

High-Level Abstractions for Simplifying Extended String Constraints in SMT

Andrew Reynolds¹, Andres Nötzli²,
Clark Barrett², and Cesare Tinelli¹

¹ Department of Computer Science, The University of Iowa

² Department of Computer Science, Stanford University



Abstract. Satisfiability Modulo Theories (SMT) solvers with support for the theory of strings have recently emerged as powerful tools for reasoning about string-manipulating programs. However, due to the complex semantics of *extended string functions*, it is challenging to develop scalable solvers for the string constraints produced by program analysis tools. We identify several classes of simplification techniques that are critical for the efficient processing of string constraints in SMT solvers. These techniques can reduce the size and complexity of input constraints by reasoning about arithmetic entailment, multisets, and string containment relationships over input terms. We provide experimental evidence that implementing them results in significant improvements over the performance of state-of-the-art SMT solvers for extended string constraints.

1 Introduction

Most programming languages support strings natively and a considerable number of programs perform some form of string manipulation. Automated reasoning about string-manipulating programs for verification and test case generation purposes is then highly relevant for these languages and programs. Applications to security, such as finding SQL injection and XSS vulnerabilities in web applications [18, 16, 23] or proving their absence, are of critical importance. String constraints have also been used to generate relational database tables from SQL queries for unit testing purposes [21]. These applications require modeling all of the string operations that appear in real programs. This is challenging since some of those operations are complex and often realized by iterative applications of simpler operations. Additionally, since strings in many programming languages have variable length, reasoning accurately about them cannot be done by a reduction to bounded types such as bit-vectors, and requires instead the development of solvers for *unbounded* strings. To make this type of reasoning more scalable, the use of dedicated theory solvers natively supporting common string operations has been proposed [5, 9]. Some string solvers are fully integrated within Satisfiability Modulo Theories (SMT) solvers [12, 4]; some are built (externally) on top of such solvers [19, 9, 16]; and others are independent of SMT solvers [23].

A major challenge in developing solvers for unbounded string constraints is the complex semantics of *extended string functions* beyond the basic operations of string concatenation and equality. Extended functions include *replace*, which replaces a string

in another string, and `indexof`, which returns the position of a string in another string. Another challenge is that constraints using extended functions are often combined with constraints over other theories, e.g. integer constraints over string lengths or applications of `indexof`, which requires the involvement of solvers for those theories. Current string solvers address these challenges by reducing constraints with extended string functions to typically more verbose constraints over basic functions. As with every reduction, some of the higher level structure of the problem may be lost, with negative repercussions on the performance and scalability.

To address this issue, we have developed new techniques that reason about constraints with extended string operators before they are reduced to simpler ones. This analysis of complex terms can often eliminate the need for expensive reductions. The techniques are based on reasoning about relationships over strings with high-level abstractions, such as their arithmetic relationships (e.g., reasoning about their length), their string containment relationships, and their relationships as multisets of characters. We have implemented these techniques in CVC4, an SMT solver with native support for string reasoning. An experimental evaluation with benchmarks from various applications shows that our new techniques allows CVC4 to significantly outperform other state-of-the-art solvers that target extended string constraints.

Our main contributions are:

- A novel procedure for proving entailments over arithmetic predicates built from the theory of strings and linear integer arithmetic.
- Extensions of this technique for showing containment relationships between strings.
- A novel simplification technique based on abstracting strings as multisets.
- Experimental evidence that the simplification techniques provide significant performance improvements over current state-of-the-art solvers.

In the remainder of this section, we discuss related work. In Section 2, we provide some background on the theory of strings and how solvers reduce extended functions. In Sections 3, 4 and 5, we describe, respectively, our arithmetic-based, containment-based, and multiset-based simplification techniques. Section 6 describes our implementation of those techniques, and Section 7 presents our evaluation.

Related Work Various approaches to solving constraints over extended string functions have been proposed. Saxena et al. [16] showed that constraints from the symbolic execution of JavaScript code contain a significant number of extended string functions, which underlines their importance. Their approach translates string constraints to bit-vector constraints, similar to other approaches based on bounded strings such as HAMPI [9]. Bjørner et al. [5] proposed native support for extended string operators in string solvers for scaling symbolic execution of .NET code. They reduce extended string functions to basic ones after getting bounds for string lengths from an integer solver. They also showed that constraints involving unbounded strings and `replace` are undecidable. PASS [11] reduces string constraints over extended functions to arrays. Z3-str and its successors [25, 24, 4] reduce extended string functions to basic functions eagerly during preprocessing. S3 [18] reduces recursive functions such as `replace` incrementally by splitting and unfolding. Its successor S3P [19] refines this reduction by pruning the resulting subproblems for better performance. CVC4 [3] reduces constraints with extended functions lazily and leverages context-dependent simplifications to simplify

$$\begin{array}{llll}
n : \text{Int} \text{ for all } n \in \mathbb{N} & + : \text{Int} \times \text{Int} \rightarrow \text{Int} & - : \text{Int} \rightarrow \text{Int} & \geq : \text{Int} \times \text{Int} \rightarrow \text{Bool} \\
l : \text{Str} \text{ for all } l \in \mathcal{A}^* & \text{con} : \text{Str} \times \cdots \times \text{Str} \rightarrow \text{Str} & & \text{len} : \text{Str} \rightarrow \text{Int} \\
\text{substr} : \text{Str} \times \text{Int} \times \text{Int} \rightarrow \text{Str} & & \text{contains} : \text{Str} \times \text{Str} \rightarrow \text{Bool} & \\
\text{indexof} : \text{Str} \times \text{Str} \times \text{Int} \rightarrow \text{Int} & & \text{replace} : \text{Str} \times \text{Str} \times \text{Str} \rightarrow \text{Str} & \\
\text{str.to.int} : \text{Str} \rightarrow \text{Int} & & \text{int.to.str} : \text{Int} \rightarrow \text{Str} &
\end{array}$$

Fig. 1. Functions in signature Σ_S . Str and Int denote strings and integers respectively.

the reductions [15]. TRAU [1] reduces certain extended functions, such as replace, to context-free membership constraints. OSTRICH [7] implements a decision procedure for a subset of constraints that include extended string functions. The simplification techniques presented in this paper are agnostic to the underlying solving procedure, so they can be combined with all of these approaches.

2 Preliminaries

We work in the context of many-sorted first-order logic with equality and assume the reader is familiar with the notions of signature, term, literal, formula, and formal interpretation of formulas. We review a few relevant definitions in the following. 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 . We assume Σ contains the equality predicate \approx , interpreted as the identity relation, and the predicates \top (for true) and \perp (for false). A Σ -formula φ is *satisfiable* (resp., *unsatisfiable*) in T if it is satisfied by some (resp., no) interpretation in \mathbf{I} . We write $\models_T \varphi$ to denote that the Σ -formula φ is *T-valid*, i.e., is satisfied in every model of T . Two Σ -terms t_1 and t_2 are *equivalent in T* if $\models_T t_1 \approx t_2$.

We consider an extended theory T_S of strings and length equations, whose signature Σ_S is given in Figure 1 and whose models differ only on how they interpret variables.³ We assume a fixed finite alphabet \mathcal{A} of characters which includes the digits $\{0, \dots, 9\}$. The signature includes the sorts Bool, Int, and Str denoting the Booleans, the integers (\mathbb{Z}), and Kleene closure of \mathcal{A} (\mathcal{A}^*), respectively. The top half of Figure 1 includes the usual symbols of *linear* integer arithmetic, interpreted as expected, a *string literal* l for each word/string of \mathcal{A}^* , a variadic function symbol con, interpreted as word concatenation, and a function symbol len, interpreted as the word length function. We write ϵ for the empty word and abbreviate $\text{len}(s)$ as $|s|$. We use words over the characters a, b, and c, as in abca, as concrete examples of string literals.

