

# Induction with Generalization in Superposition Reasoning

Márton Hajdú $^1$ , Petra Hozzová $^1$ , Laura Kovács $^{1,2(\boxtimes)}$ , Johannes Schoisswohl $^{1,3}$ , and Andrei Voronkov $^{3,4}$ 

TU Wien, Vienna, Austria
 laura.kovacs@tuwien.ac.at
 Chalmers University of Technology, Gothenburg, Sweden
 University of Manchester, Manchester, UK
 EasyChair, Manchester, UK

Abstract. We describe an extension of automating induction in superposition-based reasoning by strengthening inductive properties and generalizing terms over which induction should be applied. We implemented our approach in the first-order theorem prover VAMPIRE and evaluated our work against state-of-the-art reasoners automating induction. We demonstrate the strength of our technique by showing that many interesting mathematical properties of natural numbers and lists can be proved automatically using this extension.

#### 1 Introduction

Automating inductive reasoning opens up new possibilities for generating and proving inductive properties, for example properties with inductive data types [4,21] or inductive invariants in program analysis and verification [13,14]. Recent advances related to automating inductive reasoning, such as first-order reasoning with inductively defined data types [16], the AVATAR architecture [26], inductive strengthening of SMT properties [22], structural induction in superposition [10] and general induction rules within saturation [19], make it possible to re-consider the grand challenge of mechanizing mathematical induction [5]. In this paper, we contribute to these advances by generalizing inductive reasoning within the saturation-based proof search of first-order theorem provers using the superposition calculus.

It is common in inductive theorem proving, that given a formula/goal F, to try to prove a more general goal instead [5]. This makes no sense in saturation-based theorem proving, which is not based on a goal-subgoal architecture. As we aim to automate and generalize inductive reasoning within saturation-based proof search, our work follows a different approach than the one used in inductive theorem provers. Namely, our methodology in Sect. 4 picks up a formula F (not necessarily the goal) in the search space and adds to the search space new induction axioms with generalization, that is, instances of generalized induction schemata, aiming at proving both  $\neg F$  and a more general formula than  $\neg F$ .

<sup>©</sup> Springer Nature Switzerland AG 2020

C. Benzmüller and B. Miller (Eds.): CICM 2020, LNAI 12236, pp. 123–137, 2020.

In Sect. 3 we give a concrete example motivating our approach, illustrating the advantage of induction with generalization in saturation-based proof search. We then present our main contributions, as follows:

- (1) We introduce a new inference rule for first-order superposition reasoning, called induction with generalization (Sect. 4). Our work extends [19] by proving properties with multiple occurrences of the same induction term and by instantiating induction axioms with logically stronger versions of the property being proved. Our approach is conceptually different from previous attempts to use induction with superposition [10,11,15], as we are not restricted to specific clause splitting algorithms and heuristics used in [10], nor are we limited to induction over term algebras with the subterm ordering in [11]. As a result, we stay within the standard saturation framework and do not have to introduce constraint clauses, additional predicates or change the notion of redundancy as in [11].
- (2) We implemented our work in the Vampire theorem prover [17] and compared it to state-of-the-art reasoners automating induction, including Acl2 [5], Cvc4 [2], Imandra [18], Zeno [24] and Zipperposition [10] (Sect. 5). We also provide a set of handcrafted mathematical problems over natural numbers and lists. We show that induction with generalization in Vampire can solve problems that existing systems, including Vampire without this rule, cannot.
- (3) We provide a new digital dataset consisting of over 3,300 inductive benchmarks, for which generalized applications of induction is needed (Sect. 5). Our dataset is formalized within the SMT-LIB format using data types [3] and available at: https://github.com/vprover/inductive\_benchmarks.

#### 2 Preliminaries

We assume familiarity with multi-sorted first-order logic and saturation-based superposition reasoning. For details, we refer to [17]. Throughout this paper we denote fresh Skolem constants by  $\sigma$ , variables by x, y, z and terms by t, all possibly with indices. We denote the equality predicate by t and consider t as part of the language. Further, we write  $t_1 \neq t_2$  for the formula  $t_1 = t_2$ .

Given a set of formulas (including a negated conjecture), superposition-based theorem provers run saturation algorithms on a set of clauses corresponding to the clausal normal form (CNF) of the input set of formulas. We denote literals by L and clauses by C, all possibly with indices. We use  $\square$  to denote the empty clause. In [16] we showed how superposition-based provers can be extended with reasoning about the theory for finite term algebras.

We will denote term algebras corresponding to natural numbers by  $\mathbb{N}$  and lists of natural numbers by  $\mathbb{L}$ . We refer to the elements of the signature of the term algebras as *constructors*. We will use the same notations  $\mathbb{N}$  and  $\mathbb{L}$  for these term algebras extended by additional function and predicate symbols shown in Fig. 1.

