

Scaling Up DPLL(T) String Solvers using Context-Dependent Simplification

Andrew Reynolds¹, Maverick Woo²,
Clark Barrett³, David Brumley², Tianyi Liang⁴, and Cesare Tinelli¹

¹ Department of Computer Science, The University of Iowa

² CyLab, Carnegie Mellon University

³ Department of Computer Science, Stanford University

⁴ Two Sigma



Abstract. Efficient reasoning about strings is essential to a growing number of security and verification applications. We describe satisfiability checking techniques in an extended theory of strings that includes operators commonly occurring in these applications, such as `contains`, `index_of` and `replace`. We introduce a novel context-dependent simplification technique that improves the scalability of string solvers on challenging constraints coming from real-world problems. Our evaluation shows that an implementation of these techniques in the SMT solver CVC4 significantly outperforms state-of-the-art string solvers on benchmarks generated using PyEx, a symbolic execution engine for Python programs. Using a test suite sampled from four popular Python packages, we show that PyEx uses only 41% of the runtime when coupled with CVC4 than when coupled with CVC4's closest competitor while achieving comparable program coverage.

1 Introduction

A growing number of applications of static analysis techniques have benefited from automated reasoning tools for string constraints. The effect of such tools on symbolic execution in particular has been transformative. At a high level, symbolic execution *runs* a program under analysis by representing its input values symbolically and tracking the program variables as expressions over these symbolic values, together with other concrete values in the program. The collected expressions are then analyzed by an automated reasoning tool to determine path feasibility at branches or other properties such as security vulnerabilities at points of interest. With the ever-increasing hardware and automated reasoning capabilities, symbolic execution has enjoyed much success in practice. A recent example in the cybersecurity realm was the DARPA Cyber Grand Challenge,⁵ which featured a number of Cyber Reasoning Systems that heavily relied on symbolic execution techniques (see, e.g., [7, 24]).

Prior to the availability of string-capable reasoning tools, developers of symbolic execution engines had to adopt various ad-hoc heuristics to cope with strings and other variable-length inputs. One popular heuristic is to impose an artificial upper-bound on the length of these inputs. Unfortunately, this not only compromises analysis accuracy,

⁵ See <http://www.darpa.mil/program/cyber-grand-challenge>.

since the chosen upper bounds may be too low in face of adversarial inputs, but it also leads to inefficiencies in solvers when practitioners, in an attempt to mitigate this problem, may end up setting the upper bounds too high. To address this issue, a number of SMT solvers have been extended recently with native support for *unbounded* strings and length constraints [30, 19, 2, 26]. These solvers have dramatically improved both in performance and robustness over the past few years, enabling a new generation of symbolic execution tools that support faithful string reasoning.

This paper revisits approaches for solving *extended string constraints*, which allow a rich language of string terms over operators such as `contains`, `index_of` and `replace`. Earlier techniques for extended string constraints [5, 18, 30, 27] often rely on eager reductions to a core language of constraints, with the effect of requiring the solver to deal with fairly large or complex combination of basic constraints. For instance, encoding the constraint $\neg \text{contains}(x, y)$ commonly involves bounded universal quantification over the integers, to state that string y does not occur at any position in x . In this work, we start with the observation that DPLL(T)-based SMT solvers [12] often reason in contexts where string variables are equated to (partially) concrete values. Based on this, we have developed a way to leverage efficient context-dependent simplification techniques and reduce extended constraints to a core language lazily instead.

Contribution and significance We extend a calculus by Liang et al. [19] to handle *extended string constraints* (i.e. constraints over `substr`, `contains`, `index_of` and `replace`) using a combination of two techniques:

- a reduction of extended string constraints to basic ones involving bounded quantification (Section 3.1), and
- an inference technique based on context-dependent simplification (Section 3.2) which supplements this reduction and in practice significantly improves the scalability of our approach on constraints coming from symbolic execution.

Additionally, we provide a new set of 25,421 publicly-available benchmarks over extended string constraints. These benchmarks were generated by running PyEx, a new symbolic executor for Python programs based on PyExZ3 [3], over a test suite of 19 target functions sampled from four popular Python packages. We discuss an experimental evaluation showing that our implementation in the SMT solver CVC4 significantly outperforms other state-of-the-art string solvers in finding models for these benchmarks.

Finally, we discuss how the superior performance of CVC4 in comparison to other solvers translates into real-life benefits for Python developers using PyEx.

Structure of the paper After some formal preliminaries, we briefly review a calculus for basic string constraints in Section 2 that is an abbreviated version of [19]. We present new techniques for extended string constraints in Section 3, and evaluate these techniques on real-world queries generated by PyEx in Section 5.

1.1 Related work

The satisfiability of word equations was proven decidable by Makanin [21] and then given a PSPACE algorithm by Plandowski [22]. The decidability of the fairly restricted

language of word equations with length constraints is an open problem [11]. In practice, a number of approaches for solving string and regular expression constraints rely on reductions to automata [10, 15, 28] or bit-vector constraints for fixed-length strings [16]. More recently, new approaches have been developed for the satisfiability problem for (unbounded) word equations and memberships with length [30, 19, 27, 2, 20] within SMT solvers. Among these, Z3-STR [30] and S3 [27] are third-party extensions of the SMT solver Z3 [9] adding support for string constraints via reductions to linear arithmetic and uninterpreted functions. This support includes extended string constraints over a signature similar to the one we consider in this paper. With respect to these solvers, our string solver is fully integrated into the architecture of CVC4, meaning it can be combined with other theories of CVC4, such as algebraic datatypes and arrays.

This paper is similar in scope to work by Bjørner et al. [5], which gives decidability results and an approach for string library functions, including contains and replace. As in that work, we reduce the satisfiability problem for extended string constraints to a core language with bounded quantification. We also incorporate simplification techniques that improve performance by completely or partially avoiding this reduction.

Our target application, symbolic execution, has a rich history, starting from the seminal work of King [17]. Common modern symbolic execution tools include SAGE [13], KLEE [6], S2E [8], and Mayhem [7], which are all designed to analyze low-level binary or source code. In contrast, this paper considers constraints generated from the symbolic executor PyEx, which is designed to analyze Python code and includes support for string variables.

1.2 Formal preliminaries

We work in the context of many-sorted first-order logic with equality and assume the reader is familiar with the notions signature, term, literal, (quantified) formula, and free variable. We consider many-sorted signatures Σ that contain an (infix) logical symbol \approx for equality—which has type $\sigma \times \sigma$ for all sorts σ in Σ and is always interpreted as the identity relation. We also assume signatures Σ contain the Boolean sort `Bool` and Boolean constant symbols \top and \perp for true and false. Without loss of generality, we assume \approx is the only predicate symbol in Σ , as all other predicates may be modeled as functions with return sort `Bool`. If P is a function with return sort `Bool`, we will commonly write $P(\mathbf{t})$ as shorthand for $P(\mathbf{t}) \approx \top$. If e is a term or a formula, we denote by $\mathcal{V}(e)$ and $\mathcal{T}(e)$ the set of free variables and subterms of e respectively, extending these notations to tuples and sets of terms or formulas as expected.

A *theory* is a pair $T = (\Sigma, \mathbf{I})$ where Σ is a signature and \mathbf{I} is a class of Σ -interpretations, the *models* of T . A Σ -formula φ is *satisfiable* (resp., *unsatisfiable*) in T if it is satisfied by some (resp., no) interpretation in \mathbf{I} . A set Γ of Σ -formulas *entails in* T a Σ -formula φ , written $\Gamma \models_T \varphi$, if every interpretation in \mathbf{I} that satisfies all formulas in Γ satisfies φ as well. We write $\Gamma \models \varphi$ to denote entailment in the (empty) theory of equality. We say Γ *propositionally entails* φ , written $\Gamma \models_p \varphi$, if Γ entails φ when considering all atoms in γ and φ as propositional variables. Two Σ -terms or Σ -formulas are *equivalent in* T if they are satisfied by the same models of T . Two formulas φ_1 and φ_2 are *equisatisfiable in* T if $\exists \mathbf{x}_1. \varphi_1$ and $\exists \mathbf{x}_2. \varphi_2$ are equivalent in T where \mathbf{x}_i collects the free variables of φ_i that do not occur free in φ_j with $i \neq j$.

Σ_A	$n : \text{Int}$ for all $n \in \mathbb{N}$	$+$: $\text{Int} \times \text{Int} \rightarrow \text{Int}$	$-$: $\text{Int} \rightarrow \text{Int}$	\geq : $\text{Int} \times \text{Int} \rightarrow \text{Bool}$
Σ_S	$l : \text{Str}$ for all $l \in \mathcal{A}^*$	con : $\text{Str} \times \dots \times \text{Str} \rightarrow \text{Str}$	len : $\text{Str} \rightarrow \text{Int}$	
Σ_X	substr : $\text{Str} \times \text{Int} \times \text{Int} \rightarrow \text{Str}$	contains : $\text{Str} \times \text{Str} \rightarrow \text{Bool}$	replace : $\text{Str} \times \text{Str} \times \text{Str} \rightarrow \text{Str}$	
	index_of : $\text{Str} \times \text{Str} \times \text{Int} \rightarrow \text{Int}$			

Fig. 1. Functions in signature Σ_{ASX} . Str and Int denote strings and integers respectively.