We refer to the function symbols in the bottom half of the figure as *extended functions* and refer to terms containing them as *extended terms*. A *position* in a string $l \in \mathcal{A}^*$ is a non-negative integer n smaller than the length of l that identifies the $(n + 1)^{\text{th}}$ character of l —with 0 identifying the first character, 1 the second, and so on. For all models \mathcal{I} of T_S , all $l, l_1, l_2 \in \mathcal{A}^*$, and $n, m \in \mathbb{Z}$, $\text{substr}^{\mathcal{I}}(l, n, m)$ (the interpretation of substr in \mathcal{I} applied to l, n, m) is the longest substring of l starting at position n with length at most

³ Our implementation supports a larger set of symbols, but for brevity, we only show the subset of the symbols used throughout this paper.

m , or ϵ if n is an invalid position or m is not positive; $\text{contains}^{\mathcal{I}}(l_1, l_2)$ is true if and only if l_2 is a substring of l_1 , with ϵ being a substring of every string; $\text{indexof}^{\mathcal{I}}(l_1, l_2, n)$ is the position of the first occurrence of l_2 in l_1 at or after position n , n if l_2 is empty and $0 \leq n \leq |l_1|$, and -1 if n is an invalid position, or if no such occurrence exists; $\text{replace}^{\mathcal{I}}(l, l_1, l_2)$ is the result of replacing the first occurrence of l_1 in l by l_2 , l if l does not contain l_1 , or the result of prepending l_2 to l if l_1 is empty; $\text{str.to.int}^{\mathcal{I}}(l)$ is the non-negative integer represented by l in decimal notation or -1 if the string contains non-digit characters; $\text{int.to.str}^{\mathcal{I}}(n)$ is the result of converting n to the corresponding string in decimal notation if n is non-negative, or ϵ otherwise. We write $\text{substr}(t, u)$ as shorthand for the term $\text{substr}(t, u, |t|)$, i.e. the suffix of t starting at position u .

Note that the semantics for `replace` and `indexof` correspond to the semantics in the current draft of the SMT-LIB standard for the theory of strings [17]; they are slightly different from the ones described in previous work [4, 15, 20].

2.1 Solving Extended String Constraints (with Simplification)

Various efficient solvers have been designed for the satisfiability problem for quantifier-free $T_{\mathcal{S}}$ -constraints, including CVC4 [3], S3# [20] and Z3STR3 [4]. In this section, we give an overview of how these solvers process extended functions in practice.

Generally speaking, constraints involving extended functions are converted to basic ones through a series of reductions performed in an incremental fashion by the solver. Operators whose reduction requires universal quantification are dealt with by guessing upper bounds on the lengths of input strings or by lazily adding constraints that block models that do not satisfy extended string constraints.

Example 1. To determine the satisfiability of $\neg\text{contains}(t, s)$, the application of `contains` is reduced to constraints that ensure that s is not a substring of t at any position. Assuming we have a fixed upper bound n on the length of t , the above constraint is equivalent to the finite conjunction $\text{substr}(t, 0, |s|) \not\approx s \wedge \dots \wedge \text{substr}(t, n, |s|) \not\approx s$. Each application of `substr` is then eliminated by introducing an equality that constrains a fresh variable x_i to have the semantics of that substring. Thus, reducing the formula above results in

$$\bigwedge_{i=0}^n |t| \geq i + |s| \Rightarrow (x_i \not\approx s \wedge t \approx \text{con}(x_i^{\text{pre}}, x_i, x_i^{\text{post}}) \wedge |x_i^{\text{pre}}| \approx i \wedge |x_i| \approx |s|)$$

where $x_i, x_i^{\text{pre}}, x_i^{\text{post}}$ are fresh string variables.⁴ The above conjunction involves only string concatenation, string length, and equality, and thus can be handled by a string solver with support for word equations with length constraints.

The reduction in Example 1 introduces $5 \cdot n$ theory literals over basic string functions and $3 \cdot n$ string variables. A full reduction accounting for all corner cases of `substr` is even more complex and thus more expensive to process, even for small values of n . These performance challenges can be addressed by aggressive simplifications that *eliminate* extended functions using high-level reasoning, as shown in the next example.

⁴ This formula is a simplified form of the general reduction. The general reduction also expresses that i is a valid position in t and that the third argument of `substr` is non-negative [15].

Example 2. Consider an instance of the previous example where $s = \text{con}(a, x)$ and $t = \text{con}(b, \text{substr}(x, 0, n))$. A full reduction of $\neg\text{contains}(t, s)$ that eliminates all applications of substr , including those in t , introduces $5 \cdot n + 5$ new theory literals and $3 \cdot n + 3$ string variables. However, based on the semantics of contains it is easy to see that $\neg\text{contains}(t, s)$ is T_5 -valid: if t were to contain s , then s would have to occur in the portion of t after its first character b , since the first character of s is a . However, $\text{con}(a, x)$ cannot be contained in $\text{substr}(x, 0, n)$, since the length of the former is at least $|x| + 1$, while the length of the latter is at most $|x|$. A solver which recognizes that $\neg\text{contains}(t, s)$ can be simplified to \top in this case can avoid the reduction altogether.

We advocate for aggressive simplification techniques to improve the performance of string solvers for extended functions. In the next sections, we describe several classes of such techniques that can be applied to inputs as a preprocessing step or during solving as part of a context-dependent solving strategy [15]. We present them as sets R of rewrite rules of the form $t \rightarrow_R s$, where s is a (simplified) term equivalent to t in T_5 . We assume a deterministic application strategy for these rules, such that each term t rewrites to a unique *simplified form*, denoted by $t \downarrow$, which is irreducible by the rules. We split our simplifications into four categories, presented in Figures 4, 6, 7 and 8.⁵

3 Arithmetic-Based String Simplification

To simplify string terms, it is useful to establish relationships between quantities such as the lengths of strings. For example, $\text{contains}(t, s)$ can be simplified to \perp for a particular s and t if it can be inferred that $|s|$ is strictly greater than $|t|$. This section defines an inference system for such arithmetic relationships and the simplifications that it enables.

We are interested in proving the T_5 -validity of formulas of the form $u \geq 0$, where u is a Σ_5 -term of integer type. We describe an inference system as a set of rules for deriving judgments of the form $\vdash u \geq 0$ and a specific rule application strategy we have implemented. The inference system is *sound* in the sense that $\models_{T_5} u \geq 0$ whenever $\vdash u \geq 0$ is derivable in it. It is, however, *incomplete* as it may fail to derive $\vdash u \geq 0$ in some cases when $\models_{T_5} u \geq 0$. This incompleteness is by design, since proving the T_5 -validity of inequalities is generally expensive due to the NP-hardness of linear integer arithmetic. Without loss of generality, we require that the term u be in a simplified form, where terms of the form $|l|$ with l a string literal of n characters are rewritten to n , terms of the form $|\text{con}(t_1, \dots, t_n)|$ are rewritten to $|t_1| + \dots + |t_n|$, and like monomials in arithmetic terms are combined in the usual way (e.g., $2 \cdot |x| + |x|$ is rewritten to $3 \cdot |x|$).