	Natural numbers $\mathbb{N}$	Natural lists $\mathbb{L}$			
Constructors	$0: \mathbb{N} \qquad s: \mathbb{N} \to \mathbb{N}$	$\mathtt{nil}: \mathbb{L} \qquad \mathtt{cons}: \mathbb{N}  imes \mathbb{L}  o \mathbb{L}$			
Symbols	$+: \mathbb{N} \times \mathbb{N} \to \mathbb{N}$	$++: \mathbb{L} \times \mathbb{L} \to \mathbb{L}$			
	$\leq : \mathbb{N} \times \mathbb{N} \to \text{bool}$	$\mathtt{prefix}: \mathbb{L}  imes \mathbb{L}  o \mathrm{bool}$			
Axioms	$\forall y.(0+y=y)$	$\forall l.(\mathtt{nil} + + l = l)$			
	$\forall x, y.(s(x) + y = s(x+y))$	$\forall x, l, k.(cons(x, l) ++ k = cons(x, l ++ k))$			
	$\forall x.0 \leq x$	$orall l.\mathtt{prefix}(\mathtt{nil},l)$			
	$\forall x. \neg s(x) \le 0$	$\forall x, l. \neg \texttt{prefix}(\texttt{cons}(x, l), \texttt{nil})$			
	$\forall x, y. (s(x) \le s(y))$	$\forall x, l, y, k. (\texttt{prefix}(\texttt{cons}(x, l), \texttt{cons}(y, k))$			
	$\leftrightarrow x \leq y$ )	$\longleftrightarrow (x = y \land \mathtt{prefix}(l,k))\big)$			

**Fig. 1.** Term algebras of  $\mathbb{N}$  and  $\mathbb{L}$ , together with additional symbols and axioms.

Specifically, we will deal with + and  $\le$  for  $\mathbb{N}$  having their standard meaning and ++ and prefix for  $\mathbb{L}$ , denoting the list concatenation and the prefix relation, respectively. These additional symbols are axiomatized by first-order formulas corresponding to their recursive definitions, shown in Fig. 1.

While we use  $\mathbb{N}$  and  $\mathbb{L}$  for illustration, we however note that our approach can be used for proving properties over any other theories with various forms of induction.

Theorem proving of first-order properties of inductively defined data types needs to handle the domain closure, injectivity, distinctness and acyclicity axioms of term algebras – a detailed definition of these axioms can be found in [16, 23]. The challenge we address in [16] is how to automate proving term algebras properties given the fact that the acyclicity axiom is not finitely axiomatizable.

Throughout this paper, we will be using the structural induction axiom and rule for  $\mathbb{N}$ , introduced in [19], for illustrating our approach. Given a literal  $\neg L[t]$ , where t is chosen as an induction term, a structural induction axiom for  $\mathbb{N}$  is:

$$(L[0] \land \forall x. (L[x] \to L[s(x)])) \to \forall y. (L[y]). \tag{1}$$

Informally, the axiom expresses that if the base case holds, and if the induction step holds, then the literal holds for all possible values. The structural induction rule for  $\mathbb{N}$ , given a clause  $\neg L[t] \lor C$ , adds a clausified form of this axiom to the search space:

$$\frac{\neg L[t] \lor C}{(\neg L[0] \lor L[\sigma] \lor L[y]) \land (\neg L[0] \lor \neg L[s(\sigma)] \lor L[y])}. \tag{2}$$

After using the rule, the L[y] in both resulting clauses can be resolved against the  $\neg L[t]$  in the premise clause.

## 3 Motivating Example

Let us now motivate our approach to induction with generalization, by considering the following formula expressing the associativity of addition over  $\mathbb{N}$ :

$$\forall x, y, z.(x + (y + z) = (x + y) + z), \quad \text{with } x, y, z \in \mathbb{N}.$$
 (3)

$$\begin{aligned} &(C_1) \ \sigma_1 + (\sigma_2 + \sigma_3) \neq (\sigma_1 + \sigma_2) + \sigma_3 & \text{[input]} \\ &(C_2) \ 0 + (\sigma_2 + \sigma_3) \neq (0 + \sigma_2) + \sigma_3 \vee \sigma + (\sigma_2 + \sigma_3) = (\sigma + \sigma_2) + \sigma_3 & \text{[induct. } C_1] \\ &(C_3) \ 0 + (\sigma_2 + \sigma_3) \neq (0 + \sigma_2) + \sigma_3 \vee s(\sigma) + (\sigma_2 + \sigma_3) \neq (s(\sigma) + \sigma_2) + \sigma_3 & \text{[induct. } C_1] \\ &(C_4) \ 0 + (\sigma_2 + \sigma_3) \neq (0 + \sigma_2) + \sigma_3 \vee s(\sigma + (\sigma_2 + \sigma_3)) \neq s((\sigma + \sigma_2) + \sigma_3) & \text{[induct. } C_1] \\ &(C_5) \ 0 + (\sigma_2 + \sigma_3) \neq (0 + \sigma_2) + \sigma_3 \vee \sigma + (\sigma_2 + \sigma_3) \neq (\sigma + \sigma_2) + \sigma_3 & \text{[injective } C_4] \\ &(C_6) \ 0 + (\sigma_2 + \sigma_3) \neq (0 + \sigma_2) + \sigma_3 & \text{[res. } C_2, C_5] \\ &(C_7) \ \sigma_2 + \sigma_3 \neq \sigma_2 + \sigma_3 & \text{[trivial ineq. } C_7] \end{aligned}$$

