

# Datatypes with Shared Selectors

Andrew Reynolds<sup>1</sup>, Arjun Viswanathan<sup>1</sup>, Haniel Barbosa<sup>1</sup>,  
Cesare Tinelli<sup>1</sup>, and Clark Barrett<sup>2</sup>

<sup>1</sup> University of Iowa, Iowa City, USA

<sup>2</sup> Department of Computer Science, Stanford University

**Abstract.** We introduce a new theory of algebraic datatypes where selector symbols can be shared between multiple constructors, thereby reducing the number of terms considered by current SMT-based solving approaches. We show that the satisfiability problem for the traditional theory of algebraic datatypes can be reduced to problems where selectors are mapped to shared symbols based on a transformation provided in this paper. The use of shared selectors addresses a key bottleneck for an SMT-based enumerative approach to the Syntax-Guided Synthesis (SyGuS) problem. Our experimental evaluation of an implementation of the new theory in the SMT solver *cvc4* on syntax-guided synthesis and other domains provides evidence that the use of shared selectors improves state-of-the-art SMT-based approaches for constraints over algebraic datatypes.

## 1 Introduction

Algebraic datatypes, also known as inductive or recursive datatypes, are composite types commonly used for expressing finite data structures in computer science applications, such as lists or trees. Reasoning efficiently about (algebraic) datatypes is thus paramount in such fields as program analysis and verification, which has led to numerous approaches for automating solving in this setting. In this paper, we follow the semantic approach introduced by Barrett et al. [10], which is generally the basis for datatype decision procedures in satisfiability modulo theories (SMT) solvers [11].

In semantic presentations of the theory of algebraic datatypes [10, 22], a datatype is an absolutely free algebra over a signature of function symbols called *constructors*; the immediate subterms of a datatype value are accessed with function symbols called *selectors*, or *projections*, which are specific for each constructor and its arguments. Datatypes also have *discriminators*, or *testers*, associated with each constructor. They are predicates indicating whether a given datatype value was built with a specific constructor.

The satisfiability of quantifier-free formulas in the theory of algebraic datatypes is decidable. A basic decision procedure for this problem [10, 22] used by a number of SMT solvers operates by *progressively unrolling* datatypes: it tries to satisfy constraints by guessing top-level constructors in order to build values for the constraint variables incrementally. Concretely, if  $x$  is a datatype variable and  $c$  is an  $n$ -ary constructor for the datatype, the procedure may guess the equality constraint  $x \approx c(x_1, \dots, x_n)$  where  $x_1, \dots, x_n$  are fresh variables. If such a choice leads to an inconsistency, the procedure backtracks and tries different constructors until it determines that the constraints are satisfiable or no more choices are possible. During this process, lemmas in the form of

quantifier-free clauses may be learned by the procedure that prevent the procedure from making guesses already shown to be infeasible. However, these lemmas may include selectors, and because each selector is associated with only a single constructor, the generality and hence the usefulness of such lemmas is limited.

To address this limitation we introduce a new (formulation of the) theory of datatypes that allows certain selectors to be shared by multiple constructors. This way, information previously acquired when reasoning with a constructor, i.e., the learned lemmas on the applications of its selectors, can be reused when an argument of the same type is considered in another constructor. We illustrate this point with the following example.

*Example 1.* Consider a binary tree whose internal nodes store either one or two integer values, and whose leaves store both a Boolean and an integer value. A datatype **Tree** modeling this data structure has three constructors: one ( $N_1$ ) taking an integer and two **Tree** elements as arguments, another ( $N_2$ ) taking two integers and two **Tree** elements as arguments, and a third ( $L$ ) taking as arguments a Boolean and an integer element. We write this datatype in the following BNF-style notation:

$$\mathbf{Tree} = N_1(\mathbf{Int}, \mathbf{Tree}, \mathbf{Tree}) \mid N_2(\mathbf{Int}, \mathbf{Int}, \mathbf{Tree}, \mathbf{Tree}) \mid L(\mathbf{Bool}, \mathbf{Int})$$

We assume each constructor has selectors associated with them. The *subfields* (i.e., the immediate subterms) of terms constructed by  $N_1$  are accessed, respectively, by the selectors  $S^{N_1,1}$ ,  $S^{N_1,2}$ , and  $S^{N_1,3}$  of type  $\mathbf{Tree} \rightarrow \mathbf{Int}$ ,  $\mathbf{Tree} \rightarrow \mathbf{Tree}$  and  $\mathbf{Tree} \rightarrow \mathbf{Tree}$ . The selectors for the other constructors are similar. We also assume each constructor is associated with a tester predicate, i.e.  $isN_1$ ,  $isN_2$ , and  $isL$ , each of which takes a **Tree** as an argument. Given term  $t$  of type **Tree**, consider the following set of clauses:

$$\{ \neg isN_1(t) \vee S^{N_1,1}(t) \geq 0, \neg isL(t) \vee S^{L,2}(t) \geq 0 \} \quad (1)$$

The first clause states that when  $t$  has top symbol  $N_1$ , its first subfield (which is of type **Int**) is non-negative. Similarly, the second says that when  $t$  has top symbol  $L$ , its second subfield is non-negative.

Consider now a different kind of selector symbol  $S^{\mathbf{Int},1}$  of type  $\mathbf{Tree} \rightarrow \mathbf{Int}$  which maps each value of type **Tree** to the first (i.e., leftmost) subfield of  $t$  of type **Int**, *regardless* of the top constructor symbol of  $t$ . We will refer to such selectors as *shared selectors*. While nine selectors in the standard sense are necessary for **Tree**, five shared selectors suffice to access all possible subfields of a value of type **Tree**: two to access the **Tree** subfields, two to access the **Int** subfields, and one to access the **Bool** subfield of  $L$ . In particular, clause set (1) can be rewritten as follows using only one shared selector:

$$\{ \neg isN_1(t) \vee S^{\mathbf{Int},1}(t) \geq 0, \neg isL(t) \vee S^{\mathbf{Int},1}(t) \geq 0 \} \quad (2)$$

stating that when  $t$  has top symbol  $N_1$  or  $L$ , its first integer child is non-negative. •

In Example 1, the second set of clauses has one unique arithmetic constraint whereas the first set has two. In practice, reducing the number of unique constraints can substantially improve the performance of SMT solvers. Our experiments show that shared selectors lead to a significant reduction in the number of unique constraints for several

classes of benchmarks from real applications, with resulting SMT solver performance improvements that are proportional to the magnitude of this reduction.

*Contributions* We introduce a conservative extension of the (generic) theory of algebraic datatypes that features shared selectors. We show how using shared selectors instead of standard (unshared) selectors can improve the performance of current satisfiability procedures for the theory and also, as a result, the performance of procedures for syntax-guided synthesis. Specifically:

1. We formalize the new theory and show that constraints in the original signature can be reduced to equisatisfiable constraints whose selectors are all shared selectors. We present a decision procedure for the satisfiability of quantifier-free formulas in this theory as a natural modification of an earlier procedure for datatypes [22].
2. We provide details on an SMT-based approach for syntax-guided synthesis [24], and demonstrate how it can significantly benefit from native support in the SMT solver for a theory of datatypes with shared selectors.
3. We present an extensive experimental evaluation of our implementation in the SMT solver `cvc4` [7] on benchmarks from SMT-LIB [8] and from the most recent edition of SyGuS-COMP [3], the syntax-guided synthesis competition. This evaluation shows that shared selectors can reduce the number of terms introduced during solving, thus leading to more solved problems with respect to the state of the art.

