

An Evaluation of Automata Algorithms for String Analysis

Pieter Hooimeijer^{1,*} and Margus Veanes²

¹ University of Virginia
pieter@cs.virginia.edu

² Microsoft Research
margus@microsoft.com

Abstract. There has been significant recent interest in automated reasoning techniques, in particular constraint solvers, for string variables. These techniques support a wide variety of clients, ranging from static analysis to automated testing. The majority of string constraint solvers rely on finite automata to support regular expression constraints. For these approaches, performance depends critically on fast automata operations such as intersection, complementation, and determinization. Existing work in this area has not yet provided conclusive results as to which core algorithms and data structures work best in practice.

In this paper, we study a comprehensive set of algorithms and data structures for performing fast automata operations. Our goal is to provide an apples-to-apples comparison between techniques that are used in current tools. To achieve this, we re-implemented a number of existing techniques. We use an established set of regular expressions benchmarks as an indicative workload. We also include several techniques that, to the best of our knowledge, have not yet been used for string constraint solving. Our results show that there is a substantial performance difference across techniques, which has implications for future tool design.

1 Introduction

There has been significant recent interest in decision procedures for string constraints. Traditionally, many analyses and automated testing approaches relied on their own built-in models to reason about string values [5,8,9,21,23,28]. Recent work on string analysis has focused on providing external decision procedures that reason about strings [2,11,12,15,24,26,27]. The separation of string decision procedures from their client analyses is desirable because it allows for the independent development of more scalable algorithms.

Existing string decision procedures vary significantly in terms of feature sets. Although most tools support regular expression constraints in one form or another, there is no strong consensus on which features are necessary. Many programming languages provide a variety of string manipulating functions. For example, the user manual for PHP lists 111 distinct string manipulation functions [20], and the `System.Text` namespace in the .NET class library contains

* This work was done during an internship at Microsoft Research.

more than a dozen classes that deal with Unicode encodings [19]. Constraint solving tools must strike a balance between expressive utility and other concerns like scalability and generality.

In this paper, we focus on a subclass of string decision procedures that support regular expression constraints and use finite automata as their underlying representation. This includes the DPRLE [11], JSA [5], and Rex [26,27] tools, as well as a prototype from our recent work [12]. These tools are implemented in different languages, they parse different subsets of common regular expression idioms, they differ in how they use automata internally, and the automata data structures themselves are different. However, each approach relies crucially on the efficiency of basic automaton operations like intersection and determinization/complementation. Existing work provides reasonable information on the relative performance of the tools as a whole, but does not give any insight into the relative efficiency of the individual data structures.

Our goal is to provide an apples-to-apples comparison between the core algorithms that underlie these solvers. To achieve this, we have performed a faithful re-implementation in C# of several often-used automaton data structures. We also include several data structures that, to the best of our knowledge, have not yet featured in existing work on string decision procedures. We conduct performance experiments on an established set of string constraint benchmarks. We use both ASCII and UTF-16 encodings where possible, to investigate the impact of alphabet size on the various approaches. Our testing harness includes a relatively full-featured regular expression parser that is based on the .NET framework's built-in parser. By fixing factors like the implementation language and the front-end parser, and by including a relatively large set of regular expression features, we aim to provide practical insight into which data structures are advisable for future string decision procedure work.

This paper makes the following contributions:

- To the best of our knowledge, we provide the first performance comparison of automata data structures (and associated algorithms) designed for string decision procedures. We pay special attention to isolating the core operations of interest.
- We present several approaches that have not yet been used for string decision procedures. We demonstrate that these approaches frequently out-perform existing approaches.

The rest of this paper is structured as follows. Section 2 gives formal definitions of the automata constructs of interest. Section 3 presents the design and implementation of the data structures under consideration. In Section 4, we provide experimental results using those data structures. We discuss closely related work in Section 5, and we conclude in Section 6.

2 Preliminaries

We assume familiarity with classical automata theory [13]. We work with a representation of automata where several transitions from a source state to a target

state may be combined into a single *symbolic move*. Symbolic moves have labels that denote *sets of characters* rather than individual characters. This representation naturally separates the structure of automata graph from the representation used for the character sets that annotate the edges. In this context, classical finite automata can be viewed as a special case in which each label is a singleton set. The following definition builds directly on the standard definition of finite automata.

Definition 1. A *Symbolic Finite Automaton* or *SFA* A is a tuple $(Q, \Sigma, \Delta, q^0, F)$, where Q is a finite set of *states*, Σ is the *input alphabet*, $q^0 \in Q$ is the *initial state*, $F \subseteq Q$ is the set of *final states* and $\Delta : Q \times 2^\Sigma \times Q$ is the *move relation*.

We indicate a component of an SFA A by using A as a subscript. The word *symbolic* stems from the intension that a move (q, ℓ, p) denotes the *set* of transitions $[(q, \ell, p)] \stackrel{\text{def}}{=} \{(q, a, p) \mid a \in \ell\}$. We let $[\Delta] \stackrel{\text{def}}{=} \cup\{[\rho] \mid \rho \in \Delta\}$. An SFA A denotes the finite automaton $[A] \stackrel{\text{def}}{=} (Q_A, \Sigma_A, [[\Delta_A], q_A^0, F_A)$. The *language* $L(A)$ accepted by A is the language $L([A])$ accepted by $[A]$.