**Fig. 2.** Proof of associativity of + in a saturation-based theorem prover with induction

The induction approach introduced in [19] is able to prove this problem. The main steps of such a proof are shown in Fig. 2 and discussed next. First, the negation of formula (3) is skolemized, yielding the (unit) clause  $C_1$  of Fig. 2. As already mentioned, the  $\sigma_i$  denote fresh Skolem constants introduced during clausification. Next, the structural induction axiom (1) is instantiated so that its conclusion can resolve against  $C_1$  using the constant  $\sigma_1$  as the induction term, resulting in the formula:

$$\begin{pmatrix}
0 + (\sigma_2 + \sigma_3) = (0 + \sigma_2) + \sigma_3 \land \\
\forall x.(x + (\sigma_2 + \sigma_3) = (x + \sigma_2) + \sigma_3 \to s(x) + (\sigma_2 + \sigma_3) = (s(x) + \sigma_2) + \sigma_3)
\end{pmatrix} (4)$$

$$\to \forall y.(y + (\sigma_2 + \sigma_3) = (y + \sigma_2) + \sigma_3).$$

Then, the CNF of the induction axiom (4) is added to the search space using the following instance of the structural induction rule (2):

$$\frac{\sigma_{1} + (\sigma_{2} + \sigma_{3}) \neq (\sigma_{1} + \sigma_{2}) + \sigma_{3}}{(0 + (\sigma_{2} + \sigma_{3}) \neq (0 + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) \neq (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + \sigma_{3} \vee \sigma + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{2}) + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{3}) + (\sigma_{2} + \sigma_{3}) = (\sigma + \sigma_{3}) + (\sigma_{3} + \sigma_{3})$$

The clauses from the inference conclusion are resolved against  $C_1$ , yielding clauses  $C_2$ ,  $C_3$  of Fig. 2. Clause  $C_4$  originates by repeated demodulation into  $C_3$  using the second axiom of Fig. 1 over  $\mathbb{N}$ . Further,  $C_5$  is derived from  $C_4$  by using the injectivity property of term algebras and  $C_6$  is a resolvent of  $C_2$  and  $C_5$ . Clause  $C_7$  is then derived by repeated demodulation into  $C_6$ , using the first axiom of Fig. 1 over  $\mathbb{N}$ . By removing the trivial inequality from  $C_7$ , we finally derive the empty clause  $C_8$ .

Consider now the following instance of the associativity property (3):

$$\forall x.(x + (x + x) = (x + x) + x). \tag{6}$$

While (6) is an instance of (3), we cannot prove it using the same approach. Let us explain why this is the case. By instantiating the induction axiom (1) using (6), we get:

$$\begin{pmatrix}
0 + (0 + 0) = (0 + 0) + 0 \land \\
\forall x.(x + (x + x) = (x + x) + x \to s(x) + (s(x) + s(x)) = (s(x) + s(x)) + s(x))
\end{pmatrix} (7)$$

$$\to \forall y.(y + (y + y) = (y + y) + y).$$

After resolving this axiom with the skolemized negation of (6), we get the following two clauses<sup>1</sup>:

$$0 + (0+0) \neq (0+0) + 0 \lor \sigma + (\sigma + \sigma) = (\sigma + \sigma) + \sigma \tag{8}$$

$$0 + (0+0) \neq (0+0) + 0 \lor s(\sigma) + (s(\sigma) + s(\sigma)) \neq (s(\sigma) + s(\sigma)) + s(\sigma)$$
 (9)

While the first literals of (8) and (9) are easily resolved using axioms of +, not much can be done with the latter literals. We can only apply repeated demodulations over the second literal of (9) using axioms of + and the injectivity property of term algebras, yielding  $\sigma + s(\sigma + s(\sigma)) \neq (\sigma + s(\sigma)) + s(\sigma)$ . No further inference over this formula can be applied, in particular it cannot be resolved against the second literal of (8). Hence, the approach of [19] fails proving (6).

The existing approaches to induction also suffer from the same problem. For example [2,5,10,18,24], can prove property (3) but fail to prove its weaker instance (6). The common recipe in inductive theorem proving [5] is to try to prove (3) in addition to trying to prove (6).

Interestingly, in saturation-based theorem proving we can do better. If we follow the common recipe, we would add a generalized goal and then an induction axiom for it. Instead, we only add the induction axiom instance corresponding to the generalized goal without adding the extra goal, which results in a smaller number of clauses. More precisely, in addition to the instance of the induction schema corresponding to (6), we also add instance (4) corresponding to (3). We call this new inference rule *induction with generalization*.

#### 4 Induction with Generalization

Following [19], we consider an induction axiom to be any valid formula of the form  $premise \to \forall y.(L[y])$ , in the underlying theory, such as the theory of term algebras. An example of an induction axiom is the structural induction axiom (4). An induction schema is a collection of induction axioms. Each induction schema we consider is the set of first-order instances of some valid higher-order formula. The work [19] introduces a rule of induction where a ground literal  $\neg L[t]$  appearing in the proof search triggers addition of the corresponding induction axioms  $premise \to \forall y.(L[y])$  to the search space:

$$\frac{\neg L[t] \lor C}{\text{CNF}(premise \to \forall y.(L[y]))}(\text{induction}), \tag{10}$$

<sup>&</sup>lt;sup>1</sup> These clauses are instances of  $C_2$  and  $C_3$  from Fig. 2.

where L[y] is obtained from L[t] by replacing *all* occurrences of t by y. An example of an instance of the induction rule is (5).

While addition of a large number of such formulas may seem to blow up the search space, in practice Vampire handles such addition with little overhead, resulting in finding proofs containing nearly 150 induction inferences [19]. The reason why the overhead of adding structural induction axioms is small is explained in [20]: the added clauses only contain one variable (the y in L[y]), and the clauses containing this literal are immediately subsumed by a ground clause. The net result is adding a small number of ground clauses, which are especially easy to handle in the AVATAR architecture implemented in VAMPIRE.