We consider an extended theory T_{ASX} of strings and length equations, whose signature Σ_{ASX} is given in Figure 1. We assume a fixed finite alphabet \mathcal{A} of characters. The signature includes the sorts Str and Int denoting strings and integers respectively. Figure 1 divides the signature Σ_{ASX} into three parts, which we denote by Σ_A , Σ_S and Σ_X . We will write Σ_{AS} to denote $\Sigma_A \cup \Sigma_S$ and so on. The subsignature Σ_A is provided on the top line of Figure 1 and includes the usual symbols of *linear* integer arithmetic, interpreted as expected. We will commonly write $t_1 \bowtie t_2$, with $\bowtie \in \{>, <, \leq\}$, as syntactic sugar for the equivalent inequality between t_1 and t_2 expressed using only \geq . The subsignature Σ_S is provided on the second line and includes: a constant symbol, or *string constant*, for each word of \mathcal{A}^* (including ϵ for the empty word), interpreted as that word; a variadic function symbol $\text{con} : \text{Str} \times \dots \times \text{Str} \rightarrow \text{Str}$, interpreted as word concatenation; and a function symbol $\text{len} : \text{Str} \rightarrow \text{Int}$, interpreted as the word length function. The subsignature Σ_X is provided in the remainder of the figure.

We refer to the function symbols in Σ_X as *extended functions*, and terms whose top symbol is in Σ_X as *extended function terms*. A *position* in a string x is a non-negative integer smaller than the length of x that identifies a character in x — with 0 identifying the first character, 1 the second, and so on. For all x, y, z, n, m , the term $\text{substr}(x, n, m)$ is interpreted as the maximal substring of x starting at position n with length at most m , or the empty string if n is an invalid position; $\text{contains}(x, y)$ is interpreted as true if and only if x contains y ; $\text{index_of}(x, y, n)$ is interpreted as the position of the first occurrence of y in x starting at position n , or -1 if y is empty, n is an invalid position, or if no such occurrence exists; $\text{replace}(x, y, z)$ is interpreted as the result of replacing the first occurrence in x of y by z , or just x if y is empty or x does not contain y .

An *atomic term* is either a constant or a variable. A *flat term* is a term of the form $f(x_1, \dots, x_n)$, where x_1, \dots, x_n are variables. A *string term* is one that contains function symbols from Σ_{SX} only. A string term is *basic* if contains function symbols from Σ_S only, and *extended* otherwise. A (*basic, extended*) *string constraint* is a (dis)equality between (basic, extended) string terms. An *arithmetic constraint* is an inequality or (dis)equality between linear combinations of atomic and/or string terms with integer sort. Notice that (dis)equalities between integer variables and constraints such as $\text{len } x \approx \text{len } y$ are both string and arithmetic constraints. A T_{ASX} -constraint is either a string or an arithmetic constraint. Without loss of generality, we consider the satisfiability problem for Σ_{ASX} -formulas composed of T_{ASX} -constraints only.

If E is a finite set of basic string constraints, the *congruence closure* of E is the set

$$\begin{array}{ll}
\text{con}(s, \epsilon, \mathbf{u}) \rightarrow \text{con}(s, \mathbf{u}) & \text{con}() \rightarrow \epsilon \quad \text{con}(s) \rightarrow s \\
\text{con}(s, \text{con}(t, \mathbf{u})) \rightarrow \text{con}(s, t, \mathbf{u}) & \text{con}(s, c_1 \cdots c_i, c_{i+1} \cdots c_n, \mathbf{u}) \rightarrow \text{con}(s, c_1 \cdots c_n, \mathbf{u}) \\
\text{len}(c_1 \cdots c_n) \rightarrow n & \text{len}(\text{con}(s_1, \dots, s_n)) \rightarrow \text{len } s_1 + \cdots + \text{len } s_n
\end{array}$$

Fig. 2. Simplification rules for Σ_{AS} -terms.

$$\begin{aligned}
\hat{E} = & \{s \approx t \mid s, t \in \mathcal{T}(E), E \models s \approx t\} \cup \\
& \{s \not\approx t \mid s, t \in \mathcal{T}(E), s' \not\approx t' \in E \cup L, E \models s \approx s' \wedge t \approx t' \text{ for some } s', t'\}
\end{aligned}$$

where $L = \{l_1 \not\approx l_2 \mid \text{for all distinct } l_1, l_2 \in \mathcal{A}^*\} \cup \{\top \not\approx \perp\}$. The congruence closure of E induces an equivalence relation over $\mathcal{T}(E)$ where two terms s, t are equivalent iff $s \approx t \in \hat{E}$. For all $t \in \mathcal{T}(E)$, we denote its equivalence class in \hat{E} by $[t]_E$ or just $[t]$ when E is clear.

Given a term t , we write $t\downarrow$ to denote its *simplified form*, where $t\downarrow$ is a term that is equivalent to t in T_{ASX} and is not simplifiable further (i.e., $(t\downarrow)\downarrow = t\downarrow$). We do not insist that simplified forms be canonical, that is, equivalent terms need not have the same simplified form. Rules for computing the simplified form of Σ_{AS} -terms are given in Figure 2. Rules for computing the simplified form of other Σ_{ASX} -terms are fairly sophisticated and will be described in detail in Section 3.2. Given a tuple of basic string terms \mathbf{t} , we write $\mathbf{t}\downarrow$ to denote a tuple of atomic string terms corresponding to the arguments of $\text{con}(\mathbf{t})\downarrow$. For example, if $c_1, c_2 \in \mathcal{A}$, $(c_1, \text{con}(c_2, x), y)\downarrow = (c_1 c_2, x, y)$ and $(x, \epsilon)\downarrow = (x)$. Given an arbitrary Σ_{ASX} -formula φ , we write $[\varphi]$ to denote an equisatisfiable (purified) formula whose atoms are T_{ASX} -constraints, and where t is in simplified form (i.e., $t = t\downarrow$) for all its subterms t . For example, $[\text{substr}(x, n_0, n_1 + 1) \approx y]$ is $\text{substr}(x, n_0, n_2) \approx y \wedge n_2 \approx n_1 + 1$ for a fresh integer variable n_2 .

2 A Calculus for Basic String Constraints

Liang et al. [19] developed a calculus for the satisfiability of finite conjunctions of constraints in a theory of strings with length and regular expressions. This section presents a modified version of that calculus. We focus on the portion of that calculus that handles string equalities and length constraints, and omit discussion of its other aspects. Furthermore, we extend that calculus with support for constraints involving extended functions. To simplify the description of this extension and make it self-contained, we also extend the calculus to model propositional reasoning as well, making it applicable to Σ_{ASX} -formulas instead of just conjunctions of T_{ASX} -constraints.

Definition 1 (Configurations). *A configuration is either the distinguished symbol unsat or a tuple of the form $\langle G, S, A \rangle$ where G is a set of Σ_{ASX} -formulas, A is a set of arithmetic constraints, and S is a tuple of the form (E, X, F, N) , where:*

- E is a set of basic string equalities;
- X is a set of equalities of the form $x \approx t$, where x is a variable and t is a flat extended function term;

$$\begin{array}{c}
\text{Prop-Assign} \frac{G = G', \varphi \quad \varphi \text{ is quantifier-free}}{\|M \models_p \varphi \quad G := G' \quad S := S \oplus M|_S \quad A := A \cup M|_A} \\
\text{A-Prop} \frac{S \models s \approx t \quad s, t : \text{Int}}{A := A, s \approx t} \quad \text{S-Prop} \frac{A \models_{\text{LIA}} s \approx t \quad s, t : \text{Int}}{S := S \oplus \{s \approx t\}} \quad \text{A-Conf} \frac{A \models_{\text{LIA}} \perp}{\text{unsat}} \\
\text{L-Eq} \frac{x \approx t \in \hat{E} \quad x : \text{Str}}{A := A, \text{len } x \approx (\text{len } t) \downarrow} \quad \text{L-Geq} \frac{x \in \mathcal{V}(S \cup A) \quad x : \text{Str}}{S := S \oplus \{x \approx \epsilon\} \quad \| \quad A := A, \text{len } x > 0}
\end{array}$$

Fig. 3. Rules modeling the interaction between the propositional, string and arithmetic subsolvers.

- F is a set of pairs $s \mapsto \mathbf{a}$ where $s \in \mathcal{T}(E)$ and \mathbf{a} is a tuple of atomic string terms;
- N is a set of pairs $e \mapsto \mathbf{a}$ where e is an equivalence class of \hat{E} and \mathbf{a} is a tuple of atomic string terms. \square

A configuration $\langle G, S, A \rangle$ models the internal state of various modules of a DPLL(T)-based solver. The component G collects the formulas being processed by the solver’s propositional satisfiability (SAT) engine; S models the state of a theory solver for strings; and A collects the constraints given to a solver for linear integer arithmetic. Initial configurations have the form $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$ where φ is a *quantifier-free* T_{ASX} -formula to be checked for satisfiability and $\overline{\emptyset}$ abbreviates the tuple $(\emptyset, \emptyset, \emptyset, \emptyset)$.

