

Minimization of Symbolic Automata

Loris D'Antoni

University of Pennsylvania
lorisdan@cis.upenn.edu

Margus Veanes

Microsoft Research
margus@microsoft.com

Abstract

Symbolic Automata extend classical automata by using symbolic alphabets instead of finite ones. Most of the classical automata algorithms rely on the alphabet being finite, and generalizing them to the symbolic setting is not a trivial task. In this paper we study the problem of minimizing symbolic automata. We formally define and prove the basic properties of minimality in the symbolic setting, and lift classical minimization algorithms (Huffman-Moore's and Hopcroft's algorithms) to symbolic automata. While Hopcroft's algorithm is the fastest known algorithm for DFA minimization, we show how, in the presence of symbolic alphabets, it can incur an exponential blowup. To address this issue, we introduce a new algorithm that fully benefits from the symbolic representation of the alphabet and does not suffer from the exponential blowup. We provide comprehensive performance evaluation of all the algorithms over large benchmarks and against existing state-of-the-art implementations. The experiments show how the new symbolic algorithm is faster than previous implementations.

Categories and Subject Descriptors F.2.2 [Theory of Computation]: Automata over infinite objects, Regular languages

Keywords Minimization, Symbolic Automata

1. Introduction

Classical automata theory builds on two basic assumptions: there is a *finite state space*; and there is a *finite alphabet*. The topic of this paper is along the line of work challenging the second assumption. *Symbolic finite automata* (SFAs) are finite state automata in which the alphabet is given by a *Boolean algebra* that may have an infinite domain, and transitions are labeled with predicates over such algebra. Symbolic automata originated from the intent to support regular expressions in the context of static and dynamic program analysis [43]. Lately, they were also used for supporting regular expressions (modulo label theories) in the context of modern logical inference engines [9, 42].

Most classical automata algorithms view the size k of the alphabet as a constant and use specialized data structures that are optimized for this view [7]. Therefore, it is not clear if or how such algorithms would work when the alphabet is infinite. Understanding how operations over the finite alphabet lift to the symbolic setting

is a challenging task. Some classical automata constructions are extended to SFAs in [26]. For example, the *product* (intersection) $M_1 \times M_2$ of two symbolic automata M_1 and M_2 is computed by building product transitions of the form $\langle p_1, p_2 \rangle \xrightarrow{\varphi_1 \wedge \varphi_2} \langle q_1, q_2 \rangle$ from transitions $p_1 \xrightarrow{\varphi_1} q_1$, $p_2 \xrightarrow{\varphi_2} q_2$, in M_1, M_2 , where the guards φ_1 and φ_2 are composed using *conjunction* and pruned when *unsatisfiable*. The complexity of such constructions clearly depends on the complexity of checking satisfiability in the label theory. In this particular case, constructing the product has complexity $\mathcal{O}(f(\ell)m^2)$, where m is the number of transitions, $f(\ell)$ is the cost of checking satisfiability of predicates of size ℓ in the label theory, and ℓ is the size of the biggest predicate in M_1 and M_2 .

This paper focuses on the problem of minimizing automata over symbolic alphabets and, to the best of our knowledge, it is the first paper that investigates this problem. Minimizing *deterministic finite automata* (DFAs) is one of the fundamental concepts in automata theory. It occurs in numerous areas, such as programming languages, text processing, computational linguistics, graphics, etc.

Before looking for an algorithm for minimizing SFAs we first need to answer a fundamental question: what does it mean for an SFA to be minimal? Intuitively, any two distinct states p and q of a minimal SFA must be *distinguishable*, where two states p and q are distinguishable if there exists an input sequence s that starting from p leads to a final (non-final) state and starting from q leads to a non-final (final) state. This notion is similar to DFA minimality.

The original algorithms for minimizing DFAs were given by Huffman [29], Moore [34], and Hopcroft [27]. Since, all such algorithms use iterations over the finite alphabet, they do not immediately extend to the symbolic setting. In the following paragraphs we briefly describe how we extended the classical algorithms to SFAs, and how we designed a new minimization algorithm that fully takes advantage of the symbolic representation of the alphabet.

Our first algorithm is called Min_{SFA}^M and takes inspiration from a reformulation of Moore's (Huffman's) algorithm described in [28]. The key idea from the algorithm in [28] is the following:

if two states p and q are distinguishable, and there exists a character a , and transitions $\delta(a, p') = p$ and $\delta(a, q') = q$, then p' and q' are distinguishable.

This idea nicely translates to the symbolic setting as:

if two states p and q are distinguishable, and there exist transitions $p' \xrightarrow{\varphi} p$ and $q' \xrightarrow{\psi} q$ such that $\varphi \wedge \psi$ is satisfiable, then p' and q' are distinguishable.

Starting with the fact that final and non-final states are distinguishable, we can use a fixpoint computation for grouping states into groups of indistinguishable states. This procedure uses a number of iterations that is quadratic in the number of states. Moreover, each iteration is linear in the number of transitions.

