typically much simpler

¹ For pragmatic reasons, such as efficiency, most of LEO’s tactics are incomplete. LEO’s philosophy is to rely on a theoretically complete calculus, but to practically provide a set of complimentary strategies so that these cover a broad range of theorems.

and often even *propositional-like* or *first-order-like* (FO-like, for short), that is, they do not contain any ‘real’ higher-order subterms (such as a λ -abstraction or embedded equations), and are therefore suitable for treatment by a first-order ATP or even a propositional logic decision procedure.

SET624+3 Sometimes, extensionality treatment is not required and the originally higher-order problem is immediately reduced to only FO-like clauses. For example, after expanding the definitions, problem SET624+3 yields the following clause (where $B_{o\alpha}, C_{o\alpha}, D_{o\alpha}$ are again Skolem constants):

$$[(\exists x_{\alpha\bullet}(Bx \wedge (Cx \vee Dx)) \Leftrightarrow ((\exists x_{\alpha\bullet}Bx \wedge Cx) \vee (\exists x_{\alpha\bullet}Bx \wedge Dx)))]^F$$

Normalisation results in 26 FO-like clauses, which present a hard problem for LEO: it needs approx. 35 seconds (see Sec. 4) to find a refutation, whereas first-order ATPs only need a fraction of a second.

SET646+3 Sometimes, problems are immediately refuted after the initial clause normalisation. For example, after definition expansion in problem SET646+3 we get the following clause (where $B_{o\alpha}, C_{o\alpha}, x_{\alpha}$ are again Skolem constants):

$$[Ax \Rightarrow (\forall y_{\beta\bullet}By \Rightarrow (\forall u_{\alpha\bullet}\forall v_{\beta\bullet}(u = x \wedge v = y) \Rightarrow ((\neg\perp) \wedge (\neg\perp))))]^F$$

Normalisation in LEO immediately generates a basic refutation (i.e., a clause $[\perp]^T \vee [\perp]^T$) without even starting proof search.

SET611+3 The examples discussed so far all essentially apply extensionality treatment and normalisation to the input problem in order to immediately generate a set of inconsistent FO-like clauses. Problem SET611+3 is more complicated as it requires several reasoning steps in LEO before the initially consistent set of available FO-like clauses grows into an inconsistent one. After definition expansion, LEO is first given the input clause:

$$[\forall A_{o\alpha}, B_{o\alpha\bullet}(\lambda x_{\alpha\bullet}(Ax \wedge Bx)) = (\lambda x_{\alpha\bullet}\perp) \Leftrightarrow (\lambda x_{\alpha\bullet}(Ax \wedge \neg Bx)) = (\lambda x_{\alpha\bullet}Ax)]^F$$

which it normalises into:

$$[(\lambda x_{\alpha\bullet}(Ax \wedge Bx)) =^? (\lambda x_{\alpha\bullet}\perp)] \vee [(\lambda x_{\alpha\bullet}(Ax \wedge \neg Bx)) =^? (\lambda x_{\alpha\bullet}Ax)] \quad (9)$$

$$[(\lambda x_{\alpha\bullet}(Ax \wedge Bx)) = (\lambda x_{\alpha\bullet}\perp)]^T \vee [(\lambda x_{\alpha\bullet}(Ax \wedge \neg Bx)) = (\lambda x_{\alpha\bullet}Ax)]^T \quad (10)$$

As mentioned before, the unification constraint (9) corresponds to:

$$[(\lambda x_{\alpha\bullet}(Ax \wedge Bx)) = (\lambda x_{\alpha\bullet}\perp)]^F \vee [(\lambda x_{\alpha\bullet}(Ax \wedge \neg Bx)) = (\lambda x_{\alpha\bullet}Ax)]^F \quad (11)$$

LEO has to apply to each of these clauses and to each of their literals appropriate extensionality rules. Thus, several rounds of LEO’s set-of-support-based reasoning procedure are required, so that all necessary extensionality reasoning steps are performed, and sufficiently many FO-like clauses are generated which can be refuted by BLIKSEM.

In summary, each of the examples discussed in this section exposes a motivation for our higher-order/first-order cooperative approach to theorem proving. In particular, they show that:

- Higher-order formulations allow for a concise problem representation which often allows easier and faster proof search than first-order formulations.
- Higher-order problems can often be reduced to a set of first-order clauses that can be more efficiently handled by a first-order ATP.
- Some problems are trivially refutable after clause normalisation.
- Some problems require in-depth higher-order reasoning before a refutable first-order clause set can be extracted.