An ϵ SFA A may in addition have moves where the label is ϵ , denoting the corresponding epsilon move in $[A]$. For $\rho = (p, \ell, q)$, let $Src(\rho) \stackrel{\text{def}}{=} p$, $Tgt(\rho) \stackrel{\text{def}}{=} q$, $Lbl(\rho) \stackrel{\text{def}}{=} \ell$. We use the following notation for the set of moves starting from a given state q in A : $\Delta_A(q) \stackrel{\text{def}}{=} \{\rho \mid \rho \in \Delta_A, Src(\rho) = q\}$ and furthermore will allow lifting functions to sets. For example, $\Delta_A(Q) \stackrel{\text{def}}{=} \cup\{\Delta_A(q) \mid q \in Q\}$. We write Δ_A^ϵ for the set of all epsilon moves in Δ_A and Δ_A^\neq for $\Delta_A \setminus \Delta_A^\epsilon$. An ϵ SFA A is *clean* if for all $\rho \in \Delta_A^\neq$, $Lbl(\rho) \neq \emptyset$; A is *normalized* if for all $q \in Q_A$, if $\rho_1, \rho_2 \in \Delta_A^\neq(q)$ and $Tgt(\rho_1) = Tgt(\rho_2)$ then $\rho_1 = \rho_2$; A is *total* if for all states $q \in Q_A$, $\Sigma_A = \bigcup_{\rho \in \Delta_A(q)} Lbl(\rho)$. Let \complement denote the complement operation over character sets. For making an SFA A total, do the following: for all $q \in Q_A$ such that $\ell = \complement(\bigcup_{\rho \in \Delta_A(q)} Lbl(\rho)) \neq \emptyset$, add $(q, \ell, sink)$ to Δ_A , and add the state $sink$ to Q_A and the move $(sink, \Sigma, sink)$ to Δ_A . Epsilon elimination from an ϵ SFA is straightforward and makes use of union [27]. Also, any SFA can easily be made clean and normalized. When we say SFA we assume that epsilon moves are not present.

We assume a translation from regex (extended regular expression) patterns to ϵ SFAs that follows closely the standard algorithm, see e.g., [13, Section 2.5]. A sample regex and corresponding ϵ SFA are illustrated in Figure 1.

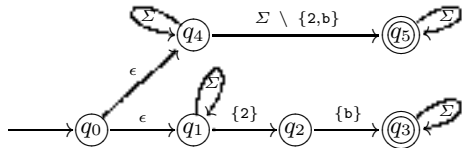


Fig. 1. ϵ SFA generated from regex $2b[\sim 2b]$

3 Automata Data Structures and Algorithms

In this section, we describe the automaton data structures and algorithms of interest. We assume a graph-based data structure for the automata; each transition edge is annotated with a data structure that represents the label of that

transition. Section 3.1 describes the data structures we use for those annotations. In Section 3.2, we discuss lazy and eager algorithms for two key operations on automata: language intersection (using the cross product construction) and language difference (using the subset construction). Later, in Section 4, we will evaluate experimentally how well each data structure/algorithm combination performs.

3.1 Representing Character Sets

We start by defining an interface for character set operations. This interface represents all the operations that the higher-level automata algorithms use to perform intersection and complementation. We then discuss several representations that can be used to implement this interface. A given representation affects both the performance of the algorithms discussed below; it also affects how effectively the algorithms can be combined with existing solvers.

Definition 2. Let I denote an ordered set of indices. A *minterm* of a finite nonempty sequence $(\ell_i)_{i \in I}$ of nonempty subsets of Σ is a sequence $(\ell'_i)_{i \in I}$ where each ℓ'_i is equal to either ℓ_i or $\mathbb{C}(\ell_i)$, and $\bigcap_{i \in I} \ell'_i \neq \emptyset$, where $\bigcap_{i \in I} \ell'_i$ is the set represented by the minterm.

Intuitively, a minterm (ℓ'_1, ℓ'_2) of a sequence (ℓ_1, ℓ_2) of two nonempty subsets $\ell_1, \ell_2 \subseteq \Sigma$ represents a minimal nonempty region $\ell'_1 \cap \ell'_2$ of a Venn diagram for ℓ_1 and ℓ_2 , see Figure 2. Note that the collection of all minterms of a given sequence of nonempty sets over Σ represents a partition of Σ .¹

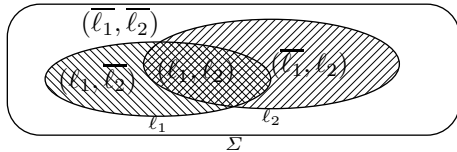


Fig. 2. Intuition behind minterms of (ℓ_1, ℓ_2)

Character set solver interface:

Boolean operations: union (\cup), intersection (\cap), and complement (\mathbb{C}).

Nonemptiness check: decide if a given set is nonempty.

Minterm generation: compute all minterms for a given sequence of character sets.

For any character set representation, we can use an algorithm similar to that of Figure 3 to perform the minterm computation in terms of repeated complementation and intersections. This algorithm is similar in spirit to the standard *Quine-McCluskey algorithm* for minimizing Boolean functions. We compute intersections combinatorially in rank order; once a given combination reaches the

¹ A partition of Σ is a collection of nonempty subsets of Σ such that every element of Σ is in exactly one of these subsets.

Input: Finite nonempty sequence of nonempty sets $(\ell_i)_{i \in I}$.
Rank propagation: Iterate over ranks $r \leq |I|$ and construct S_r :

$$S_0 = \{\emptyset\}$$

$$S_{r+1} = \{J \cup \{i\} \mid J \in S_r, i \in I \setminus J, \ell_i \cap \bigcap_{j \in J} \ell_j \neq \emptyset\}$$

Output: Set of all $J \in \bigcup_r S_r$ such that $(\bigcap_{i \in J} \ell_i) \cap (\bigcap_{i \in I \setminus J} \mathcal{C}(\ell_i)) \neq \emptyset$.

Fig. 3. Minterm generation algorithm

empty set at rank i , we know that derived combinations of rank $i' > i$ are necessarily empty (and thus uninteresting). In practice, we use this algorithm for every character set representation.

For each character-set representation we discuss how the above operations are supported, and indicate how the minterm generation is implemented.

Character sets as BDDs. In the general case, regexes, and string operations involving regexes, assume Unicode UTF-16 character encoding. This means that characters correspond to 16-bit bit-vectors. For the special case of ASCII range (resp. extended ASCII range) the number of bits is 7 (resp. 8). Independent of the number of bits, the general idea behind the BDD encoding is that each bit of a character corresponds to a variable of the BDD. The crucial point is to determine a good variable order for the BDD encoding. In order to do so, we considered typical uses of regexes, and in a separate evaluation determined an order which yielded BDDs of small size.

