

# Automated Theorem Proving for Prolog Verification

Fred Mesnard and Thierry Marianne

LIM, université de La Réunion, France

`frederic.mesnard@univ-reunion.fr` `thierry.marianne@univ-reunion.fr`

## Abstract

LPTP (Logic Program Theorem Prover) is an interactive natural-deduction-based theorem prover for pure Prolog programs with negation as failure, unification with the occurs check, and a restricted but extensible set of built-in predicates. With LPTP, one can formally prove termination and partial correctness of such Prolog programs. LPTP was designed in the mid 90's by Robert F. Stärk. It is written in ISO-Prolog and comes with an Emacs user-interface.

From a theoretical point of view, in his publications about LPTP, Stärk associates a set of first-order axioms  $\text{IND}(P)$  to the considered Prolog program  $P$ .  $\text{IND}(P)$  contains the Clark's equality theory for  $P$ , definitions of success, failure and termination for each user-defined logic procedure in  $P$ , axioms relating these three points of view, and an axiom schema for proving inductive properties. LPTP is thus a dedicated proof editor where these axioms are hard-wired.

We propose to explicit these axioms as first-order formulas, and apply automated theorem provers to check the property of interest. Using FOF as an intermediary language, we experiment the use of automated theorem provers for Prolog program verification. We evaluate the approach over about a benchmark of 400 properties of Prolog programs from the library available with LPTP. Both the logic compiler from a given input Prolog program  $P$  to the FOF version of  $\text{IND}(P)$  and the benchmark are publicly available.

## 1 Introduction

In the mid 90's, Robert F. Stärk defined a framework for Prolog verification in [23, 25]. He considered a subset of ISO-Prolog [11]: pure Prolog programs with negation as failure, unification with the occurs check, and allowed a restricted but extensible set of built-in predicates. He presented a first-order formalisation with axiom schemas of the usual operational semantics of Prolog. A safeness condition included in termination condition imposes groundness before evaluation of negated goals. He showed soundness and completeness for termination, success, and failure. The framework also allows partial correctness properties to be proved by induction wrt. the clauses defining predicates, considered as inductive definitions. The logical theory was hard-wired in an interactive dedicated natural-deduction-based theorem prover called LPTP (Logic Program Theorem Prover). Stärk implemented<sup>1</sup> LPTP in ISO-Prolog, together with an Emacs user-interface, an HTML and  $\text{\TeX}$  manager, a detailed user-manual, and a library of predicates for Peano numbers, lists, sorting algorithms, etc. with numerous proven properties.

Thirty years later, LPTP is still running on any ISO-Prolog processor, with its initial interface. Today, formal verification of computer programs is an established discipline within computer science. Nonetheless, program verification by interactive theorem proving is still a slow process and requires non-trivial skills. On the other hand, during the last three decades, the increase in computing power and the advances in automated theorem proving have been notable. For instance, the TPTP (Thousands of Problems for Theorem Provers) [26] is a library of test problems for automated theorem proving. It provides online tools to check syntax of

---

<sup>1</sup>available at <https://github.com/FredMesnard/lptp>

input problems and apply bunch of user selected automated theorem provers. Among them, E [22] and Vampire [14] are two powerful freely available automated theorem provers, performing very well in many international competitions over the years. Interactive theorem prover implementers starting with [20, 3] were able to take advantage of these progress by implementing so called *hammers* for their tools.

This evolution raises the following questions: can we also use the TPTP *Esperanto* to formulate the logic theory Stärk associates to a logic program? Can we use *off the shelf* TPTP provers and obtain automatic proofs in reasonable time? Can we get an acceptable success rate with such an approach?

The main contribution of this paper is the following. Using FOF (*first-order form*, one of the logic languages proposed by TPTP [27]) as an intermediary language, we describe the first – to the best of our knowledge – experiment of the use of automated theorem provers, namely E and Vampire, for Prolog program verification, including termination and partial correctness. We evaluate the approach over about 400 properties of Prolog programs. Both the compiler from Stärk’s theory applied to a given input Prolog program to FOF and the benchmark are publicly available<sup>2</sup>.

We organize the paper as follows. The next section presents a brief summary of the LPTP system. The third section describes step by step how to compile a Prolog program and its associated LPTP axioms to FOF. Then we present an experimental evaluation, related work and we conclude.

## 2 Notation

FOF (*First Order Form*) is a well-known logic language from TPTP for expressing first-order logic (FOL) axioms and conjectures. A formula is written `fof(name, role, formula)`, where *name* is the name of the formula, *role* is either `axiom` or `conjecture` and *formula* is informally defined as:

FOL	FOF	FOL	FOF
$A \wedge B$	<code>A &amp; B</code>	$\neg p(x)$	<code>~ p(X)</code>
$A \vee B$	<code>A   B</code>	$\exists x.A$	<code>?[X] : A</code>
$A \rightarrow B$	<code>A =&gt; B</code>	$\forall x.A$	<code>![X] : A</code>

Numerous examples will appear in the next section.