3 Higher-Order/First-Order Cooperation via OANTS

The cooperation between higher-order and first-order reasoners, which we investigate in this paper, is realised in the concurrent hierarchical blackboard architecture OANTS [7]. We first describe in Sec. 3.1 the existing OANTS architecture. In order to overcome some of its problems, in particular efficiency problems, we devised within OANTS a new and improved cooperation method for the higher-order ATP LEO and first-order provers (in particular, BLIKSEM) – we describe this in Sec. 3.2. We address the question of how to generate the necessary clauses in Sec. 3.3, and discuss soundness and completeness of our implementation of the higher-order/first-order cooperation in Sec. 3.4.

3.1 OANTS

OANTS was originally conceived to support interactive theorem proving but was later extended to a fully automated proving system [22, 8]. Its basic idea is to compose a *central proof object* by generating, in each proof situation, a ranked list of potentially applicable inference steps. In this process, all inference rules, such as calculus rules or tactics, are uniformly viewed with respect to three sets: premises, conclusions, and additional parameters. The elements of these three sets are called *arguments* of the inference rule and they usually depend on each other. An inference rule is applicable if at least some of its arguments can be instantiated with respect to the given proof context. The task of the OANTS architecture is now to determine the applicability of inference rules by computing instantiations for their arguments.

The architecture consists of two layers. On the lower layer, possible instantiations of the arguments of individual inference rules are computed. In particular, each inference rule is associated with its own blackboard and concurrent processes, one for each argument of the inference rule. The role of every process is to compute possible instantiations for its designated argument of the inference rule, and to record these on the blackboard. The computations are carried out with respect to the given proof context and by exploiting information already present on the blackboard, that is, argument instantiations computed by other

processes. On the upper layer, the information from the lower layer is used for computing and heuristically ranking the inference rules that are applicable in the current proof state. The most promising rule is then applied to the central proof object and the data on the blackboards is cleared for the next round of computations.

OANTS employs resource reasoning to guide search.² This enables the controlled integration (e.g., by specifying time-outs) of full-fledged external reasoning systems such as automated theorem provers, computer algebra systems, or model generators into the architecture. The use of the external systems is modelled by inference rules, usually one for each system. Their corresponding computations are encapsulated in one of the independent processes in the architecture. For example, an inference rule modelling the application of an ATP has its conclusion argument set to be an open goal. A process can then place an open goal on the blackboard, where it is picked up by a process that applies the prover to it. Any computed proof or partial-proof from the external system is again written to the blackboard from where it is subsequently inserted into the proof object when the inference rule is applied. While this setup enables proof construction by a collaborative effort of diverse reasoning systems, the cooperation can only be achieved via the central proof object. This means that all partial results have to be translated back and forth between the syntaxes of the integrated systems and the language of the proof object. Since there are many types of integrated systems, the language of the proof object — a higher-order language even richer than LEO's, together with a natural deduction calculus — is expressive but also cumbersome. This leads not only to a large communication overhead, but also means that complex proof objects have to be created (large clause sets need to be transformed into large single formulae to represent them in the proof object; the support for this in OANTS to date is inefficient), even if the reasoning of all systems involved is clause-based. Consequently, the cooperation between external systems is typically rather inefficient [5].

3.2 Cooperation via a single inference rule

In order to overcome the problem of the communication bottleneck described above, we devised a new method for the cooperation between a higher-order and a first-order theorem prover within OANTS. Rather than modelling each theorem prover as a separate inference rule (and hence needing to translate the communication via the language of the central proof object), we model the cooperation between a higher-order (concretely, LEO) and a first-order theorem prover (in our case study BLIKSEM) in OANTS as a single inference rule. The cooperation between these two theorem provers is carried out directly and not via the central proof object. This avoids translating clause sets into single formulae and back. While in our previous approach the cooperation between LEO and an FO-ATP was modelled at the upper layer of the OANTS architecture, our

² OANTS provides facilities to define and modify the processes at run-time. But notice that we do not use these advanced features in the case study presented in this paper.

new approach presented in this paper models their cooperation by exploiting the lower layer of the OANTS blackboard architecture. This is not an ad hoc solution, but rather, it demonstrates OANTS's flexibility in modelling the integration of cooperative reasoning systems.