**Induction with Generalization.** In a nutshell, given a goal, we add an induction axiom corresponding to a more general one. The rule can be formulated in the same way as (10), yet with a different conclusion:

$$\frac{\neg L[t] \lor C}{\text{CNF}(premise' \to \forall y.(L'[y]))} (\texttt{IndGen}), \tag{11}$$

where L'[y] is obtained from L[t] by replacing *some* occurrences of t by y, and premise' is the premise corresponding to L'[y]. Both induction rules are obviously sound because their conclusions are constructed such that they are valid in the underlying theory.

To implement IndGen, if a clause selected for inferences contains a ground literal  $\neg L[t]$  having more than one occurrence of t, we should select a non-empty subset of occurrences of t in L[t], select an induction axiom corresponding to this subset, and then apply the rule.

Motivating Example, Continued. Suppose that t is  $\sigma_1$  and  $\neg L[t]$  is  $\sigma_1 + (\sigma_1 + \sigma_1) \neq (\sigma_1 + \sigma_1) + \sigma_1$ , which is obtained by negating and skolemizing (6). Then by applying IndGen we can add the following induction axiom:

$$\begin{pmatrix}
0 + (\sigma_1 + \sigma_1) = (0 + \sigma_1) + \sigma_1 \land \\
\forall x.(x + (\sigma_1 + \sigma_1) = (x + \sigma_1) + \sigma_1 \to s(x) + (\sigma_1 + \sigma_1) = (s(x) + \sigma_1) + \sigma_1)
\end{pmatrix} (12)$$

$$\to \forall y.(y + (\sigma_1 + \sigma_1) = (y + \sigma_1) + \sigma_1),$$

which is different from (7). When we add this formula, we can derive the empty clause in the same way as in Fig. 2.

Saturation with Induction with Generalization. The main questions to answer when applying induction with generalization is which occurrences of the induction term in the induction literal we should choose.

Generally, if the subterm t occurs n times in the premise, there are  $2^n - 1$  ways of applying the rule, all potentially resulting in formulas not implying each other. Thus, an obvious heuristic to use all non-empty subsets may result in too many formulas. For example,  $\sigma_1 + (\sigma_1 + \sigma_1) \neq (\sigma_1 + \sigma_1) + \sigma_1$  would result in adding 63 induction formulas.

Another simple heuristic is to restrict the number of occurrences selected as induction term to a fixed number. This strategy reduces the number of applications of induction at the cost of losing proofs that would need subsets of

cardinality larger than the limit. Finding possible heuristics for selecting specific subsets for common cases of literals can be subject of future work, especially interesting in proof assistants in mathematics.

Note that some of the conclusions of (11) can, in turn, have many children obtained by induction with generalization. Our experiments in Sect. 5 show that, even when we generate all possible children, VAMPIRE can still solve large examples with more than 10 occurrences of the same induction variable, again thanks to the effect that, for each application of induction, only a small number of ground clauses turn out to be added to the search space.

We therefore believe that our work can potentially be also useful for larger examples, and even in cases when the inductive property to be proved is embedded in a larger context.

### 5 Experiments

**Implementation.** We implemented induction with generalization in VAMPIRE, with two new options:

- boolean-valued option indgen, which turns on/off the application of induction with generalization, with the default value being off, and
- integer-valued option indgenss, which sets the maximum size of the subset of occurrences used for induction, with the default value 3. This option is ignored if indgen is off.

Our implementation of induction with generalization is available at: https://github.com/vprover/vampire.

In experiments described here, if indgen is off, VAMPIRE performs induction on all occurrences of a term in a literal as in [19]. In this section

- VAMPIRE refers to the (default) version of VAMPIRE with induction rule (10) (i.e., the option -ind struct)
- VAMPIRE\* additionally uses the IndGen rule of induction with generalization (11) (i.e., the options -ind struct -indgen on).
- VAMPIRE\*\* uses the same options as VAMPIRE\* plus the option -indoct on, which applies induction to arbitrary ground terms, not just to constants as in VAMPIRE or in VAMPIRE\*.

SMT-LIB Experiments. We evaluated our work using the UFDT and UFDTLIA problem sets from SMT-LIB [3], yielding all together 4854 problems. Many of these problems come from program analysis and verification and contain large numbers of axioms, so they are different from standard mathematical examples used in many other papers on automation of induction. Given the nature of the benchmarks, we were interested in two questions:

 What is the overhead incurred by using induction with generalization in large search spaces, especially when it is not used in proofs? If the new rule is prohibitively expensive, this means it could probably only be used in smaller examples used in interactive theorem proving. 2. Is the new rule useful at all for this kind of benchmarks? While the new rule can be used in principle, should it (or can it) be used in program analysis and verification?

Our results show that the overhead is relatively small but we could not solve problems not solvable without the use of the new rule.

Induction (10) in Vampire was already evaluated in [19] against other solvers on these examples. Hence, we only compare how Vampire\*/Vampire\*\* performs against Vampire, using both the default and the portfolio modes. (In the default mode, Vampire/Vampire\*/Vampire\*\* uses default values for all parameters except the ones specified by the user; in the portfolio mode, Vampire/Vampire\*/Vampire\*\* sequentially tries different configurations for parameters not specified by the user.) Together, we ran 18 instances: Vampire, Vampire\* with indgenss set to 2, 3, 4 and unlimited, and Vampire\*\* with the same four variants of indgenss; each of them in both default and portfolio mode. We ran our experiments on the StarExec cluster [25].