Let  $P$  be a pure logic program where negative literals may appear in the body of clauses (also called *normal program* in [17]). For sake of conciseness, we do not consider built-in predicates (see [25] for a full treatment) other than the equality `=/2`. We start with  $\mathcal{L}$ , the first-order language associated to  $P$ . The *goals* of  $\mathcal{L}$  are:

$$G, H ::= \text{true} \mid \text{fail} \mid s = t \mid A \mid \backslash +G \mid (G, H) \mid (G; H) \mid \text{some } x \, G$$

where  $s$  and  $t$  are two terms,  $x$  is a variable and  $A$  is an atomic goal. The goals of  $\mathcal{L}$  have the operational semantics specified by ISO-Prolog [11] with the occurs check.

$\hat{\mathcal{L}}$  is the specification language of LPTP. For each user-defined predicate symbol  $R$ ,  $\hat{\mathcal{L}}$  does not contain  $R$ , but instead it contains three predicate symbols  $R^s$ ,  $R^f$ ,  $R^t$  of the same arity as  $R$ , which express success, failure and termination of  $R$ .  $\hat{\mathcal{L}}$  also contains a unary constraint for groundness  $gr$ , expressing that its argument is ground. The *formulas* of  $\hat{\mathcal{L}}$  are:

$$\phi, \psi ::= \top \mid \perp \mid s = t \mid R(\vec{t}) \mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \phi \rightarrow \psi \mid \forall x\phi \mid \exists x\phi$$

<sup>2</sup>Somewhere on GitHub

where  $\vec{t}$  is a sequence of  $n$  terms and  $R$  denotes a  $n$ -ary predicate symbol of  $\hat{\mathcal{L}}$ . The semantics of  $\hat{\mathcal{L}}$  is the first-order predicate calculus of classical logic.

For any of the user-defined logic procedure  $R$  in a logic program  $P$ ,  $D_R^P(\vec{x})$  denotes its Clark's *if-and-only-if* completed definition [4, 17],

For defining the declarative semantics of logic programs, Stärk uses three syntactic operators **S**, **F** and **T** which map goals of  $\mathcal{L}$  into  $\hat{\mathcal{L}}$ -formulas. Intuitively, **SG** means  $G$  succeeds (any breadth-first evaluation of  $G$  succeeds), **FG** means  $G$  fails (the ISO-Prolog evaluation stops without any answer), and **TG** means  $G$  terminates (the ISO-Prolog evaluation produces a finite number of answers then stops). The definition of the operators follows:

$$\begin{array}{llll}
\mathbf{S}R(\vec{t}) := R^s(\vec{t}) & \mathbf{S} \text{ true} := \top & \mathbf{S} \text{ fail} := \perp & \mathbf{S}(s = t) := (s = t) \\
\mathbf{S} \backslash +G := \mathbf{F}G & \mathbf{S}(G, H) := \mathbf{S}G \wedge \mathbf{S}H & \mathbf{S}(G; H) := \mathbf{S}G \vee \mathbf{S}H & \mathbf{S}(\text{some } x \ G) := \exists x \mathbf{S}G \\
\\
\mathbf{F}R(\vec{t}) := R^f(\vec{t}) & \mathbf{F} \text{ true} := \perp & \mathbf{F} \text{ fail} := \top & \mathbf{F}(s = t) := \neg(s = t) \\
\mathbf{F} \backslash +G := \mathbf{S}G & \mathbf{F}(G, H) := \mathbf{F}G \vee \mathbf{F}H & \mathbf{F}(G; H) := \mathbf{F}G \wedge \mathbf{F}H & \mathbf{F}(\text{some } x \ G) := \forall x \mathbf{F}G \\
\\
\mathbf{T}R(\vec{t}) := R^t(\vec{t}) & \mathbf{T} \text{ true} := \top & & \\
\mathbf{T} \text{ fail} := \top & \mathbf{T}(s = t) := \top & & \\
\mathbf{T} \backslash +G := \mathbf{T}G \wedge gr(G) & \mathbf{T}(G, H) := \mathbf{T}G \wedge (\mathbf{F}G \vee \mathbf{T}H) & & \\
\mathbf{T}(G; H) := \mathbf{T}G \wedge \mathbf{T}H & \mathbf{T}(\text{some } x \ G) := \forall x \mathbf{T}G & & 
\end{array}$$

Finally, we add the definition of  $gr$ , which is a constraint of the specification language and is needed for defining  $\mathbf{T} \backslash +G$ :

$$\begin{array}{ll}
gr(\text{true}) := \top & gr((G, H)) := gr(G) \wedge gr(H) \\
gr(\text{fail}) := \top & gr((G; H)) := gr(G) \wedge gr(H) \\
gr(s = t) := gr(s) \wedge gr(t) & gr(\text{some } x \ G) := \exists x \ gr(G) \\
gr(R(t_1, \dots, t_n)) := gr(t_1) \wedge \dots \wedge gr(t_n) & gr(\backslash +G) := gr(G)
\end{array}$$

### 3 Compiling LPTP axioms to FOF

With LPTP, we prove properties of a logic program  $P$  wrt. its *inductive extension*  $\text{IND}(P)$  which includes Clark's completion [4] and induction along the definition of the predicates. Stärk shows that the inductive extension is always consistent and proves various correctness and completeness results wrt. the operational semantics of Prolog in [25]. The first-order theory  $\text{IND}(P)$  is defined by nine kinds of axioms which we describe now, along with their translation in FOF.