## 2 Preliminaries

Our setting is a many-sorted classical first-order logic similar in essence to the one adopted by the SMT-LIB standard [9]. A signature  $\Sigma = (\mathcal{Y}, \mathcal{F})$  consists of a set  $\mathcal{Y}$  of first-order types, or *sorts*, and a set  $\mathcal{F}$  of first-order function symbols over these types. Each symbol  $f \in \mathcal{F}$  is associated with a list  $\tau_1, \dots, \tau_n$  of argument types and a return type  $\tau$ , written  $f : \tau_1 \times \dots \times \tau_n \rightarrow \tau$  or just  $f : \tau$  if  $n = 0$ . The function  $\text{arity}(f)$  returns  $n$ . We assume that any signature contains a **Bool** type and constants `true`, `false`; a family  $(\approx : \tau \times \tau \rightarrow \mathbf{Bool})_{\tau \in \mathcal{Y}}$  of equality symbols; a family  $(\text{ite} : \mathbf{Bool} \times \tau \times \tau \rightarrow \tau)_{\tau \in \mathcal{Y}}$  of *if-then-else* symbols; and the Boolean connectives  $\neg, \wedge, \vee$  with their expected types. Function symbols of **Bool** return type play the role of predicate symbols.

Typed terms are built as usual over function symbols from  $\mathcal{F}$  and typed variables from a fixed family  $(\mathcal{V}_\tau)_{\tau \in \mathcal{Y}}$  of pairwise-disjoint infinite sets. Formulas are terms of type **Bool**. The syntax  $t \neq u$  is short for  $\neg(t \approx u)$ . We reserve the names  $x, y, z$  for variables;  $r, s, t, u$  for terms (which may be formulas); and  $\varphi, \psi$  for formulas. We use the symbol  $\approx$  for equality at the meta-level. The *set of all terms* occurring in a term  $t$  is denoted by  $\mathbf{T}(t)$ . When convenient, we write an enumeration of (meta)symbols  $a_1, \dots, a_n$  as  $\bar{a}$ . If  $b_1, \dots, b_k$  is another enumeration,  $\bar{a}\bar{b}$  denotes the enumeration  $a_1, \dots, a_n, b_1, \dots, b_k$ .

Given a signature  $\Sigma = (\mathcal{Y}, \mathcal{F})$ , a  $\Sigma$ -interpretation  $\mathcal{I}$  maps: each  $\tau \in \mathcal{Y}$  to a non-empty set  $\tau^{\mathcal{I}}$ , the *domain* of  $\tau$  in  $\mathcal{I}$ , with  $\mathbf{Bool}^{\mathcal{I}} = \{\top, \perp\}$ ; each  $x \in \mathcal{V}_\tau$  to an element of  $\tau^{\mathcal{I}}$ ; each  $f \in \mathcal{F}$  s.t.  $f : \tau_1 \times \dots \times \tau_n \rightarrow \tau$  to a total function  $u^{\mathcal{I}} : \tau_1^{\mathcal{I}} \times \dots \times \tau_n^{\mathcal{I}} \rightarrow \tau^{\mathcal{I}}$  when  $n > 0$  and to an element of  $\tau^{\mathcal{I}}$  when  $n = 0$ . A satisfiability relation between  $\Sigma$ -interpretations and  $\Sigma$ -formulas is defined inductively as usual.

A *theory* is a pair  $\mathcal{T} = (\Sigma, \mathbf{I})$  where  $\Sigma$  is a signature and  $\mathbf{I}$  is a non-empty class of  $\Sigma$ -interpretations, the *models of  $\mathcal{T}$* , that is closed under variable reassignment (i.e., every  $\Sigma$ -interpretation that differs from one in  $\mathbf{I}$  only in how it interprets the variables is also in  $\mathbf{I}$ ) and isomorphism. A  $\Sigma$ -formula  $\varphi$  is  $\mathcal{T}$ -satisfiable (respectively  $\mathcal{T}$ -unsatisfiable) if it is satisfied by some (resp., no) interpretation in  $\mathbf{I}$ . A satisfying interpretation for  $\varphi$  *models (or is a model of)  $\varphi$* . A formula  $\varphi$  is *valid in  $\mathcal{T}$*  (or  $\mathcal{T}$ -*valid*) if every model of  $\mathcal{T}$  is a model of  $\varphi$ .

### 3 Theory of Datatypes With Shared Selectors

In this section, we consider a theory  $\mathcal{D}$  of algebraic datatypes over some signature  $\Sigma = (\mathcal{Y}, \mathcal{F})$  and then extend it conservatively to an expanded signature *with shared selectors*. The terms of  $\mathcal{D}$  are quantifier-free. As a technical convenience, we treat free variables as constants in a suitable expansion of  $\Sigma$ . The types of  $\mathcal{D}$  are partitioned into a set of *datatypes*  $\mathcal{Y}_{\text{dt}}$ , and a set of other types  $\mathcal{Y}_{\text{ord}}$ . We use the metavariables  $\delta, \epsilon$  to refer to datatypes and  $\tau, \nu$  for arbitrary first-order types. Each datatype  $\delta$  is equipped with one or more *constructors*, distinguished function symbols from  $\mathcal{F}$  with return type  $\delta$ . For every argument  $k$  of a constructor  $C : \tau_1 \dots, \tau_n \rightarrow \delta$  for  $\delta$ , we assume  $\mathcal{F}$  contains a (*standard*) *selector*  $S_{\delta}^{C,k} : \delta \rightarrow \tau_k$ . We omit  $\delta$  from the selector name when it is understood or not important. We refer the reader to the SMT-LIB 2 reference document [9] or Barrett et al. [10] for a formal definition of this theory.<sup>1</sup> We recall salient properties of its symbols as needed.

To start, each model of the theory, when reduced to the constructors of a datatype in the theory, is isomorphic to a term (or *Herbrand*) algebra. Concretely, this means that if  $\delta \in \mathcal{Y}_{\text{dt}}$  is a datatype whose constructors are  $\{C_1, \dots, C_m\}$ , then the following formulas are all  $\mathcal{D}$ -valid for all distinct  $i, j \in \{1, \dots, m\}$

$$\forall x_1, \dots, x_{p_i}, z_1, \dots, z_{q_i}. C_i(x_1, \dots, x_{p_i}) \neq C_j(z_1, \dots, z_{q_i}) \quad (\textit{Distinctness})$$

$$\forall x_1, \dots, x_{p_i}, z_1, \dots, z_{p_i}. C_i(x_1, \dots, x_{p_i}) \approx C_i(z_1, \dots, z_{p_i}) \rightarrow x_1 \approx z_1 \wedge \dots \wedge x_{p_i} \approx z_{p_i} \quad (\textit{Injectivity})$$

$$\forall x. \text{is}C_1(x) \vee \dots \vee \text{is}C_m(x) \quad (\textit{Exhaustiveness})$$

Above, we write  $\text{is}C_i(t)$  to denote the predicate that holds if and only if the top symbol of  $t$  is  $C_i$ . Strictly speaking, we do not need to extend our signature with the tester symbols  $\text{is}C$  since a term of the form  $\text{is}C(t)$  can be considered an abbreviation for the equality  $t \approx C(S^{C,1}(t), \dots, S^{C,n}(t))$  where  $n = \text{arity}(C)$ .

Interpretations must also respect *acyclicity*, which states that constructor terms cannot be equal to any of their proper subterms.

Since all models of  $\mathcal{D}$  interpret a datatype  $\delta$  in the same way modulo isomorphism, we will say that  $\delta$  is *finite* if its interpretation is a finite set. For simplicity, we will assume that *every type  $\tau$  in  $\mathcal{D}$  that is not a datatype is interpreted as an infinite set in every model of  $\mathcal{D}$* . This is not a strong restriction in practice, since types with some fixed, finite cardinality  $k$  can be treated as datatypes with  $k$  nullary constructors.

<sup>1</sup> The two references differ on how they make selectors (which are naturally partial functions) total. We follow the SMT-LIB 2 standard here.