Based on the evaluation, a good choice of order turned out to be one where the highest bit has the lowest ordinal, in particular, for 16 bits, the MSB has ordinal 0 and the LSB has ordinal 15. This decision was, in part, based on the observation that typical regexes make extensive use of predefined character patterns called *character classes* where several bits are clustered and yield compact BDDs, such as $\backslash w$ (matches any word character), $\backslash W$ (matches any non-word character), $\backslash d$ (matches any decimal digit), $\mathfrak{p}\{\text{Lu}\}$ (matches any single character in the Unicode general category Lu of uppercase letters), etc. There are a total of 30 general Unicode categories and the character class $\backslash w$ is a union of 7 of those categories (namely categories 0 (Lu), 1 (Ll), 2 (Lt), 3 (Lm), 4 (Lo), 8 (Nd) and 18 (Pc)). In total, the class $\backslash w$ alone denotes a set of characters consisting of 323 nonoverlapping ranges, totaling 47057 characters. Let β_p denote the BDD of a character pattern p . The BDD $\beta_{\backslash w}$ has 1656 nodes. Note that Boolean operations with the above set interpretation correspond directly to Boolean operations over BDDs.

As an example, consider the pattern $[\backslash w - [\backslash da - d]]$, which denotes the set of all word characters that are neither decimal digits nor characters a through d , i.e., $\beta_{[\backslash w - [\backslash da - d]]} = \beta_{\backslash w} \cap \mathcal{C}(\beta_{\backslash d} \cup \beta_{a-d})$. We write $\llbracket \beta \rrbracket$ for the set of characters represented by β , thus for example $\llbracket \beta_1 \cap \beta_2 \rrbracket = \llbracket \beta_1 \rrbracket \cap \llbracket \beta_2 \rrbracket$. It is easy to write infeasible character patterns. For example, $[\backslash W - [\backslash D]]$ denotes an empty set since $\llbracket \beta_{\backslash d} \rrbracket \subset \llbracket \beta_{\backslash w} \rrbracket$ and

$$\llbracket \beta_{\backslash w} \cap \mathcal{C}(\beta_{\backslash d}) \rrbracket = \llbracket \mathcal{C}(\beta_{\backslash w}) \cap \mathcal{C}(\mathcal{C}(\beta_{\backslash d})) \rrbracket = \mathcal{C}(\llbracket \beta_{\backslash w} \rrbracket) \cap \llbracket \beta_{\backslash d} \rrbracket.$$

Performing nonemptiness checking for BDDs is trivial because the empty BDD β^\perp is unique. Similarly, the BDD of all characters β^\top is unique. Except for β^\perp and β^\top , two BDDs are not guaranteed to be identical by construction even though they are isomorphic when representing the same sets. However, checking isomorphism (equivalence) of BDDs is linear. Equivalence checking is used prior to calling the minterm algorithm in order to eliminate duplicate sets from the input sequence.

Character sets as bitvector predicates. A common alternative representation for character sets is to use interval arithmetic over bitvectors or integers. Here we assume that we are working in the context of a constraint solver that provides built-in support for bit vectors. We write BV^n for the sort of characters used by the solver, which is assumed to be a sort of n -bit vectors for a given fixed $n \in \{7, 8, 16\}$. Standard logic operations as well as standard arithmetic operations over BV^n , such as ‘ \leq ’, are assumed to be built-in and can be used to form predicates for expressing character ranges. Let $\varphi_p(\chi)$ denote a predicate (with a single *fixed* free variable $\chi:BV^n$) corresponding to the regex character pattern p and let $\llbracket \varphi_p \rrbracket$ denote the set of all characters a for which $\varphi_p(a)$ is true modulo the built-in theories. For example, consider BV^7 and the character pattern $\setminus w$, the predicate $\varphi_{\setminus w}$ is as follows where each disjunct corresponds to a Unicode category (the Unicode categories 2, 3 and 4 are empty for the ASCII character range):

$$('A' \leq \chi \wedge \chi \leq 'Z') \vee ('a' \leq \chi \wedge \chi \leq 'z') \vee ('0' \leq \chi \wedge \chi \leq '9') \vee \chi = '-'$$

where ‘ \dots ’ is the bitvector representation of a character. The Boolean operations are directly supported by the corresponding built-in logical operators. For example $\llbracket \varphi_{\setminus w - \setminus da - a} \rrbracket = \llbracket \varphi_{\setminus w} \wedge \neg(\varphi_{\setminus d} \vee \varphi_{a - a}) \rrbracket = \llbracket \varphi_{\setminus w} \rrbracket \cap \mathbb{C}(\llbracket \varphi_{\setminus d} \rrbracket \cup \llbracket \varphi_{a - a} \rrbracket)$

For the ASCII range (or extended ASCII range), the direct range representation has several advantages by being succinct and taking advantage of the built-in optimizations of the underlying solver. Based on our experiments, however, for full Unicode the representation produces predicates that do not scale. Instead, we use *if-then-else* (*Ite*) terms, supported in all state-of-the-art constraint solvers, to encode *Shannon expansions*, thus mimicking BDDs. An *Ite*-term is a term $Ite(\psi, t_1, t_2)$ that equals t_1 , if ψ is true; equals t_2 , otherwise. We explain the main idea. Given a BDD β the corresponding predicate $Ite[\beta]$ is constructed as follows where all shared subterms are constructed only once (and cached) and are thus maximally shared in the resulting term of the solver. Given a BDD β (other than β^\perp or β^\top) we write $BitIs0(\beta)$ for the predicate over $\chi:BV^n$, that is true if and only if the i ’th bit of χ is 0, where $i = n - ordinal(\beta) - 1$ (recall that $ordinal(\beta)$ is the reverse bit position).

$$Ite[\beta] \stackrel{\text{def}}{=} \begin{cases} true, & \text{if } \beta = \beta^\top; \\ false, & \text{if } \beta = \beta^\perp; \\ Ite(BitIs0(\beta), Ite[Left(\beta)], Ite[Right(\beta)]), & \text{otherwise.} \end{cases}$$

It follows from the definition that $\llbracket \beta \rrbracket = \llbracket Ite[\beta] \rrbracket$.