Concretely, the single inference rule modelling the cooperation between LEO and a first-order theorem prover needs four arguments to be applicable: (1) an open proof goal, (2) a partial LEO proof, (3) a set of FO-like clauses in the partial proof, (4) a first-order refutation proof for the set of FO-like clauses. Each of these arguments is computed, that is, its instantiation is found, by an independent process. The first process finds open goals in the central proof object and posts them on the blackboard associated with the new rule. The second process starts an instance of the LEO theorem prover for each new open goal on the blackboard. Each LEO instance maintains its own set of FO-like clauses. The third process monitors these clauses, and as soon as it detects a change in this set, that is, if new FO-like clauses are added by LEO, it writes the entire set of clauses to the blackboard. Once FO-like clauses are posted, the fourth process first translates each of the clauses directly into a corresponding one in the format of the first-order theorem prover, and then starts the first-order theorem prover on them. Note that writing FO-like clauses on the blackboard is by far not as time consuming as generating higher-order proof objects. As soon as either LEO or the first-order prover finds a refutation, the second process reports LEO's proof or partial proof to the blackboard, that is, it instantiates argument (2). Once all four arguments of our inference rule are instantiated, the rule can be applied and the open proof goal can be closed in the central proof object. That is, the open goal can be proved by the cooperation between LEO and a first-order theorem prover. When computing applicability of the inference rule, the second and the fourth process concurrently spawn processes running LEO or a first-order prover on a different set of FO-like clauses. Thus, when actually applying the inference rule, all these instances of provers working on the same open subgoal are stopped.

The cooperation can be carried out between any first-order theorem prover and LEO instantiated with any strategy, thus resulting in different instantiations of the inference rule discussed above. While several first-order provers are integrated in OANTS and could be used, BLIKSEM was sufficient for the case study reported in this paper (see Sec. 4). In most cases, more than one BLIKSEM process was necessary. But as the problems were always concerned with only one subgoal, only one LEO process had to be started.

Our approach to the cooperation between a higher-order and a first-order theorem prover has many advantages. The main one is that the communication is restricted to the transmission of clauses, and thus it avoids intermediate translation into the language of the central proof object. This significantly reduces the communication overhead and makes effective proving of more involved theorems feasible. A disadvantage of this approach is that we cannot easily translate and integrate the two proof objects produced by LEO and BLIKSEM into the central proof object maintained by OANTS, as is possible when applying only

one prover per open subgoal. Providing such translation remains future work. The repercussions will be discussed in more detail in Sec. 3.4.

3.3 Extracting FO-like clauses from LEO

Crucial to a successful cooperation between LEO and a first-order ATP is obviously the generation of FO-like clauses. LEO always maintains a heap of FO-like clauses. In the current LEO system this heap remains rather small since LEO's standard calculus intrinsically avoids primitive equality and instead provides a rule that replaces occurrences of primitive equality with their corresponding Leibniz definitions which are higher-order. The Leibniz principle defines equality as follows $=_{o\alpha\alpha} := \lambda x_{\alpha}. \lambda y_{\alpha}. [\forall P_{o\alpha}. Px \Rightarrow Py]$. LEO also provides a rule which replaces syntactically non-unifiable unification constraints between terms of non-Boolean base type by their respective representations that use Leibniz equality. While the clauses resulting from these rules are still refutable in LEO, they are not refutable by BLIKSEM without adding set theory axioms. We illustrate the effect by the following simple example, where a_i , b_i , and f_i are constants:

$$a = b \Rightarrow f(a) = f(b)$$

Depending on whether we work with primitive equality or Leibniz equality this problem is reduced to the clause sets in either (12) or (13) respectively (in the latter \mathbf{P}_{o_i} is a new free variable, and Q_{o_i} is a new Skolem constant):

$$[\mathbf{P}a]^F \vee [\mathbf{P}b]^T \quad [a = b]^T \quad [f(a) =? f(b)] \quad (12)$$

$$[\mathbf{P}a]^F \vee [\mathbf{P}b]^T \quad [Q(f(a))]^T \quad [Q(f(b))]^F \quad (13)$$

While the former is obviously refutable in BLIKSEM, the latter is not. LEO, however, still finds a refutation for the latter and generates the crucial substitution $\mathbf{P} \leftarrow \lambda x_{\alpha}. Q(f(x))$ by higher-order pre-unification.

To circumvent this problem, we adapted the relevant rules in LEO. Instead of immediately constructing Leibniz representation of clauses, an intermediate representation containing primitive equality is generated and dumped on the heap of FO-like clauses. As a consequence, additional useful FO-like clauses are accumulated and the heap can become quite large, in particular, since we do not apply any subsumption to the set of FO-like clauses (this is generally done more efficiently by a first-order ATP anyway). Recent research has shown that Leibniz equality is generally very bad for automating higher-order proof search. Thus, future work in LEO includes providing support for full primitive equality and avoiding Leibniz equations.

3.4 Soundness and completeness of the cooperation

Clearly, soundness and completeness properties depend on the corresponding properties of the systems involved, in our case, of LEO and BLIKSEM.