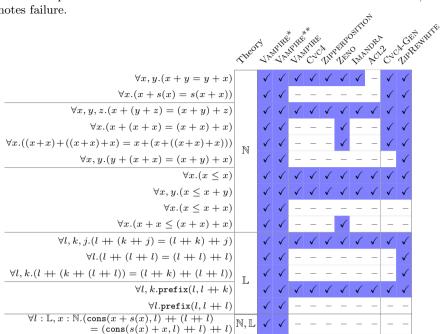
The best Vampire\*/Vampire\*\* solved 5 problems in the portfolio mode and 1 problem in the default mode not solved by Vampire. However, the proofs found by them did not use induction with generalization. This is a common problem in experiments with saturation theorem proving: new rules change the direction of the proof search and may result in new simplifications that also drastically affect the search space. As a result, new proofs may be found, yet these proofs do not actually use the new rule. There were no problems solved by Vampire that were not solved by any Vampire\*/Vampire\*\*.

The maximum number of IndGen applications in proofs was 3 and the maximum depth of induction was 4. Vampire\*/Vampire\*\* used generalized induction in proofs of 10 problems. However, these problems are also solvable by Vampire (without generalized induction). Thus, we conclude that SMT-LIB problems (probably as well as other typical program analysis and verification benchmarks) typically do not gain from using generalization.

Experiments with Mathematical Problems. We handcrafted a number of natural problems over natural numbers and lists and tested the new rule on these problems. Our benchmarks are available at: https://github.com/vprover/inductive\_benchmarks.

Table 1 lists 16 of such examples using the functions defined in Fig. 1. Some examples were taken from or inspired by the TIP benchmark library [9]: e.g., the seventh benchmark in Table 1 is adapted from the TIP library and the second problem is inspired by a symmetric problem from the TIP library,  $\forall x.(s(x)+x=s(x+x))$ . While they are handcrafted, we believe they are representative since no attempt was done to exclude problems not solvable by VAMPIRE using induction with generalization.

We evaluated and compared several state-of-the-art reasoners supporting standard input formats and, due to the nature of our work, either superposition-based approaches or approaches to generalization. It was not easy to make these experiments since provers use different input syntaxes (see Table 2). As a result, we also had to design translations of our benchmarks.



**Table 1.** Experiments with 16 handcrafted benchmarks. "✓" denotes success, "\_' denotes failure.

Except for IMANDRA (which is a cloud-based service), we ran our experiments on a 2,9 GHz Quad-Core Intel Core i7 machine. We ran each solver as a single-threaded process with a 5s time limit. Our results are summarized in Table 1, where Cvc4-Gen refers to the solver Cvc4 with the automatic lemma discovery enabled. ZIPREWRITE refers to ZIPPERPOSITION with function and predicate definitions encoded as rewrite rules instead of ordinary logical formulas, in order to trigger its generalization heuristics [10]. Configurations used for running all solvers are listed in Table 2.

Table 1 shows that Vampire\*/Vampire\*\* (with indgenss=3) outperforms all solvers, including Vampire itself. When considering solvers without fine-tuned heuristics, such as in Ziprewrite and Cvc4-Gen, Vampire\*\* solves many more problems. Interestingly, Ziprewrite heuristics work well with addition and list concatenation, but not with orders. Further, Cvc4-Gen heuristics prove associativity of addition, but not the list counterpart for concatenation. We believe our experiments show the potential of using induction with generalization as a new inference rule since it outperforms heuristic-driven approaches with no special heuristics or fine-tuning added to Vampire.

Experiments with Problems Requiring Associativity and Commutativity. The (x + x) + x = x + (x + x) is a special case of a family of problems over natural numbers. The problems can be formulated as follows.

Let  $t_1$  and  $t_2$  be two terms built using variables, + and the successor function. Then the equality  $t_1 = t_2$  is valid over natural numbers if and only if they have the same number of occurrences of the successor function and each variable of this equality has the same number of occurrences in  $t_1$  and  $t_2$ . For example, the following equality is valid over natural numbers:

$$s(x + (x + s(y + z))) + s(z) = (z + s(x)) + (x + s(s((z + y)))).$$

To prove such problems over  $\mathbb{N}$ , one needs both induction and generalization. Without the successor function, they can be easily proved using associativity and commutativity of +, but associativity and commutativity are not included in the axioms of  $\mathbb{N}$ . When the terms are large, the problems become highly challenging.

Solver	Configuration	Input format		
Vampire	-ind struct	SMT-LIB		
Vampire*	-ind struct -indgen on	SMT-LIB		
Vampire**	-ind struct -indgen on -indoct on	SMT-LIB		
Cvc4	quant-ind	SMT-LIB		
CVC4-GEN	quant-indconjecture-gen	SMT-LIB		
ZIPPERPOSITION	default mode	.zf (native input format)		
ZIPREWRITE	default mode	.zf with definitions as rewrite rules		
ZENO	default mode	functional program encoding		
Imandra	default mode	functional program encoding		
ACL2	default mode	functional program encoding		

Table 2. Configurations and input format of solvers for the mathematical problems.

We generated a set of instances of these problems (with and without the successor function, and also other functions and predicates) by increasing term sizes. We also generated similar problems for lists using concatenation and reverse functions, and prefix predicate. Some of the terms were, e.g., variations of (6) with 20 occurrences of x. Our entire dataset, containing over 3,300 examples, is available at the previously mentioned URL.

We were again interested in evaluating and comparing various reasoners and approaches on these problems. The interesting feature of these problems is that they are natural yet we can generate problems of almost arbitrary complexity.