Definition 1 (Polynomial Form). *An arithmetic term u is in polynomial form if $u = m_1 \cdot u_1 + \dots + m_n \cdot u_n + m$, where m_1, \dots, m_n are non-zero integer constants, m is an integer constant, and each u_1, \dots, u_n is a unique term and one of the following:*

1. *an integer variable,*
2. *an application of length to a string variable, e.g. $|x|$,*
3. *an application of length to an extended function, e.g. $|\text{substr}(t, v, w)|$, or*
4. *an application of an extended function of integer type, e.g. $\text{indexof}(t, s, v)$.*

⁵ Some specialized rules have been omitted for space reasons.

$$\begin{array}{c}
\frac{u \rightarrow_U^* n \quad n \geq 0}{\vdash u \geq 0} \quad \text{where} \\
|t| \rightarrow_U 0 \\
|\text{substr}(t, v, w)| \rightarrow_U \begin{cases} w & \text{if } \vdash v \geq 0 \text{ and } \vdash |t| \geq v + w \\ |t| - v & \text{if } \vdash v \geq 0 \text{ and } \vdash v + w \geq |t| \end{cases} \\
|\text{replace}(t, s, r)| \rightarrow_U \begin{cases} |t| & \text{if } \vdash |r| \geq |s| \text{ or } \vdash |r| \geq |t| \\ |t| - |s| & \end{cases} \\
|\text{int.to.str}(v)| \rightarrow_U 1 \quad \text{if } \vdash v \geq 0 \\
|\text{indexof}(t, s, v)| \rightarrow_U -1 \\
|\text{str.to.int}(t)| \rightarrow_U -1 \\
m \cdot v + u' \rightarrow_U m \cdot w + u' \text{ if } v \rightarrow_U w \text{ and } m > 0 \text{ or } v \rightarrow_O w \text{ and } m < 0 \\
\\
|\text{substr}(t, v, w)| \rightarrow_O w \quad \text{if } \vdash w \geq 0 \\
|\text{substr}(t, v, w)| \rightarrow_O \begin{cases} |t| - v & \text{if } \vdash |t| \geq v \\ |t| & \end{cases} \\
|\text{replace}(t, s, r)| \rightarrow_O \begin{cases} |t| & \text{if } \vdash |s| \geq |r| \\ |t| + |r| & \end{cases} \\
|\text{int.to.str}(v)| \rightarrow_O \begin{cases} v & \text{if } \vdash v > 0 \\ v + 1 & \text{if } \vdash v \geq 0 \end{cases} \\
|\text{indexof}(t, s, v)| \rightarrow_O \begin{cases} |t| - |s| & \text{if } \vdash |t| \geq |s| \\ |t| & \end{cases} \\
m \cdot v + u' \rightarrow_O m \cdot w + u' \text{ if } v \rightarrow_O w \text{ and } m > 0 \text{ or } v \rightarrow_U w \text{ and } m < 0
\end{array}$$

Fig. 2. Rules for arithmetic entailment based on under- and over-approximations computed for arithmetic terms containing extended string operators. We write t, s, r to denote string terms, u, u', v, w to denote integer terms and m, n to denote integer constants.

Given u in polynomial form, our inference system uses a set of over- and under-approximations for showing that $u \geq 0$ holds in all models of T_S . We define two auxiliary rewrite systems, denoted \rightarrow_O and \rightarrow_U . If u rewrites to v (in zero or more steps) in \rightarrow_O , written $u \rightarrow_O^* v$, we say that v is an *over-approximation* of u . We can prove in that case that $\models_{T_S} v \geq u$. Dually, if u rewrites to v in \rightarrow_U , written $u \rightarrow_U^* v$, we say that v is an *under-approximation* of u and can prove that $\models_{T_S} u \geq v$. Based on these definitions, the core of our inference system can be summarized by the single inference rule schema provided in Figure 2 together with the conditional rewrite systems \rightarrow_O and \rightarrow_U which are defined inductively in terms of the inference system and each other.

A majority of the rewrite rules have side conditions requiring the derivability of certain judgments in the same inference system. To improve their readability we take some liberties with the notation and write $\vdash u_1 \geq u_2$, say, instead of $\vdash u_1 - u_2 \geq 0$. For example, $|\text{substr}(t, v, w)|$ is under-approximated by w if it can be inferred that the interval from v to $v + w$ is a valid range of positions in string t , which is expressed by the side conditions $\vdash v \geq 0$ and $\vdash |t| \geq v + w$. Note that some arithmetic terms, such as $|\text{substr}(t, v, w)|$, can be approximated in *multiple* ways—hence the need for a strategy for choosing the best approximation for arithmetic string terms, described later. The rules for polynomials are written modulo associativity of $+$ and state that a monomial $m \cdot v$ in them can be over- or under-approximated based on the sign of the coefficient m .

STR-ARITH-APPROX(u), where $u = u_x + u_\ell + u_s + m$ and:

- $u_x = m_1^y \cdot y_1 + \dots + m_n^y \cdot y_n$,
- $u_\ell = m_1^\ell \cdot |x_1| + \dots + m_p^\ell \cdot |x_p|$,
- $u_s = m_1^v \cdot v_1 + \dots + m_q^v \cdot v_q$.

for variables $x_1, \dots, x_p, y_1, \dots, y_n$ and extended terms v_1, \dots, v_q :

1. If $q > 0$, choose a v_i and v_i^a that maximize the following criteria (in descending order), where $u' = (u[m_i^v \cdot v_i \mapsto m_i^v \cdot v_i^a]) \downarrow$:
 - (a) (Soundness) $v_i \rightarrow_U v_i^a$ if $m_i^v > 0$ and $v_i \rightarrow_O v_i^a$ if $m_i^v < 0$;
 - (b) (Avoids new terms) Minimizes the size of $\text{negcoeff}(u') \setminus \text{negcoeff}(u)$;
 - (c) (Cancels existing terms) Maximizes the size of $\text{negcoeff}(u) \setminus \text{negcoeff}(u')$.
 Return $u \rightarrow_U u'$.
2. If $p > 0$ and $m_j^\ell > 0$ for some j , return $u \rightarrow_U (u[m_j^\ell \cdot |x_j| \mapsto 0]) \downarrow$.

Fig. 3. A greedy strategy for showing arithmetic entailments in the theory T_5 . We write $\text{negcoeff}(u)$ to denote the set of terms whose coefficient is negative in u .

For simplicity, we silently assume in the figure that basic arithmetic simplifications are applied after each rewrite step to put the right-hand side in polynomial form.

Example 3. Let u be $|\text{replace}(x, aa, b)|$. Because $\vdash |aa| \geq |b|$, the first case of the over-approximation rule for `replace` applies, and we get that $u \rightarrow_O |x|$. This reflects that the result of replacing the first occurrence, if any, of `aa` in x with `b` is no longer than x .

Example 4. Let u be the same as in the previous example and let v be $-1 \cdot u + 2 \cdot |x|$. Since $u \rightarrow_O |x|$ and the coefficient of u in v is negative, we have that $v \rightarrow_U -1 \cdot |x| + 2 \cdot |x|$, which simplifies to $|x|$; moreover, $|x| \rightarrow_U 0$. Thus, $v \rightarrow_U^* 0$ and so $\vdash v \geq 0$. In other words, we can use the approximations to show that u is at most $2 \cdot |x|$.

3.1 A Strategy for Approximation

The rewrite systems \rightarrow_O and \rightarrow_U allow for many possible derivations. Thus, it is important to devise a strategy that is efficient and succeeds often in practice. We use a greedy rule application strategy that favors rule applications leading to the cancellation of monomials. For example, consider the term $|x| - |\text{substr}(y, 0, |x|)|$, and observe that the subtrahend can be over-approximated either by $|y|$ or by $|x|$. However, proving the T_5 -validity of $|x| - |\text{substr}(y, 0, |x|)| \geq 0$ with the former over-approximation is impossible since $|x| - |y| \geq 0$ does not hold in all models of T_5 . In contrast, the latter approximation produces $|x| - |x| \geq 0$ which is trivially T_5 -valid.

Recall that, given an arithmetic inequality $u \geq 0$, our goal is to find a reduction $u \rightarrow_U^* n$ where n is a non-negative constant. Our strategy for choosing which rule of \rightarrow_U to apply to u is given in Figure 3. We decompose u into three parts: the portion u_x consisting of a sum of integer variables, the portion u_ℓ consisting of a sum of lengths of string variables, and the remaining portion u_s which is a sum of monomials involving extended terms v_1, \dots, v_q as defined in Definition 1.

Since there are multiple choices for how terms in u_s are approximated, the strategy focuses primarily on this portion. In particular, we apply an approximation for one of

the terms v_i , under-approximating or over-approximating depending on the sign of its coefficient, and replace the monomial in t by its corresponding approximation. The choice of v_i and v_i^a is based on maximizing the likelihood that the overall derivation will produce a non-negative constant.