#### 3.1 First steps

Let us consider the following logic program ADD as our running example.

```

nat(0).
nat(s(X)) :- nat(X).
add(0, Y, Y).
add(s(X), Y, s(Z)) :- add(X, Y, Z).

```

We discuss the axioms proposed by Stärk and apply them to the ADD program.

The axioms of Clark's equality theory

1.  $f(x_1, \dots, x_n) = f(y_1, \dots, y_n) \rightarrow x_i = y_i$  [if  $f$  is  $n$ -ary and  $1 \leq i \leq n$ ]
2.  $f(x_1, \dots, x_n) \neq g(y_1, \dots, y_m)$  [if  $n \neq m$  or  $f \neq g$ ]
3.  $t \neq x$  [if  $x$  occurs in  $t$  and  $t \neq x$ ]

These first axioms specify that the universe is the set of trees built from the symbol functions extracted from the program under consideration. The third axiom forbids infinite rational trees. Note that it is an axiom schema, i.e., an infinite set of first order axioms. We will omit it and it can be a source of imprecision, but we stay sound. Here is the FOF version:

```
fof(id1,axiom,! [Xx4] : ! [Xx5] : (s(Xx4) = s(Xx5) => Xx4 = Xx5)).
fof(id2,axiom,! [Xx3] : ~ ('0' = s(Xx3))).
```

#### Axioms for *gr/1*

4.  $\text{gr}(c)$  [if  $c$  is a constant]
5.  $\text{gr}(x_1) \wedge \dots \wedge \text{gr}(x_m) \leftrightarrow \text{gr}(f(x_1, \dots, x_m))$  [ $f$  is  $m$ -ary]

Actually, LPTP deals with *non-ground* terms, as any ISO-Prolog processor does. LPTP offers a predefined predicate *gr/1* that we can consider as a constraint. This relation is useful for instance for dealing with negation as failure as LPTP only allows negation by failure for *ground* goals. Back to our example, here is the FOF version:

```
fof(id4,axiom,gr('0')).
fof(id5,axiom,! [Xx6] : (gr(Xx6) <=> gr(s(Xx6)))).
```

The ADD program contains two user-defined predicates, *add/3* and *nat/1*. LPTP considers each user-defined predicates through three points of view: finite failure, success and termination. So LPTP will create the following predicates: *add\_fails/3*, *add\_succeeds/3*, *add\_terminates/3*, and similarly for *nat/1*. These three viewpoints are linked with the following axioms, where  $R^s$  (resp.  $R^f$  and  $R^t$ ) denotes *R\_succeeds* (resp. *R\_fails* and *R\_terminates*).

#### Uniqueness axioms and totality axioms

6.  $\neg(R^s(\vec{x}) \wedge R^f(\vec{x}))$  [if  $R$  is a user-defined predicate]
7.  $R^t(\vec{x}) \rightarrow (R^s(\vec{x}) \vee R^f(\vec{x}))$  [if  $R$  is a user-defined predicate]

Axiom 6 says that for any tuple of (possibly non-ground) terms, we cannot have at the same time success and failure of  $R$ . Axiom 7 states that given termination, we have success or failure. Altogether, it means that for any tuple of terms  $\vec{x}$ , assuming termination, either  $R(\vec{x})$  succeeds or (exclusively)  $R(\vec{x})$  finitely fails. So for our example, we get:

```
fof(ida6,axiom,! [Xx7,Xx8,Xx9] :
  ~ ((add_succeeds(Xx7,Xx8,Xx9) & add_fails(Xx7,Xx8,Xx9)))).
fof(ida7,axiom,! [Xx7,Xx8,Xx9] :
  (add_terminates(Xx7,Xx8,Xx9) =>
    (add_succeeds(Xx7,Xx8,Xx9) | add_fails(Xx7,Xx8,Xx9)))).
fof(idn6,axiom,! [Xx10] :
  ~ ((nat_succeeds(Xx10) & nat_fails(Xx10)))).
fof(idn7,axiom,! [Xx10] :
  (nat_terminates(Xx10) =>
    (nat_succeeds(Xx10) | nat_fails(Xx10)))).
```

Fixed point axioms for user-defined predicates  $R$ 

$$8. R^s(\vec{x}) \leftrightarrow \mathbf{SD}_R^P(\vec{x}), R^f(\vec{x}) \leftrightarrow \mathbf{FD}_R^P(\vec{x}), R^t(\vec{x}) \leftrightarrow \mathbf{TD}_R^P(\vec{x})$$

We recall that  $D_R^P(\vec{x})$  denotes the definition of the completion [4] of the user-defined procedure  $R(\vec{x})$  in the logic program  $P$ . In the previous section, we saw how to apply the operator  $\mathbf{S}$ ,  $\mathbf{F}$ , and  $\mathbf{T}$  to formulas. So for instance, the first equivalence  $R^s(\vec{x}) \leftrightarrow \mathbf{SD}_R^P(\vec{x})$  defines  $R^s(\vec{x})$ . Back to our running example, we get:

```
fof(idns8,axiom,! [Xx1] : (nat_succeeds(Xx1) <=>
  (? [Xx2] : (Xx1 = s(Xx2) & nat_succeeds(Xx2)) | Xx1 = '0')))).
fof(idnf8,axiom,! [Xx1] : (nat_fails(Xx1) <=>
  (! [Xx2] : (~ (Xx1 = s(Xx2)) |
    nat_fails(Xx2)) & ~ (Xx1 = '0')))).
fof(idnt8,axiom,! [Xx1] : (nat_terminates(Xx1) <=>
  (! [Xx2] : ((~ (Xx1 = s(Xx2)) | nat_terminates(Xx2)))))).
```

and similarly for `add/3`.