Soundness: The general philosophy of OANTS is to ensure the correctness of proofs by the generation of explicit proof objects, which can be checked independently from the proof generation. In particular, reasoning steps of ATPs have to be translated into OANTS’s natural deduction calculus via the TRAMP proof transformation system [17] to be machine-checkable. Since the cooperative proof result of LEO-BLIKSEM cannot yet be directly inserted into the centralised proof object, the generation of a machine-checkable proof object is not yet supported. One possible solution is to insert BLIKSEM proofs into LEO proofs at the right places. Then, the modified LEO proofs can be inserted into the centralised proof object, and hence, explicit proof objects can be generated by OANTS. In principle, there is no problem with this, however, it is not yet implemented.

While there are many advantages in guaranteeing correctness of proofs by checking them, it is worth noting that the combination of LEO and BLIKSEM is sound under the assumption that the two systems are sound. Namely, to prove a theorem it is sufficient to show that a subset of clauses generated in the proof is inconsistent. If LEO generates an inconsistent set of clauses, then it does so correctly by assumption, be it a FO-like set or not. Assuming that the translation from FO-like clauses to truly first-order clauses preserves consistency/inconsistency, then a set of clauses that is given to BLIKSEM is inconsistent only if LEO generated an inconsistent set of clauses in the first place. By the assumption that BLIKSEM is sound follows that BLIKSEM will only generate the empty clause when the original clause set was inconsistent.

Thus, soundness of our cooperative approach critically relies only on the soundness of the selected transformational mapping from FO-like clauses to proper first-order clauses. We use the mapping from TRAMP, which has been previously shown to be sound and is based on [16]. Essentially, it injectively maps expressions such as $P(f(a))$ to expressions such as $@_{\text{pred}}^1(P, @_{\text{fun}}^1(f, a))$, where the @ are new first-order operators describing function and predicate application for particular types and arities. The injectivity of the mapping guarantees soundness, since it allows each proof step to be mapped back from first-order to higher-order. Hence, our higher-order/first-order cooperative approach between LEO and BLIKSEM is sound.

Completeness: Completeness (in the sense of Henkin completeness) can in principle be achieved in higher-order systems, but practically, the strategies used are typically not complete for efficiency reasons. Let us assume that we use a complete strategy in LEO. All that our procedure does is pass FO-like clauses to BLIKSEM. Hence, no proofs can be lost in this process. That is, completeness follows trivially from the completeness of LEO.

The more interesting question is whether particular cooperation strategies will be complete as well. For instance, in LEO we may want to give higher preference to real higher-order steps which guarantee the generation of first-order clauses.

4 Experiments and Results

We conducted several experiments to evaluate our hybrid reasoning approach. In particular, we concentrated on problems given in Table 1. We investigated several LEO strategies in order to compare LEO's individual performance with the performance of the LEO-BLIKSEM cooperation. Our example set differs from the one in [14] in that it contains some additional problems, and it also omits an entry for problem SET108+1. This problem addresses the universal class and can therefore not be formalised in type theory in the same concise way as the other examples, but only in a way very similar to the one given in TPTP.

Table 4 presents the results of our experiments. All timings given in the table are in seconds. The first column contains the TPTP identifier of the problem. The second column relates some of the problems to their counterparts in the Journal of Formalized Mathematics (JFM; see mizar.org/JFM) where they originally stem from. This eases the comparison with the results in [6, 2], where the problems from the JFM article *Boolean Properties of Sets* were already solved: the problems are named with prefix 'B:'. Prefix 'RS1:' stands for the JFM article *Relations Defined on Sets*. The third column lists the TPTP (v3.0.1 as of 20 January 2005, see <http://www.tptp.org>) difficulty rating of the problem, which indicates how hard the problem is for first-order ATPs (difficulty rating 1.00 indicates that no TPTP prover can solve the problem).

The fourth, fifth and sixth columns list whether SATURATE, MUSCADET (v2.4) and E-SETHEO (csp04), respectively, can (+) or cannot (-) solve a problem. The seventh column lists the timing results for VAMPIRE (v7). The results for SATURATE are taken from [14] (a '?' in Table 4 indicates that the result was not listed in [14] and is thus unavailable). The results for MUSCADET and E-SETHEO are taken from the on-line version of the solutions provided with the TPTP. Since the listed results were obtained from different experiments on different platforms, their run-time comparison would be unfair, and was thus not carried out. The timings for VAMPIRE, on the other hand, are based on private communication with A. Voronkov and they were obtained on a computer with a very similar specification as we used for the LEO-BLIKSEM timings. Note, that the results for VAMPIRE and E-SETHEO reported in [14] differ for some of the problems to the ones in TPTP. This is probably due to different versions of the systems tested, for instance, the TPTP uses VAMPIRE version 7, while the results reported in [14] are based on version 5. The results in columns four through to seven show that some problems are still very hard for first-order ATPs, as well as for the special purpose theorem prover MUSCADET. Column eight and nine in Table 4 list the results for LEO alone and LEO-BLIKSEM, respectively. Each of these two columns is further divided into sub-columns to allow for a detailed comparison.