We evaluated and compared Vampire\*, Vampire\*\*, Cvc4-Gen, Zeno and Ziprewrite, that is the best performing solvers on inductive reasoning with generalization according to Table 1, using the same experimental setting as already described for Table 1. Table 3 lists a partial summary of our experiments, displaying results for 2,007 large instances of four simple properties with one variable, corresponding to the fourth, ninth, twelfth and fifteenth problem from Table 1. (Due to space constraints, we chose these problems as a representative subset of our large benchmarks, since the solvers' performance was very similar for the whole benchmark set.)

In Table 3, we use the following notation. By nx = nx we denote formulas of the form  $x \circ ... \circ x = x \circ ... \circ x$  with n occurrences of x on both sides of the equality, and parentheses on various places in the expressions, with  $\circ$  being +, or

++ for the datatypes  $\mathbb{N}$  and  $\mathbb{L}$ , respectively. By  $mx \leq nx$  and  $\mathtt{prefix}(mx,nx)$  we denote formulas of the form  $x+\ldots+x\leq x+\ldots+x$  and  $\mathtt{prefix}(x++\ldots+x,x++\ldots+x)$ , respectively, with m occurrences of x on the left and n occurrences of x on the right hand side of the  $\leq$  or  $\mathtt{prefix}$  predicates, and with parentheses on various places in the expressions. Result N%(M) means that the solver solved M of the problems from this category, which corresponds to N%.

From Table 3, we conclude that Vampire\*\* scales better than CVC4-Gen on a large majority of benchmarks, and scales comparably to Zeno. While Ziprewrite can solve more problems than Vampire\*\*, Vampire\*\* is more consistent in solving at least some problems from each category. Ziprewrite can solve many problems thanks to treatment of equalities as rewrite rules. We are planning to add an option of using recursive definitions as rewrite rules in Vampire in the future too.

#### 6 Related Work

Research into automating induction has a long history with a number of techniques developed, including for example approaches based on semi-automatic inductive theorem proving [5,7,8,18], specialized rewriting procedures [12], SMT reasoning [22] and superposition reasoning [10,11,15,19].

Previous works on automating induction mainly focus on inductive theorem proving [7,8,24]: deciding when induction should be applied and what induction axiom should be used. Further restrictions are made on the logical expressiveness, for example induction only over universal properties [5,24] and without uninterpreted symbols [18], or only over term algebras [11,15]. Inductive proofs usually rely on auxiliary lemmas to help proving an inductive property. In [8] heuristics for finding such lemmas are introduced, for example by randomly generating equational formulas over random inputs and using these formulas if they hold reasonably often. The use of [8] is therefore limited to the underlining heuristics. Other approaches to automating induction circumvent the need for auxiliary lemmas by using uncommon cut-free proof systems for inductive reasoning, such as a restricted  $\omega$ -rule [1], or cyclic reasoning [6].

The work presented in this paper automates induction by integrating it directly in superposition-based proof search, without relying on rewrite rules and external heuristics for generating auxiliary inductive lemmas/subgoals as in [5,7,8,18,24]. Our new inference rule IndGen for induction with generalization adds new formulas to the search space and can replace lemma discovery heuristics used in [7,8,22]. Our work also extends [19] by using and instantiating induction axioms with logically stronger versions of the property being proved. Unlike [10], our methods do not necessarily depend on AVATAR [26], can be used with any (inductive) data type and target induction rules different than structural induction. Contrarily to [11], we are not limited to induction over term algebras with the subterm ordering and we stay in a standard saturation framework. Moreover, compared to [5,7,8,22], one of the main advantages of our approach is that it does not use a goal-subgoal architecture and can, as a result, combine superposition-based equational reasoning with inductive reasoning.