Finally, for any property of the form  $\forall \vec{x}[R^s(\vec{x}) \rightarrow \phi(\vec{x})]$ , where  $R(\vec{x})$  is a user-defined procedure and  $\phi(\vec{x})$  an  $\hat{\mathcal{L}}$ -formula, we have an induction schema. The interactive prover LPTP is able to *dynamically* generate an induction axiom on demand while the user interacts with it. In our approach, we statically generate the induction axiom *once* from the form of the conjecture to be proved. This is a potential source of imprecision, but again we stay sound. Let us examine a simple case. It is exactly what happens using LPTP, which slightly generalizes [25]. By *directly recursive user-defined predicate* in the box below, we forbid mutual recursive definitions. Of course, LPTP is able to handle mutually recursive properties, see [23] for some examples.

A (simplified) induction schema for a user-defined predicate  $R$ 

Let  $R$  be a directly recursive user-defined predicate and let  $\phi(\vec{x})$  be an  $\hat{\mathcal{L}}$ -formula such that the length of  $\vec{x}$  is equal to the arity of  $R$ .

Let  $sub(\phi(\vec{x})/R)$  be the formula to be proven  $\forall \vec{x}(R^s(\vec{x}) \rightarrow \phi(\vec{x}))$ .

Let  $closed(\phi(\vec{x})/R)$  be the formula obtained from  $\forall \vec{x}(\mathbf{SD}_R^P(\vec{x}) \rightarrow R^s(\vec{x}))$  by replacing

- $R^s(\vec{x})$  by  $\phi(\vec{x})$  on the right of  $\rightarrow$ ,
- all occurrences of  $R(\vec{t})$  appearing on the left of  $\rightarrow$  by  $\phi(\vec{t}) \wedge R(\vec{t})$ .

Then the induction axiom is the following formula:

$$9. closed(\phi(\vec{x})/R) \rightarrow sub(\phi(\vec{x})/R)$$

Let us apply this axiom to the following property, informally stated as: for any term  $x$ , if  $\mathbf{nat}(x)$  then  $\mathbf{add}(x, 0, x)$ . Expressed in LPTP, it gives: for any term  $x$ , if  $\mathbf{nat\_succeeds}(x)$  then  $\mathbf{add\_succeeds}(x, 0, x)$ , which is exactly the formula  $sub(\phi(\vec{x})/R)$  of axiom 9. So  $R \equiv \mathbf{nat}$ ,  $R^s \equiv \mathbf{nat\_succeeds}$  and  $\phi(\vec{x}) \equiv \mathbf{add\_succeeds}(x, 0, x)$ .

For the left hand side of axiom 9, we start from

$$\forall x(\mathbf{SD}_{\mathbf{nat}}^{ADD}(x) \rightarrow \mathbf{nat\_succeeds}(x))$$

We have  $D_{nat}^{ADD}(x) \equiv x = 0 \vee \exists y(x = s(y) \wedge \text{nat}(y))$ . We replace  $\text{nat}(y)$  by  $\text{nat}(y) \wedge \text{add\_succeeds}(y, 0, y)$ . We replace  $\text{nat\_succeeds}(x)$  by  $\text{add\_succeeds}(x, 0, x)$ . We get:  $\forall x(\mathbf{S}[x = 0 \vee \exists y(x = s(y) \wedge \text{nat}(y) \wedge \text{add\_succeeds}(y, 0, y))] \rightarrow \text{add\_succeeds}(x, 0, x))$ . We apply **S** and obtain:  $\forall x([x = 0 \vee \exists y(x = s(y) \wedge \text{nat\_succeeds}(y) \wedge \text{add\_succeeds}(y, 0, y))] \rightarrow \text{add\_succeeds}(x, 0, x))$ .

Summarizing, in FOF, associated to the property to be proved:

```
fof(lemma,conjecture,
! [Xx] : (nat_succeeds(Xx) => add_succeeds(Xx,'0',Xx))).
```

we obtain the following induction axiom:

```
fof(induction,axiom,(
! [Xx] :
  ((? [Xx2] : (Xx = s(Xx2) & (nat_succeeds(Xx2)
                                & add_succeeds(Xx2,'0',Xx2)))
   | Xx = '0') => add_succeeds(Xx,'0',Xx))
=>
! [Xx] : (nat_succeeds(Xx) => add_succeeds(Xx,'0',Xx))).
```

We can gather all the 15 axioms, including the axioms defining `add_success/3`, `add_fails/3`, and `add_terminates/3` and the conjecture plus its induction axiom in a file, say `test.fof` and submit it<sup>3</sup> to the E prover or to Vampire. Both systems will find a refutation in a fraction of second<sup>4</sup> on a standard laptop.