We describe a calculus for the satisfiability of string constraints by a set of derivation rules that modify configurations. The rules are given in *guarded assignment form*, where the top of the rule describes the conditions under which the rule can be applied, and the bottom of the rule either is *unsat*, or otherwise describes the resulting modifications to the components of our configuration. A rule may have multiple, alternative conclusions separated by $\|$. An application of a rule is *redundant* if it has a conclusion where each component in the derived configuration is a subset of the corresponding component in the premise configuration. A configuration other than *unsat* is *saturated* if every possible application of a derivation rule to it is redundant. A *derivation tree* is a tree where each node is a configuration whose children, if any, are obtained by a non-redundant application of a rule of the calculus. A *closed* derivation tree (where all terminal nodes are *unsat*) with root node $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$ is a proof that φ is unsatisfiable in T_{ASX} . A derivation tree with root node $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$ and a saturated leaf is, under certain assumptions (see Theorem 1) a witness that φ is satisfiable in T_{ASX} .

To discuss the rules, we first introduce the notation for updating the internal state $S = (E, X, F, N)$ of our string solver. Let M be a set of T_{ASX} -constraints. By introducing enough fresh variables, we can construct an equisatisfiable set $M_E \cup M_X$ where M_E is a set of Σ_{AS} -literals and M_X is a set of equalities of the form $y \approx t$, with y a variable and t a flat extended function term. For simplicity, we assume here that M_E does not contain string disequalities.⁶ We define the *external update of S with M* , written $S \oplus M$, to be the tuple $(E \cup M_E, X \cup M_X, \emptyset, \emptyset)$. For Σ_{AS} -terms t_1, t_2 ,

⁶ String disequalities can be reduced to a finite disjunction of equalities, e.g. see [1]. More sophisticated and efficient methods for handling string disequalities are given in [19].

we define the *internal update of S with $t_1 \approx t_2$* , written $S \odot t_1 \approx t_2$, to be the tuple $(E \cup \{t_1 \approx t_2\}, X, F\sigma, \emptyset)$, where σ is the substitution $\{t_1 \mapsto t_2\}$ if t_1 is a string variable and is the empty substitution otherwise, and $F\sigma$ is the result of replacing the right hand side of all pairs $y \mapsto a$ in F with $(a\sigma)\downarrow$.

The rules in Figure 3 model the basic interaction between the various subsolvers in our approach. The rule Prop-Assign considers each *propositional satisfying assignment* M for a quantifier-free formula $\varphi \in G$, i.e., each set M of literals such as every atom of φ occurs either positively or negatively (but not both) in M , all the atoms in M occur in φ , and $M \models_p \varphi$. For each such M , Prop-Assign has a conclusion where the string constraints in M (denoted as $M|_S$) are given to the string subsolver and the arithmetic constraints (denoted as $M|_A$) are given to the arithmetic subsolver. The arithmetic and string solvers use rules A-Prop and S-Prop to share equalities between (shared) arithmetic terms, and use A-Conf and S-Conf to report that their respective set of constraints is unsatisfiable. In those rules and in the rest of the paper, \models_{LIA} denotes entailment in linear integer arithmetic. The rules L-Eq and L-Geq respectively infer and guess arithmetic constraints involving string length.

Example 1. Consider the formula φ of the form $\text{len } x > \text{len } y \wedge y \approx \text{con}(x, a)$. Starting from configuration $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$, we may apply Prop-Assign to remove φ from G and update S and A based the propositional satisfying assignment for φ . We obtain a configuration where A is $\{\text{len } x > \text{len } y\}$ and the E component of S is $\{y \approx \text{con}(x, a)\}$. Since $(\text{len}(\text{con}(x, a)))\downarrow = \text{len } x + 1$, we may add $\text{len } y \approx \text{len } x + 1$ to A by the rule L-Eq. Since $\text{len } x > \text{len } y$, $\text{len } y \approx \text{len } x + 1 \models_{LIA} \perp$, we may apply A-Conf to derive the unsat configuration, establishing that φ is unsatisfiable in T_{ASX} . \square

The rules in Figure 4 are used by the string solver for building the mappings F and N . We defer discussion of the X component of configurations until Section 3. In the rules, we assume a total ordering $<$ on string terms, whose only restriction is that $t_1 < t_2$ if t_1 is variable and t_2 is not. For a term t , we call $F t$ the *flat form* of t , and for an equivalence class e , we call $N e$ the *normal form* of e . We construct the mappings F and N using the rules F-Form1, F-Form2, N-Form1 and N-Form2 in a mutually recursive fashion. The remaining three rules apply to cases where the above rules do not result in complete mappings for F and N . In the case where we compute flat forms for two terms s and t in the same equivalence class that share a common prefix w followed by two distinct variables u and v , we apply F-Unify to infer $u \approx v$ in the case that the lengths of u and v are equal, otherwise we apply F-Split to infer that u is a prefix of v or vice versa in the case that the lengths of u and v are disequal. The rule L-Split splits the derivation based on equality between the lengths of string variables x and y , which we use to derive configurations where one of these two rules applies.

Example 2. Consider a configuration where E is $\{y \approx \text{con}(a, x), x \approx \text{con}(u, z), y \approx \text{con}(a, v, z), w \approx a, \text{len } u \approx \text{len } v\}$. The equivalence classes of \hat{E} are

$$\{y, \text{con}(a, x), \text{con}(a, v, a)\} \quad \{x, \text{con}(u, a)\} \quad \{u\} \quad \{v\} \quad \{w, a\}$$

Using N-Form2, we obtain $N[u] = (u)$ and $N[v] = (v)$. Using F-Form1 and N-Form1, we obtain $F \text{con}(u, a) = N[x] = (u, z)$ and $F a = N[w] = (a)$. For $[y]$, we use F-Form1 to obtain $F \text{con}(a, x) = (a, u, a)$ and $F \text{con}(a, v, a) = (a, v, a)$. Since the flat

$$\begin{array}{c}
\text{F-Form1} \frac{\text{con}(t_1, \dots, t_n) \in \mathcal{T}(\mathbb{E}) \quad \mathbb{N}[t_i] = s_i \text{ for } i = 1, \dots, n}{\mathbb{F} := \mathbb{F}, \text{con}(t_1, \dots, t_n) \mapsto (s_1, \dots, s_n) \downarrow} \quad \text{F-Form2} \frac{l \in \mathcal{T}(\mathbb{E})}{\mathbb{F} := \mathbb{F}, l \mapsto (l) \downarrow} \\
\\
\text{N-Form1} \frac{e \notin \mathcal{V}(\mathbb{E}) \quad \mathbb{F}t = \mathbf{a} \text{ for all } t \in e \setminus \mathcal{V}(\mathbb{E})}{\mathbb{N} := \mathbb{N}, e \mapsto \mathbf{a}} \quad \text{N-Form2} \frac{[x] \subseteq \mathcal{V}(\mathbb{E}) \quad x \leq y \text{ for all } y \in [x]}{\mathbb{N} := \mathbb{N}, [x] \mapsto (x)} \\
\\
\text{F-Unify} \frac{\mathbb{F}s = (\mathbf{w}, u, \mathbf{u}_1) \quad \mathbb{F}t = (\mathbf{w}, v, \mathbf{v}_1) \quad s \approx t \in \hat{\mathbb{E}} \quad u < v \quad \mathbb{E} \models \text{len } u \approx \text{len } v}{\mathbb{S} := \mathbb{S} \odot v \approx u} \\
\\
\text{F-Split} \frac{\mathbb{F}s = (\mathbf{w}, u, \mathbf{u}_1) \quad \mathbb{F}t = (\mathbf{w}, v, \mathbf{v}_1) \quad s \approx t \in \hat{\mathbb{E}} \quad \mathbb{E} \models \text{len } u \not\approx \text{len } v \quad u \notin \mathcal{V}(\mathbf{v}_1) \quad v \notin \mathcal{V}(\mathbf{u}_1)}{\mathbb{S} := \mathbb{S} \odot u \approx \text{con}(v, z) \quad \parallel \quad \mathbb{S} := \mathbb{S} \odot v \approx \text{con}(u, z)} \\
\\
\text{L-Split} \frac{x, y \in \mathcal{V}(\mathbb{E}) \quad x, y : \text{Str}}{\mathbb{E} := \mathbb{E}, \text{len } x \approx \text{len } y \quad \parallel \quad \mathbb{E} := \mathbb{E}, \text{len } x \not\approx \text{len } y}
\end{array}$$

Fig. 4. String derivation rules. The rules construct flat forms \mathbb{F} and normal forms \mathbb{N} for string terms. The letter l denotes a string constant, and z denotes a fresh string variable.

forms of these terms are not the same and $\text{len } u \approx \text{len } v \in \mathbb{E}$, we may apply F-Unify to conclude $v \approx u$. We update \mathbb{S} to $\mathbb{S} \odot u \approx v$, after which the equivalence classes are:

$$\{y, \text{con}(\mathbf{a}, x), \text{con}(\mathbf{a}, v, \mathbf{a})\} \quad \{x, \text{con}(u, \mathbf{a})\} \quad \{u, v\} \quad \{w, \mathbf{a}\}$$

and $\mathbb{F} \text{con}(\mathbf{a}, v, \mathbf{a})$ is now $(\mathbf{a}, u, \mathbf{a})$. Then we can use N-Form1 and N-Form2 to reconstruct \mathbb{N} . Since $\mathbb{F} \text{con}(\mathbf{a}, x) = \mathbb{F} \text{con}(\mathbf{a}, v, \mathbf{a}) = (\mathbf{a}, u, \mathbf{a})$, we can obtain $\mathbb{N}[y] = (\mathbf{a}, u, \mathbf{a})$. This results in a configuration where \mathbb{N} is a complete mapping over the equivalence classes and no more rules apply, indicating that \mathbb{E} is satisfiable in T_{ASX} . \square

We say a configuration is *cyclic* if $\hat{\mathbb{E}}$ either contains a chain of equalities of the form $s \approx \text{con}(t_1), s_1 \approx \text{con}(t_2), \dots, s_{n-1} \approx \text{con}(t_n), s_n \approx s$ where s_i is a term from t_i for each i , or an equality of the form $s \approx t$ where $\mathbb{F}s = (\mathbf{w}, \mathbf{u})$, $\mathbb{F}t = (\mathbf{w}, \mathbf{v})$, \mathbf{w} is the maximal common prefix of $\mathbb{F}s$, and $\mathbb{F}t$ and $\mathcal{V}(\mathbf{u}) \cap \mathcal{V}(\mathbf{v})$ is non-empty. Recent techniques have been proposed for cyclic string constraints [19, 29]. We instead focus primarily on acyclic string constraints in the following result.