**Table 3.** Experiments on 2,007 arithmetical problems.

	THEORY	Vampire*	Vampire**	CVC4-GEN	Zeno	ZipRewrite
3x = 3x		100% (1)	100% (1)	100% (1)	100% (1)	100% (1)
4x = 4x	И	90% (9)	100% (10)	100% (10)	20% (2)	100% (10)
5x = 5x		30% (15)	50% (25)	100% (50)	12% (6)	100% (50)
6x = 6x		8% (4)	18% (9)	100% (50)	22% (11)	100% (50)
7x = 7x		- '	10% (5)	100% (50)	2% (1)	100% (50)
8x = 8x		_	2% (1)	100% (50)	4% (2)	100% (50)
9x = 9x		_	2% (1)	100% (50)	8% (4)	84% (42)
10x = 10x		_	_	100% (50)	8% (4)	90% (45)
3x = 3x	L	100% (1)	100% (1)	_	_	100% (1)
4x = 4x		70% (7)	90% (9)	_	_	100% (10)
5x = 5x		46% (23)	48% (24)	_	_	100% (50)
6x = 6x		6% (3)	26% (13)	_	6% (3)	100% (50)
7x = 7x		2% (1)	6% (3)	_	_	100% (50)
8x = 8x		-	_	_	_	90% (45)
9x = 9x		_	_	_	_	88% (44)
10x = 10x		_	_	_	_	68% (34)
$3x \leq 3x$		100% (2)	100% (2)	100% (2)	100% (2)	100% (2)
$4x \le 4x$		_ ( <u>-</u> )	15% (3)	100% (20)	20% (4)	100% (20)
$5x \le 5x$		_	4% (2)	100% (50)	12% (6)	100% (50)
$1x \le 2x$		100% (1)	100% (1)	_	-	-
$2x \le 3x$		50% (1)	50% (1)	_	100% (2)	_
$3x \le 4x$		-	30% (3)	_	40% (4)	_
$4x \le 5x$		_	8% (4)	_	16% (8)	_
$5x \le 6x$		_	6% (3)	_	10% (5)	_
$1x \leq 3x$		100% (2)	100% (2)	_	100% (2)	100% (2)
$2x \le 4x$		-	40% (2)	_	40% (2)	100% (5)
$3x \le 5x$		_	14% (4)	_	28% (8)	100% (28)
$4x \leq 6x$	N	_	10% (5)	_	18% (9)	100% (50)
$5x \leq 7x$		_	4% (2)	_	18% (9)	100% (50)
$1x \leq 4x$		100% (5)	100% (5)	_	80% (4)	100% (5)
$2x \leq 5x$		_	35% (5)	_	42% (6)	100% (14)
$3x \leq 6x$		-	18% (9)	_	38% (19)	100% (50)
$4x \leq 7x$		_	6% (3)	_	16% (8)	100% (50)
$5x \leq 8x$		_	- '	_	6% (3)	100% (50)
$1x \leq 5x$		100% (14)	100% (14)	_	85% (12)	100% (14)
$2x \leq 6x$		-	33% (14)	_	26% (11)	100% (42)
$3x \leq 7x$		_	14% (7)	_	32% (16)	100% (50)
$4x \leq 8x$		_	4% (2)	_	18% (9)	100% (50)
$5x \leq 9x$		-	-	-	14% (7)	100% (50)
prefix(3x,3x)		100% (2)	50% (1)	_		100% (2)
prefix(4x, 4x)		-	25% (5)	_	-	100% (20)
prefix(5x, 5x)		_	2% (1)	_	4% (2)	100% (50)
prefix(1x, 2x)		100% (1)	100% (1)	-	-	_
prefix(2x, 3x)		-	50% (1)	_	50% (1)	_
prefix(3x, 4x)		-	20% (2)	-	20% (2)	-
prefix(4x, 5x)		-	8% (4)	-	8% (4)	_
prefix(5x, 6x)		-	-	-	-	-
prefix(1x, 3x)		100% (2)	100% (2)	-	50% (1)	100% (2)
prefix(2x, 4x)		20% (1)	40% (2)	-	20% (1)	100% (5)
prefix(3x, 5x)		-	14% (4)	-	14% (4)	100% (28)
prefix(4x, 6x)	L	-	6% (3)	-	8% (4)	100% (50)
prefix(5x,7x)		-	2% (1)	-	2% (1)	100% (50)
prefix(1x, 4x)		100% (5)	100% (5)	-	40% (2)	100% (5)
prefix(2x, 5x)		-	35% (5)	-	21% (3)	100% (14)
prefix(3x, 6x)		-	14% (7)	-	12% (6)	100% (50)
prefix(4x,7x)		_	4% (2)	_	4% (2)	100% (50)
prefix(5x, 8x)		-	-	-	4% (2)	100% (50)
prefix(1x, 5x)		100% (14)	100% (14)	-	42% (6)	100% (14)
prefix(2x, 6x)		-	33% (14)	-	21% (9)	100% (42)
prefix(3x,7x)		-	16% (8)	-	16% (8)	100% (50)
prefix(4x, 8x)		_	10% (5)	-	12% (6)	100% (50)
prefix(5x, 9x)		_	- '	_	- 1	100% (50)
- ' ' ''/						

Normally, generalization in theorem proving means that given a goal F, we try to prove a more general goal. In logic, a statement F' is more general than F if F' implies F. Thus, by proving F' we also prove F. One way to generalize is to replace one or more occurrences of a subterm by a fresh variable, using the fact that  $\forall x.(F[x])$  implies F[t]. This is essentially the idea behind approaches to generalization in all systems we compared with. While our approach is superficially similar, it does something fundamentally different. Instead of (or in addition to) adding an instance I of the induction schema that can be used to prove F[t], we add an instance I' that can be used to prove  $\forall x.(F[x])$ . An interesting observation is that, in general, neither I implies I', nor I' implies I, so neither of I and I' is more general.

The second fundamental difference is that, because induction in VAMPIRE is not based on a goal-subgoal architecture, we can add both induction formulas I and I' at the same time. While this may seem inefficient, for some induction schemata, including structural induction, the overhead is very small (as also confirmed by our experiments).

#### 7 Conclusions

We introduced a new rule for induction with generalization in saturation-based reasoning based on adding induction axioms for proving generalizations of the goals appearing during proof-search. Our experiments show that we solve many problems that other existing systems cannot solve. Future work includes designing heuristics to guide proof search, using rewriting approaches, and performing other kinds of generalization and induction.

Acknowledgments. We thank Giles Reger for discussions related to the work. We acknowledge funding supporting this work, in particular the ERC starting grant 2014 SYMCAR 639270, the EPSRC grant EP/P03408X/1, the ERC proof of concept grant 2018 SYMELS 842066, the Wallenberg Academy fellowship 2014 TheProSE, the Austrian FWF research project W1255-N23, and the Hungarian-Austrian project 101öu8.