It allows us to conclude for any term  $x$ , if  $\text{nat}(x)$  then  $\text{add}(x, 0, x)$  is true. Operationally, for any natural number  $n$ , in the Prolog search tree corresponding to the goal  $\text{add}(s^n(0), 0, s^n(0))$ , the empty clause appears. Assuming termination which can also be shown with LPTP, it means that the user will get (at least) one positive answer for the query  $:- \text{add}(s^n(0), 0, s^n(0))$ . when executed with any ISO-Prolog system.

Here's the manual proof of the same property in its LPTP version (a Prolog file), followed by its T<sub>E</sub>X version produced by LPTP. We began this proof by invoking the `ind` tactic, asking for an inductive proof. Both the base case and the inductive were generated and automatically completed by LPTP.

```
:- lemma(add:x_0_x,
all [x]: succeeds nat(?x) => succeeds add(?x,0,?x),
induction([all x: succeeds nat(?x) => succeeds add(?x,0,?x)],
[step([],[],[],succeeds add(0,0,0)),
step([x],
[succeeds add(?x,0,?x),
succeeds nat(?x)],
[],
succeeds add(s(?x),0,s(?x)))]))).
```

<sup>3</sup>E.g., in the terminal `eprover --auto test.fof` and `vampire test.fof`

<sup>4</sup>0.05s for Vampire, 0.017s for E on a MacBook Air, Apple M2, 16Go, macOS Sonoma.

**Lemma**  $[add:x\_0\_x] \forall x (\mathbf{S} \mathbf{nat}(x) \rightarrow \mathbf{S} \mathbf{add}(x, 0, x)).$

**Proof.**

Induction<sub>0</sub>:  $\forall x (\mathbf{S} \mathbf{nat}(x) \rightarrow \mathbf{S} \mathbf{add}(x, 0, x)).$

Hypothesis<sub>1</sub>: none.

Conclusion<sub>1</sub>:  $\mathbf{S} \mathbf{add}(0, 0, 0).$

Hypothesis<sub>1</sub>:  $\mathbf{S} \mathbf{add}(x, 0, x)$  and  $\mathbf{S} \mathbf{nat}(x).$

Conclusion<sub>1</sub>:  $\mathbf{S} \mathbf{add}(s(x), 0, s(x)). \quad \square$

### 3.2 A second property

Now let us consider the following property: for any  $x, y$  and  $z$  such that  $\mathbf{nat}(x)$ ,  $\mathbf{nat}(y)$  and  $\mathbf{add}(s(x), y, z)$ , we have  $\mathbf{add}(x, s(y), z).$

Let us consider the previous property as an axiom, which can be freely used by the automated prover. We have our new conjecture:

```
fof('lemma-(add:x_0_x)', axiom,
    ! [Xx] : (nat_succeeds(Xx) => add_succeeds(Xx, '0', Xx)).
```

```
fof('lemma-(add:succ)', conjecture,
    ! [Xx, Xy, Xz] : (((nat_succeeds(Xx) & nat_succeeds(Xy))
                        & add_succeeds(s(Xx), Xy, Xz))
                      => add_succeeds(Xx, s(Xy), Xz))).
```

In order to generate an induction axiom for this property, we first rewrite it in the form  $\forall \vec{x} [R^s(\vec{x}) \rightarrow \phi(\vec{x})]$  and we apply the simplified induction schema for user-defined predicate. It gives:

```
fof(induction, axiom, (
    ! [Xx] :
      ((? [Xy25] :
        (Xx = s(Xy25) & (nat_succeeds(Xy25)
          & ! [Xy, Xz] : ((add_succeeds(s(Xy25), Xy, Xz)
            & nat_succeeds(Xy))
          => add_succeeds(Xy25, s(Xy), Xz))))
      | Xx = '0') =>
      ! [Xy, Xz] : ((add_succeeds(s(Xx), Xy, Xz) & nat_succeeds(Xy))
        => add_succeeds(Xx, s(Xy), Xz)))
    => ! [Xx] : (nat_succeeds(Xx)
      => ! [Xy, Xz] : ((add_succeeds(s(Xx), Xy, Xz) & nat_succeeds(Xy))
        => add_succeeds(Xx, s(Xy), Xz)))).
```

Again, we can gather all axioms, the conjecture and its induction axiom in a file and submit it to Vampire, which will find a refutation in about one minute.

Here's a manual LPTP proof of the same property in its  $\text{\TeX}$  version. We began this proof by invoking the **ind** tactic, asking for an inductive proof. Both the base case and the inductive case were generated and manually completed. Clearly, the proof is more complex than the previous one.

**Lemma**  $[add:succ] \forall x, y, z (\mathbf{S}nat(x) \wedge \mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z) \rightarrow \mathbf{S}add(x, s(y), z)).$

**Proof.**

Induction<sub>0</sub>:  $\forall x (\mathbf{S}nat(x) \rightarrow \forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z) \rightarrow \mathbf{S}add(x, s(y), z))).$

Hypothesis<sub>1</sub>: none.

Assumption<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(0), y, z). \mathbf{DS}add(s(0), y, z)$  by completion.