The relationship between an  $n$ -ary constructor  $C$  and each of its selectors  $S^{C,k}$  with  $k = 1, \dots, n$  is captured by the following  $\mathcal{D}$ -valid formula:

$$\forall x_1, \dots, x_{n_i}. S_{\delta}^{C,k}(C(x_1, \dots, x_{n_i})) \approx x_k \quad (\text{Standard selection})$$

### 3.1 Shared Selectors

We extend the signature of  $\mathcal{D}$  with additional selectors which we call *shared selectors* and denote as  $S_{\delta}^{\tau,k}$ , for each datatype  $\delta$  and type  $\tau$  in  $\mathcal{D}$  and each natural number  $k$ . Intuitively, a shared selector  $S_{\delta}^{\tau,k}$  for  $\delta$ , when applied to a  $\delta$ -term  $C(t_1, \dots, t_n)$  returns the  $k$ -th argument of  $C$  that has type  $\tau$ , if one exists.

*Example 2.* Consider again the **Tree** datatype introduced in Example 1. For term

$$t = N_1(1, N_2(2, 3, L(\text{true}, 4), L(\text{false}, 5)), L(\text{true}, 6))$$

the equalities  $S^{\text{Int},1}(t) \approx 1$ ,  $S^{\text{Int},2}(S^{\text{Tree},1}(t)) \approx 3$ , and  $S^{\text{Int},1}(S^{\text{Tree},2}(S^{\text{Tree},1}(t))) \approx 6$  are all valid in our extension of  $\mathcal{D}$  to shared selectors. •

To define shared selectors formally, let us first define a partial function *stoa* (for *selector to argument*) that takes as input a natural number  $k$ , a type  $\tau$ , and a constructor  $C$ , and returns the index of the  $k$ -th argument of  $C$  of type  $\tau$ . We leave *stoa* undefined if  $C$  has fewer than  $k$  arguments of type  $\tau$ .

*Example 3.* For the **Tree** datatype,  $\text{stoa}(1, \text{Int}, N_1) = 1$ ,  $\text{stoa}(2, \text{Tree}, N_1) = 3$  and  $\text{stoa}(1, \text{Int}, L) = 2$ , whereas  $\text{stoa}(2, \text{Int}, N_1)$ ,  $\text{stoa}(1, \text{Bool}, N_2)$ , and  $\text{stoa}(1, \text{Tree}, L)$  are undefined. •

More formally, in our extension of theory  $\mathcal{D}$  with shared selectors, which we also refer to as  $\mathcal{D}$  for convenience, the following holds for all datatypes  $\delta$ , constructors  $C$  of  $\delta$ , and shared selectors  $S^{\tau,k}$ , whenever  $\text{stoa}(k, \tau, C)$  is defined:

$$\forall x_1, \dots, x_n. S_{\delta}^{\tau,k}(C(x_1, \dots, x_n)) \approx x_i, \text{ where } i = \text{stoa}(k, \tau, C) \quad (\text{Shared selection})$$

It is not difficult to argue that every  $\Sigma$ -formula  $\varphi$  without shared selectors is valid in the extended theory if and only if it is valid in the original theory.

### 3.2 From Standard Selectors to Shared Selectors

The satisfiability problem for *constraints*, i.e., finite sets of literals, over the original theory of datatypes (without shared selectors) is decidable [10]. In this section, we introduce a transformation  $\mathcal{H}$  that reduces arbitrary constraints in our extended theory  $\mathcal{D}$ , which may have both standard and shared selectors, to constraints with no standard selectors. Applying this transformation as an initial step allows us to determine the satisfiability of arbitrary  $\Sigma$ -constraints by means of a decision procedure for  $\Sigma$ -constraints without standard selectors.

To define this transformation, let  $\max_{\Sigma}$  denote some natural number that is greater than the arity of all constructors in  $\Sigma$ . We define the dual of the *stoa* function from

$$\begin{aligned}
\mathcal{H}(t, M) &= \text{match } t \text{ with} \\
&\quad x \rightarrow x \\
\mathbf{C}(t_1, \dots, t_n) &\rightarrow \mathbf{C}(\mathcal{H}(t_1, M), \dots, \mathcal{H}(t_n, M)) \\
\mathbb{S}_\delta^{\tau, k}(t_1) &\rightarrow \mathbb{S}_\delta^{\tau, k}(\mathcal{H}(t_1, M)) \\
\mathbb{S}_\delta^{\mathbf{C}, k}(t_1) &\rightarrow \begin{cases} \mathbb{S}_\delta^{\tau, \text{atos}(\tau, \mathbf{C}, k)}(\mathcal{H}(t_1, M)) & \text{if } M(t_1) = \mathbf{C} \\ \mathbb{S}_\delta^{\tau, \text{err}(\mathbf{C}, k)}(\mathcal{H}(t_1, M)) & \text{otherwise} \end{cases} \\
&\quad \text{where } \mathbb{S}_\delta^{\mathbf{C}, k} : \delta \rightarrow \tau
\end{aligned}$$

Fig. 1: Definition of  $\mathcal{H}(t, M)$

Subsection 3.1 as the partial function  $\text{atos}$  (for *argument to selector*) that takes as input a type  $\tau$ , a constructor  $\mathbf{C} : \tau_1 \times \dots \times \tau_n \rightarrow \delta$ , and a natural number  $k \leq n$ , and returns the number of times  $\tau$  occurs in  $\tau_1, \dots, \tau_k$ .

Figure 1 defines the transformation  $\mathcal{H}$ , which takes as arguments a  $\Sigma$ -term  $t$  and a mapping  $M$ . The latter consists of one entry of the form  $s \mapsto \mathbf{C}$  for each datatype term  $s$  in  $\mathbf{T}(t)$  where  $\mathbf{C}$  is one of the constructors for the type of  $s$ . Without loss of generality, we assume that all applications of shared selectors  $\mathbb{S}_\delta^{\tau, k}$  occurring in  $t$  are such that  $k < \max_\Sigma$ . The transformation  $\mathcal{H}$  leaves variables unchanged; for terms whose top symbol is a constructor or a shared selector,  $\mathcal{H}$  behaves homorphically. For terms  $t$  with a standard selector  $\mathbb{S}_\delta^{\mathbf{C}, k} : \delta \rightarrow \tau$  as top symbol, we distinguish whether the argument  $t_1$  is mapped to  $\mathbf{C}$  by  $M$  or not. In the first case, we replace  $\mathbb{S}_\delta^{\mathbf{C}, k}$  by the shared selector  $\mathbb{S}_\delta^{\tau, \text{atos}(\tau, \mathbf{C}, k)}$ . In the second case, we replace  $\mathbb{S}_\delta^{\mathbf{C}, k}$  by the shared selector  $\mathbb{S}_\delta^{\tau, \text{err}(\mathbf{C}, k)}$ , where  $\text{err}$  is a function that takes as arguments a constructor and a  $k$  such that  $1 \leq k \leq \text{arity}(\mathbf{C})$ , and returns a natural number. Additionally,  $\text{err}$  has the properties:

1. If  $\mathbf{C}_1 \neq \mathbf{C}_2$  or  $k_1 \neq k_2$ , then  $\text{err}(\mathbf{C}_1, k_1) \neq \text{err}(\mathbf{C}_2, k_2)$ , and
2.  $\text{err}(\mathbf{C}_1, k) \geq \max_\Sigma$ .

We use the function  $\text{err}$  in this transformation to introduce shared selectors that are unique to the pair  $(\mathbf{C}, k)$ , as guaranteed by Property 1 above, and whose return value is undefined, as guaranteed by Property 2. In either case,  $\mathcal{H}$  is applied recursively to  $t_1$ .

We extend  $\mathcal{H}$  to sets of equalities and disequalities  $E$  as follows:

$$\begin{aligned}
\mathcal{H}(E, M) &= \{ \mathcal{H}(t_1, M) \approx \mathcal{H}(t_2, M) \mid t_1 \approx t_2 \in E \} \cup \\
&\quad \{ \mathcal{H}(t_1, M) \not\approx \mathcal{H}(t_2, M) \mid t_1 \not\approx t_2 \in E \} \cup \{ \text{isC}(t) \mid t \mapsto \mathbf{C} \in M \}
\end{aligned}$$

In other words, for each (dis)equality, we include the corresponding constraint where the transformation is applied to both its terms. We add to this set an application of the discriminator for  $\mathbf{C}$  to  $t$  for each  $t \mapsto \mathbf{C}$  in the mapping  $M$ .