# References

- Baker, S., Ireland, A., Smaill, A.: On the use of the constructive omega-rule within automated deduction. In: Voronkov, A. (ed.) LPAR 1992. LNCS, vol. 624, pp. 214–225. Springer, Heidelberg (1992). https://doi.org/10.1007/BFb0013063
- Barrett, C., et al.: CVC4. In: Gopalakrishnan, G., Qadeer, S. (eds.) CAV 2011.
   LNCS, vol. 6806, pp. 171–177. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-22110-1\_14
- 3. Barrett, C., Fontaine, P., Tinelli, C.: The Satisfiability Modulo Theories Library (SMT-LIB) (2016). www.SMT-LIB.org
- Blanchette, J.C., Peltier, N., Robillard, S.: Superposition with datatypes and codatatypes. In: Galmiche, D., Schulz, S., Sebastiani, R. (eds.) IJCAR 2018. LNCS (LNAI), vol. 10900, pp. 370–387. Springer, Cham (2018). https://doi.org/10.1007/ 978-3-319-94205-6\_25

- 5. Boyer, R.S., Moore, J.S.: A Computational Logic Handbook, Perspectives incomputing, vol. 23. Academic Press, Cambridge (1979)
- Brotherston, J., Simpson, A.: Sequent calculi for induction and infinite descent. J. Log. Comput. 21(6), 1177–1216 (2011)
- 7. Bundy, A., Stevens, A., van Harmelen, F., Ireland, A., Smaill, A.: Rippling: a heuristic for guiding inductive proofs. Artif. Intell. **62**(2), 185–253 (1993)
- 8. Claessen, K., Johansson, M., Rosén, D., Smallbone, N.: HipSpec: automating inductive proofs of program properties. In: Proceedings of the ATx/WInG, pp. 16–25 (2012)
- Claessen, K., Johansson, M., Rosén, D., Smallbone, N.: TIP: tons of inductive problems. In: Kerber, M., Carette, J., Kaliszyk, C., Rabe, F., Sorge, V. (eds.) CICM 2015. LNCS (LNAI), vol. 9150, pp. 333–337. Springer, Cham (2015). https://doi. org/10.1007/978-3-319-20615-8-23
- Cruanes, S.: Superposition with structural induction. In: Dixon, C., Finger, M. (eds.) FroCoS 2017. LNCS (LNAI), vol. 10483, pp. 172–188. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66167-4\_10
- Echenim, M., Peltier, N.: Combining induction and saturation-based theorem proving. J. Autom. Reason. 64(2), 253–294 (2019). https://doi.org/10.1007/s10817-019-09519-x
- Falke, S., Kapur, D.: Rewriting induction + linear arithmetic = decision procedure.
   In: Gramlich, B., Miller, D., Sattler, U. (eds.) IJCAR 2012. LNCS (LNAI), vol. 7364, pp. 241–255. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-31365-3-20
- 13. Gleiss, B., Kovács, L., Robillard, S.: Loop analysis by quantification over iterations. In: Proceedings of the LPAR, pp. 381–399 (2018)
- Gurfinkel, A., Shoham, S., Vizel, Y.: Quantifiers on demand. In: Lahiri, S.K., Wang, C. (eds.) ATVA 2018. LNCS, vol. 11138, pp. 248–266. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01090-4\_15
- Kersani, A., Peltier, N.: Combining superposition and induction: a practical realization. In: Fontaine, P., Ringeissen, C., Schmidt, R.A. (eds.) FroCoS 2013. LNCS (LNAI), vol. 8152, pp. 7–22. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-40885-4\_2
- 16. Kovács, L., Robillard, S., Voronkov, A.: Coming to terms with quantified reasoning. In: Proceedings of POPL, pp. 260–270 (2017)
- Kovács, L., Voronkov, A.: First-order theorem proving and VAMPIRE. In: Shary-gina, N., Veith, H. (eds.) CAV 2013. LNCS, vol. 8044, pp. 1–35. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-39799-8\_1
- 18. Passmore, G., et al.: The Imandra automated reasoning system. In: Proceedings of the IJCAR (2020, to appear)
- Reger, G., Voronkov, A.: Induction in saturation-based proof search. In: Fontaine,
   P. (ed.) CADE 2019. LNCS (LNAI), vol. 11716, pp. 477–494. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-29436-6\_28
- 20. Reger, G., Voronkov, A.: Induction in Saturation-Based Proof Search. EasyChair Smart Slide (2020). https://easychair.org/smart-slide/slide/hXmP
- Reynolds, A., Blanchette, J.C.: A decision procedure for (Co)datatypes in SMT solvers. In: Felty, A.P., Middeldorp, A. (eds.) CADE 2015. LNCS (LNAI), vol. 9195, pp. 197–213. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-21401-6\_13
- 22. Reynolds, A., Kuncak, V.: Induction for SMT solvers. In: D'Souza, D., Lal, A., Larsen, K.G. (eds.) VMCAI 2015. LNCS, vol. 8931, pp. 80–98. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-46081-8-5

- 23. Rybina, T., Voronkov, A.: A decision procedure for term algebras with queues. ACM Trans. Comput. Log. **2**(2), 155–181 (2001)
- 24. Sonnex, W., Drossopoulou, S., Eisenbach, S.: Zeno: an automated prover for properties of recursive data structures. In: Flanagan, C., König, B. (eds.) TACAS 2012. LNCS, vol. 7214, pp. 407–421. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-28756-5\_28
- Stump, A., Sutcliffe, G., Tinelli, C.: StarExec: a cross-community infrastructure for logic solving. In: Demri, S., Kapur, D., Weidenbach, C. (eds.) IJCAR 2014. LNCS (LNAI), vol. 8562, pp. 367–373. Springer, Cham (2014). https://doi.org/10. 1007/978-3-319-08587-6\_28
- 26. Voronkov, A.: AVATAR: the architecture for first-order theorem provers. In: Biere, A., Bloem, R. (eds.) CAV 2014. LNCS, vol. 8559, pp. 696–710. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-08867-9\_46