$\exists z_1 (z = s(z_1) \wedge \mathbf{S}add(0, y, z_1)).$

Let<sub>3</sub>  $z_1$  with  $z = s(z_1) \wedge \mathbf{S}add(0, y, z_1). \mathbf{DS}add(0, y, z_1)$  by completion.  $y = z_1.$

$z = s(y). \mathbf{S}add(0, s(y), z).$

Thus<sub>3</sub>:  $\mathbf{S}add(0, s(y), z).$

Thus<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(0), y, z) \rightarrow \mathbf{S}add(0, s(y), z).$

Conclusion<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(s(0), y, z) \rightarrow \mathbf{S}add(0, s(y), z)).$

Hypothesis<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z) \rightarrow \mathbf{S}add(x, s(y), z))$  and  $\mathbf{S}nat(x).$

Assumption<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(s(x)), y, z). \mathbf{DS}add(s(s(x)), y, z)$  by completion.

$\exists z_2 (z = s(z_2) \wedge \mathbf{S}add(s(x), y, z_2)).$

Let<sub>3</sub>  $z_2$  with  $z = s(z_2) \wedge \mathbf{S}add(s(x), y, z_2). \mathbf{S}add(x, s(y), z_2). z = s(z_2).$

$\mathbf{S}add(s(x), s(y), z).$

Thus<sub>3</sub>:  $\mathbf{S}add(s(x), s(y), z).$

Thus<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(s(x)), y, z) \rightarrow \mathbf{S}add(s(x), s(y), z).$

Conclusion<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(s(s(x)), y, z) \rightarrow \mathbf{S}add(s(x), s(y), z)). \quad \square$

### 3.3 Commutativity of Peano addition

Finally let us consider commutativity of Peano addition: for any  $x, y, z$ , if  $\mathbf{add}(x, y, z)$  then  $\mathbf{add}(y, x, z)$ . Of course, stated this way, the property is false. We need to enforce that  $x$  and  $y$  are Peano numbers. So first we add our two previous properties as axioms. Here is our new conjecture, associated to its induction axiom:

```
fof('theorem-(add:commutative)', conjecture,
    ! [Xx,Xy,Xz] : (((nat_succeeds(Xx) & nat_succeeds(Xy))
                      & add_succeeds(Xx,Xy,Xz))
                    => add_succeeds(Xy,Xx,Xz))).

fof(induction, axiom,
    (! [Xx] :
      ((? [Xy26] : (Xx = s(Xy26) & (nat_succeeds(Xy26)
        & ! [Xy,Xz] : ((add_succeeds(Xy26,Xy,Xz) & nat_succeeds(Xy))
          => add_succeeds(Xy,Xy26,Xz))))
      | Xx = '0') =>
      ! [Xy,Xz] : ((add_succeeds(Xx,Xy,Xz) & nat_succeeds(Xy))
        => add_succeeds(Xy,Xx,Xz)))
    =>
    ! [Xx] : (nat_succeeds(Xx) =>
      ! [Xy,Xz] : ((add_succeeds(Xx,Xy,Xz) & nat_succeeds(Xy))
        => add_succeeds(Xy,Xx,Xz)))).
```

The conjecture is proved in less than one second by one<sup>5</sup> of the automated

---

<sup>5</sup>0.6s for Vampire.



provers. The LPTP proof explicitly uses the two previous lemmas:

**Theorem**  $[add:commutative] \forall x, y, z (\mathbf{S}nat(x) \wedge \mathbf{S}nat(y) \wedge \mathbf{S}add(x, y, z) \rightarrow \mathbf{S}add(y, x, z)).$

**Proof.**

Induction<sub>0</sub>:  $\forall x (\mathbf{S}nat(x) \rightarrow \forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(x, y, z) \rightarrow \mathbf{S}add(y, x, z))).$

Hypothesis<sub>1</sub>: none.

Assumption<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(0, y, z).$   $\mathbf{D} \mathbf{S}add(0, y, z)$  by completion.  $y = z.$

$\mathbf{S}add(y, 0, y)$  by Lemma  $add:x-0-x$   $[add:x_0x].$   $\mathbf{S}add(y, 0, z).$

Thus<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(0, y, z) \rightarrow \mathbf{S}add(y, 0, z).$

Conclusion<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(0, y, z) \rightarrow \mathbf{S}add(y, 0, z)).$

Hypothesis<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(x, y, z) \rightarrow \mathbf{S}add(y, x, z))$  and  $\mathbf{S}nat(x).$

Assumption<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z).$   $\mathbf{S}add(x, s(y), z)$  by Lemma  $add:succ$   $[add:succ].$

$\mathbf{S}add(s(y), x, z).$   $\mathbf{S}add(y, s(x), z)$  by Lemma  $add:succ$   $[add:succ].$

Thus<sub>2</sub>:  $\mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z) \rightarrow \mathbf{S}add(y, s(x), z).$

Conclusion<sub>1</sub>:  $\forall y, z (\mathbf{S}nat(y) \wedge \mathbf{S}add(s(x), y, z) \rightarrow \mathbf{S}add(y, s(x), z)).$   $\square$

## 4 Experimental Results

We applied the schema explained in the previous section to various libraries available with LPTP which we summarize now. The library **nat** defines some basic Peano relations with the expected properties. The library **gcd** defines a version of the greatest common divisor relation, with its full correctness proof. The library **list** proposes some elementary relations about lists with their properties. The library **suffix** defines two versions of the sublist relation, one as the prefix of a suffix, the other as the suffix of a prefix, and shows that the two versions are equivalent wrt. termination, success and failure. Similarly, the library **reverse** defines the two classical versions of the reverse relation, one with the append relation, the other with an accumulator and shows their equivalence. The library **permutation** defines the permutation relation with some properties useful for the correctness proofs of the sorting algorithms defined in the libraries **sort** and **mergesort**. The library **taut** defines a tautology checker for propositional formulas, together with its correctness proof [24].