For a term u in polynomial form, let $\text{negcoeff}(u)$ be a set of integer terms whose coefficient is negative in u , e.g. $\text{negcoeff}(y_1 + -1 \cdot y_2) = \{y_2\}$. Terms in this set can be seen as *obligations* for proving entailments in our derivations since if $y_2 \in \text{negcoeff}(u)$, it must be the case that our derivation applies a rule that introduces a term with a positive coefficient for y_2 . In Figure 3, we say that our choice of $v_i \rightarrow_U v_i^a$ *avoids new terms* if it does not have the effect of adding any new terms to $\text{negcoeff}(u)$, and *cancel existing terms* if it has the effect of removing terms from this set. If the portion u_s is empty, we apply the rule $|x_j| \rightarrow_U 0$ if there exists a monomial $m_j^\ell \cdot |x_j|$ where m_j^ℓ is positive. This rule is applied with lowest priority because these monomials may help to cancel negative terms introduced by the other steps.

Step 1 depends on knowing the set of possible one-step approximations $v_i \rightarrow_U v_i^a$ and $v_i \rightarrow_O v_i^a$ for terms from u . These are determined using the rules of Figure 2. Whenever applicable, we break ties between rewrites in Step 1 by considering a fixed arbitrary ordering over extended terms.

Example 5. Let u be $1 + |t_1| + |t_2| - |x_1|$, where t_1 is $\text{substr}(x_2, 1, |x_2| + |x_4|)$ and t_2 is $\text{replace}(x_1, x_2, x_3)$. Step 1 of STR-ARITH-APPROX considers the possible approximations $|t_1| \rightarrow_U |x_2| - 1$ and $|t_2| \rightarrow_U |x_1| - |x_2|$. Note that under-approximations are needed because the coefficients of $|t_1|$ and $|t_2|$ are positive. The first approximation is an instance of the third rule in Figure 2, noting that both $\vdash 1 \geq 0$ and $\vdash 1 + |x_2| + |x_4| \geq |x_2|$ are derivable by a *basic* strategy that, wherever applicable, under-approximates string length terms as zero. Our strategy chooses the first approximation since it introduces no new negative coefficient terms, thus obtaining: $u \rightarrow_U |x_2| + |t_2| - |x_1|$. We now choose the approximation $|t_2| \rightarrow_U |x_1| - |x_2|$, noting that it introduces no new negative coefficient terms and cancels an existing one, $|x_1|$. After arithmetic simplification, we have derived $u \rightarrow_U^* 0$, and hence $\vdash u \geq 0$.

One can show that our strategy is sound, terminating, and deterministic. This means that applying STR-ARITH-APPROX to completion produces a unique rewrite chain of the form $t \rightarrow_U u_1 \rightarrow_U \dots \rightarrow_U u_n$ for a finite n , where each step is an application of one of the rewrite rules from Figure 2.

3.2 Simplification Rules with Arithmetic Side Conditions

We use the inference system from the previous section for simplifications of string terms with arithmetic side conditions. Figure 4 summarizes those simplifications.

The first rule rewrites a string equality to \perp if one of the two sides can be inferred to be strictly longer than the other. In the second rule, if one side of an equality, $\text{con}(s, r, q)$, is such that the sum of lengths of s and q alone can be shown to be greater than or equal to the length of the other side, then r must be empty. The third rule recognizes that string containment reduces to string equality when it can be inferred that string s is at least as long as the string t that must contain it. The next rule captures the fact

$t \approx s \rightarrow$	\perp	$\text{if } \vdash t \geq s + 1$
$t \approx \text{con}(s, r, q) \rightarrow$	$t \approx \text{con}(s, q) \wedge r \approx \epsilon$	$\text{if } \vdash s + q \geq t $
$\text{contains}(t, s) \rightarrow$	$t \approx s$	$\text{if } \vdash s \geq t $
$\text{substr}(t, v, w) \rightarrow$	ϵ	$\text{if } \vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(\text{con}(t, s), v, w) \rightarrow$	$\text{substr}(s, v - t , w)$	$\text{if } \vdash v \geq t $
$\text{substr}(\text{con}(s, t), v, w) \rightarrow$	$\text{substr}(s, v, w)$	$\text{if } \vdash s \geq v + w$
$\text{substr}(\text{con}(t, s), 0, w) \rightarrow$	$\text{con}(t, \text{substr}(s, 0, w - t))$	$\text{if } \vdash w \geq t $
$\text{indexof}(t, s, v) \rightarrow$	$\text{ite}(\text{substr}(t, v) \approx s, v, -1)$	$\text{if } \vdash v + s \geq t $

Fig. 4. String simplification rules. Letters t, s, r, q denote string terms; v, w denote integer terms.

that substr simplifies to the empty string if it can be inferred that its position v is not within bounds, or its length w is not positive. In the figure, we write that rule with a disjunctive side condition; this is a shorthand to denote that we can pick any disjunct and show that it holds assuming the negation of the other disjuncts. We can use those assumptions to perform substitutions to simplify the derivation. Concretely, to show $\vdash u_1 \geq u_2 \vee \dots \vee u \not\approx u'$ it is sufficient to infer $\vdash (u_1 \geq u_2)[u \mapsto u']$. We demonstrate this with an example.

Example 6. Consider the term $\text{substr}(t, |t| + w, w)$. Our rules may simplify this term to ϵ by inferring that its start position ($|t| + w$) is not within the bounds of t if we assume that its size (w) is positive. In detail, assume that $w > 0$ (the negation of the last disjunct in the side condition of the fourth rule), which is equivalent to $w \approx |x| + 1$ where x is a fresh string variable and $|x|$ denotes an unknown non-negative quantity. It is sufficient to derive the formula obtained by replacing all occurrences of w by $|x| + 1$ in the disjunct $|t| + w \geq |t|$ to show that the start position of our term is out of bounds. After simplification, we obtain $|x| + 1 \geq 0$, which is trivial to derive.

The next two rules in Figure 4 apply if we can infer respectively that the start position of the substring comes strictly after a prefix t or that the end position of the substring comes strictly before a suffix t of the first argument string. In either case, t can be dropped.

Example 7. Let t be $\text{substr}(\text{con}(x_1, \text{replace}(x_2, x_3, x_4)), 0, w)$, where w is $|x_1| - |x_2|$. We have that $t \rightarrow \text{substr}(x_1, 0, w)$, noting that $\vdash |x_1| \geq 0 + |x_1| - |x_2|$. In other words, only the first component x_1 of the string concatenation is relevant to the substring since its end point must occur before the end of x_1 .

The final rule for substr shows that a prefix of a substring can be pulled upwards if the start position is zero and we can infer that the substring is guaranteed to include at least a prefix string t . Finally, if we can infer that the last position of s in t starting from position v is at or beyond the end of t , then the indexof term can be rewritten as an if-then-else (ite) term that checks whether s is a suffix of t .

4 Containment-Based String Simplification

This section provides an overview of simplifications that are based on reasoning about the containment relationship between strings. We describe an inference system for

$$\begin{array}{c}
\frac{l_1 \text{ contains } l_2}{\vdash \text{con}(l_1, t) \ni l_2} \quad \frac{\vdash s \ni r}{\vdash \text{con}(t, s) \ni r} \quad \frac{\vdash t \ni^s r \quad \vdash s \ni^p q}{\vdash \text{con}(t, s) \ni \text{con}(r, q)} \quad \frac{}{\vdash t \ni \text{substr}(t, v, w)} \\
\frac{l_1 \text{ does not contain } l_2}{\vdash l_1 \not\ni \text{con}(l_2, t)} \quad \frac{\vdash r \not\ni t}{\vdash r \not\ni \text{con}(s, t)} \quad \frac{\vdash l_1 \setminus l_2 \not\ni t}{\vdash l_1 \not\ni \text{con}(l_2, t)} \\
\frac{l_2 \text{ is a prefix of } l_1}{\vdash \text{con}(l_1, t) \ni^p l_2} \quad \frac{\vdash s \ni^p r}{\vdash \text{con}(t, s) \ni^p \text{con}(t, r)} \quad \frac{}{\vdash t \ni^p t} \quad \frac{\vdash v \leq 0}{\vdash t \ni^p \text{substr}(t, v, w)} \\
\frac{l_2 \text{ is a suffix of } l_1}{\vdash \text{con}(t, l_1) \ni^s l_2} \quad \frac{\vdash s \ni^s r}{\vdash \text{con}(s, t) \ni^s \text{con}(r, t)} \quad \frac{}{\vdash t \ni^s t} \quad \frac{\vdash v + w \geq |t|}{\vdash t \ni^s \text{substr}(t, v, w)}
\end{array}$$

Fig. 5. Inferences for string containment \ni , is-prefix \ni^p and is-suffix \ni^s . We write l_1, l_2 to denote string literals, v, w to denote integer terms and t, s, r, q to denote string terms. Although not shown here, we also assume rules that reflect the transitivity of \ni , \ni^p and \ni^s .

deriving when one string is definitely contained or not contained in another. Following the notation from the last section, we write $\vdash t \ni s$ to denote the judgment of our inference system, denoting that string t contains string s in all models of T_S . Conversely, we write $\vdash t \not\ni s$ to denote string t does not contain string s . We write $\vdash t \ni^p s$ (resp., $\vdash t \ni^s s$) to denote the judgment indicating that s must be a prefix (resp., suffix) of t .