Theorem 1. *For all quantifier-free Σ_{AS} -formulas φ , the following hold.*

1. *There is a closed derivation tree with root $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$ only if φ is unsat in T_{ASX} .*
2. *There is a derivation tree with root $\langle \{\varphi\}, \overline{\emptyset}, \emptyset \rangle$ containing an acyclic saturated configuration only if φ is sat in T_{ASX} .*

3 Techniques for Extended String Constraints

This section gives a novel extension of the calculus in the previous section for determining the satisfiability of T_{ASX} -constraints. While the decidability of this problem is

$$\begin{array}{c}
\text{Ext-Expand} \frac{x \approx t \in X}{G := G, \llbracket x \approx t \rrbracket \quad X := X \setminus \{x \approx t\}} \quad \text{where} \\
\llbracket x \approx \text{substr}(y, n, m) \rrbracket = \text{ite}(0 \leq n < \text{len } y \wedge 0 < m, y \approx \text{con}(z_1, x, z_2) \wedge \text{len } z_1 \approx n \wedge \\
\text{len } z_2 \approx \text{len } y \dot{-} m, x \approx \epsilon) \\
\llbracket x \approx \text{contains}(y, z) \rrbracket = (x \not\approx \top) \Leftrightarrow \forall k. 0 \leq k \leq \text{len } y - \text{len } z \Rightarrow \text{substr}(y, k, \text{len } z) \not\approx z \\
\llbracket x \approx \text{index_of}(y, z, n) \rrbracket = \text{ite}(0 \leq n \wedge z \not\approx \epsilon \wedge \text{contains}(y', z), \text{substr}(y', x', \text{len } z) \approx z \wedge \\
\neg \text{contains}(\text{substr}(y', 0, x' + \text{len } z - 1), z), x \approx -1) \\
\text{with } y' = \text{substr}(y, n, \text{len } y - n) \text{ and } x' = x - n \\
\llbracket x \approx \text{replace}(y, z, w) \rrbracket = \text{ite}(\text{contains}(y, z) \wedge z \not\approx \epsilon, x \approx \text{con}(z_1, w, z_2) \wedge \\
y \approx \text{con}(z_1, z, z_2) \wedge \text{index_of}(y, z, 0) \approx \text{len } z_1, x \approx y) \\
\text{B-Val} \frac{t : \text{Int} \quad n \text{ is a numeral}}{A := A, t \leq n \quad \parallel \quad A := A, t > n} \quad \text{B-Inst} \frac{G = G', \varphi[\forall k. 0 \leq k \leq t \Rightarrow \psi] \\
A \models_{\text{LIA}} t \leq n \text{ for some numeral } n}{G := G', [\varphi[\wedge_{i=0}^n \psi\{k \mapsto i\}]]}
\end{array}$$

Fig. 5. Rules for reducing Σ_{ASX} -constraints to Σ_{AS} -constraints, where z_1, z_2 are fresh variables, $n_1 \dot{-} n_2$ denotes the maximum of $n_1 - n_2$ and 0, and ite is the if-then-else connective.

not known [5], we focus on techniques that are both refutation-sound and model-sound but guaranteed to terminate in general. We introduce two techniques for establishing the (un)satisfiability of T_{ASX} -constraints S , described by the additional derivation rules in Figures 5 and 6 which supplement those from Section 2.

3.1 Expanding Extended Function Terms to Bounded Integer Quantification

The satisfiability problem for equalities over Σ_{ASX} -terms can be reduced to the satisfiability problem for possibly quantified Σ_{AS} -formulas. This reduction is provided by rule Ext-Expand in Figure 5 which, given $x \approx t \in X$, adds an equisatisfiable formula $\llbracket x \approx t \rrbracket$ to G , the *expanded form* of $x \approx t$. The rule also removes $x \approx t$ from X , effectively marking it as processed. Since $\llbracket x \approx t \rrbracket$ keeps the (free) variables of $x \approx t$, any interpretation later found to satisfy $\llbracket x \approx t \rrbracket$ will also satisfy $x \approx t$.

The definition of expanded form for the possible cases of t are given below rule Ext-Expand.⁷ We remark that this rule can be applied only finitely many times. For an intuition of why, consider any ordering $<$ on function symbols such that $f < g$ if g is an extended function and f is not, and $\text{substr} < \text{contains} < \text{index_of} < \text{replace}$. In all cases, all function symbols of $\llbracket x \approx f(t) \rrbracket$ are smaller than f in this ordering.

Note that the reduction introduces formulas with *bounded integer quantification*, that is, formulas of the form $\forall k. 0 \leq k \leq t \Rightarrow \varphi$, where t does not contain k and φ is quantifier-free. Special consideration is needed to handle formulas of this form. We employ a pragmatic approach, modeled by the other two rules in Figure 5, which guesses upper bounds on certain arithmetic terms and eliminates those quantified formulas based on these bounds. Specifically, rule B-Val splits the search into two cases

⁷ We use a number of optimizations of this encoding in our implementation.

$t \leq n$ and $t > n$, where t is an integer term and n is a numeral. In the rule B-Inst, if G contains a formula φ having a subformula $\forall k. 0 \leq k \leq t \Rightarrow \psi$ and A entails a (concrete) upper bound on t , then that subformula is replaced by a finite conjunction. Since all quantifiers introduced in a configuration are bounded, these two rules used in combination suffice to eliminate them.

Example 3. Consider the formula $\varphi = \text{contains}(y, z) \wedge 0 < \text{len } y \leq 3 \wedge 0 < \text{len } z$. Applying Prop-Assign to φ results in a configuration where E , X and A respectively are $\{x \approx \top\}$, $\{x \approx \text{contains}(y, z)\}$, and $\{0 < \text{len } y \leq 3, 0 < \text{len } z\}$. Using Ext-Expand, we remove $x \approx \top$ from X and add to G , $\llbracket x \approx \text{contains}(y, z) \rrbracket$ which is:

$$(x \approx \top) \Leftrightarrow \forall k. 0 \leq k \leq \text{len } y - \text{len } z \Rightarrow \text{substr}(y, k, \text{len } z) \approx z.$$

Since $A \models_{\text{LIA}} \text{len } y - \text{len } z \leq 2$, by B-Inst we can replace this formula with:

$$((x \approx \top) \Leftrightarrow \bigwedge_{j=0}^2 \text{substr}(y, n_j, \text{len } z) \approx z) \wedge \bigwedge_{j=0}^2 n_j \approx j$$

where n_j is a fresh integer variable for $j = 0, 1, 2$. Applying Prop-Assign to this formula gives a branch where E and X are updated to $\{x \approx \top, m_0 \approx \text{len } z, x_0 \approx z\}$ and $\{x_0 \approx \text{substr}(y, n_0, m_0)\}$ for fresh variables x_0 and m_0 . Using the rule Ext-Expand, we remove the equality from X and add $\llbracket x_0 \approx \text{substr}(y, n_0, m_0) \rrbracket$ to G , which is:

$$\text{ite}(0 \leq n_0 < \text{len } y \wedge 0 \leq m_0, \\ y \approx \text{con}(z_1, x_0, z_2) \wedge \text{len } z_1 \approx 0 \wedge \text{len } z_2 \approx \text{len } y - m_0, x_0 \approx \epsilon)$$

with z_1, z_2 fresh string variables. Applying Prop-Assign again produces a branch with $E = \{x \approx \top, x_0 \approx z, m_0 \approx \text{len } z, y \approx \text{con}(z_1, x_0, z_2)\}$ and X empty. The set E is satisfiable in T_{ASX} . Deriving a saturated configuration from this point indicates that φ is satisfiable in T_{ASX} as well. Indeed, all interpretations that satisfy both $y \approx \text{con}(z_1, x_0, z_2)$ and $x_0 \approx z$ also satisfy $\text{contains}(y, z)$. \square

Although these rules give the general idea of our approach, our implementation actually handles the bounded quantifiers in a more sophisticated way, using model-based quantifier instantiation [23]. In a nutshell, we avoid generating all the instances of a quantified formula either by backtracking when a subset of them are unsatisfiable, or by determining that (sets of) instances are already satisfied by a candidate model.