A nonemptiness check of $\llbracket \varphi \rrbracket$ for a character predicate φ is conducted by checking whether the formula φ is *satisfiable* by using the constraint solver. For minterm generation we use an incremental constraint solving technique that is sometimes called *cube-formula* solving [26]. The core idea behind the technique is as follows:

Given a finite sequence of character formulas $\varphi = (\varphi_i)_{i < m}$ and distinct Boolean variables $\mathbf{b} = (b_i)_{i < m}$ define $Cube(\varphi, \mathbf{b}) \stackrel{\text{def}}{=} \bigwedge_{i < m} \varphi_i \Leftrightarrow b_i$. A *solution* of $Cube(\varphi, \mathbf{b})$ is an assignment M of truth values to \mathbf{b} such that $M \models Cube(\varphi, \mathbf{b})$. Given a solution M of $Cube(\varphi, \mathbf{b})$, let $\varphi_M \stackrel{\text{def}}{=} \bigwedge_{i < m, M \models b_i} b_i \wedge \bigwedge_{i < m, M \models \neg b_i} \neg b_i$.

We use the following iterative model generation procedure to generate the set S of all solutions of $Cube(\varphi, \mathbf{b})$. Although the procedure is exponential in (n in) the *worst case* due to the inherent complexity of the problem, it seems to work well in practice [26]. The algorithm assumes that the solver is capable of model generation, rather than just checking existence of a model, which is the case for most state-of-the-art constraint solvers and in particular SMT solvers.

Minterm generation with cubes: Input is a nonempty sequence of satisfiable character predicates φ . Output is the set \mathbf{M} of solutions of $Cube(\varphi, \mathbf{b})$.

Initialize: Let $\mathbf{M}_0 = \emptyset$ and $\psi_0 = Cube(\varphi, \mathbf{b})$.

Iterate: If ψ_i has a model M let $\mathbf{M}_{i+1} = \mathbf{M}_i \cup \{M\}$ and $\psi_{i+1} = \psi_i \wedge \neg \varphi_M$, otherwise let $\mathbf{M} = \mathbf{M}_i$ and stop.

3.2 Primitive Automata Algorithms

The algorithms on automata that we are most interested in, and that we will focus on here, are *product* $A \times B$ and *difference* $A - B$, where difference provides a way of checking *subset* constraints between regular expressions as well as doing *complementation* \bar{B} such that $L(\bar{B}) = \Sigma^* \setminus L(B)$ and *determinization* of B , $Det(B)$, as special cases. The algorithms assume a representation of SFAs where move labels are symbolic and use a *character set solver* that provides the functionality discussed in Section 3.1. Both algorithms also have *lazy* versions that do not construct the full automaton: lazy difference is *subset checking* $L(A) \subseteq L(B)$ (with witness in $L(A) \setminus L(B)$ if $L(A) \not\subseteq L(B)$), and lazy product is *disjointness checking* $L(A) \cap L(B) = \emptyset$ (with witness in $L(A) \cap L(B)$ if $L(A) \cap L(B) \neq \emptyset$). We do not describe the lazy versions explicitly, since, apart from not constructing the resulting automaton, they are similar to the eager versions.

The product algorithm is shown in Figure 4. Note that the character set solver is used for performing intersection and nonemptiness check on labels. The difference algorithm, shown in Figure 5, makes, in addition, essential use of minterm generation during an implicit determinization of the subtracted automaton B in $A - B$. Intuitively, the algorithm is a combined product and complementation algorithm. The main advantage of the combination is *early pruning of search stack* S by keeping the resulting automaton *clean*. In our implementation, the difference algorithm uses a standard set implementation to represent state sets, denoted by $State\{\dots\}$ in Figure 5, to represent subsets of automata states. The difference algorithm uses three different kinds of set data structures, each with

Input: ϵ SFAs A and B over Σ .
Initialize: Eliminate epsilons from A and B .
 Let S be a stack, initially $S = (\langle q_A^0, q_B^0 \rangle)$.
 Let V be a hashtable, initially $V = \{\langle q_A^0, q_B^0 \rangle\}$.
 Let Δ be a hashtable, initially $\Delta = \emptyset$.
Search: While S is nonempty, repeat the following.
 Pop $\langle q_1, q_2 \rangle$ from S .
For each $\rho_1 \in \Delta_A(q_1)$ and $\rho_2 \in \Delta_B(q_2)$, let $\ell = Lbl(\rho_1) \cap Lbl(\rho_2)$, if $\ell \neq \emptyset$:
 Let $p_1 = Tgt(\rho_1)$, $p_2 = Tgt(\rho_2)$ and add $(\langle q_1, q_2 \rangle, \ell, \langle p_1, p_2 \rangle)$ to Δ .
 If $\langle p_1, p_2 \rangle \notin V$ then add $\langle p_1, p_2 \rangle$ to V and push $\langle p_1, p_2 \rangle$ to S .
Output: $A \times B = (\langle q_A^0, q_B^0 \rangle, \Sigma, V, \{q \in V \mid q \in F_A \times F_B\}, \Delta)$.

Fig. 4. Product algorithm for symbolic automata. Constructs $A \times B$ such that $L(A \times B) = L(A) \cap L(B)$.