Rules for inferring judgments of these forms are given in Figure 5. Like our rules for arithmetic, these rules are solely based on the syntactic structure of terms, so inferences in this system can be computed statically. Both the assumptions and conclusions of the rules assume associativity of string concatenation with identity element ϵ , that is, $\text{con}(t, s)$ may refer to a term of the form $\text{con}(\text{con}(t_1, t_2), s) = \text{con}(t_1, t_2, s)$ or alternatively to $\text{con}(\epsilon, s) = s$. Most of the rules are straightforward. The inference system has special rules for substrings $\text{substr}(t, v, w)$, using arithmetic entailments from Section 3 to show prefix and suffix relationships with the base string t . For negative containment, the rules of the inference system together can show a (possibly non-constant) string cannot occur in a constant string by reasoning that its characters cannot appear in order in that string. We write $l_1 \setminus l_2$ to denote the empty string if l_1 does not contain l_2 , or the result of removing the smallest prefix of l_1 that contains l_2 from l_1 otherwise.

Example 8. Let t be `abca` and let s be `con(b, x, a, y, c)`. String s is not contained in t for any value of x, y . We derive $\vdash t \not\ni s$ using two applications of the rightmost rule for negative containment in Figure 5, noting `abca` \setminus `b` = `ca`, `ca` \setminus `a` = `c`, and `c` does not contain `c`. In other words, the containment does not hold since the characters `b`, `a` and `c` cannot be found in order in the constant `abca`.

4.1 Simplification Rules based on String Containment

Figure 6 gives rules for simplifying extended function terms based on the aforementioned judgments pertaining to string containment. First, equalities can be rewritten to false and applications of `contains` can be rewritten to a constant based on the appropriate judgment of our inference system. Applications of `indexof` can be simplified to -1 if it can be shown that the second argument does not appear in the suffix of the first

$t \approx s \rightarrow \perp$	\perp	if $\vdash t \not\approx s$
$\text{contains}(t, s) \rightarrow \perp$	\perp	if $\vdash t \not\approx s$
$\text{contains}(t, s) \rightarrow \top$	\top	if $\vdash t \approx s$
$\text{indexof}(t, s, v) \rightarrow -1$	-1	if $\vdash \text{substr}(t, v) \not\approx s$
$\text{indexof}(\text{con}(t, r), s, v) \rightarrow \text{indexof}(t, s, v)$	$\text{indexof}(t, s, v)$	if $\vdash \text{substr}(t, v) \approx s$
$\text{indexof}(\text{con}(t, r), s, v) \rightarrow \text{indexof}(r, s, v - t) + t $	$\text{indexof}(r, s, v - t) + t $	if $\vdash \text{substr}(\text{con}(t, r), v) \approx s$ and $\vdash v \geq t $
$\text{indexof}(t, s, v) \rightarrow v$	v	if $\vdash \text{substr}(t, v) \approx^p s$ and $\vdash v < t $
$\text{replace}(t, s, r) \rightarrow t$	t	if $\vdash t \not\approx s$
$\text{replace}(\text{con}(t, q), s, r) \rightarrow \text{con}(\text{replace}(t, s, r), q)$	$\text{con}(\text{replace}(t, s, r), q)$	if $\vdash t \approx s$
$\text{replace}(t, s, r) \rightarrow \text{con}(r, \text{substr}(t, s))$	$\text{con}(r, \text{substr}(t, s))$	if $\vdash t \approx^p s$

Fig. 6. Simplification rules based on string containment.

argument starting at the position given by the third argument. The next two rules reason about cases where the second argument s definitely occurs in the first argument starting from position v . In this case, if we additionally know that s occurs within (beyond) a prefix t of the first argument, then the suffix r (prefix t) can be dropped, where the start position and the return value of the result are modified accordingly. If we know s is a prefix of the first argument at position v , then the result is v if indeed v is in the bounds of t . Notice that the latter condition is necessary to handle the case where s is the empty string. The three rules for `replace` are analogous. First, the `replace` rewrites to the first argument if we know it does not contain the second argument s . If we know s is definitely contained in a prefix of the first argument, then we can pull the remainder of that string upwards. Finally, if we know s is a prefix of the first argument, then we can replace that prefix with r while concatenating the remainder. We use the term $\text{substr}(t, |s|)$ to denote the remainder after the replacement for the sake of brevity, although this term typically does not involve extended functions after simplification, e.g. $\text{replace}(\text{con}(x, y), x, z) \rightarrow \text{con}(z, y)$ noting that $(\text{substr}(\text{con}(x, y), |x|)) \downarrow = y$, or $\text{replace}(\text{ab}, a, x) \rightarrow \text{con}(x, \text{b})$ noting that $(\text{substr}(\text{ab}, |a|)) \downarrow = \text{b}$.

4.2 Simplifications based on Equivalence of String Containment

We further refine our approach based on inferring when one containment is *equivalent* to another one. For example, $\text{con}(a, x)$ is contained in $\text{con}(b, y)$ if and only if $\text{con}(a, x)$ is contained in y alone. We introduce simplifications for such equivalences by reasoning about the maximal overlap between two strings.

We adapt and extend the notation given in previous work [15]. Given string literals l_1 and l_2 , the *sufficient left overlap* of l_1 and l_2 , written $l_1 \sqcup_l l_2$, is the largest suffix of l_1 that is a prefix of l_2 or has l_2 as a prefix. For example, we have $\text{abc} \sqcup_l \text{cd} = \text{c}$, $\text{abc} \sqcup_l \text{b} = \text{bc}$, and $\text{abc} \sqcup_l \text{ba} = \epsilon$. We extend this definition to arbitrary strings s such that $l_1 \sqcup_l s$ is equivalent to $l_1 \sqcup_l l_2$ for the largest constant prefix l_2 of s , where notice that l_2 is the empty string if s does not have a constant prefix. For example, we have $\text{abc} \sqcup_l \text{con}(\text{cde}, y) = \text{c}$, $\text{abc} \sqcup_l \text{con}(\text{b}, y) = \text{bc}$, and $\text{abc} \sqcup_l \text{con}(\text{a}, y) = \text{abc}$. We define the dual operator *sufficient right overlap*, written $l_1 \sqcup_r l_2$, which is the largest prefix of l_1 that is a suffix of l_2 or has l_2 as a suffix, e.g. $\text{abc} \sqcup_r \text{b} = \text{ab}$, and extend this to

$$\begin{array}{lcl}
\text{contains}(\text{con}(t, l), s) & \rightarrow & \text{contains}(\text{con}(t, l \sqcup_r s), s) \\
\text{contains}(\text{con}(l, t), s) & \rightarrow & \text{contains}(\text{con}(l \sqcup_l s, t), s) \\
\text{indexof}(\text{con}(t, l), s, v) & \rightarrow & \text{indexof}(\text{con}(t, l \sqcup_r s), s, v) \\
\text{indexof}(\text{con}(l, t), s, v) & \rightarrow & \text{indexof}(\text{con}(l_2, t), s, v - |l_1|) \quad \text{if } l = l_1 \cdot l_2 \text{ and } l_2 = l \sqcup_l s \\
& & + |l_1| \quad \vdash \text{substr}(\text{con}(l, t), v) \ni s \\
\text{replace}(\text{con}(t, l), s, r) & \rightarrow & \text{con}(\text{replace}(\text{con}(t, l_1), s, r), l_2) \quad \text{if } l = l_1 \cdot l_2 \text{ and } l_1 = l \sqcup_r s \\
\text{replace}(\text{con}(l, t), s, r) & \rightarrow & \text{con}(l_1, \text{replace}(\text{con}(l_2, t), s, r)) \quad \text{if } l = l_1 \cdot l_2 \text{ and } l_2 = l \sqcup_l s
\end{array}$$