3.2 Context-Dependent Simplification of Extended Function Terms

The reductions described above may be impractical due to the size and complexity of the formulas they introduce. For this reason, we have developed techniques for recognizing when the interpretation of an extended function term can be deduced based on the constraints in the current context. As a simple example, if $\text{contains}(\text{abc}, x)$ is a term belonging to the string theory (with abc a string constant), and the string solver has inferred that x is equal to a concrete value (e.g., d) or even a partially concrete value (e.g., $\text{con}(d, y)$), then we can already infer that $\text{contains}(\text{abc}, x)$ is equivalent to \perp , thereby avoiding the construction of its expanded form. We present next a generic technique for inferring such facts that has a substantial performance impact in practice.

$$\begin{array}{c}
\text{Ext-Simplify} \frac{x \approx t \in X \quad E \models \mathbf{y} \approx \mathbf{s} \quad (t\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow \text{ is a } \Sigma_{AS}\text{-term}}{G := G, [x \approx (t\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow] \quad X := X \setminus \{x \approx t\}} \\
\\
\text{Ext-Simplify-Pred} \frac{x \approx t \in X \quad E \models \mathbf{y} \approx \mathbf{s} \quad t : \text{Bool}}{G := G, [(x \approx \top) \Leftrightarrow ((t \approx \top)\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow] \quad X := X \setminus \{x \approx t\}} \\
\\
\text{Ext-Eq} \frac{x_1 \approx t_1, x_2 \approx t_2 \in X \quad E \models \mathbf{y} \approx \mathbf{s} \quad (t_1\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow = (t_2\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow}{S := S \oplus \{x_1 \approx x_2\}}
\end{array}$$

Fig. 6. Rules for context-dependent simplification of extended functions terms.

The rule Ext-Simplify from Figure 6 applies to configurations in which an extended function term t can be simplified to an equivalent form, modulo the current set of constraints, that does not involve extended functions. In this rule, we derive a set of equalities $\mathbf{y} \approx \mathbf{s}$ that are consequences of our current set of string constraints E , where typically \mathbf{y} are the (free) variables of t . We will refer to $\{\mathbf{y} \mapsto \mathbf{s}\}$ as a *derivable substitution (in E)*. If the simplified form of t under this substitution is a Σ_{AS} -term, then we add the equality $x \approx (t\{\mathbf{y} \mapsto \mathbf{s}\})\downarrow$ to G , and remove $x \approx t$ from X . Similarly when t is of sort `Bool`, the rule Ext-Simplify-Pred adds an equivalence to G based on the result of simplifying the formula $t \approx \top$ under a derivable substitution, and removes $x \approx t$ from X . The rule Ext-Eq is used to deduce equalities between extended terms that are syntactically identical after simplification under a derivable substitution.

These rules require methods for computing the simplified form $t\downarrow$ of Σ_{ASX} -terms t , as well as for choosing substitutions $\{\mathbf{y} \mapsto \mathbf{s}\}$. We describe these methods in the following, and give several examples.

Simplification rules for extended string functions Recall that by construction, a term t and its simplified form $t\downarrow$ are equivalent in T_{ASX} . It is generally advantageous to use techniques that often simplify Σ_{ASX} -terms t to Σ_{AS} -terms $t\downarrow$, since this eliminates the need to apply Ext-Expand to compute the expanded form of t . For this reason, we use aggressive and non-trivial simplification techniques when considering Σ_{ASX} -terms.

Examples of some of the simplification rules for `contains` are given in Figure 7. There, for string constants l_1, l_2 , we write $l_1 \setminus l_2$ to denote the empty string if l_1 does not contain l_2 , and the remainder obtained from removing the largest prefix of l_1 containing l_2 otherwise. We use $l_1 \sqcup_l l_2$ (resp. $l_1 \sqcup_r l_2$) to denote l_2 if l_1 contains l_2 , and the largest suffix (resp. prefix) of l_1 that is a prefix (resp. suffix) of l_2 otherwise. For example, $(abcde \setminus cd) = e$, $(abcde \setminus ba) = \epsilon$, $(abcde \sqcup_l def) = de$, $(abcde \sqcup_r def) = \epsilon$, and $(abcdc \sqcup_l cd) = cd$. Also, $s \rightarrow^* t$ indicates that t can be obtained from s by zero or more applications of the rules in the figure. One can prove that his rewrite system is terminating by noting that all conditions and right hand side of each non-trivial rule involve only concatenation terms with strictly fewer arguments than its left hand side.

In practice, the rules are implemented by a handful of recursive passes over the arguments of `contains` terms. Computing the simplified form of other operators is also fairly sophisticated and not shown here. (Our simplifier is around 2000 lines of C++ code.) Despite its complexity, simplification often results in significant performance

$\text{contains}(l_1, l_2) \rightarrow \top$		if l_1 contains l_2
$\text{contains}(l_1, l_2) \rightarrow \perp$		if l_1 does not contain l_2
$\text{contains}(l_1, \text{con}(l_2, \mathbf{t})) \rightarrow \perp$		if l_1 does not contain l_2
$\text{contains}(l_1, \text{con}(l_2, \mathbf{t})) \rightarrow \perp$		if $\text{contains}(l_1 \setminus l_2, \text{con}(\mathbf{t})) \rightarrow^* \perp$
$\text{contains}(l_1, \text{con}(x, \mathbf{t})) \rightarrow \perp$		if $\text{contains}(l_1, \text{con}(\mathbf{t})) \rightarrow^* \perp$
$\text{contains}(\text{con}(l_1, \mathbf{t}), l_2) \rightarrow \top$		if l_1 contains l_2
$\text{contains}(\text{con}(x, \mathbf{t}), s) \rightarrow \top$		if $\text{contains}(\text{con}(\mathbf{t}), s) \rightarrow^* \top$
$\text{contains}(\text{con}(t, s), \text{con}(t, \mathbf{u})) \rightarrow \top$		if $\text{contains}(\text{con}(s), \text{con}(\mathbf{u})) \rightarrow^* \top$
$\text{contains}(\text{con}(l_1, \mathbf{t}), l_2) \rightarrow \text{contains}(\text{con}(\mathbf{t}), l_2)$		if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\text{con}(t, l_1), l_2) \rightarrow \text{contains}(\text{con}(t), l_2)$		if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) \approx \top \rightarrow \epsilon \approx t$		
$\text{contains}(\text{con}(t_1, l_1, t_2), l_2) \approx \top \rightarrow \bigvee_{i=1}^2 \text{contains}(\text{con}(t_i), l_2) \approx \top$ if $l_1 \sqcup_r l_2 = l_1 \sqcup_l l_2 = \epsilon$		

Fig. 7. Examples of simplification rules for contains.

improvements, by eliminating the need to generate the expanded form of Σ_{ASX} -terms. We illustrate this in the following examples.

Example 4. Given input $y \approx bc \wedge \text{contains}(\text{con}(a, y), \text{con}(b, z, a))$, our calculus considers a configuration where E and X respectively are

$$\{y \approx bc, x_1 \approx \top, z_1 \approx \text{con}(a, y), z_2 \approx \text{con}(b, z, a)\} \text{ and } \{x_1 \approx \text{contains}(z_1, z_2)\}$$

where x_1, z_1, z_2 are fresh string variables. We have that $E \models z_1 \approx abc \wedge z_2 \approx \text{con}(b, z, a)$. Hence the substitution $\sigma = \{z_1 \mapsto abc, z_2 \mapsto \text{con}(b, z, a)\}$ is derivable in this configuration. Since $(\text{contains}(z_1, z_2)\sigma) \downarrow = \text{contains}(abc, \text{con}(b, z, a)) \downarrow = \perp$, we may apply Ext-Simplify to remove $x_1 \approx \text{contains}(z_1, z_2)$ from X and add $x_1 \approx \perp$ to E, after which unsat may be derived, since $x_1 \approx \top \in E$. In this example, we have avoided expanding the input formula by reasoning that $\text{con}(a, y)$ does not contain $\text{con}(b, z, a)$ in the context where y is bc . \square

Example 5. Given input $y \approx \text{con}(a, z) \wedge \text{contains}(\text{con}(x, y), bc)$, our calculus considers the configuration where E and X respectively are

$$\{y \approx \text{con}(a, z), x_1 \approx \top, z_1 \approx \text{con}(x, y), z_2 \approx bc\} \text{ and } \{x_1 \approx \text{contains}(z_1, z_2)\}$$

The substitution $\sigma = \{z_1 \approx \text{con}(x, a, z), z_2 \mapsto bc\}$ is derivable in this configuration. Computing $(\text{contains}(z_1, z_2)\sigma) \downarrow$ results in $\text{contains}(x, bc) \vee \text{contains}(z, bc)$. We may apply Ext-Simplify-Pred to remove $x_1 \approx \text{contains}(z_1, z_2)$ from X and add this formula to G, after which we consider the two disjuncts independently. \square

Example 6. Given input $y \approx ab \wedge \text{contains}(\text{con}(b, z), y) \wedge \neg \text{contains}(z, y)$, our calculus considers the configuration where E and X respectively are

$$\{y \approx ab, x_1 \approx \top, x_2 \approx \perp, z_1 \approx \text{con}(b, z)\} \text{ and } \{x_1 \approx \text{contains}(z_1, y), x_2 \approx \text{contains}(z, y)\}$$