*Example 4.* Consider again the **Tree** datatype from Example 1. Let:

$$E = \{x \approx \mathbf{N}_1(2, y), \mathbf{S}^{\mathbf{N}_1, 2}(x), \mathbf{S}^{\mathbf{N}_1, 1}(x) \approx 2, \mathbf{S}^{\mathbf{L}, 2}(x) \not\approx 0\} \text{ and } M = \{x \mapsto \mathbf{N}_1, y \mapsto \mathbf{L}\}$$

Then,  $\mathcal{H}(E, M)$  is the set:

$$\{x \approx N_1(2, y, \mathbf{S}^{\text{Tree},1}(x)), \mathbf{S}^{\text{Int},1}(x) \approx 2\} \cup \{\mathbf{S}^{\text{Int},\text{err}(L,2)}(x) \neq 0\} \cup \{\text{isN}_1(x), \text{isL}(y)\}$$

Since  $M$  maps  $x$  to  $N_1$ , the standard selector application  $\mathbf{S}^{\text{N}_1,2}(x)$  is converted to the shared selector application  $\mathbf{S}^{\text{Tree},1}(x)$ , whereas  $\mathbf{S}^{\text{L},2}(x)$  is converted to  $\mathbf{S}^{\text{Int},\text{err}(L,2)}(x)$ . •

The following theorem states the key property of the transformation  $\mathcal{H}$ , namely that a set of arbitrary  $\Sigma$ -constraints  $E$  is satisfiable if and only if there exists some mapping  $M$  for which  $\mathcal{H}(E, M)$  is satisfiable. The full proof of this statement is available in an extended version of this paper [26].

**Theorem 1.**  *$E$  is  $\mathcal{D}$ -satisfiable iff  $\mathcal{H}(E, M)$  is  $\mathcal{D}$ -satisfiable for some  $M$ .*

*Proof. (Sketch)* We split the statement into its two implications. The proof relies on the construction of a mapping  $M$  from a model of  $E$ .

“ $\Rightarrow$ ”: If  $E$  is satisfied by some  $\Sigma$ -model  $\mathcal{I}$  of  $\mathcal{D}$ , there exists a mapping  $M_{\mathcal{I}}$  and  $\Sigma$ -model  $\mathcal{J}$  of  $\mathcal{D}$  such that  $\mathcal{H}(E, M_{\mathcal{I}})$  is satisfied by  $\mathcal{J}$ . We show this by a particular construction for  $M_{\mathcal{I}}$  and  $\mathcal{J}$ . Let the mapping  $M_{\mathcal{I}}$  be  $\{t \mapsto C \mid \mathcal{I} \models \text{isC}(t), t \in \mathbf{T}(E)\}$ . Construct  $\mathcal{J}$  as follows. First, all types  $\tau$  and constructors are interpreted by  $\mathcal{J}$  the same way as in  $\mathcal{I}$ . Furthermore, we interpret all variables and standard selectors in  $\mathcal{J}$  the same as in  $\mathcal{I}$ . It remains to state how shared selectors are interpreted in  $\mathcal{J}$ . Notice that our transformation generates shared selectors of the form  $\mathbf{S}_{\delta}^{\tau, \text{err}(C, k)}$ . We distinguish these in the following construction.

$$\mathbf{S}_{\delta}^{\tau, \text{err}(C, k)\mathcal{J}} = \mathbf{S}_{\delta}^{\text{C}, k\mathcal{I}}, \text{ and } \mathbf{S}_{\delta}^{\tau, k\mathcal{J}} = \mathbf{S}_{\delta}^{\tau, k\mathcal{I}} \text{ for all other shared selectors.}$$

The above construction is well-defined due to our definition of  $\text{err}$ . In particular,  $\text{err}(C, k)$  is defined uniquely for each (constructor, natural number) pair. By this construction, it can be shown that  $\mathcal{H}(t, M_{\mathcal{I}})^{\mathcal{J}} = t^{\mathcal{I}}$  by structural induction on  $t$  for all  $t \in \mathbf{T}(E)$ . Since  $\mathcal{I}$  satisfies  $E$  and since  $\mathcal{H}(t, M_{\mathcal{I}})^{\mathcal{J}} = t^{\mathcal{I}}$  for all terms  $t \in \mathbf{T}(E)$ , we have that  $\mathcal{J}$  satisfies the equalities and disequalities in  $\mathcal{H}(E, M_{\mathcal{I}})$  of the form  $(\neg)\mathcal{H}(t_1, M_{\mathcal{I}}) \approx \mathcal{H}(t_2, M_{\mathcal{I}})$ . By construction of  $M_{\mathcal{I}}$ , we have that  $\mathcal{J}$  satisfies the constraints in  $\mathcal{H}(E, M_{\mathcal{I}})$  of the form  $\text{isC}(t)$  where  $t \mapsto C \in M_{\mathcal{I}}$ . Hence,  $\mathcal{J}$  satisfies  $\mathcal{H}(E, M_{\mathcal{I}})$ .

“ $\Leftarrow$ ”: If  $\mathcal{H}(E, M)$  is satisfied by some  $\Sigma$ -model  $\mathcal{J}$  of  $\mathcal{D}$  for some mapping  $M$ , then  $E$  is satisfied by some  $\Sigma$ -model  $\mathcal{I}$  of  $\mathcal{D}$ . We show this by constructing  $\mathcal{I}$  as follows. First, all types, constructors, variables and have the same interpretation in  $\mathcal{I}$  as in  $\mathcal{J}$ . Furthermore, all shared selectors have the same interpretation in  $\mathcal{I}$  as in  $\mathcal{J}$ . We interpret standard selectors in  $\mathcal{I}$  as follows.

$$\mathbf{S}_{\delta}^{\text{C}, k}\mathcal{I} = \begin{cases} \mathbf{S}_{\delta}^{\tau, \text{atos}(\tau, \text{C}, k)}(t)^{\mathcal{I}} & \text{if } M(t) = \text{C} \\ \mathbf{S}_{\delta}^{\tau, \text{err}(\text{C}, k)}(t)^{\mathcal{I}} & \text{otherwise} \end{cases}$$

Similar to the first part, it can be shown that  $t^{\mathcal{I}} = \mathcal{H}(t, M)^{\mathcal{J}}$  by structural induction on  $t$  for all  $t \in \mathbf{T}(E)$ . Since  $\mathcal{J}$  satisfies  $\mathcal{H}(E, M)$  and  $t^{\mathcal{I}} = \mathcal{H}(t, M)^{\mathcal{J}}$  for all  $t \in \mathbf{T}(E)$ , we have that  $\mathcal{I}$  satisfies the equalities and disequalities in  $\mathcal{H}(E, M)$  of the form  $(\neg)\mathcal{H}(t_1, M) \approx \mathcal{H}(t_2, M)$ . Furthermore, since  $\mathcal{J}$  satisfies the constraints  $\text{isC}(t)$  for all  $t \mapsto C \in M$  and since  $t^{\mathcal{I}} = t^{\mathcal{J}}$ , we have that  $\mathcal{I}$  satisfies these constraints as well. Thus,  $\mathcal{I}$  satisfies  $\mathcal{H}(E, M)$ .  $\square$

**Corollary 1.** For some index sets  $I$  and  $J$ , and set  $E$  of  $\Sigma$ -literals without standard selectors, let

$$E_0 = E \cup \{ \mathbf{S}^{C_{j_i}, k_i}(x_i) \approx y_i \mid i \in I, j \in J \} \quad \text{and}$$

$$E_1 = E \wedge \{ \text{ite}(\text{isC}_{j_i}(x_i), \mathbf{S}^{\tau, \text{atos}(\tau, C_{j_i}, k_i)}(x_i), \mathbf{S}^{\tau, \text{err}(C_{j_i}, k_i)}(x_i)) \approx y_i \mid i \in I, j \in J \}.$$

The sets  $E_0$  and  $E_1$  are equisatisfiable in  $\mathcal{D}$ .