Fig. 7. Simplification rules based on equivalence of string containment. We write l, l_1, l_2 to denote string literals, v, w to denote integer terms and t, s to denote string terms.

arbitrary strings in an analogous way. The sufficient left (resp., right) overlap operator can be used to determine how much of a constant string prefix l_1 (resp., suffix) can be safely removed from a string without impacting whether it contains another string.

The rules in Figure 7 simplify extended terms by considering string overlaps. The first two rules drop parts of string literals from the suffix or prefix of their first arguments. The two rules for `indexof` are similar: a suffix of the first argument can be dropped if it does not contribute to whether it contains the second argument. A prefix of an `indexof` term can be dropped if it does not contribute to containment, but only in the case where we know the second argument is definitely contained in the first argument. This is to guard against the case where the entire `indexof` term returns -1 . The rules for `replace` are similar to those for `contains`, except that the suffix (resp., prefix) of the first argument is pulled upwards instead of being dropped.

5 Multiset-Based String Simplification

Next, we introduce simplifications based on reasoning about strings as multisets, i.e. collections of unordered characters. Such reasoning is sufficient for showing that equalities like $\text{con}(a, x) \approx \text{con}(x, b)$ are equivalent to \perp , since the left side of the equality contains exactly one more occurrence of character a than the right-hand side. Similar to arithmetic reasoning from Section 3, we use approximations when reasoning about strings as multisets. We define the *multiset abstraction* of t , written \mathcal{M}_t , as the multiset $\{t_1, \dots, t_n\}$ where t is equivalent to $\text{con}(t_1, \dots, t_n)$ and all constants in this set are characters. For example, $\mathcal{M}_{\text{con}(aba, x)} = \{a, a, b, x\}$. We define a rewrite system $\rightarrow_{\mathcal{O}}^{\mathcal{M}}$ over strings where a rewritten string over-approximates the original string in the following sense: if $t \rightarrow_{\mathcal{O}}^{\mathcal{M}} s$, then for all models of T_5 and any character c , the number of occurrences of c in the strings in \mathcal{M}_s is greater than or equal to the number of occurrences in the strings in \mathcal{M}_t .

Figure 8 lists the rules for the rewrite system $\rightarrow_{\mathcal{O}}^{\mathcal{M}}$ and the simplifications based on multiset reasoning. Given a predicate $\text{contains}(t, s)$, if over-approximating t with respect to the rules of $\rightarrow_{\mathcal{O}}^{\mathcal{M}}$ results in a string r , and it can be determined that s contains strictly more occurrences of some character c than r , then it cannot be the case that s is contained in t . To establish this, we check whether the multiset difference of \mathcal{M}_s and \mathcal{M}_r contains c , and conversely the difference of \mathcal{M}_r and \mathcal{M}_s contains only character constants which are distinct from c . In the second rule, if one side of an equality can be determined to contain *only* a character c , then one occurrence of that character can be dropped from both sides of the equality, since the relative position of that character does not

$$\begin{array}{l}
\text{contains}(t, s) \rightarrow \perp \\
\text{con}(t, c, s) \approx \text{con}(q, c, r) \rightarrow \text{con}(t, s) \approx \text{con}(q, r) \text{ if }
\end{array}
\begin{array}{l}
\text{if } t \rightarrow_{\mathcal{O}}^{\mathcal{M}} *r, \\
\mathcal{M}_s \setminus \mathcal{M}_r = \{c, s_1, \dots, s_n\} \text{ and} \\
\mathcal{M}_r \setminus \mathcal{M}_s = \{c_1, \dots, c_m\} \\
\mathcal{M}_{\text{con}(t, c, s)} \rightarrow_{\mathcal{O}}^{\mathcal{M}} *p \text{ and} \\
\mathcal{M}_p = \{c, \dots, c\}
\end{array}$$

$$\begin{array}{l}
\text{substr}(t, v, w) \rightarrow_{\mathcal{O}}^{\mathcal{M}} t \\
\text{where replace}(t, s, r) \rightarrow_{\mathcal{O}}^{\mathcal{M}} \text{con}(t, r) \\
\text{con}(t, s, r) \rightarrow_{\mathcal{O}}^{\mathcal{M}} \text{con}(t, q, r) \text{ if } s \rightarrow_{\mathcal{O}}^{\mathcal{M}} q
\end{array}$$

Fig. 8. Simplification rules based on multiset reasoning. We write c, c_1, \dots to denote characters, v, w to denote integer terms, and t, s, r, q, p to denote string terms.

matter. The three rules for $\rightarrow_{\mathcal{O}}^{\mathcal{M}}$ state that the multiset abstraction of a term of the form $\text{substr}(t, v, w)$ can be over-approximated as the entire string t ; a term $\text{replace}(t, s, r)$ can be over-approximated as a string having both t and r ; and over-approximation can be applied to the children of con terms.

Example 9. We have that $\text{con}(\text{aaa}, \text{substr}(x, y_1, y_2)) \approx \text{con}(x, \text{b}) \rightarrow \perp$ by noting that $\text{con}(\text{aaa}, \text{substr}(x, y_1, y_2)) \rightarrow_{\mathcal{O}}^{\mathcal{M}} * \text{con}(\text{aaa}, x)$, $\mathcal{M}_{\text{con}(\text{aaa}, x)} = \{\text{a}, \text{a}, \text{a}, x\}$ and $\mathcal{M}_{\text{con}(x, \text{b})} = \{\text{b}, x\}$. The difference of the latter with the former is $\{\text{b}\}$, and the former with the latter is $\{\text{a}, \text{a}, \text{a}\}$. Thus, the right side of the equality contains at least one more occurrence of b than the left side; hence, the equality is equivalent to false.

6 Implementation

We implemented the above simplification rules and others in the DPLL-based SMT solver CVC4, which implements a theory solver for a basic fragment of word equations with length, several other theory solvers, and reduction techniques for extended string functions as described in Section 2.1. Our simplification rules are run in a *preprocessing* pass as well as an *inprocessing* pass during solving. For the latter, we use a context-dependent simplification strategy that infers when an extended string constraint, e.g., $\neg \text{contains}(t, s)$, simplifies to \perp based on other assertions, e.g., $s \approx \epsilon$. Our simplification techniques do not affect the core procedure for the theory of strings, nor the compatibility of the string solver with other theories. In total, our implementation is about 3,500 lines of C++ code. We cache the results of the simplifications and the approximation-based arithmetic entailments to amortize their costs.

Additional Simplification Rules The simplification rules in this paper are a subset of the rules in the implementation. We omit other uncategorized rules for lack of space. Many of these apply to specific term patterns, such as cases where two nested applications of substr can be combined; cases where an application of replace can be eliminated by case splitting; and other cases like $\text{con}(t, t) \approx a \rightarrow \perp$. An example of such rules is $\text{contains}(\text{replace}(t, w_1, w_2), w_3) \rightarrow \text{contains}(t, w_3)$ if w_3 does not overlap with either w_1 or w_2 , because the replace does not change whether t contains w_3 or not. Another class of rules only applies to strings of length one because they cannot span multiple components of a concatenations, e.g. $\text{contains}(\text{con}(t, s), c) \rightarrow \text{contains}(t, c) \vee \text{contains}(s, c)$

where c is a character. Finally, there are rewrites that benefit from multiple techniques presented in this paper. For example, we have a rewrite that splits string equations into multiple smaller equations if it can determine that prefixes must have the same length: $\text{con}(a, t, s) \approx \text{con}(t, b, r) \rightarrow \text{con}(a, t) \approx \text{con}(t, b) \wedge s \approx r \rightarrow \perp$.