The substitution $\sigma = \{y \approx ab, z_1 \approx \text{con}(b, z)\}$ is derivable in this configuration, and $(\text{contains}(z_1, y)\sigma) \downarrow = \text{contains}(z, ab) = (\text{contains}(z, y)\sigma) \downarrow$. Hence we can apply Ext-Eq to add $x_1 \approx x_2$ to E, after which unsat can be derived. \square

Choosing Substitutions A simple and general heuristic for choosing substitutions $\{y \mapsto s\}$ for terms t in the rules from Figure 6 is to map each variable y in t to some representative of its equivalence class $[y]$. We assume string constants are chosen as representatives whenever possible. We call this the *representative substitution* for t (in E). Representative substitutions are both easy to compute and often enough for reducing Σ_{ASX} -terms. A more powerful method for choosing substitutions is to consider substitutions that map each free variable y in t to $\text{con}(a_1, \dots, a_n) \downarrow$ where $N[y] = (a_1, \dots, a_n)$. We call this the *normal form substitution* for t . Intuitively, the normal form of t is a schema representing all known information about t . In this sense, a substitution mapping variables to their normal forms gives the highest likelihood of enabling our simplification techniques. In practice, our implementation takes advantage of both of these heuristics for choosing substitutions.

Example 7. Say we are in a configuration where the equivalence classes of \hat{E} are:

$$\{y, \text{con}(a, x)\} \quad \{x, \text{con}(z, c), u\} \quad \{z, w, b\}$$

and the normal forms are $N[y] = abc$, $N[x] = bc$ and $N[z] = b$. If y , x , and b are chosen as the representatives of these classes, the representative substitution σ_r for this configuration is $\{y \mapsto y, x \mapsto x, u \mapsto x, z \mapsto b, w \mapsto b\}$, whereas the normal form substitution σ_n is $\{y \mapsto abc, x \mapsto bc, u \mapsto bc, z \mapsto b, w \mapsto b\}$. Only the latter substitution suffices to show that $\text{contains}(y, \text{con}(z, z))$ is false in the current context, noting $(\text{contains}(y, \text{con}(z, z))\sigma_r) \downarrow = \text{contains}(y, \text{con}(z, z))$ and $(\text{contains}(y, \text{con}(z, z))\sigma_n) \downarrow = \perp$. \square

4 Implementation

We have implemented all of these techniques in the DPLL(T)-based SMT solver CVC4 [4]. At a high level, our implementation can be summarized as a particular strategy for applying the rules of the calculus, which we outline in the following.

Strategy 1 Start with a derivation tree consisting of (root) node $\langle \{\varphi\}, \emptyset, \emptyset \rangle$. Let t_{len} be $\text{len } x_1 + \dots + \text{len } x_m$ where x_1, \dots, x_m are the string variables of φ .

While the tree is not closed, consider as current configuration the left-most leaf in the tree that is not unsat and apply to it a derivation rule to that configuration, based on the steps below.

1. Let n be the smallest numeral such that $t_{\text{len}} > n \notin A$. If $t_{\text{len}} \leq n \notin A$, apply B-Val for t_{len} and n .
2. If G contains a formula φ with subformula $\forall k. 0 \leq k \leq t \Rightarrow \psi$, then let n be the smallest numeral such that $t > n \notin A$. If $t \leq n \in A$, apply B-Inst for φ and n . Otherwise, apply B-Val for t and n .
3. If possible, apply a rule from Figure 3, giving priority to A-Conf and S-Conf.
4. If possible, apply a rule from Figure 6 based on representative substitutions.
5. If possible, apply a rule from Figure 4.
6. If possible, apply a rule from Figure 6 based on normal form substitutions.

7. If X is non-empty, apply the rule *Ext-Expand* for some equality $x \approx t$ in X .

If no rule applies and the current configuration is acyclic, return *sat*. If the tree is closed, return *unsat*. \square

The strategy above is sound both for refutations and models, although it is not terminating in general.

Theorem 2. For all initial configurations $\langle \{\varphi\}, \emptyset, \emptyset \rangle$ where φ is a quantifier-free Σ_{ASX} -formula:

1. Strategy 1 returns *unsat* only if φ is unsatisfiable in T_{ASX} .
2. Strategy 1 returns *sat* only if φ is satisfiable in T_{ASX} .

Implementation While a comprehensive description of our implementation is beyond the scope of this work, we mention a few salient implementation details. The rule *Prop-Assign* is implemented by converting G to clausal normal form and giving the resulting clauses to a SAT solver with support for conflict-driven clause learning. The rule *A-Conf* is implemented by a standard theory solver for linear integer arithmetic. The rules of the calculus that modify the S component of our configuration are implemented in a dedicated $\text{DPLL}(T)$ *theory solver* for strings which generates *conflict clauses* when branches of a derivation tree are closed, and *theory lemmas* for rules that add formulas to G or A and those that have multiple conclusions. Conflict clauses are generated by tracking *explanations* so that each literal internally added to S can be justified in terms of input literals. Finally, we do not explicitly introduce fresh variables when constructing the set X , and instead record the set of extended terms that occur in E , which are implicitly treated as variables. We now revisit a few of the examples, giving concrete details on the operation of the solver.

Example 8. In Example 1, our input was $\text{len } x > \text{len } y \wedge y \approx \text{con}(x, a)$. For this input, the SAT solver finds a propositionally satisfying assignment that assigns both conjuncts to true, which causes the literal $\text{len } x > \text{len } y$ to be given to the theory solver for linear integer arithmetic, and $y \approx \text{con}(x, a)$ to be given to the theory solver for strings. This corresponds to an application of the rule *Prop-Assign*. The string solver sends $(\neg y \approx \text{con}(x, a) \vee \text{len } y \approx \text{len } x + 1)$ as a theory lemma to the SAT solver, corresponding to an application of the rule *L-Eq*. After that, the SAT solver assigns $\text{len } y \approx \text{len } x + 1$ to true, causing that literal to be asserted to the arithmetic solver, which subsequently generates a conflict clause of the form $(\neg \text{len } x > \text{len } y \vee \neg \text{len } y \approx \text{len } x + 1)$ corresponding to an application of *A-Conf*. After receiving this clause, the SAT solver is unable to find another satisfying assignment and causes the system to terminate with “unsat.” \square

Example 9. In Example 4, the string solver is given as input the literals $y \approx bc$ and $\text{contains}(\text{con}(a, y), \text{con}(b, z, a)) \approx \top$. The intermediate variables z_1 and z_2 are not explicitly introduced. Instead, using the substitution $\sigma = \{y \mapsto bc\}$, the solver directly infers that $\text{contains}(\text{con}(a, y), \text{con}(b, z, a))\sigma \downarrow = \perp$. Based on this simplification, it infers $\text{contains}(\text{con}(a, y), \text{con}(b, z, a)) \approx \perp$ with the explanation $y \approx bc$. Since the inferred literal conflicts with the second input literal, the string solver reports the conflict clause $\neg y \approx bc \vee \text{contains}(\text{con}(a, y), \text{con}(b, z, a)) \approx \top$. \square

Example 10. In Example 6, the equality $y \approx ab$ is the explanation for the substitution $\sigma = \{y \mapsto ab\}$ under which $\text{contains}(\text{con}(b, z), y)\sigma \downarrow = \text{contains}(z, ab) = \text{contains}(\text{con}(z, y))\sigma \downarrow$. Hence, the solver reports $(\neg y \approx ab \vee \neg \text{contains}(\text{con}(b, z), y) \vee \text{contains}(z, y))$ as a conflict clause in this example. \square

Example 11. Explanations are tracked for normal form substitutions as well. In Example 7, a possible explanation for the substitution σ_n is $y \approx \text{con}(a, x) \wedge x \approx \text{con}(z, c) \wedge u \approx x \wedge z \approx b \wedge w \approx z$. Explanations for simplifications that occur under the substitution σ_n must include these equalities. \square

In practice, we minimize explanations by only including the variables in substitutions that are relevant for certain inferences. In particular, the domain of derivable substitutions is restricted to the free variables of the terms they apply to. We further reduce this set based on dependency analysis. For example, $\text{contains}(abc, \text{con}(x, y)) \approx \perp$ can be explained by $x \approx d \wedge y \approx a$. However, $x \approx d$ alone is enough.

5 Evaluation

This section reports on our experimental evaluation of our approach for extended string constraints as implemented in the SMT solver CVC4.⁸ We used benchmark queries generated by running PyEx, a symbolic executor for Python programs, over a test suite that mimics the usage by a real-world Python developer. The technical details of our benchmark generation process are provided in Section 5.2.

We considered several configurations of CVC4 that differ in the subset of steps from Strategy 1 they apply. The default configuration, denoted **cvc4**, performs Steps 2, 3, 5 and 7 only. Configurations with suffix **f** (for “finite model finding”) perform Step 1, and configurations with suffix **s** (for “simplification”) perform Steps 4 and 6. For example, configuration **cvc4+fs** performs all seven steps. We consider other solvers for string constraints, including Z3-STR [30] (git revision e398f81) and Z3 [9] (version 4.5, git revision 24eae3f) which was recently extended with native support for strings.