All our experiments (for the values of LEO and LEO-BLIKSEM) were conducted on a 2.4 GHz Xenon machine with 1GB of memory and an overall time limit of 100 seconds. For our experiments with LEO alone in column eight in Table 4 we tested four different strategies. Mainly, they differ in their treatment of equality and extensionality. This ranges from immediate expansion of

| TPTP-Problem | Mizar Problem | Diff-culty | Saturate | Muscadet | E-Se-theo | Vamp-ire 7 | LEO | | | LEO-BLIKSEM | | | | |
|--------------|---------------|------------|----------|----------|-----------|------------|--------|------|-------|-------------|-------|------|------|------|
| | | | | | | | Strat. | Cl. | Time | Cl. | Time | FOcl | FOtm | GnCl |
| SET014+4 | | .67 | + | + | + | .01 | ST | 41 | .16 | 34 | 6.76 | 19 | .01 | 7 |
| SET017+1 | | .56 | - | - | + | .03 | EXT | 3906 | 57.52 | 25 | 8.54 | 16 | .01 | 74 |
| SET066+1 | | 1.00 | ? | - | - | - | - | - | - | 26 | 6.80 | 20 | 10 | 56 |
| SET067+1 | | .56 | + | + | + | .04 | ST | 6 | .02 | 13 | .32 | 16 | .01 | 12 |
| SET076+1 | | .67 | + | - | + | .00 | - | - | - | 10 | .47 | 18 | .01 | 35 |
| SET086+1 | | .22 | + | - | + | .04 | ST | 4 | .01 | 4 | .01 | N/A | N/A | N/A |
| SET096+1 | | .56 | + | - | + | .03 | - | - | - | 27 | 7.99 | 14 | .01 | 25 |
| SET143+3 | B:67 | .67 | + | + | + | 68.71 | EIR | 37 | .38 | 33 | 7.93 | 18 | .01 | 19 |
| SET171+3 | B:71 | .67 | + | + | - | 108.31 | EIR | 36 | .56 | 25 | 4.75 | 19 | .01 | 20 |
| SET580+3 | B:23 | .44 | + | + | + | 14.71 | EIR | 25 | .19 | 6 | 2.73 | 8 | .01 | 13 |
| SET601+3 | B:72 | .22 | + | + | + | 168.40 | EIR | 145 | 2.20 | 55 | 4.96 | 8 | .01 | 13 |
| SET606+3 | B:77 | .78 | + | - | + | 62.02 | EIR | 21 | .33 | 17 | 10.8 | 15 | .01 | 5 |
| SET607+3 | B:79 | .67 | + | + | + | 65.57 | EIR | 22 | .31 | 17 | 7.79 | 15 | .01 | 6 |
| SET609+3 | B:81 | .89 | + | + | - | 161.78 | EIR | 37 | .60 | 26 | 6.50 | 19 | 10 | 17 |
| SET611+3 | B:84 | .44 | + | - | + | 60.20 | EIR | 996 | 12.69 | 72 | 32.14 | 38 | .01 | 101 |
| SET612+3 | B:85 | .89 | + | - | - | 113.33 | EIR | 41 | .54 | 18 | 3.95 | 6 | .01 | 7 |
| SET614+3 | B:88 | .67 | + | + | - | 157.88 | EIR | 38 | .46 | 19 | 4.34 | 16 | .01 | 17 |
| SET615+3 | B:89 | .67 | + | + | - | 109.01 | EIR | 38 | .57 | 17 | 3.59 | 6 | .01 | 9 |