Input: ϵ SFAs A and B over Σ .
Initialize: Eliminate epsilons from A and B and make B total.
 Let $q^0 = \langle q_A^0, State\{q_B^0\} \rangle$.
 Let S be a stack, initially $S = (q^0)$.
 Let V , F , and Δ be hashtables. Initially $V = \{q^0\}$ and $\Delta = \emptyset$.
 Initially $F = \{q^0\}$, if $q_A^0 \in F_A$ and $q_B^0 \notin F_B$; $F = \emptyset$, otherwise.
Search: While S is nonempty repeat the following.
 Pop $\langle p, Q \rangle$ from S . Let $\Delta_B(Q) = \{\rho_i\}_{i \in I}$ and let $\ell_i = Lbl(\rho_i)$ for $i \in I$.
Compute minterms: Compute \mathbf{M} as the set of all minterms of the sequence $(\ell_i)_{i \in I}$. Here minterms are given as subsets of I . For $J \in \mathbf{M}$ let $\ell_J = \bigcap_{i \in J} \ell_i \cap \bigcap_{i \in I \setminus J} \bar{\ell}_i$.
For each minterm $J \in \mathbf{M}$:
For each A move $\rho \in \Delta_A(p)$ such that $Lbl(\rho) \cap \ell_J \neq \emptyset$:
 Let $P = State\{Tgt(\rho_i) \mid i \in J\}$.
 Let $q = \langle Tgt(\rho), P \rangle$.
 Add $(\langle p, Q \rangle, Lbl(\rho) \cap \ell_J, q)$ to Δ .
 If $q \notin V$ then add q to V and push q to S .
 If $Tgt(\rho) \in F_A$ and $P \cap F_B = \emptyset$ then add q to F .
Output: $A - B = (q^0, \Sigma, V, F, \Delta)$.

Fig. 5. Difference algorithm. Constructs $A - B$ such that $L(A - B) = L(A) \setminus L(B)$.

different strengths: character sets, state sets, and hashtables for algorithm variables, where the character set representation depends on the character set solver.

Note that the difference algorithm reduces to complementation of B when $L(A) = \Sigma^*$, e.g., when $A = (q_A^0, \Sigma, \{q_A^0\}, \{q_A^0\}, \{(q_A^0, \Sigma, q_A^0)\})$. Then the condition $Lbl(\rho) \cap \ell_J \neq \emptyset$ above is trivially true, since $Lbl(\rho) \cap \ell_J = \ell_J$, for $\rho \in \Delta_A$. Consequently, there is no pruning of the search stack S then with respect to A , and the full complement \bar{B} is constructed. Moreover, \bar{B} is *deterministic*, since for any two distinct moves $(\langle p, Q \rangle, \ell_J, q)$ and $(\langle p, Q \rangle, \ell_{J'}, q')$ that are added to Δ above, $\ell_J \cap \ell_{J'} = \emptyset$ by definition of minterms. It follows that the difference algorithm also provides a *determinization* algorithm for B : construct \bar{B} as above and let $Det(B) = (q_B^0, \Sigma_{\bar{B}}, Q_{\bar{B}}, Q_{\bar{B}} \setminus F_{\bar{B}}, \Delta_{\bar{B}})$.

4 Experiments

In this section we first compare the performance of the product and difference algorithms with respect to different character set representations and eager vs. lazy versions of the algorithms. For the predicate representation of character

	A_0	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8	A_9
$ Q $	36	30	31	17	11	18	4573	104	2228	42
$ \Delta^\epsilon $	25	15	10	11	4	14	8852	99	3570	32
$ \Delta^f $	36	24	28	15	11	13	148	96	524	40

Fig. 6. Sizes of ϵ SFAs A_{i-1} constructed for regexes $\#i$ for $1 \leq i \leq 10$ in [27, Table I], where each regex is assumed to have an implicit start-anchor \sim and end-anchor $\$$

sets we use Z3 as the constraint solver. Integration with Z3 uses the .NET API that is publicly available [31]. All experiments were conducted on a laptop with an Intel dual core T7500 2.2GHz processor. We then compare the performance of our implementation with the `dk.brics.automaton` library that underlies the Java String Analyzer (JSA) [5].

Our experiment uses a set of ten benchmark regexes that have previously been used to evaluate string constraint solving tools [12,27]. The regexes are representative for various practical usages and originate from a case study [18]. The sizes of the automata constructed from the regexes A_0 are shown in Figure 6.

For each pair $(A_i, A_j)_{i,j < 10}$ we conducted the following experiments to compare different character set representations, algorithmic choices, and the effect of the size of the alphabet. For product we ignored the order of the arguments due to commutativity, thus there are 55 pairs in total, and for difference there are 100 pairs in total. Figure 7 shows the evaluation results for the difference algorithm and the product algorithm, respectively. The times exclude the construction time of ϵ SFAs from the regexes but *include* the epsilon elimination time to convert ϵ SFAs to SFAs that is a preprocessing step in the algorithms. The total time to construct the ϵ SFAs for the 10 regexes (including parsing) was 0.33 seconds (for both UTF16 as well as ASCII). For parsing the regexes we use the built-in regex parser in .NET.

The top columns correspond to the eager vs. lazy versions of the algorithms and the secondary columns correspond to whether the result is an empty automaton or a nonempty automaton. The rows correspond to the different algorithmic choices: *BDD-X* denotes the use of the BDD based solver where X is ASCII or UTF16; *Pred-ASCII* denotes the use of predicate encoding of character sets using Z3 predicates over BV^7 ; *Pred-UTF16* denotes the use of predicate encoding of character sets using Z3 *Ite*-term encodings of BDDs over BV^{16} ; *Range-ASCII* denotes the use of standard range representation for ASCII character sets as *hashsets of character pairs*; *Hash-ASCII* denotes the use of *hashsets of individual characters* as character sets.

We excluded the evaluation results for most of the *Range-UTF16* and *Hash-UTF16* cases, since too many instances either timed out or caused an out-of-memory exception. The only case that completed for the same number of instances was lazy difference with *Range-UTF16* character set representation, but is approximately 70 times slower in total experiment time compared to *BDD-UTF16*. In the difference experiment, 5 instances either timed out or caused an out-of-memory exception, involving the automata A_6 and A_8 in all cases. One of the hardest instances that terminated was checking $L(A_8) \subseteq L(A_9)$ that took