Validating Simplification Rules The correctness of our simplification techniques is critical to the soundness of the overall solver. Due to the sophistication and breadth of those techniques, it is challenging to formally verify our implementation. As a pragmatic alternative, we periodically test our implementation using a testing infrastructure we developed for this purpose. We found this to be critical in our development process. Our testing infrastructure allows the developer to specify a context-free grammar in the syntax-guided synthesis format [2]. We generate all terms t in this grammar up to a fixed size and test the equivalence of t and its simplified form $t \downarrow$ on a set of randomly generated points. The most recent run of this system on two grammars (one for extended string terms and another for string predicates) up to a term size of three, validated 319,867 simplifications of string terms and 188,428 simplifications of string predicates on 1,000 sample points. This run took 924s for string terms and 971s for the string predicates using the same hardware as in Section 7.

7 Evaluation

We evaluate the impact of each simplification technique as implemented in CVC4 on three benchmark sets that use extended string operators: CMU, a dataset obtained from symbolic execution of Python code [15]; TERMEQ, a benchmark set consisting of the verification of term equivalences over strings [14]; and SLOG, a benchmark set extracted from vulnerability testing of web applications [22]. The SLOG set uses the replace function extensively but does not contain other extended functions. We also evaluate the impact on APLAS, a set of handcrafted benchmarks involving looping word equations [10] (string equalities whose left and right sides have variables in common).

We compare CVC4 with Z3 commit 9cb1a0f [8],⁶ a state-of-the-art string solver. Additionally, we compare against OSTRICH on the SLOG benchmarks but not other sets because it does not support some functions such as contains and indexof. We omit a comparison with Z3STR3 4.8.4 because we found multiple issues in its latest release including wrong answers, which we have reported to the authors. We also omit a comparison with S3# due to differing semantics. We compare four configurations of CVC4: **all**, which enables all optimizations; **-arith**, which disables arithmetic-based simplification techniques (discussed in Section 3); **-contain**, which disables containment-based simplification techniques (discussed in Section 4); and **-msets**, which disables multiset-based simplification techniques (discussed in Section 5). Additionally, to test the applicability of our techniques to other solvers, we test the effect of our simplifications on Z3 by using CVC4 to generate simplified benchmarks and then running Z3 on those benchmarks. We generate a set of simplified benchmarks that are simplified with CVC4 with ($Z3_f$) and without ($Z3_b$) the simplification techniques presented in this paper.

We ran all benchmarks on a cluster equipped with Intel E5-2637 v4 CPUs running Ubuntu 16.04 and dedicated one core, 8GB RAM, and 600 seconds for each job. Table 1

⁶ 9cb1a0f is newer than the current release 4.8.4 and includes several fixes for critical issues.

Set		all	-arith	-contain	-msets	z3	z3 _b	z3 _f	OSTRICH	R%	
CMU	sat	5703	5535	5703	5703	2343	3923	3943	■	32%	
	unsat	65	29	65	65	50	58	61			
	×	154	358	154	154	3529	1941	1918			
TERMEQ	sat	10	10	10	10	4	5	5	■	68%	
	unsat	51	37	28	51	35	40	60			
	×	19	33	42	19	41	35	15			
SLOG	sat	1302	1302	1302	1302	1133	1225	1225	■	27%	
	unsat	2082	2082	2082	2082	2080	2080	2080			1304
	×	7	7	7	7	178	86	86			5
APLAS	sat	135	135	135	135	9	51	46	■	n/a	
	unsat	292	292	171	171	94	129	292			
	×	159	159	280	280	483	406	248			
Total	sat	7150	6982	7150	7150	3489	5204	5219	■	1304	
	unsat	2490	2440	2346	2369	2259	2307	2493			2082
	×	339	557	483	460	4231	2468	2267			5

Table 1. Number of solved problems per benchmark set. Best results are in **bold**. Gray cells indicate benchmark sets not supported by a solver. “R%” indicates the reduction of extended string functions during preprocessing. All benchmarks ran with a timeout of 600s.

summarizes the number of solved instances for each configuration and the baseline solvers grouped by benchmark sets. We remark that the average reduction of extended string functions (with all simplification techniques enabled) shown in column “R%” is significant on all benchmark sets. The scatter plots in Figure 9 detail the effects of disabling each family of simplifications. They distinguish between satisfiable and unsatisfiable instances. To emphasize non-trivial benchmarks, we omit the benchmarks that are solved in less than a second by all solvers.

The arithmetic-based simplification techniques have the most significant performance impact on the symbolic execution benchmarks CMU. The number of solved benchmarks is significantly lower when disabling those techniques. The scatter plot shows that for longer running satisfiable queries there is a large portion of the benchmarks that are solved up to an order of magnitude faster with the simplifications. These improvements in runtime on the CMU set are particularly compelling because they come from a symbolic execution application, which involves a large number of queries with a short timeout. The improvements are more pronounced for unsatisfiable benchmarks, where our results show that simplifications often give the solver the ability to derive a refutation in a matter of seconds, something that is infeasible with configurations without these techniques. The APLAS set contains no extended string operators and hence our arithmetic-based simplification techniques have little impact on this set.

In contrast, both containment and multiset-based rewrites have a high impact on the APLAS set, as **-contain** and **-msets** both solve 121 fewer benchmarks. Additionally, **-contain** has a high impact on the TERMEQ set, where the simplifications enable the best configuration to solve 61 out of 80 benchmarks. Since these techniques apply most

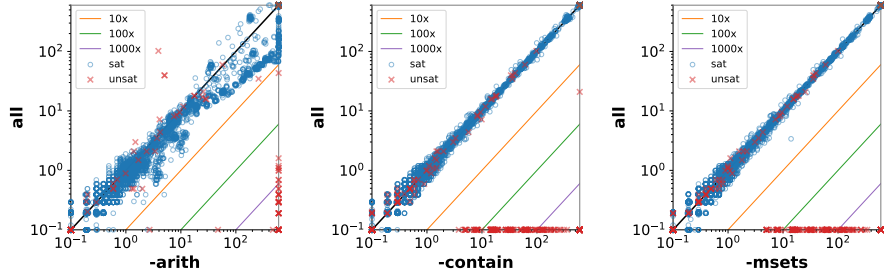


Fig. 9. Scatter plots showing the impact of disabling simplification techniques in CVC4 on both satisfiable and unsatisfiable benchmarks. All benchmarks ran with a timeout of 600s.

frequently to looping word equations, they are less important for the CMU set, which does not have such equations. The containment-based and multiset-based techniques primarily help on unsatisfiable benchmarks, as shown in the scatter plots. On TERMEQ benchmarks, it tends to be easier to find counterexamples, i.e. to solve the satisfiable ones, so there is more to gain on unsatisfiable benchmarks.

On SLOG, OSTRICH solves two more instances than CVC4 but CVC4 is over 50 times faster on commonly solved instances while supporting a richer set of string operators. On all benchmark sets, CVC4 solves at least as many benchmarks as Z3 and CVC4 has $12\times$ fewer timeouts than Z3. On the simplified benchmarks, Z3 performs significantly better. On the CMU and the APLAS benchmarks, Z3_b outperforms Z3 by a large margin. Additionally simplifying the benchmarks with the techniques presented in this paper improves performance further on most benchmark sets and allows Z3_f to solve the most unsatisfiable benchmarks overall. These results indicate that Z3 could benefit from additional simplifications, and they underscore the importance of curating and publishing simplification techniques in order to improve the state-of-the-art.