| SET623+3 | B:99 | 1.00 | ? | - | - | - | EXT | 43 | 8.84 | 23 | 9.54 | 10 | .01 | 14 |
| SET624+3 | B:100 | .67 | + | - | + | .04 | ST | 4942 | 34.71 | 54 | 9.61 | 46 | .01 | 212 |
| SET630+3 | B:112 | .44 | + | - | + | 60.39 | EIR | 11 | .07 | 6 | .08 | 8 | 10 | 4 |
| SET640+3 | RS1:2 | .22 | + | - | + | 70.41 | EIR | 2 | .01 | 2 | .01 | N/A | N/A | N/A |
| SET646+3 | RS1:8 | .56 | + | - | + | 59.63 | EIR | 2 | .01 | 2 | .01 | N/A | N/A | N/A |
| SET647+3 | RS1:9 | .56 | + | - | + | 64.21 | EIR | 26 | .15 | 13 | .30 | 13 | .01 | 15 |
| SET648+3 | RS1:10 | .56 | + | - | + | 64.22 | EIR | 26 | .15 | 14 | .30 | 13 | .01 | 16 |
| SET649+3 | RS1:11 | .33 | - | - | + | 63.77 | EIR | 45 | .30 | 29 | 5.49 | 12 | .01 | 16 |
| SET651+3 | RS1:13 | .44 | - | - | + | 63.88 | EIR | 20 | .10 | 11 | .16 | 10 | 10 | 11 |
| SET657+3 | RS1:19 | .67 | + | - | + | 1.44 | EIR | 2 | .01 | 2 | .01 | N/A | N/A | N/A |
| SET669+3 | RS1:19 | .22 | - | - | + | .34 | EI | 35 | .22 | 35 | .23 | N/A | N/A | N/A |
| SET670+3 | RS1:33 | 1.00 | ? | - | - | - | EXT | 15 | .17 | 17 | .36 | 16 | .01 | 6 |
| SET671+3 | RS1:34 | .78 | - | - | + | 218.02 | EIR | 78 | .64 | 7 | 2.71 | 10 | .01 | 14 |
| SET672+3 | RS1:35 | 1.00 | ? | - | - | - | EXT | 27 | .4 | 30 | .70 | 21 | .01 | 11 |
| SET673+3 | RS1:36 | .78 | - | - | + | 47.86 | EIR | 78 | .65 | 14 | 5.66 | 14 | .01 | 16 |
| SET680+3 | RS1:47 | .33 | + | - | + | .07 | ST | 185 | .88 | 29 | 4.61 | 18 | .01 | 24 |
| SET683+3 | RS1:50 | .22 | + | - | + | .06 | ST | 46 | .20 | 35 | 8.90 | 18 | 10 | 24 |
| SET684+3 | RS1:51 | .78 | - | - | + | .33 | ST | 275 | 2.45 | 46 | 5.95 | 26 | .01 | 47 |
| SET686+3 | RS1:53 | .56 | - | - | + | .11 | ST | 274 | 2.36 | 46 | 5.37 | 26 | .01 | 46 |
| SET716+4 | | .89 | + | + | - | - | ST | 39 | .45 | 18 | 3.81 | 18 | .01 | 118 |
| SET724+4 | | .89 | + | + | - | - | EXT | 154 | 2.75 | 18 | 7.21 | 15 | 10 | 23 |
| SET741+4 | | 1.00 | ? | - | - | - | - | - | - | - | - | - | - | - |
| SET747+4 | | .89 | - | + | - | - | ST | 34 | .46 | 25 | 1.11 | 18 | 10 | 10 |
| SET752+4 | | .89 | ? | + | - | - | - | - | - | 50 | 6.60 | 48 | .01 | 4363 |
| SET753+4 | | .89 | ? | + | - | - | - | - | - | 15 | 3.07 | 12 | 10 | 19 |
| SET764+4 | | .56 | + | + | + | .02 | EI | 9 | .05 | 8 | .04 | N/A | N/A | N/A |
| SET770+4 | | .89 | + | + | - | - | - | - | - | - | - | - | - | - |