	Eager				Lazy				
	Empty		Nonempty		Empty		Nonempty		
	Difference	BDD-ASCII	(8).33/.007	(78)36/.008	(8).33/.006	(87)41/.001			
BDD-UTF16		(8).56/.02	(78)38/.02	(8).52/.012	(87)41/.01				
Pred-ASCII		(8)1.6/.06	(78)72/.08	(8)1.6/.06	(87)58/.02				
Pred-UTF16		(8)5.5/.12	(78)179/.24	(8)5.3/.12	(87)66/.11				
Range-ASCII		(8).9/.03	(78)67/.03	(8)1/.03	(87)44/.003				
Hash-ASCII		(8).9/.03	(78)67/.03	(8)1/.03	(87)45/.003				
brics-ASCII		(9)32/.015	(72)273/.016						
brics-UTF16		(9)39/.11	(72)341/.44						
Product		Eager				Lazy			
		Empty		Nonempty		Empty		Nonempty	
	BDD-ASCII	(29)9.7/.001	(26)90/.002	(29)9.7/.001	(26)19/.001				
	BDD-UTF16	(29)9.7/.001	(26)92/.003	(29)9.7/.001	(26)19/.001				
	Pred-ASCII	(29)10/.003	(26)142/.003	(29)10/.004	(26)25/.007				
	Pred-UTF16	(29)10/.01	(26)150/.05	(29)10/.01	(26)25/.03				
	Range-ASCII	(29)10/.002	(23)16/.005	(29)10/.002	(26)69/.001				
	Hash-ASCII	(29)10/.002	(23)16/.005	(29)10/.002	(26)70/.001				
	brics-ASCII	(25)66/.015	(19)65/.016						
	brics-UTF16	(25)66/.03	(19)70/.09						

Fig. 7. Experimental evaluation of the SFA algorithms. Each entry in the tables has the form $(n)t/m$ where n is the number of combinations solved, t is the total time it took to solve the n instances, m is the median. Time is in seconds. For product, the eager experiment constructs $A_i \times A_j$ for $0 \leq i \leq j < 10$; the lazy experiment tests emptiness of $L(A_i) \cap L(A_j)$ for $0 \leq i \leq j < 10$. For difference, the eager experiment constructs $A_i - A_j$ for $0 \leq i, j < 10$; the lazy experiment tests if $L(A_i) \subseteq L(A_j)$ ($L(A_i - A_j) = \emptyset$) for $0 \leq i, j < 10$. Timeout for each instance (i, j) was 20 min.

14 minutes with *Range-UTF16* and 1.3 seconds with *BDD-UTF16*, where the minterm generation algorithm is used heavily.

With the BDD representation, the size of the alphabet turns out not to play a major role in the algorithms, although the BDD sizes are typically considerably larger for UTF16 (hundreds or even thousands of nodes) compared to ASCII (tens of nodes). This is a useful indication that BDDs work surprisingly well as character sets and may even work for UTF32 encoding (which we have not tried). In several non-BDD cases, many individual experiments timed out. The BDD character set based algorithms are an *order of magnitude faster in the median for UTF16*. As expected, all the lazy versions of the algorithms are faster in the nonempty case. The eager constructions are however needed when converting the resulting SFAs to *symbolic language acceptors* for combination with other theories [27]. Symbolic language acceptors for SFAs are made use of for example in the parameterized unit testing framework Pex [23], where regexes are handled in branch conditions of .NET programs using SFAs with Pred-UTF16 representation of character sets, integration of the algorithms in this paper into Pex is currently ongoing work.²

² <http://www.rise4fun.com/>

Comparison with brics. For this comparison we first serialized the automata A_i , $i < 10$, in textual format. This is for two reasons: 1) to provide a fair comparison *at the level of automata algorithms*, i.e., to exclude differences in the automata constructed from the regexes; 2) to avoid semantic differences regarding the meaning of the notations used in the regexes. In the measurements we bootstrapped the automata by excluding the time to deserialize and to reconstruct the automata in *brics*, but included the time to add epsilon transitions (that add extra character interval transitions), as this is effectively equivalent to epsilon elimination using the *brics* library.

The Java code responsible for the lazy difference experiment is

```
... construct Ai, Aj, epsAi, epsAj ...
long t = System.currentTimeMillis();
boolean empty = (Ai.addEpsilons(epsAi)).subsetOf(Aj.addEpsilons(epsAj));
t = System.currentTimeMillis() - t;
```

using the `brics.Automaton` method `subsetOf`. The eager experiment used the `brics.Automaton` method `minus`. The code for the product experiment is similar (using `intersection`). For running each instance we assigned 1.5GB to the java runtime (which was the maximum possible). For all cases involving A_6 , the product and difference experiments timed out or caused an out-of-memory exception. For the product experiment all cases involving A_8 also timed out. The instance $L(A_8) \subseteq L(A_8)$ did not time out using *brics*, while caused an out-of-memory exception using our tool during minterm generation.

	Product		Difference	
	Empty	Nonempty	Empty	Nonempty
BDD-ASCII	.022/.001	.024/.001	.32/.005	.33/.001
BDD-UTF16	.022/.001	.024/.001	.54/.015	.78/.002
Pred-ASCII	.07/.003	.07/.004	1.6/.06	2/.01
Pred-UTF16	.2/.01	.2/.01	5.2/.12	8.2/.05
brics-ASCII	.14/.001	.2/.015	.1/.016	.8/.016
brics-UTF16	.6/.03	1.4/.05	3.6/.05	24.8/.08

Fig. 8. Lazy difference and product experiment times with A_6 and A_8 excluded. Each entry is *total/median* in seconds. For product, 21 instances are empty and 15 instances are nonempty. For difference, 8 instances are empty and 56 instances are nonempty.

For a better comparison of the easy cases we conducted a separate experiment where both A_6 and A_8 are excluded from the experiments and considered the lazy versions only. Thus, there are total of 36 product instances and 64 difference instances. The outcome of this experiment is shown in Figure 8. Let t_X^{op} denote the total time for experiment row X and operation op in Figure 8. Then:

$$\begin{aligned}
 t_{brics-ASCII}^{prod} / t_{BDD-ASCII}^{prod} &\approx 7, & t_{brics-UTF16}^{prod} / t_{BDD-UTF16}^{prod} &\approx 43, \\
 t_{brics-ASCII}^{diff} / t_{BDD-ASCII}^{diff} &\approx 1.4, & t_{brics-UTF16}^{diff} / t_{BDD-UTF16}^{diff} &\approx 21.
 \end{aligned}$$

5 Related Work

In this section we discuss related work, focusing on automated reasoning techniques; Figure 9 provides a brief overview. DPRLE [11] has two associated implementations: a fully verified core algorithm written in Gallina [6], and the OCaml implementation used for the experiments. The Gallina specification is designed primarily to help discharge proof obligations (using Coq’s built-in inversion lemmas for cons lists, for example). This specification is of limited interest, since it is not actually designed to be executed. The OCaml implementation relies heavily on OCaml’s built-in hashtable data structure. The transition function is modeled using two separate mappings: one for character-consuming transitions (`state -> state -> character set`) and one for ϵ -transitions (`state -> state set`). This representation reflects the fact that the main DPRLE algorithms rely heavily on tracking states across automaton operations [11].