5.1 Comparison with Other String Solvers

We first evaluated the raw performance of CVC4, Z3, and Z3-STR on the string benchmarks we collected with PyEx. We considered three sets of benchmarks produced by PyEx using **cvc4+fs**, Z3, and Z3-STR as the path constraint solver during program exploration. We denote these sets, which collectively consisted of 25,421 benchmark problems, as **PyEx-cfs**, **PyEx-z3** and **PyEx-z32**, respectively. We omit a small number (35) of these benchmarks for the following reasons: for 13 of them, at least one of the solvers produced a parse error; for the other 22, one solver returned a model (i.e., a satisfying assignment for the variables in the input problem) and agreed with its own model, but another solver answered “unsat” and disagreed with the model of the first solver.⁹ We

⁸ For details, see <http://cvc4.cs.stanford.edu/papers/CAV2017-strings/>.

⁹ We say solver *A* (*dis*)agrees with solver *B*’s model for input formula φ if *A* finds that $\varphi \wedge \mathcal{M}_B$ is (un)sat, where \mathcal{M}_B is a conjunction of equalities encoding *B*’s.

	PyEx-cfs (5557)			PyEx-z3 (8399)			PyEx-z32 (11430)			Total (25386)			
Solver	sat	unsat	×	sat	unsat	×	sat	unsat	×	sat	time	unsat	time
cvc4+fs	4229	1256	72	5694	1325	1380	10104	1194	132	20027	5h1m	3775	7m
cvc4+s	4133	1270	154	5461	1325	1613	9884	1193	353	19478	8h40m	3788	6m
cvc4+f	4160	1217	180	5571	1308	1520	9210	1145	1075	18941	6h38m	3670	6m
cvc4	4160	1213	184	5570	1308	1521	9211	1145	1074	18941	6h32m	3666	6m
z3	3421	1274	862	4925	1333	2141	7219	1196	3015	15565	11h30m	3803	3m
z3str2	2013	1278	2266	2803	1333	4263	4726	1182	5522	9542	15h47m	3793	13m

Fig. 8. Results of running each solver over benchmarks generated by PyEx over our test suite. All benchmarks run with a 30 second timeout.

attribute the parse errors to how the solvers process certain escape sequences in string constants, and the model discrepancies to minor differences in the semantics of substr when input indices are out of bounds.

All results were produced on StarExec [25], a public execution service for running comparative evaluations of logical solvers. The results for the three solvers on the three benchmark sets are shown in Figure 8 based on a 30 second timeout. The columns show the number of benchmarks that were determined to be satisfiable and unsatisfiable by each solver. The column with heading \times indicates the number of times the solver either timed out or terminated with an inconclusive response such as “unknown.” The best configuration of CVC4 (**cvc4+fs**) had a factor of 3.8 fewer timeouts than Z3, and a factor of 7.6 fewer timeouts or failures than Z3-STR in total over all benchmark sets. In particular, we note that **cvc4+fs** solved 1,451 unique benchmarks with respect to Z3 among those generated during a symbolic execution run using Z3 as the solver (**PyEx-z3**). Since PyEx supports concurrent solver invocation, this suggests a mixed-solver strategy that employs both **cvc4+fs** and Z3 would likely have reduced the number of failed queries for that run. For unsatisfiable benchmarks, the solvers were relatively closer in performance, where Z3 solved 24 more unsatisfiable benchmarks than **cvc4+fs**, which is not tuned for the unsatisfiable case due to its use of finite model finding. This further suggests a mixed-solver strategy would likely be beneficial for symbolic execution since it is often used for both program exploration (where sat leads to progress) and vulnerability checking (where unsat implies safety).

In addition to solving more benchmarks, **cvc4+fs** was significantly faster over them. Figure 9 plots the cumulative run time of the three solvers on benchmarks that each solves. With respect to Z3, which took 11 hours and 33 minutes on the 19,368 benchmarks it solves, **cvc4+fs** solved its first 19,368 benchmarks in 1 hour and 23 minutes, and overall took only 5 hours and 8 minutes on the 23,802 benchmarks it solves.

Using the context-dependent simplification techniques from Section 3.2, **cvc4+s** was able to solve 663 more benchmarks than **cvc4**, which does not apply simplification. By incorporating finite model finding, **cvc4+fs** was able to solve 536 more benchmarks in significantly less time, with a cumulative difference of more than 3 hours on solved benchmarks compared to **cvc4+s**. Taking the virtual best configuration of CVC4, our techniques find 20,594 benchmarks to be satisfiable, 567 more than **cvc4+fs**, indicating that a portfolio approach for the various configurations would be advantageous.

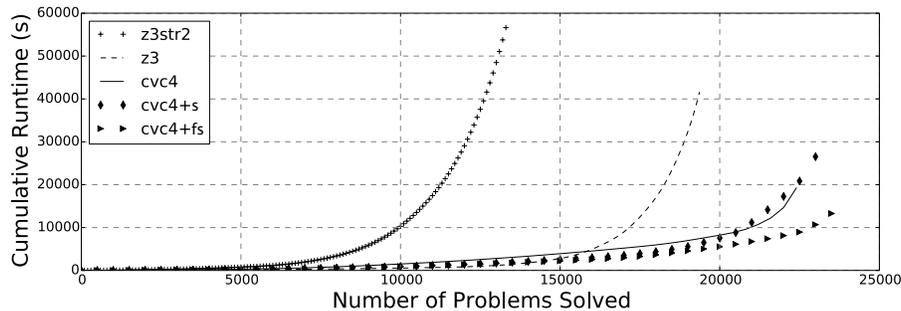


Fig. 9. Cactus plot of configurations of CVC4, Z3 and Z3-STR on solved benchmarks across all three benchmark sets.

We also measured how often CVC4 resorted to expanding extended function terms. We considered a modified configuration **cvc4+fs'** that is identical to **cvc4+fs** except that it does not use the rule Ext-Simplify. The 23,738 benchmarks solved by both **cvc4+fs** and **cvc4+fs'** had 619.2 extended function terms on average. On average over these benchmarks, the configuration **cvc4+fs'** found that 63.5 unique extended functions terms were relevant to satisfiability (e.g., were added to a configuration), and of these 24.3 were expanded (38%). Likewise, **cvc4+fs** found that 66.4 unique extended functions terms were relevant to satisfiability, and of these 12.6 were expanded (19%). With Ext-Simplify, it inferred 405.2 equalities per benchmark on average based on context-dependent simplification, showing that simplification is possible in a majority of contexts (97%). Limited to expansions that introduce universal quantification, which excludes expansions of substr and positively asserted contains, **cvc4+fs'** considered 11.4 expansions on average compared to 2.7 considered by **cvc4+fs**. This means that approximately 4 times fewer quantified formulas were introduced thanks to context-dependent simplification in **cvc4+fs**.

5.2 Symbolic Execution for Python

Our benchmarks were generated by PyEx, which is a symbolic executor designed to assist Python developers achieve high-coverage testing in their daily development routine, e.g., as part of the nightly tests run on the most recent version of a code base. To demonstrate the relative performance of CVC4, Z3, and Z3-STR in our nightly tests scenario, we ran PyEx on a test suite of 19 functions sampled from 4 popular Python packages: httpplib2, pip, pymongo, and requests. The set of queries generated during this experiment was used for our evaluation of the raw solver performance in Section 5.1. In this section, we show that the superior raw performance of CVC4 over other current solvers also translates into *real-life benefits* for PyEx users.

Our experiment was conducted on a developer machine featuring an Intel E3-1275 v3 quad-core processor (3.5GHz) and 32GB of memory. PyEx was run on each of the 19 functions for a maximum CPU time of 2 hours. By design, PyEx issues concurrent queries using Python multiprocessing when multiple queries are pending. To reflect the

configuration of our test machine, we capped the number of concurrent processes to 8. Note that, due to the nature of our test infrastructure, each of the 19 functions was tested in sequence and thus concurrency happened only within the testing of each individual function. In addition, PyEx has a notion of a per-path timeout, which is a heuristic to steer code exploration away from code paths that are stuck with hard queries. For this experiment, that timeout was set to 10 minutes.

We argue that the most important metrics for a developer are (i) the wall-clock time to run PyEx over the test suite and (ii) the coverage achieved over this time. In our experiments, PyEx with Z3 and with Z3-STR finished in 717 minutes and 829 minutes, respectively. By comparison, PyEx with the recommended configuration of CVC4 (**cvc4+fs**) finished in 295 minutes, which represents a speedup of 59% and 64% respectively over the other solvers. To compare coverage, we used the Python coverage library to measure both *line coverage*, the percentage of executed source lines, and *branch coverage*, the number of witnessed branch outcomes, of the test suite during symbolic execution. The line coverage of PyEx with **cvc4+fs**, Z3, and Z3-STR was respectively 8.48%, 8.41%, and 8.34%,¹⁰ whereas the branch coverage was 3612, 3895, and 3500. Taking both metrics into account, we conclude that **cvc4+fs** is highly effective for PyEx since it achieves similar coverage as the other tools while running significantly faster.