Using this corollary, we can reduce (possibly after some literal flattening) the satisfiability of an arbitrary set of  $\Sigma$ -constraints  $E_0$  to a set of  $\Sigma$ -constraints  $E_1$  not containing standard selectors. In particular, our implementation in `cvc4` replaces each application of the form  $\mathbf{S}^{C_{j_i}, k_i}(x_i)$  by the term  $\text{ite}(\text{isC}_{j_i}(x_i), \mathbf{S}^{\tau, \text{atos}(\tau, C_{j_i}, k_i)}(x_i), \mathbf{S}^{\tau, \text{err}(C_{j_i}, k_i)}(x_i))$  during a preprocessing pass on the input formula.

## 4 Decision Procedure for Datatypes with Shared Selectors

This section describes a tableau-like calculus for deciding constraint satisfiability in  $\mathcal{D}$ , with constraint variables interpreted existentially. The calculus is parametrized by the theory's signature  $\Sigma$ . By the results of the previous section, we can restrict with no loss of generality the input language to sets of equalities and disequalities between  $\Sigma$ -terms with no standard selectors and no discriminators. Since our calculus is based on similar calculi for datatypes that have been presented in detail in previous work [10, 22], we focus on our modifications to accommodate shared selectors.

The derivation rules of the calculus operate on a current set  $E$  of constraints as specified in Figure 2. A derivation rule can be applied to  $E$  if its premises are met. Some of those premises check membership in the *congruence closure*  $E^*$  of  $E$ , the smallest superset of  $E$  that is closed under entailment in the theory of equality.<sup>2</sup> A rule's conclusion either modifies  $E$  or replaces it by  $\perp$  to indicate unsatisfiability. There, the notation  $E, t \approx s$  abbreviates  $E \cup \{t \approx s\}$ ; the notation  $\bar{t} \approx \bar{u}$  stands for the set of equalities between the corresponding elements of  $\bar{t}$  and  $\bar{u}$ . The `SPLIT` rule has multiple alternative conclusions, denoting branching.

A rule application is *redundant* if (one of) its conclusion(s) leaves  $E$  unchanged. The rules are applied to build a *derivation tree*, i.e., a tree whose nodes are finite sets of (dis)equalities, with an *initial constraint set*  $E_0$  as its root and child nodes obtained by a non-redundant rule application to their parent. We say that  $E_0$  has a derivation tree  $D$  if  $D$  is a derivation tree with root  $E_0$ . A node is *saturated* if it admits only redundant rule applications. A derivation tree is *closed* if all of its leaf nodes are  $\perp$ . Intuitively, a derivation tree is generated progressively from  $E_0$  by applying a derivation rule to a leaf node. The rules are applied until the derivation tree becomes closed (indicating that the initial set  $E_0$  is  $\mathcal{D}$ -unsat) or contains a saturated leaf node (indicating that  $E_0$  is  $\mathcal{D}$ -sat).

In the calculus, all reasoning based on the general properties of equality is encapsulated in the rule `CONFLICT`, which detects that congruent terms are forced to be distinct. The remaining rules perform datatype reasoning proper, with `DECOMPOSE` computing a downward equality closure based on the injectivity of constructors and `CLASH` detecting

<sup>2</sup> Such tests are effective by well-known results about the theory of equality [6].



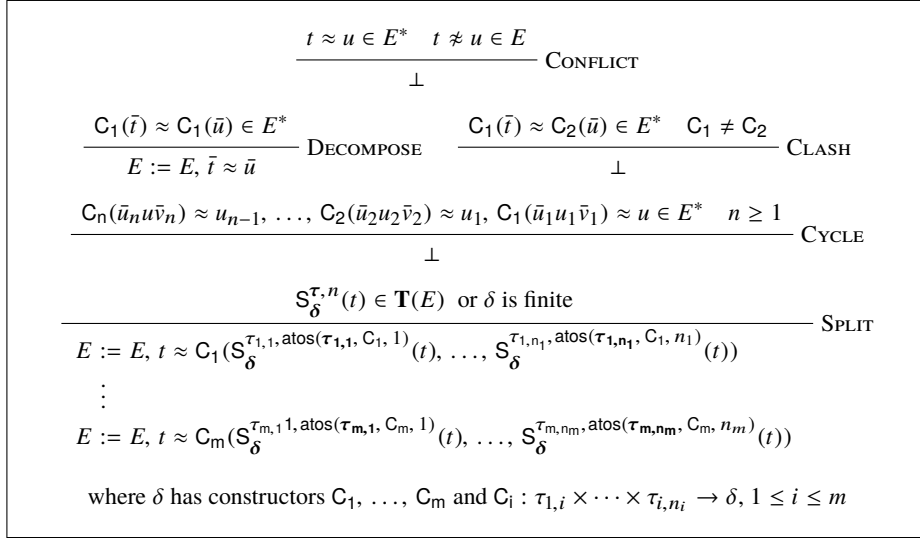


Fig. 2: Derivation rules.

failures based on their distinctness. The CYCLE rule recognizes when a constructor term must be equivalent to one of its subterms, which is forbidden in all models of the theory.

The calculus also incrementally unrolls terms by branching on different constructors, with the SPLIT rule performing case distinctions on constructors for various terms occurring in  $E$ . The main modification from the previous calculi for the theory of datatypes is that this SPLIT rule operates on shared selectors. Its application can be seen as an on-the-fly transformation from standard to shared selectors as described in Section 3.2. Indeed, for each constructor  $C_i$  in its conclusion, the following holds with a mapping  $M$  such that  $M(t) = C_i$ :

$$S_\delta^{\tau_i, \text{atos}(\tau_i, C_i, 1)}(t) = \mathcal{H}(S_\delta^{C_i, 1}(t), M), \dots, S_\delta^{\tau_{n_i}, \text{atos}(\tau_{n_i}, C_i, n_i)}(t) = \mathcal{H}(S_\delta^{C_i, n_i}(t), M)$$

Any derivation strategy for the calculus that does not stop until it generates a closed tree or a saturated node yields a decision procedure for the  $\mathcal{D}$ -satisfiability of sets of  $\Sigma$ -literals. We prove this similarly to previous work [10, 22], but using shared selectors and in the simpler setting obtained by assuming the availability of a congruence closure procedure. The full proofs are available in an extended version of this paper [26].

**Proposition 1 (Termination).** *All derivation trees in the calculus are finite.*

**Proposition 2 (Refutation Soundness).** *If a constraint set  $E_0$  has a closed derivation tree, then it is  $\mathcal{D}$ -unsatisfiable.*

**Proposition 3 (Solution Soundness).** *If a constraint set  $E_0$  has a derivation tree with a saturated node, then it is  $\mathcal{D}$ -satisfiable.*

**Theorem 2.** *Constraint satisfiability in the theory  $\mathcal{D}$  of datatypes with (standard and) shared selectors is decidable.*

## 5 Using Shared Selectors for Syntax-Guided Synthesis

In this section, we show how the theory of datatypes with shared selectors can substantially improve the performance of an approach by Reynolds et al. [24] for performing *syntax-guided synthesis* (SyGuS) [1] directly within an SMT solver.

Syntax-guided synthesis is the problem of automatically synthesizing a function that satisfies a given specification, but with the addition of explicit syntactic restrictions on the solution space. These restrictions specify that the function must be built with selected operators over basic types (such as arithmetic and Boolean operators) and belong to the language generated by a given grammar. Grammars allow users to specify formally a set of candidates for the desired function, thus reducing the search effort of a SyGuS solver.

More technically, a syntax-guided synthesis problem for a function  $f$  in a background theory  $T$  of the basic types consists of:

1. a set of semantic restrictions, or specification, given by a (second-order)  $T$ -formula of the form  $\exists f. \forall \bar{x}. \varphi[f, \bar{x}]$ , and
2. a set of syntactic restrictions on the solutions for  $f$ , given by a grammar  $R$ .