The Rex tool provides a SFA representation that is similar to the formal definition given in Section 2. The core idea, based on work by van Noord and Gerdeman [22], is to represent automaton transitions using logical predicates. Rex works in the context of *symbolic language acceptors*, which are first-order encodings of symbolic automata into the theory of algebraic datatypes. The Rex implementation uses the Z3 satisfiability modulo theories (SMT) solver [7] to solve the produced constraints. The encoding process and solving process are completely disjoint. This means that many operations, like automaton intersection, can be offloaded to the underlying solver. For example, to find a string w that matches two regular expressions, r_1 and r_2 , Rex can simply assert the existence of w , generate symbolic automaton encodings for r_1 and r_2 , and assert that s is accepted by both those automata. We refer to this as a *Boolean encoding* of the string constraints.

The initial Rex work [27] explores various optimizations, such as minimizing the symbolic automata prior to encoding them. These optimizations make use of the underlying SMT solver to find combinations of edges that have internally-consistent move conditions. Subsequent work [26] explored the trade-off between the Boolean encoding and the use of automata-specific algorithms for language intersection and language difference. In this case, the automata-specific algorithms make repeated calls to Z3 to solve *cube formulae* to enumerate edges

Technique Name	Data structure	Corresponds to
Christensen, Moller, et al. (JSA) [5]	Char. ranges (Java, Unicode)	Eager Range
Hooimeijer and Weimer (DPRLE) [11]	Single-char. hashset (OCaml; ASCII)	Eager Hashset
Hooimeijer and Weimer [12]	Char. ranges (C++; ASCII)	Lazy Range
Minamide [21,28]	Single-char. functional set (OCaml; ASCII)	Eager Hashset
Veanes, Halleux, and Tillmann [27]	Unary pred. (C# and Z3[31]; Unicode)	Lazy predicate
Henriksen, Jensen, et al. [10,29]	BDDs C	Eager BDD

Fig. 9. Existing automata-based string analyses and, the data structures they use, and the closest-matching experimental prototype tested in this paper

with co-satisfiable constraints. In practice, this approach is not consistently faster than the Boolean encoding.

The Java String Analyzer (JSA) by Christensen *et al.* [5] uses finite automata internally to represent strings. The underlying `dk.brics.automaton` library is not presented as a constraint solving tool, but it is used in that way. The library represents automata using a pointer-based graph representation of node and edge objects; edges represent contiguous character ranges. The library, which is written in Java, includes a deterministic automaton representation that is specialized for matching given strings efficiently. This is not a common use case for constraint solving, since the goal there is to efficiently *find* string assignments rather than verify them. Several other approaches include a built-in model of string operations; Minamide [21] and Wassermann and Su [28] rely on an ad-hoc OCaml implementation that uses that language’s built-in applicative data structures.

In previous work, we present a lazy backtracking search algorithm for solving regular inclusion constraints [12]. The underlying automaton representation, written in C++, is based on the Boost Graph Library [25], which allows for a variety of adjacency list representations. We annotate transitions with integer ranges, similar to JSA. The implementation pays special attention to memory management, using fast pool allocation for small objects such as the abstract syntax tree for regular expressions. This implementation uses lazy algorithms intersection and determinization algorithms, allowing for significant performance benefits relative to DPRLE [11] and the original Rex [27] implementation.

Other Approaches. Bjørner *et al.* describe a set of string constraints based on common string library functions [2]. The approach is based on a direct encoding to a combination of theories provided by the Z3 SMT solver [7]. Supported operations include substring extraction (and, in general, the combination of integer constraints with string constraints), but the work does not provide an encoding for regular sets.

The Hampi tool [15] uses an eager bitvector encoding from regular expressions to bitvector logic. The encoding does not use quantification, and requires the enumeration of all positional shifts for every subexpression. The Kaluza tool extends this approach to systems of constraints with multiple variables and concatenation [24].

The MONA implementation [10] provides decision procedures for several varieties of monadic second-order logic (M2L). MONA relies on a highly-optimized BDD-based representation for automata. The implementation has seen extensive engineering effort; many of the optimizations are described in a separate paper [16]. We are not aware of any work that investigates the use of multi-terminal BDDs for nondeterministic finite automata directly. In this paper, we restrict our attention to a class of graph-based automata representations.

There is a wide range of application domains that, at some level, rely on implicit or explicit automata representations. Fang *et al.* [30] use an automaton-based method to model string constraints and length bounds for abstract interpretation. There is significant work on using automata for arithmetic constraint solving; Boigelot and Wolper provide an overview [4]. A subset of

work on explicit state space model checking (e.g., [3]) has focused on general algorithms to perform fast state space reduction. Similarly, there is theoretical work on word equations that involves automata [1,17], as well as algorithms work on the efficient construction of automata from regular expressions (e.g., [14]).

6 Conclusions

In this paper we presented a study of automata representations for efficient intersection and difference constructions. We conducted this study to evaluate which combination of data structures and algorithms is most effective in the context of string constraint solving. Existing work in this area has consistently included performance comparisons at the tool level, but has been largely inconclusive regarding which automata representations work best in general. To answer this question, we re-implemented a number of data structures in the same language (C#) using a front-end parser that correctly handles a large subset of .NET's regular expression language, and using a UTF-16 alphabet. Our experiments showed that, over the sample inputs under consideration, the BDD-based representation provides the best performance when paired with the lazy versions of the intersection and difference algorithms. This suggests that future string decision procedure work can achieve significant direct benefits by using this combination of data structure (BDD) and algorithm (lazy) as a starting point. To the best of our knowledge, no existing tools currently use this combination.