6 Concluding Remarks and Future Work

We have presented a calculus for extended string constraints that relies on both bounded quantifier elimination and context-dependent simplification techniques. The latter led to significant performance benefits for constraints coming from symbolic execution of Python programs. An implementation of these techniques in CVC4 has 3.8 times fewer timeouts and enables the PyEx symbolic executor to achieve comparable program coverage on our test suite while using only 41% of the runtime compared to other solvers.

Our analysis on program coverage indicates that an interesting research avenue for future work would be to determine correlations between certain features of models generated by string solvers and their utility in a symbolic executor, since different models may lead to different symbolic executions and hence different overall analyses.

We plan to develop context-dependent simplification techniques for other string functions, including conversion functions `str_to_int` and `int_to_str`, and to adapt these techniques to other SMT theories. Notably, we would like to use context-dependent simplification to optimize lazy approaches for fixed-width bit-vectors [14] where it is beneficial to avoid bit-blasting bit-vector operators, such as multiplication, that require elaborate encodings.

Acknowledgments This work was supported in part by the National Science Foundation under grants CNS-1228765, CNS-1228768, and CNS-1228827. We express our immense gratitude to Peter Chapman, who served as the first lead developer of PyEx.

¹⁰ Overall coverage appears to be low because we tested only some functions from each library.

References

- [1] P. A. Abdulla, M. F. Atig, Y. Chen, L. Holík, A. Rezine, P. Rümmer, and J. Stenman. String constraints for verification. In *Computer Aided Verification - 26th International Conference, CAV 2014, Proceedings*, pages 150–166, 2014.
- [2] P. A. Abdulla, M. F. Atig, Y. Chen, L. Holík, A. Rezine, P. Rümmer, and J. Stenman. Norn: An SMT solver for string constraints. In *Computer Aided Verification - 27th International Conference, CAV 2015, Proceedings*, pages 462–469, 2015.
- [3] T. Ball and J. Daniel. Deconstructing dynamic symbolic execution. In *Proceedings of the 2014 Marktoberdorf Summer School on Dependable Software Systems Engineering*. IOS Press, 2014.
- [4] C. Barrett, C. L. Conway, M. Deters, L. Hadarean, D. Jovanović, T. King, A. Reynolds, and C. Tinelli. CVC4. In *Proceedings of the 23rd International Conference on Computer Aided Verification, CAV'11*, pages 171–177. Springer-Verlag, 2011.
- [5] N. Bjørner, N. Tillmann, and A. Voronkov. Path feasibility analysis for string-manipulating programs. In *Proceedings of the 15th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS 2009*, pages 307–321. Springer-Verlag, 2009.
- [6] C. Cadar, D. Dunbar, and D. Engler. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proceedings of the 8th USENIX Symposium on Operating System Design and Implementation*, pages 209–224. USENIX, 2008.
- [7] S. K. Cha, T. Avgerinos, A. Rebert, and D. Brumley. Unleashing Mayhem on binary code. In *Proceedings of the 2012 IEEE Symposium on Security and Privacy*, pages 380–394. IEEE, 2012.
- [8] V. Chipounov, V. Kuznetsov, and G. Candea. S2E: A platform for in-vivo multi-path analysis of software systems. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 265–278. ACM, 2011.
- [9] L. De Moura and N. Bjørner. Z3: An efficient smt solver. In *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS'08/ETAPS'08*, pages 337–340. Springer-Verlag, 2008.
- [10] X. Fu and C. Li. A string constraint solver for detecting web application vulnerability. In *Proceedings of the 22nd International Conference on Software Engineering and Knowledge Engineering, SEKE'2010*. Knowledge Systems Institute Graduate School, 2010.
- [11] V. Ganesh, M. Minnes, A. Solar-Lezama, and M. Rinard. Word equations with length constraints: What's decidable? In *Proceedings of the 8th International Conference on Hardware and Software: Verification and Testing, HVC'12*, pages 209–226. Springer-Verlag, 2013.
- [12] H. Ganzinger, G. Hagen, R. Nieuwenhuis, A. Oliveras, and C. Tinelli. DPPL(T): fast decision procedures. In *Computer Aided Verification, 16th International Conference, CAV 2004, Proceedings*, pages 175–188, 2004.
- [13] P. Godefroid, M. Y. Levin, and D. Molnar. Automated whitebox fuzz testing. In *Proceedings of the 16th Annual Network and Distributed System Security Symposium*. Internet Society, 2008.
- [14] L. Hadarean, K. Bansal, D. Jovanovic, C. Barrett, and C. Tinelli. A tale of two solvers: Eager and lazy approaches to bit-vectors. In *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings*, pages 680–695, 2014.

- [15] P. Hooimeijer and M. Veanes. An evaluation of automata algorithms for string analysis. In *Proceedings of the 12th international conference on Verification, model checking, and abstract interpretation*, pages 248–262. Springer-Verlag, 2011.
- [16] A. Kiezun, V. Ganesh, P. J. Guo, P. Hooimeijer, and M. D. Ernst. HAMPI: a solver for string constraints. In *Proceedings of the eighteenth international symposium on Software testing and analysis*, pages 105–116. ACM, 2009.
- [17] J. C. King. Symbolic execution and program testing. *Communications of the ACM*, 19(7):385–394, 1976.
- [18] G. Li and I. Ghosh. PASS: string solving with parameterized array and interval automaton. In *9th International Haifa Verification Conference, HVC 2013, Haifa, Israel, November 5-7, 2013, Proceedings*, pages 15–31, 2013.
- [19] T. Liang, A. Reynolds, C. Tinelli, C. Barrett, and M. Deters. A DPLL(T) theory solver for a theory of strings and regular expressions. In *Proceedings of the 26th International Conference on Computer Aided Verification*, volume 8559 of *Lecture Notes in Computer Science*. Springer, 2014.
- [20] T. Liang, N. Tsiskaridze, A. Reynolds, C. Tinelli, and C. Barrett. A decision procedure for regular membership and length constraints over unbounded strings. In *Proceedings of the 10th International Symposium on Frontiers of Combining Systems*, volume 9322 of *Lecture Notes in Computer Science*, pages 135–150. Springer, 2015.
- [21] G. S. Makanin. The problem of solvability of equations in a free semigroup. *English transl. in Math USSR Sbornik*, 32:147–236, 1977.
- [22] W. Plandowski. Satisfiability of word equations with constants is in PSPACE. *Journal of the ACM*, 51(3):483–496, May 2004.
- [23] A. Reynolds, C. Tinelli, A. Goel, S. Krstic, M. Deters, and C. Barrett. Quantifier instantiation techniques for finite model finding in SMT. In *Automated Deduction - CADE-24 - 24th International Conference on Automated Deduction, Lake Placid, NY, USA, June 9-14, 2013. Proceedings*, pages 377–391, 2013.
- [24] N. Stephens, J. Grosen, C. Salls, A. Dutcher, R. Wang, J. Corbetta, Y. Shoshitaishvili, C. Kruegel, and G. Vigna. Driller: Augmenting fuzzing through selective symbolic execution. In *Proceedings of the Network and Distributed System Security Symposium*, 2016.
- [25] A. Stump, G. Sutcliffe, and C. Tinelli. Starexec: a cross-community infrastructure for logic solving. In S. Demri, D. Kapur, and C. Weidenbach, editors, *Proceedings of the 7th International Joint Conference on Automated Reasoning*, *Lecture Notes in Artificial Intelligence*. Springer-Verlag, 2014.
- [26] M. Trinh, D. Chu, and J. Jaffar. Progressive reasoning over recursively-defined strings. In *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part I*, pages 218–240, 2016.
- [27] M.-T. Trinh, D.-H. Chu, and J. Jaffar. S3: A symbolic string solver for vulnerability detection in web applications. In M. Yung and N. Li, editors, *Proceedings of the 21st ACM Conference on Computer and Communications Security*, 2014.
- [28] M. Veanes, N. Bjørner, and L. De Moura. Symbolic automata constraint solving. In *Proceedings of the 17th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning, LPAR’10*, pages 640–654. Springer-Verlag, 2010.
- [29] Y. Zheng, V. Ganesh, S. Subramanian, O. Tripp, J. Dolby, and X. Zhang. Effective search-space pruning for solvers of string equations, regular expressions and length constraints. In *Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I*, pages 235–254, 2015.
- [30] Y. Zheng, X. Zhang, and V. Ganesh. Z3-str: A z3-based string solver for web application analysis. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2013*, pages 114–124. ACM, 2013.

This document is being distributed for informational and educational purposes only and is not an offer to sell or the solicitation of an offer to buy any securities or other instruments. The information contained herein is not intended to provide, and should not be relied upon for investment advice. The views expressed herein are not necessarily the views of Two Sigma Investments, LP or any of its affiliates (collectively, "Two Sigma"). Such views reflect significant assumptions and subjective of the author(s) of the document and are subject to change without notice. The document may employ data derived from third-party sources. No representation is made as to the accuracy of such information and the use of such information in no way implies an endorsement of the source of such information or its validity.

The copyrights and/or trademarks in some of the images, logos or other material used herein may be owned by entities other than Two Sigma. If so, such copyrights and/or trademarks are most likely owned by the entity that created the material and are used purely for identification and comment as fair use under international copyright and/or trademark laws. Use of such image, copyright or trademark does not imply any association with such organization (or endorsement of such organization) by Two Sigma, nor vice versa.