8 Conclusion

We have presented a set of aggressive simplification techniques for reasoning about extended string constraints. Our results suggest that such techniques are key to advancing the state of the art in SMT string solving. Arithmetic-based simplifications lead to significant speedups in benchmarks from a symbolic execution application, while containment and multiset-based simplifications improve the performance on problems consisting of difficult term equivalences and looping word equations. Our approach is not limited to CVC4 and can be adapted to other solvers.

Given the encouraging results for each of the simplification techniques in our evaluation, we plan to extend them to other types of abstraction and make them context-aware. The latter extension involves taking into account other assertions when checking whether a side condition of a rule is fulfilled.

Acknowledgements This work was partially supported by the National Science Foundation under award 1656926, the Defense Advanced Research Projects Agency under award FA8650-18-2-7854, and Amazon Web Services.

References

- [1] P. A. Abdulla, M. F. Atig, Y. Chen, B. P. Diep, L. Holík, A. Rezine, and P. Rümmer. Trau: SMT solver for string constraints. In N. Bjørner and A. Gurfinkel, editors, *2018 Formal Methods in Computer Aided Design, FMCAD 2018, Austin, TX, USA, October 30 - November 2, 2018*, pages 1–5. IEEE, 2018.
- [2] R. Alur, R. Bodík, E. Dallal, D. Fisman, P. Garg, G. Juniwal, H. Kress-Gazit, P. Madhusudan, M. M. K. Martin, M. Raghothaman, S. Saha, S. A. Seshia, R. Singh, A. Solar-Lezama, E. Torlak, and A. Udupa. Syntax-guided synthesis. In M. Irlbeck, D. A. Peled, and A. Pretschner, editors, *Dependable Software Systems Engineering*, volume 40 of *NATO Science for Peace and Security Series, D: Information and Communication Security*, pages 1–25. IOS Press, 2015.
- [3] C. Barrett, C. L. Conway, M. Deters, L. Hadarean, D. Jovanovic, T. King, A. Reynolds, and C. Tinelli. CVC4. In G. Gopalakrishnan and S. Qadeer, editors, *Computer Aided Verification - 23rd International Conference, CAV 2011, Snowbird, UT, USA, July 14-20, 2011. Proceedings*, volume 6806 of *Lecture Notes in Computer Science*, pages 171–177. Springer, 2011.
- [4] M. Berzish, V. Ganesh, and Y. Zheng. Z3str3: A string solver with theory-aware heuristics. In D. Stewart and G. Weissenbacher, editors, *2017 Formal Methods in Computer Aided Design, FMCAD 2017, Vienna, Austria, October 2-6, 2017*, pages 55–59. IEEE, 2017.
- [5] N. Bjørner, N. Tillmann, and A. Voronkov. Path feasibility analysis for string-manipulating programs. In S. Kowalewski and A. Philippou, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 15th International Conference, TACAS 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009. Proceedings*, volume 5505 of *Lecture Notes in Computer Science*, pages 307–321. Springer, 2009.
- [6] S. Chaudhuri and A. Farzan, editors. *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part I*, volume 9779 of *Lecture Notes in Computer Science*. Springer, 2016.
- [7] T. Chen, M. Hague, A. W. Lin, P. Rümmer, and Z. Wu. Decision procedures for path feasibility of string-manipulating programs with complex operations. *PACMPL*, 3(POPL):49:1–49:30, 2019.
- [8] L. M. de Moura and N. Bjørner. Z3: an efficient SMT solver. In C. R. Ramakrishnan and J. Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer, 2008.
- [9] A. Kiezun, V. Ganesh, S. Artzi, P. J. Guo, P. Hooimeijer, and M. D. Ernst. HAMPI: A solver for word equations over strings, regular expressions, and context-free grammars. *ACM Trans. Softw. Eng. Methodol.*, 21(4):25:1–25:28, 2012.
- [10] Q. L. Le and M. He. A decision procedure for string logic with quadratic equations, regular expressions and length constraints. In S. Ryu, editor, *Programming Languages and Systems - 16th Asian Symposium, APLAS 2018, Wellington, New Zealand, December 2-6, 2018, Proceedings*, volume 11275 of *Lecture Notes in Computer Science*, pages 350–372. Springer, 2018.
- [11] G. Li and I. Ghosh. PASS: string solving with parameterized array and interval automaton. In V. Bertacco and A. Legay, editors, *Hardware and Software: Verification and Testing - 9th International Haifa Verification Conference, HVC 2013, Haifa, Israel, November 5-7, 2013, Proceedings*, volume 8244 of *Lecture Notes in Computer Science*, pages 15–31. Springer, 2013.

- [12] 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 A. Biere and R. Bloem, editors, *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*, volume 8559 of *Lecture Notes in Computer Science*, pages 646–662. Springer, 2014.
- [13] R. Majumdar and V. Kuncak, editors. *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part II*, volume 10427 of *Lecture Notes in Computer Science*. Springer, 2017.
- [14] A. Reynolds, H. Barbosa, A. Niemetz, A. Nötzli, M. Preiner, C. Barrett, and C. Tinelli. Rewrites for SMT solvers using syntax-guided enumeration. *SMT*, 2018.
- [15] A. Reynolds, M. Woo, C. Barrett, D. Brumley, T. Liang, and C. Tinelli. Scaling up DPLL(T) string solvers using context-dependent simplification. In Majumdar and Kuncak [13], pages 453–474.
- [16] P. Saxena, D. Akhawe, S. Hanna, F. Mao, S. McCamant, and D. Song. A symbolic execution framework for javascript. In *31st IEEE Symposium on Security and Privacy, S&P 2010, 16-19 May 2010, Berkeley/Oakland, California, USA*, pages 513–528. IEEE Computer Society, 2010.
- [17] C. Tinelli, C. Barrett, and P. Fontaine. Unicode Strings (Draft 1.0). <http://smtlib.cs.uiowa.edu/theories-UnicodeStrings.shtml>, 2018.
- [18] M. Trinh, D. Chu, and J. Jaffar. S3: A symbolic string solver for vulnerability detection in web applications. In G. Ahn, M. Yung, and N. Li, editors, *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, Scottsdale, AZ, USA, November 3-7, 2014*, pages 1232–1243. ACM, 2014.
- [19] M. Trinh, D. Chu, and J. Jaffar. Progressive reasoning over recursively-defined strings. In Chaudhuri and Farzan [6], pages 218–240.
- [20] M. Trinh, D. Chu, and J. Jaffar. Model counting for recursively-defined strings. In Majumdar and Kuncak [13], pages 399–418.
- [21] M. Veanes, N. Tillmann, and J. de Halleux. Qex: Symbolic SQL query explorer. In E. M. Clarke and A. Voronkov, editors, *Logic for Programming, Artificial Intelligence, and Reasoning - 16th International Conference, LPAR-16, Dakar, Senegal, April 25-May 1, 2010, Revised Selected Papers*, volume 6355 of *Lecture Notes in Computer Science*, pages 425–446. Springer, 2010.
- [22] H. Wang, T. Tsai, C. Lin, F. Yu, and J. R. Jiang. String analysis via automata manipulation with logic circuit representation. In Chaudhuri and Farzan [6], pages 241–260.
- [23] F. Yu, M. Alkhalaf, and T. Bultan. Stranger: An automata-based string analysis tool for PHP. In J. Esparza and R. Majumdar, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 16th International Conference, TACAS 2010, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2010, Paphos, Cyprus, March 20-28, 2010. Proceedings*, volume 6015 of *Lecture Notes in Computer Science*, pages 154–157. Springer, 2010.
- [24] Y. Zheng, V. Ganesh, S. Subramanian, O. Tripp, M. Berzish, J. Dolby, and X. Zhang. Z3str2: an efficient solver for strings, regular expressions, and length constraints. *Formal Methods in System Design*, 50(2-3):249–288, 2017.
- [25] Y. Zheng, X. Zhang, and V. Ganesh. Z3-str: a z3-based string solver for web application analysis. In B. Meyer, L. Baresi, and M. Mezini, editors, *Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, ESEC/FSE’13, Saint Petersburg, Russian Federation, August 18-26, 2013*, pages 114–124. ACM, 2013.