A solution for  $f$  is a lambda term  $\lambda \bar{y}. e$  of the same type as  $f$ , such that (i)  $\forall \bar{x}. \varphi[\lambda \bar{y}. e, \bar{x}]$  is valid in  $T$  (modulo beta-reductions) and (ii)  $e$  is in the language generated by  $R$ .

cvc4 incorporates a SyGuS solver that automatically encodes the solution space of a SyGuS problem as a set of algebraic datatypes mirroring the problem's syntactic restrictions [24]. A deep embedding of the datatypes in the problem's background theory  $T$ , realized as a set of automatically generated axioms, provides a semantics for datatype values in terms of the semantic values in  $T$ .

*Example 5.* Consider the problem of synthesizing a binary function  $f$  over the integers such that  $f$  is commutative (i.e.,  $\exists f \forall xy. f(x, y) \approx f(y, x)$ ), and with the solution space for  $f$  defined by a context-free grammar  $R$  with start symbol  $A$  and production rules:

$$A \rightarrow x \mid y \mid 0 \mid 1 \mid A + A \mid A - A \mid \text{ite}(B, A, A) \qquad B \rightarrow A \geq A \mid A \approx A \mid \neg B$$

The following mutually recursive datatypes capture the grammar  $R$ . The datatypes themselves correspond to  $R$ 's non-terminals (e.g.,  $\mathbf{a}$  corresponds to  $A$ ), their constructors correspond to production rules (e.g.,  $X$  corresponds to  $A \rightarrow x$ ):

$$\begin{aligned} \mathbf{a} &= X \mid Y \mid \text{Zero} \mid \text{One} \mid \text{Plus}(\mathbf{a}, \mathbf{a}) \mid \text{Minus}(\mathbf{a}, \mathbf{a}) \mid \text{lte}(\mathbf{b}, \mathbf{a}, \mathbf{a}) \\ \mathbf{b} &= \text{Geq}(\mathbf{a}, \mathbf{a}) \mid \text{Eq}(\mathbf{a}, \mathbf{a}) \mid \text{Neg}(\mathbf{b}) \end{aligned}$$

Datatypes like the ones above are associated with the programs they represent through *evaluation functions* that map datatype values, expressed as variable-free constructor terms, to expressions over the basic types. For example, the evaluation function for  $\mathbf{a}$  is denoted by a function symbol  $\text{eval}_{\mathbf{a}} : \mathbf{a} \times \mathbf{Int} \times \mathbf{Int} \rightarrow \mathbf{Int}$ , and the specific term  $\text{eval}_{\mathbf{a}}(\text{Plus}(X, X), 2, 3)$  is interpreted as  $(x + x)\{x \mapsto 2, y \mapsto 3\} = 2 + 2 = 4$ . The evaluation functions are defined axiomatically by a set of quantified formulas that, in this case, can be handled by any SMT solver that, like cvc4, supports the combined

theory of datatypes, linear arithmetic, and uninterpreted functions. The SyGuS problem for  $f$  in this example can then be stated as the *first-order* formula:

$$\forall xy. \text{eval}_{\mathbf{a}}(d, x, y) \approx \text{eval}_{\mathbf{a}}(d, y, x) \quad (3)$$

where  $d$  is a fresh constant of type  $\mathbf{a}$ . This formula has models in which  $d$  is interpreted as `Zero` or `Plus(X, Y)`, which correspond to solutions  $f = \lambda xy. 0$  and  $f = \lambda xy. x + y$  for the original problem, respectively.<sup>3</sup> •

Since `cvc4` is a  $\text{DPLL}(T)$ -based solver [11], for a problem like the one in the example above, it will find a possible solution for  $d$  by first guessing its top constructor symbol with an application of the `SPLIT` rule from Figure 2. The effect of the rule is achieved in practice with the generation of *splitting lemmas* such as the following, which we write here with discriminators and standard selectors for simplicity:

$$\text{isX}(d) \vee \text{isY}(d) \vee \dots \vee \text{islte}(d) \quad (4)$$

$$\text{isX}(\text{S}^{\text{Plus},1}(d)) \vee \text{isY}(\text{S}^{\text{Plus},1}(d)) \vee \dots \vee \text{islte}(\text{S}^{\text{Plus},1}(d)) \quad (5)$$

$$\text{isGeq}(\text{S}^{\text{lte},1}(d)) \vee \text{isEq}(\text{S}^{\text{lte},1}(d)) \vee \text{isNeg}(\text{S}^{\text{lte},1}(d)) \quad (6)$$

$$\text{isX}(\text{S}^{\text{lte},2}(d)) \vee \text{isY}(\text{S}^{\text{lte},2}(d)) \vee \dots \vee \text{islte}(\text{S}^{\text{lte},2}(d)) \quad (7)$$

The solver will subsequently guess the top constructor for other subterms of  $d$ 's value. These guesses are represented symbolically by *selector chains*, i.e. zero or more applications of selectors to  $d$ ; for example,  $\text{S}^{\text{Plus},1}(d)$  is a selector chain that corresponds to the first *child* of  $d$  (if we think of the value of  $d$  as a tree) when  $d$  is an application of `Plus`;  $\text{S}^{\text{Plus},1}(\text{S}^{\text{Plus},1}(d))$  is a selector chain that corresponds to the first child of the first child of  $d$  when  $d$  and its first child are both applications of `Plus`; and so on.

The bottleneck in solving (3) is the large number of splitting lemmas for selector chains introduced during search which, depending on the datatypes involved, is often highly exponential. Our key observation is that datatypes generated by the SyGuS approach sketched above very often include constructors with arguments of the same type. In Example 5, both  $\mathbf{a}$  and  $\mathbf{b}$  have multiple constructors with arguments of type  $\mathbf{a}$ . Using shared selectors, we can reduce the number of selectors in the example from 7 to 3 for  $\mathbf{a}$  and from 5 to 3 for  $\mathbf{b}$ . Moreover, *using shared selectors in selector chains makes splitting lemmas relevant in multiple contexts*. For example, a splitting lemma for a selector chain  $\text{S}^{\mathbf{a},1}(d)$  is relevant when  $d$  is either `Plus`, `Minus` or `lte`; likewise  $\text{S}^{\mathbf{a},1}(\text{S}^{\mathbf{a},1}(d))$  is relevant when  $d$  and its first child of type  $\mathbf{a}$  are applications of either `Plus`, `Minus` or `lte`. Notice that by using the decision procedure for shared selectors from Section 4, lemmas (5) and (7) would be instead both provided to the SAT engine as:

$$\text{isX}(\text{S}^{\text{Int},1}(d)) \vee \text{isY}(\text{S}^{\text{Int},1}(d)) \vee \dots \vee \text{islte}(\text{S}^{\text{Int},1}(d))$$

Using shared selectors can lead to a reduction in the number of other kinds of lemmas as well. For instance, during synthesis `cvc4` implements *symmetry breaking* techniques to avoid spending time on multiple candidates that are all equivalent in  $T$  [24, 25]. Redundant candidates are avoided by adding blocking clauses to the SAT engine that are also expressed in terms of discriminators applied to selector chains.

<sup>3</sup> For a thorough description of this approach, see [24].