We gather the results in table 1. The first column gives the name of library. The second column gives the number of lemma/corollary/theorem of the associated proof file. The remaining nine columns can be divided in three groups. On a MacBook Air, Apple M2, 16Go, macOS Sonoma, the first group gives the success rate for a 1 second timeout for the E prover (column E-1s), Vampire (column V-1s) and for the combination of the two provers (column EV-1s). The second group (resp. third group) gives the success rate for a timeout of 10 seconds (resp. 60 seconds).

<i>lib</i>	#	E-1s	V-1s	EV-1s	E-10s	V-10s	EV-10s	E-60s	V-60s	EV-60s
nat	91	54%	72%	78%	58%	77%	80%	61%	81%	85%
gcd	11	45%	45%	45%	45%	45%	45%	45%	45%	45%
list	84	56%	73%	83%	67%	87%	90%	68%	89%	92%
suffix	31	74%	81%	93%	81%	94%	97%	81%	97%	100%
reverse	25	52%	64%	64%	64%	80%	84%	68%	84%	88%
permut.	42	45%	50%	55%	52%	59%	64%	52%	64%	67%
sort	42	33%	33%	40%	43%	57%	57%	48%	59%	62%
merges.	24	50%	71%	71%	62%	79%	79%	67%	79%	79%
taut	43	0%	67%	67%	65%	70%	70%	65%	74%	74%

Table 1: Success rate

## 5 Related Work

There is quite a few Prolog verification frameworks, see e.g. [6, 9, 2, 21] and more recently [7]. Most of them aim at *paper and pencil* proofs. Although they may offer interesting and elegant methods, the validity of the proofs relies on the usual mathematical writing in natural language, and proofs are not automatically checked. In our opinion, verifying such hand-written proofs can be a time consuming process compared to a push-button approach.

For Answer Set Programming (a declarative specification language with a Prolog syntax, oriented towards knowledge representation and search problems), [8] describes an approach toward verification in which Vampire checks the equivalence of Answer Set programs.

Some other programming languages include automated verification tools *by design*. For example, Dafny [15] makes heavy use of SMT solving and Why3 [10] allows to export verification conditions to many automatic and interactive theorem provers.

An earlier account of the integration of automated and interactive theorem proving is described in [1]. As already announced in the introduction, most interactive theorem provers now include the possibility to run some automated theorem provers. Starting with Isabelle [18, 20, 3, 19], *hammers* can found in e.g., ACL2 [12], Coq [5] and Lean [16].

## 6 Conclusion

Let us recall the questions we ask in the introduction and propose our answers after this initial experiment:

- Can we also use the TPTP *Esperanto* to formulate the logic theory Stärk associates to a logic program? Yes. One axiom schema was not implemented: Axiom 3 which forbids rational terms. Another one was partially implemented: Axiom 9 for induction. Actually an inductive argument inside an inductive proof is not possible with our approach. We loose precision but in both cases we stay sound.
- Can we use off the shelf TPTP provers and obtain automatic proofs in reasonable time? Yes. We use Vampire and the E prover with their most basic options, essentially a timeout. Although Vampire seems to find a refutation faster, the E prover can regularly find proofs where Vampire seems to fail. Hence the two provers are complementary. For the moment, we did not try the advanced features offered by the provers like the one proposed in [13] for directly dealing with finite trees.

- Can we get an acceptable success rate with such an approach? Yes. With the E prover and Vampire running in parallel, the average success rate we get from our benchmark is about 77% for a one minute timeout on a standard laptop, which we find quite acceptable.

Compared to the efforts one spends while manually, laboriously elaborating certain proofs, such a tool is clearly a time-saver. We did not expect such a good success rate. We think there are various reasons that can explain it. The computing power of our current laptops is huge and automated theorem provers have been largely improved. Stärk’s art of proving, by slicing the proof of most theorems in more manageable lemmas, may also have an important impact.

Finally, there is room for improvement of the presented work, which can be considered as a first approach towards a hammer for LPTP according to [3]. In particular, the first step – the *premise selector*, which could select subparts of the LPTP library potentially useful for a proof – and the third step – the *proof reconstruction module*, which could rewrite the proof found by the automatic prover in the LPTP proof format – are yet to be investigated.