Table 4. Experimental data for the benchmark problems given in Table 1.

primitive equality with Leibniz equality and limited extensionality reasoning, STANDARD (ST), to immediate expansion of primitive equality and moderate extensionality reasoning, EXT, to delayed expansion of primitive equality and moderate extensionality reasoning, EXT-INPUT (EI), and finally to delayed expansion of primitive equality and advanced recursive extensionality reasoning, EXT-INPUT-RECURSIVE (EIR). Column eight in Table 4 presents the fastest strategy for a respective problem (Strat.), the number of clauses generated by LEO (Cl.), and the total runtime (Time). While occasionally there were more than one LEO strategy that could solve a problem, it should be noted that none of the strategies was successful for all the problems solved by LEO.

In contrast to the experiments with LEO alone, we used only the EXT-INPUT strategy for our experiments with the LEO-BLIKSEM cooperation. Column nine in

Table 4 presents the number of clauses generated by LEO (Cl.) together with the time (Time), and in addition, the number of first-order clauses sent to BLIKSEM (FOcl), the time used by BLIKSEM (FOtm), and the number of clauses generated by BLIKSEM (GnCl). Note, that we give the data only for the first instance that BLIKSEM actually succeeded in solving the problem. This time also includes the time needed to write and process input and output files over the network. While LEO and instances of BLIKSEM were running in separate threads (each run of BLIKSEM was given a 50 second time limit), the figures given in the ‘Time’ column reflect the overall time needed for a successful proof. That is, it contains the time needed by all concurrent processes: LEO’s own process as well as those processes administering the various instances of BLIKSEM. Since all these processes ran on a single processor, there is potential to ameliorate the overall runtimes by using real multiprocessing.

Note also, that the number of clauses in LEO’s search space is typically low since subsumption is enabled. Subsumption, however, was not enabled for the accumulation of FO-like clauses in LEO’s bag of FO-like clauses. This is why there are usually more clauses in this bag (which is sent to BLIKSEM) than there are available in LEO’s search space. Finally, observe that for some problems a refutation was found after LEO’s clausal normalisation, and therefore BLIKSEM was not applicable (N/A).

While LEO itself can solve a majority of the considered problems with some strategy, the LEO-BLIKSEM cooperation can solve more problems and, moreover, needs only a single LEO strategy. We can also observe that for many problems that appear to be relatively hard for LEO alone (e.g., SET017+1, SET611+3, SET624+3), the LEO-BLIKSEM cooperation solves them not only more quickly, but also it sometimes reduces the problems to relatively small higher-order pre-processing steps with subsequent easy first-order proofs, as for instance, in the case of SET017+1.

From a mathematical viewpoint the investigated problems are trivial and, hence, they should ideally be reliably and very efficiently solvable within a proof assistant. This has been achieved for the examples in Table 4 (except for SET741+4 and SET770+4) by our hybrid approach. While some of the proof attempts now require slightly more time than when using LEO alone with a specialised strategy, they are, in most cases, still faster than when proving with a first-order system.

5 Related Work and Conclusion

Related to our approach is the TECHS system [12], which realises a cooperation between a set of heterogeneous first-order theorem provers. Similarly to our approach, partial results in TECHS are exchanged between the different theorem provers in form of clauses. The main difference to the work of Denzinger *et al.* (and other related architectures like [13]) is that our system bridges between higher-order and first-order automated theorem proving. Also, unlike in TECHS, we provide a declarative specification framework for modelling exter-

nal systems as cooperating, concurrent processes that can be (re-)configured at run-time. Related is also the work of Hurd [15] which realises a generic interface between HOL and first-order theorem provers. It is similar to the solution previously achieved by TRAMP [17] in OMEGA, which serves as a basis for the sound integration of ATPs into OANTS. Both approaches pass essentially first-order clauses to first-order theorem provers and then translate their results back into HOL resp. OMEGA. Some further related work on the cooperation of Isabelle with VAMPIRE is presented in [18]. The main difference of our work to the related systems is that while our system calls first-order provers from within higher-order proof search, this is not the case for [15, 17, 18].

One of the motivations for our work is to show that the cooperation of higher-order and first-order automated theorem provers can be very successful and effective. The results of our case study provide evidence for this: our non-optimised system outperforms related work on state-of-the-art first-order theorem provers and their ad hoc extensions such as SATURATE [14] on 45 mathematical problems chosen from the TPTP SET category. Among them are four problems which cannot be solved by any TPTP system to date. In contrast to the first-order situation, these problems can in fact be proved in our approach reliably from first principles, that is, without avoiding relevant base axioms of the underlying set theory, and moreover, without the need to provide *relevant* lemmata and definitions by hand.

The results of our case study motivate further research in the automation of higher-order theorem proving and the experimentation with different higher-order to first-order transformation mappings (such as the ones used by Hurd) that support our hybrid reasoning approach. They also provide further evidence for the usefulness of the OANTS approach as described in [8, 5] for flexibly modelling the cooperation of reasoning systems.

Our results also motivate the need for a higher-order extension of the TPTP library in which alternative higher-order problem formalisations are linked with their first-order counterparts so that first-order theorem provers could also be evaluated against higher-order systems (and vice versa).

Future work is to investigate how far our approach scales up to more complex problems and more advanced mathematical theories. In less trivial settings as discussed in this paper, we will face the problem of selecting and adding relevant lemmata to avoid immediate reduction to first principles and to appropriately instantiate set variables. Relevant related work for this setting is Bishop's approach to *selectively expand definitions* as presented in [9] and Brown's PhD thesis on *set comprehension in Church's type theory* [10].