*Example 6.* Consider again the function  $f$ , grammar  $R$ , and datatypes  $\mathbf{a}$  and  $\mathbf{b}$  from Example 5. Assume that the solver considers  $X$  as a candidate solution for  $d$ , and later considers another candidate solution,  $\text{Plus}(X, \text{Zero})$ . Since the corresponding arithmetic terms  $x$  and  $x + 0$  are equivalent in integer arithmetic, the solver infers a *lemma template* of the form:

$$\neg \text{isPlus}(z) \vee \neg \text{isX}(\text{S}^{\text{Int},1}(z)) \vee \neg \text{isZero}(\text{S}^{\text{Int},2}(z))$$

to block a redundant candidate solution like (the one corresponding to)  $x + 0$ . This is achieved by instantiating the template with the substitution  $\{z \mapsto d\}$  for variable  $z$ . More interestingly,  $z$  can be instantiated with other selector chains to rule out *entire families* of redundant candidate solutions. For instance, the lemma obtained with  $\{z \mapsto \text{S}^{\text{Int},1}(d)\}$  rules out all terms that have  $x + 0$  as their first child of type  $\mathbf{a}$ , such as the terms  $(x + 0) + y$ ,  $\text{ite}(x \geq y, x + 0, y)$  and  $(x + 0) - 1$ , which are equivalent to the smaller expressions  $x + y$ ,  $\text{ite}(x \geq y, x, y)$  and  $x - 1$ , respectively, and hence redundant as candidate solutions. Sharing selectors allows the same blocking clause to be reused for the different constructors, whereas standard selectors would require three different clauses in this case, with  $z \mapsto \text{S}^{\text{Plus},1}(d)$ ,  $z \mapsto \text{S}^{\text{Ite},2}(d)$ , and  $z \mapsto \text{S}^{\text{Minus},1}(d)$ , respectively. •

A majority of SyGuS problems can be encoded as datatypes that have significant sharing of selectors across multiple constructors, thus making the use of shared selectors particularly effective in this domain. The next section measures the impact of shared selections when solving SyGuS problems in `cvc4`.

## 6 Experiments

We implemented our calculus for the theory of datatypes with shared selectors in `cvc4` Version 1.5, together with a preprocessing pass to convert standard selectors in input formulas to shared ones and other modifications to the existing decision procedure for datatypes, as described in Sections 3.2 and 4. We discuss here our evaluation of two configurations of `cvc4`, one with and one without support for shared selectors, on two different sets of benchmarks: the SyGuS benchmark suite from the 2017 SyGuS competition [4]; and SMT-LIB [8] benchmarks containing datatypes. Our experiments<sup>4</sup> were performed on the StarExec logic solving service [28].

### 6.1 Syntax-guided Synthesis Benchmarks

The benchmarks from the 2017 SyGuS competition are divided into five families across four tracks: (i) the General track, with problems over the theories of linear integer arithmetic (LIA) or bit-vectors; (ii) the conditional linear integer arithmetic track (CLIA), with problems over LIA; (iii) the Invariant synthesis track, also over LIA; and (iv) the Programming-by-examples track [17, 18], with a family over bit-vectors and another

<sup>4</sup> The data and details on how to reproduce our results are available at <https://cvc4.cs.stanford.edu/papers/IJCAR2018-shsel/>.

Family	#	Solved: sh / std	Time	SAT Decs	Terms	Sels
General	535	319 / 235 (232)	15.4 / 144.9	67k / 151k	189k / 284k	5.8 / 16.8
CLIA	73	18 / 17 (17)	25.1 / 142.4	158k / 405k	25k / 60k	9.6 / 22.2
Invariant	67	46 / 46 (46)	49.1 / 114.6	374k / 896k	37k / 61k	5.7 / 13.1
PBE_BV	750	665 / 253 (253)	27.4 / 211.9	54k / 3873k	14k / 202k	3.0 / 16.0
PBE_Strings	108	93 / 64 (64)	13.3 / 39.9	90k / 334k	14k / 41k	8.6 / 18.7

Fig. 3: Performance of cvc4 on benchmarks from five families of SyGuS Comp 2017.

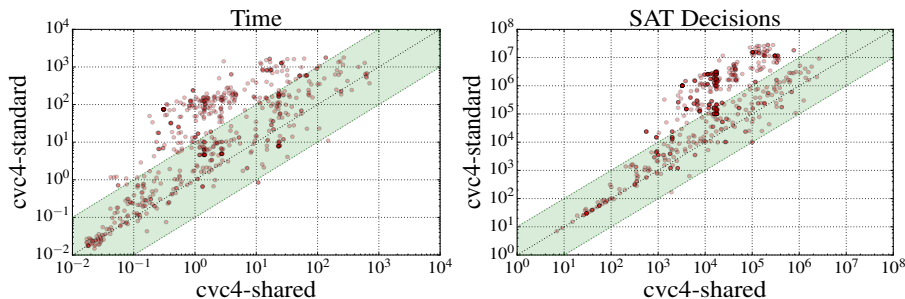


Fig. 4: Impact of shared selectors on solving time and number of SAT decisions.

over strings. We measured the impact of shared selectors by comparing for the two configurations of cvc4 the total number of solved problems and the average solving time, number of decisions performed by the SAT engine, quantifier-free terms generated, and number of selectors in the signature. Averages were computed over the set of problems solved by both configurations. We used a timeout of 30 minutes per benchmark.

A summary of the results is given in Figure 3. The first two columns show the evaluated family and the number of benchmarks in it, while the other columns present the statistics listed above, with average times expressed in seconds. The number of problems solved by both configurations is given in parentheses in the third column. The results clearly show that sharing selectors reduces the number of selectors in the signature, which generally leads to fewer terms and SAT decisions, with a positive impact on solving speed and number of problems solved. Except for the invariant family, the cvc4 configuration with shared selectors solves more problems than the one without. The impact of shared selectors is particularly significant for the bit-vector benchmark suite (PBE\_BV), with a reduction of over 80% in the average number of selectors. In that case, cvc4 is over eight times faster with shared selectors than without, solving 412 more problems, thus reducing the percentage of unsolved problems in this category from over 65% to less than 12%. Significant improvements can also be observed in the PBE\_Strings and General families, with the percentages of unsolved problems being reduced from over 40% to almost 13% and from over 55% to almost 40%, respectively.

We present a per-problem comparison in the scatter plots of Figure 4, which clearly shows that for the vast majority of the benchmarks, sharing selectors reduces the number of SAT decisions and improves the solving time, often by orders of magnitude.

Family	#	Solved: sh / std	Time	Decs	Terms	Sels
Leon	410	179 / 175 (175)	0.96 / 0.75	9920 / 9925	718 / 929	8.67 / 23.10
Sledgehammer	321	113 / 112 (112)	0.47 / 0.47	6949 / 6942	185 / 185	10.50 / 12.76
Nunchaku	158	67 / 67 (67)	0.49 / 0.44	7149 / 6653	1373 / 1297	6.22 / 7.22

Fig. 5: Performance of `cvc4` on benchmarks from three families of SMT LIB.

*Comparison against other SyGuS solvers* We also compared `cvc4`'s performance with the state-of-the-art SyGuS solver EUSOLVER [2, 5]. For fairness, in this comparison we combine the results of the above configurations of `cvc4` with its other approach for solving single-invocation synthesis problems (see [24] for details), which impacts the CLIA and General families of benchmarks. We obtained the following results for the problems solved by EUSOLVER and `cvc4` with and without shared selectors: 71/73/73 for CLIA, 404/391/334 for General, 42/46/46 for Invariant, 739/665/253 for PBE\_BV, and 68/93/64 for PBE\_Strings. These numbers show that overall `cvc4` is significantly more competitive with shared selectors than without, surpassing EUSOLVER's performance in three of the five families.

## 6.2 Datatype benchmarks from SMT LIB

We also considered all SMT-LIB benchmarks containing datatypes. Among these, we excluded from consideration 14,387 benchmarks that do not have any shareable selectors, as `cvc4` with and without shared selectors perform the same on these benchmarks. The remaining 889 benchmarks are divided into three families: (i) the Leon set contains benchmarks generated by Leon [12] for verification of Scala programs (AUFBVDTLIA logic); (ii) the Sledgehammer set has benchmarks from Isabelle [21] generated by Sledgehammer [14] (UFDT logic); and (iii) the Nunchaku set has benchmarks generated for higher order theorem provers by Nunchaku [23] (UFDT logic).

We summarize our results over the two configurations of `cvc4`, with and without shared selectors, in Figure 5, following the same schema as in Figure 3. We used a timeout of 60 seconds, since in this setting we evaluate SMT solvers as backends of verification and ITP tools, which require fast answers. The configuration with shared selectors solved at least all the benchmarks as the one without. The Leon benchmark set shows the most significant impact of sharing selectors, with a reduction of over 60% in the average number of selectors, and 4 more problems solved.