## References

- [1] W. Ahrendt, B. Beckert, R. Hähnle, W. Menzel, W. Reif, G. Schellhorn, and P. Schmitt. *Integrating Automated and Interactive Theorem Proving. Automated Deduction — A Basis for Applications: Volume II: Systems and Implementation Techniques*, pages 97–116. Springer, 1998.
- [2] K. R. Apt and E. Marchiori. Reasoning about Prolog programs: from modes through types to assertions. *Formal Aspects of Computing*, 6(6):743–765, 1994.
- [3] J. C. Blanchette, C. Kaliszyk, L. C. Paulson, and J. Urban. Hammering towards QED. *J. Formaliz. Reason.*, 9(1):101–148, 2016.
- [4] K. L. Clark. Negation as failure. In H. Gallaire and J. Minker, editors, *Logic and Databases*, pages 293–322. Plenum Press, New York, 1978.
- [5] L. Czajka, B. Ekici, and C. Kaliszyk. Concrete semantics with Coq and CoqHammer. In F. Rabe, W. M. Farmer, G. O. Passmore, and A. Youssef, editors, *CICM*, volume 11006 of *LNCS*, pages 53–59. Springer, 2018.
- [6] P. Deransart. Proof methods of declarative properties of definite programs. *Theoretical Computer Science*, pages 99–166, 1993.
- [7] W. Drabent. Correctness and completeness of logic programs. *ACM Trans. Comput. Log.*, 17(3):18, 2016.
- [8] J. Fandinno, V. Lifschitz, P. Lühne, and T. Schaub. Verifying Tight Logic Programs with Anthem and Vampire. *Theory Pract. Log. Program.*, 20(5):735–750, 2020.
- [9] G. Ferrand and P. Deransart. Proof method of partial correctness and weak completeness for normal logic programs. *J. Log. Program.*, 17(2/3&4):265–278, 1993.
- [10] J.-C. Filliâtre and A. Paskevich. Why3 - where programs meet provers. In M. Felleisen and P. Gardner, editors, *ESOP*, volume 7792 of *LNCS*, pages 125–128. Springer, 2013.
- [11] ISO/IEC 13211-1. Information Technology – Programming Languages – Prolog – Part 1: General Core. 1995.
- [12] S. J. C. Joosten, C. Kaliszyk, and J. Urban. Initial experiments with TPTP-style automated theorem provers on ACL2 problems. In F. Verbeek and J. Schmaltz, editors, *International Workshop on ACL2*, volume 152 of *EPTCS*, pages 77–85, 2014.
- [13] L. Kovács, S. Robillard, and A. Voronkov. Coming to terms with quantified reasoning. In G. Castagna and A. D. Gordon, editors, *POPL 2017*, pages 260–270. ACM, 2017.
- [14] L. Kovács and A. Voronkov. First-order Theorem Proving and Vampire. In N. Sharygina and H. Veith, editors, *CAV 2013*, volume 8044 of *LNCS*, pages 1–35. Springer, 2013.

- [15] K. R. M. Leino. Developing Verified Programs with Dafny. In R. Joshi, P. Müller, and A. Podelski, editors, *VSTTE*, volume 7152 of *LNCS*, page 82. Springer, 2012.
- [16] P. Lippe. Lean Hammer. [https://github.com/phlippe/Lean\\_hammer](https://github.com/phlippe/Lean_hammer), 2019. Accessed: 2024-01.
- [17] J. W. Lloyd. *Foundations of Logic Programming*. Springer-Verlag, 1987.
- [18] J. Meng and L. C. Paulson. Experiments on supporting interactive proof using resolution. In D. Basin and M. Rusinowitch, editors, *IJCAR*, volume 3097 of *LNCS*, pages 372–384. Springer, 2004.
- [19] L. C. Paulson. Sledgehammer: some history, some tips. <https://lawrencecpaulson.github.io/2022/04/13/Sledgehammer.html>, 2022. Accessed: 2024-02-01.
- [20] L. C. Paulson and J. C. Blanchette. Three Years of Experience with Sledgehammer, a Practical Link Between Automatic and Interactive Theorem Provers. In G. Sutcliffe, S. Schulz, and E. Ternovska, editors, *IWIL*, volume 2 of *EPiC Series in Computing*, pages 1–11. EasyChair, 2010.
- [21] D. Pedreschi and S. Ruggieri. Verification of Logic Programs. *J. Log. Program.*, 39(1-3):125–176, 1999.
- [22] S. Schulz, S. Cruanes, and P. Vukmirović. Faster, higher, stronger: E 2.3. In P. Fontaine, editor, *Proc. of the 27th CADE, Natal, Brasil*, number 11716 in *LNAI*, pages 495–507. Springer, 2019.
- [23] R. F. Stärk. First-order theories for pure Prolog programs with negation. *Arch. Math. Log.*, 34(2):113–144, 1995.
- [24] R. F. Stärk. Total correctness of logic programs: A formal approach. In R. Dyckhoff, H. Herre, and P. Schroeder-Heister, editors, *ELP’96*, volume 1050 of *LNCS*, pages 237–254. Springer, 1996.
- [25] R. F. Stärk. The theoretical foundations of LPTP (a logic program theorem prover). *Journal of Logic Programming*, 36(3):241–269, 1998.
- [26] G. Sutcliffe. The TPTP Problem Library and Associated Infrastructure. From CNF to TH0, TPTP v6.4.0. *Journal of Automated Reasoning*, 59(4):483–502, 2017.
- [27] G. Sutcliffe. The Logic Languages of the TPTP World. *Logic Journal of the IGPL*, 2022.