Acknowledgment. We thank Nikolaј Bjørner for discussions that lead to the BDD character set representation idea and for continuous support with Z3.

References

1. Bala, S.: Regular language matching and other decidable cases of the satisfiability problem for constraints between regular open terms. In: Diekert, V., Habib, M. (eds.) STACS 2004. LNCS, vol. 2996, pp. 596–607. Springer, Heidelberg (2004)
2. Bjørner, N., Tillmann, N., Voronkov, A.: Path feasibility analysis for string-manipulating programs. In: Kowalewski, S., Philippou, A. (eds.) TACAS 2009. LNCS, vol. 5505, pp. 307–321. Springer, Heidelberg (2009)
3. Blom, S., Orzan, S.: Distributed state space minimization. *J. Software Tools for Technology Transfer* 7(3), 280–291 (2005)
4. Boigelot, B., Wolper, P.: Representing arithmetic constraints with finite automata: An overview. In: Stuckey, P.J. (ed.) ICLP 2002. LNCS, vol. 2401, pp. 1–19. Springer, Heidelberg (2002)
5. Christensen, A.S., Møller, A., Schwartzbach, M.I.: Precise Analysis of String Expressions. In: Cousot, R. (ed.) SAS 2003. LNCS, vol. 2694, pp. 1–18. Springer, Heidelberg (2003)
6. Coquand, T., Huet, G.P.: The calculus of constructions. *Information and Computation* 76(2/3), 95–120 (1988)
7. de Moura, L., Bjørner, N.S.: Z3: An Efficient SMT Solver. In: Ramakrishnan, C.R., Rehof, J. (eds.) TACAS 2008. LNCS, vol. 4963, pp. 337–340. Springer, Heidelberg (2008)

8. Godefroid, P., Kiežun, A., Levin, M.Y.: Grammar-based whitebox fuzzing. In: PLDI 2008, Tucson, AZ, USA, June 9-11 (2008)
9. Godefroid, P., Klarlund, N., Sen, K.: DART: directed automated random testing. In: PLDI 2005, pp. 213–223 (2005)
10. Henriksen, J.G., Jensen, J., Jørgensen, M., Klarlund, N., Paige, B., Rauhe, T., Sandholm, A.: Mona: Monadic second-order logic in practice. In: Brinksma, E., Steffen, B., Cleaveland, W.R., Larsen, K.G., Margaria, T. (eds.) TACAS 1995. LNCS, vol. 1019. Springer, Heidelberg (1995)
11. Hooimeijer, P., Weimer, W.: A decision procedure for subset constraints over regular languages. In: PLDI, pp. 188–198 (2009)
12. Hooimeijer, P., Weimer, W.: Solving string constraints lazily. In: ASE 2010 (2010)
13. Hopcroft, J.E., Ullman, J.D.: Introduction to Automata Theory, Languages, and Computation. Addison-Wesley, Reading (1979)
14. Ilie, L., Yu, S.: Follow automata. *Information and Computation* 186(1), 140–162 (2003)
15. Kiežun, A., Ganesh, V., Guo, P.J., Hooimeijer, P., Ernst, M.D.: HAMPI: a solver for string constraints. In: ISSTA 2009, pp. 105–116. ACM, New York (2009)
16. Klarlund, N., Møller, A., Schwartzbach, M.I.: MONA implementation secrets. *International Journal of Foundations of Computer Science* 13(4), 571–586 (2002)
17. Kunc, M.: What do we know about language equations? In: Harju, T., Karhumäki, J., Lepistö, A. (eds.) DLT 2007. LNCS, vol. 4588, pp. 23–27. Springer, Heidelberg (2007)
18. Li, N., Xie, T., Tillmann, N., de Halleux, P., Schulte, W.: Reggae: Automated test generation for programs using complex regular expressions. In: ASE 2009 (2009)
19. MSDN Library. System.text namespace, <http://msdn.microsoft.com/en-us/library/system.text.aspx>
20. PHP Manual. Pcre; posix regex; strings, <http://php.net/manual/en/book.strings.php>
21. Minamide, Y.: Static approximation of dynamically generated web pages. In: WWW 2005, pp. 432–441 (2005)
22. Van Noord, G., Gerdemann, D.: Finite state transducers with predicates and identities. *Grammars* 4 (2001)
23. Pex, <http://research.microsoft.com/projects/pex>
24. Saxena, P., Akhawe, D., Hanna, S., Mao, F., McCamant, S., Song, D.: A symbolic execution framework for javascript. In: SP 2010. IEEE Computer Society, Los Alamitos (2010)
25. Siek, J.G., Lee, L.-Q., Lumsdaine, A.: The Boost Graph Library: User Guide and Reference Manual. Addison-Wesley Professional, Reading (2001)
26. Veanes, M., Bjørner, N., de Moura, L.: Symbolic automata constraint solving. In: Fermüller, C.G., Voronkov, A. (eds.) LPAR-17. LNCS, vol. 6397, pp. 640–654. Springer, Heidelberg (2010)
27. Veanes, M., de Halleux, P., Tillmann, N.: Rex: Symbolic Regular Expression Explorer. In: ICST 2010. IEEE, Los Alamitos (2010)
28. Wassermann, G., Su, Z.: Sound and precise analysis of web applications for injection vulnerabilities. In: PLDI 2007, pp. 32–41. ACM, New York (2007)
29. Yu, F., Alkhalaf, M., Bultan, T.: An automata-based string analysis tool for PHP. In: Esparza, J., Majumdar, R. (eds.) TACAS 2010. LNCS, vol. 6015, pp. 154–157. Springer, Heidelberg (2010)
30. Yu, F., Bultan, T., Ibarra, O.H.: Symbolic string verification: Combining string analysis and size analysis. In: Kowalewski, S., Philippou, A. (eds.) TACAS 2009. LNCS, vol. 5505, pp. 322–336. Springer, Heidelberg (2009)
31. Z3, <http://research.microsoft.com/projects/z3>