*Comparison against other SMT solvers* To put the shared selector version of `cvc4` in context with the state of the art, we also compared it with the only two provers that can reason about datatypes and support the SMT-LIB format: `z3` [16] and Vampire [19]. On the Nunchaku and Sledgehammer benchmarks, the number of problems solved by `cvc4/z3/Vampire` is 67/29/30 and 113/119/138, respectively. The comparison on the Leon set excludes Vampire, since it does not support the theory of bit-vectors; the split between `cvc4` and `z3` is 179/173 on that set. The results show that `cvc4` compares favorably with the other tools.

## 7 Related Work

The motivation of our work is to reduce the number of terms considered by a decision procedure for the theory of algebraic of datatypes, based on procedures introduced in previous work [10, 22]. Thus, our contributions apply to other systems that handle datatypes semantically, such as `smbc` [15] and the SMT solver `z3` [16]. On systems that reason about datatypes axiomatically, such as the first-order theorem prover and SMT solver Vampire [19], and the higher-order systems Isabelle [21] and Dafny [20], whether to share selectors and how to handle them is simply a matter of axiomatizing the datatypes theory accordingly. For example, the axiomatization in Vampire avoids selectors altogether [19, Sect. 4.3], while in Isabelle users are encouraged to write specifications directly with shared selectors [13, Sec. 3].

Most SyGuS solvers employ a variation of counter-example guided inductive synthesis (CEGIS), introduced by Solar-Lezama [27]. While `cvc4` benefits from sharing selectors by representing syntax restrictions with datatypes, other systems use an outer layer with an underlying reasoning engine, for instance using an SMT solver to verify the correctness of candidate solutions, but not for performing the enumerative search [5].

## 8 Conclusion

We have presented an extension of the theory of algebraic datatypes that adds shared selectors. We have discussed and proved correct a calculus for deciding the constraint satisfiability problem in the new theory. Moreover, we have described how algebraic datatypes can be leveraged in an SMT solver to solve syntax-guided synthesis problems and explained how the use of shared selectors in this setting can lead to significant performance gains. Our experiments demonstrate that an implementation of the new calculus in the `cvc4` solver significantly enhances its performance on syntax-guided synthesis problems and is responsible for making `cvc4` the best known solver for certain classes of problems.

In future work, we plan to generalize our approach so that distinct selector *chains* can be compressed to a single application of the same selector symbol. This requires more sophisticated criteria for recognizing when two selector chains for a datatype cannot be simultaneously constrained for arbitrary values of that datatype. We believe that this further extension can be done in a manner similar to the one presented here and expect that this will lead to further performance improvements.

## References

1. R. Alur, R. Bodík, G. Juniwal, M. M. K. Martin, M. Raghothaman, S. A. Seshia, R. Singh, A. Solar-Lezama, E. Torlak, and A. Udupa. Syntax-guided synthesis. In *FMCAD*, pages 1–8. IEEE, 2013.
2. R. Alur, P. Cerný, and A. Radhakrishna. Synthesis through unification. In *CAV*, volume 9207 of *LNCIS*, pages 163–179. Springer, 2015.
3. R. Alur, D. Fisman, R. Singh, and A. Solar-Lezama. Sygus-comp 2017: Results and analysis. In D. Fisman and S. Jacobs, editors, *Proceedings Sixth Workshop on Synthesis (SYNT)*, volume 260 of *EPTCS*, pages 97–115, 2017.

4. R. Alur, D. Fisman, R. Singh, and A. Solar-Lezama. Sygus-comp 2017: Results and analysis. *CoRR*, abs/1711.11438, 2017.
5. R. Alur, A. Radhakrishna, and A. Udupa. Scaling enumerative program synthesis via divide and conquer. In *TACAS*, volume 10205 of *LNCS*, pages 319–336, 2017.
6. F. Baader and T. Nipkow. *Term Rewriting and All That*. Cambridge University Press, New York, NY, USA, 1998.
7. C. Barrett, C. L. Conway, M. Deters, L. Hadarean, D. Jovanović, T. King, A. Reynolds, and C. Tinelli. CVC4. In *CAV*, pages 171–177. Springer, 2011.
8. C. Barrett, P. Fontaine, and C. Tinelli. The Satisfiability Modulo Theories Library (SMT-LIB). [www.SMT-LIB.org](http://www.SMT-LIB.org), 2016.
9. C. Barrett, P. Fontaine, and C. Tinelli. The SMT-LIB Standard: Version 2.6. Technical report, Department of Computer Science, The University of Iowa, 2017.
10. C. Barrett, I. Shikanian, and C. Tinelli. An abstract decision procedure for a theory of inductive data types. *JSAT*, 3(1-2):21–46, 2007.
11. C. Barrett and C. Tinelli. Satisfiability modulo theories. In *Handbook of Model Checking*. Springer, 2014. (to appear).
12. R. Blanc, V. Kuncak, E. Kneuss, and P. Suter. An overview of the leon verification system: Verification by translation to recursive functions. In *Proceedings of the 4th Workshop on Scala*, pages 1:1–1:10. ACM, 2013.
13. J. C. Blanchette, J. Hölzl, A. Lochbihler, L. Panny, A. Popescu, and D. Traytel. Truly modular (co)datatypes for isabelle/hol. In *ITP*, volume 8558 of *LNCS*, pages 93–110. Springer, 2014.
14. S. Böhme and T. Nipkow. Sledgehammer: Judgement day. In *IJCAR*, volume 6173 of *LNCS*, pages 107–121. Springer-Verlag, 2010.
15. S. Cruanes. Satisfiability modulo bounded checking. In *CADE*, volume 10395 of *LNCS*, pages 114–129. Springer, 2017.
16. L. de Moura and N. Bjørner. Z3: An Efficient SMT Solver. In *TACAS*, volume 4963 of *LNCS*, pages 337–340. Springer, 2008.
17. S. Gulwani. Automating string processing in spreadsheets using input-output examples. In *POPL*, pages 317–330. ACM, 2011.
18. S. Gulwani. Programming by examples: Applications, algorithms, and ambiguity resolution. In *IJCAR*, volume 9706 of *LNCS*, pages 9–14. Springer, 2016.
19. L. Kovács, S. Robillard, and A. Voronkov. Coming to terms with quantified reasoning. In *POPL*, pages 260–270. ACM, 2017.
20. K. R. M. Leino. Developing verified programs with dafny. *Ada Lett.*, 32(3):9–10, Dec. 2012.
21. T. Nipkow, L. C. Paulson, and M. Wenzel. *Isabelle/HOL - A Proof Assistant for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002.
22. A. Reynolds and J. C. Blanchette. A decision procedure for (co)datatypes in smt solvers. *Journal of Automated Reasoning*, 58(3):341–362, Mar 2017.
23. A. Reynolds, J. C. Blanchette, S. Cruanes, and C. Tinelli. Model finding for recursive functions in SMT. In *IJCAR*, volume 9706 of *LNCS*, pages 133–151. Springer, 2016.
24. A. Reynolds, V. Kuncak, C. Tinelli, C. Barrett, and M. Deters. Refutation-based synthesis in smt. *Formal Methods in System Design*, Feb 2017.
25. A. Reynolds and C. Tinelli. Sygus techniques in the core of an smt solver. *arXiv preprint arXiv:1711.10641*, 2017.
26. A. Reynolds, A. Viswanathan, H. Barbosa, C. Tinelli, and C. Barrett. Datatypes with Shared Selectors. Technical report, The University of Iowa, 2018. <http://cvc4.cs.stanford.edu/papers/IJCAR2018-shsel/>.
27. A. Solar-Lezama, L. Tancau, R. Bodík, S. A. Seshia, and V. A. Saraswat. Combinatorial sketching for finite programs. In *ASPLOS*, pages 404–415. ACM, 2006.
28. A. Stump, G. Sutcliffe, and C. Tinelli. Starexec: A cross-community infrastructure for logic solving. In *IJCAR*, volume 8562 of *LNCS*, pages 367–373. Springer, 2014.