Acknowledgements For advice and help we thank Chad Brown, Andreas Meier, Andrei Voronkov, and Claus-Peter Wirth.

References

1. P. Andrews. *An Introduction to mathematical logic and Type Theory: To Truth through Proof*. Number 27 in Applied Logic Series. Kluwer, 2002.

2. C. Benzmüller. *Equality and Extensionality in Higher-Order Theorem Proving*. PhD thesis, Universität des Saarlandes, Germany, 1999.
3. C. Benzmüller. Extensional higher-order paramodulation and RUE-resolution. *Proc. of CADE-16, LNAI 1632*, p. 399–413. Springer, 1999.
4. C. Benzmüller. Comparing approaches to resolution based higher-order theorem proving. *Synthese*, 133(1-2):203–235, 2002.
5. C. Benzmüller, M. Jamnik, M. Kerber, and V. Sorge. Experiments with an Agent-Oriented Reasoning System. *Proc. of KI 2001, LNAI 2174*, p.409--424. Springer, 2001.
6. C. Benzmüller and M. Kohlhase. LEO – a higher-order theorem prover. *Proc. of CADE-15, LNAI 1421*. Springer, 1998.
7. C. Benzmüller and V. Sorge. A Blackboard Architecture for Guiding Interactive Proofs. *Proc. of AIMS’98, LNAI 1480*, p. 102–114. Springer, 1998.
8. C. Benzmüller and V. Sorge. OANTS – An open approach at combining Interactive and Automated Theorem Proving. *Proc. of Calculemus-2000*. AK Peters, 2001.
9. M. Bishop and P. Andrews. Selectively instantiating definitions. *Proc. of CADE-15, LNAI 1421*. Springer, 1998.
10. C. E. Brown. *Set Comprehension in Church’s Type Theory*. PhD thesis, Dept. of Mathematical Sciences, Carnegie Mellon University, USA, 2004.
11. H. de Nivelle. *The Bliksem Theorem Prover, Version 1.12*. Max-Planck-Institut, Saarbrücken, Germany, 1999. <http://www.mpi-sb.mpg.de/~bliksem/manual.ps>.
12. J. Denzinger and D. Fuchs. Cooperation of Heterogeneous Provers. *Proc. IJCAI-16*, p. 10–15. Morgan Kaufmann, 1999.
13. M. Fisher and A. Ireland. Multi-agent proof-planning. CADE-15 Workshop “Using AI methods in Deduction”, 1998.
14. H. Ganzinger and J. Stuber. Superposition with equivalence reasoning and delayed clause normal form transformation. *Proc. of CADE-19, LNAI 2741*. Springer, 2003.
15. J. Hurd. An LCF-style interface between HOL and first-order logic. *Automated Deduction — CADE-18, LNAI 2392*, p. 134–138. Springer, 2002.
16. M. Kerber. *On the Representation of Mathematical Concepts and their Translation into First Order Logic*. PhD thesis, Universität Kaiserslautern, Germany, 1992.
17. A. Meier. TRAMP: Transformation of Machine-Found Proofs into Natural Deduction Proofs at the Assertion Level. *Proc. of CADE-17, LNAI 1831*. Springer, 2000.
18. J. Meng and L. C. Paulson. Experiments on supporting interactive proof using resolution. *Proc. of IJCAR 2004, LNCS 3097*, p. 372–384. Springer, 2004.
19. R. Nieuwenhuis, Th. Hillenbrand, A. Riazanov, and A. Voronkov. On the evaluation of indexing techniques for theorem proving. *Proc. of IJCAR-01, LNAI 2083*, p. 257–271. Springer, 2001.
20. D. Pastre. Muscadet2.3 : A knowledge-based theorem prover based on natural deduction. *Proc. of IJCAR-01, LNAI 2083*, p. 685–689. Springer, 2001.
21. A. Riazanov and A. Voronkov. Vampire 1.1 (system description). *Proc. of IJCAR-01, LNAI 2083*, p. 376–380. Springer, 2001.
22. V. Sorge. *OANTS: A Blackboard Architecture for the Integration of Reasoning Techniques into Proof Planning*. PhD thesis, Universität des Saarlandes, Germany, 2001.
23. G. Stenz and A. Wolf. E-SETHEO: An Automated³ Theorem Prover – System Abstract. *Proc. of the TABLEAUX’2000, LNAI 1847*, p. 436–440. Springer, 2000.
24. G. Sutcliffe and C. Suttner. The TPTP Problem Library: CNF Release v1.2.1. *Journal of Automated Reasoning*, 21(2):177–203, 1998.