Unfortunately, Min_{SFA}^M does not scale in the case of a performance critical application described in Section 7.1, where SFAs

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

POPL '14, January 22 - 24 2014, San Diego, CA, USA.
Copyright © 2014 ACM 978-1-4503-2544-8/14/01...\$15.00.
<http://dx.doi.org/10.1145/2535838.2535849>

	Min_{SFA}^M	Min_{SFA}^H	Min_{SFA}^N
Satisfiability checking, \wedge, \vee	✓	✓	✓
Predicate negation		✓	✓
Minterm generation		✓	

Table 1. Summary of operations needed over a given label theory in the respective minimization algorithms.

end up having thousands of states, and the complexity of Min_{SFA}^M was not acceptable. To this end, we studied a generalization of Hopcroft’s algorithm, called Min_{SFA}^H . Hopcroft’s algorithm is based on a technique called partition refinement of states, and, with worst case complexity $O(kn \log n)$ (n = number of states, and k = number of alphabet symbols), it is the most efficient algorithm for minimizing DFAs. The main idea behind Min_{SFA}^H is to subdivide all the labels in the SFA into non-overlapping predicates called minterms. Such minterms can then be seen as atomic characters in the sense of classical DFAs, allowing us to use the classical Hopcroft’s algorithm over the finite set of minterms. Even though Min_{SFA}^H performs well in the aforementioned application, it suffers from another problem: in the worst case, the number of minterms can be exponential in the number of edges of the SFA. Moreover, our experiments showed that this factor can indeed be observed when considering more complex label theories, such as the theory of pairs over $string \times int$ (§ 6.4).

This leads to our main algorithmic contribution in the paper. We designed Min_{SFA}^N , a new algorithm for SFA minimization which takes full advantage of the symbolic representation of the input alphabet. Min_{SFA}^N is inspired by the idea of refining the partition of the state space used by Hopcroft’s algorithm, but does not require the pre-computation of minterms as in Min_{SFA}^H . The key observation is that the set of relevant predicates can be computed locally, rather than using all the transitions of the SFA. While Min_{SFA}^N is similar to Min_{SFA}^H in terms of state complexity, it does not suffer from the exponential complexity related to the minterm computation. In fact, in all our experiments, both Min_{SFA}^H and Min_{SFA}^M are outperformed by Min_{SFA}^N .

Table 1 presents a summary of the minimization algorithms, illustrating their dependency on operations over the label theory.

We compared the performance of the three algorithms using: 1) the benchmark of randomly generated DFAs presented in [3], 2) the SFAs generated by common regular expressions taken from the web, 3) a set of SFAs aimed at showing the exponential minterm explosion, 4) a randomly generated set of SFAs over a complex alphabet theory, and 5) the SFAs generated during the transformation from Monadic Second Order logic to DFAs [38]. In experiments 2 and 3 we also compared the performance of our implementation of Min_{SFA}^H against the implementation of Hopcroft’s algorithm in [3] and in the `brics.automaton` [1] library, and we observed similar performance characteristics that validated our implementation. In the fifth experiment we compared our results against `Mona` [22], the state of the art tool for deciding Monadic Second Order logic.

Contributions. In summary, our contributions are:

- a formal study of the notion of minimality of SFAs (§2);
- two algorithms for minimizing SFAs based on classical DFA algorithms (§3 and §4);
- a completely new algorithm for minimizing SFAs together with a proof of its correctness (§5);
- a comprehensive evaluation of the algorithms using a variety of different benchmarks (§6); and
- a description of several concrete applications of such minimization algorithms (§7).

2. Effective Boolean algebras and SFAs

We first formally define the notion of effective Boolean algebra and symbolic finite automata. Next, we develop a theory which explains what it means for a symbolic finite automata to be minimal.

An *effective Boolean algebra* \mathcal{A} has components $(\mathcal{D}, \Psi, \llbracket _ \rrbracket, \perp, \top, \vee, \wedge, \neg)$. \mathcal{D} is an r.e. (recursively enumerable) set of *domain elements*. Ψ is an r.e. set of *predicates* closed under the Boolean connectives and $\perp, \top \in \Psi$. The *denotation function* $\llbracket _ \rrbracket : \Psi \rightarrow 2^{\mathcal{D}}$ is r.e. and is such that, $\llbracket \perp \rrbracket = \emptyset$, $\llbracket \top \rrbracket = \mathcal{D}$, for all $\varphi, \psi \in \Psi$, $\llbracket \varphi \vee \psi \rrbracket = \llbracket \varphi \rrbracket \cup \llbracket \psi \rrbracket$, $\llbracket \varphi \wedge \psi \rrbracket = \llbracket \varphi \rrbracket \cap \llbracket \psi \rrbracket$, and $\llbracket \neg \varphi \rrbracket = \mathcal{D} \setminus \llbracket \varphi \rrbracket$. For $\varphi \in \Psi$, we write $IsSat(\varphi)$ when $\llbracket \varphi \rrbracket \neq \emptyset$ and say that φ is *satisfiable*. \mathcal{A} is *decidable* if $IsSat$ is decidable.

The intuition is that such an algebra is represented programmatically as an API with corresponding methods implementing the Boolean operations and the denotation function. We are primarily going to use the following two effective Boolean algebras in the examples, but the techniques in the paper are fully generic.

2^{Bv^k} is the powerset algebra whose domain is the finite set Bv^k , for some $k > 0$, consisting of all nonnegative integers smaller than 2^k , or equivalently, all k -bit bit-vectors. A predicate is represented by a BDD of depth k .¹ The Boolean operations correspond directly to the BDD operations, \perp is the BDD representing the empty set. The denotation $\llbracket \beta \rrbracket$ of a BDD β is the set of all integers n such that a binary representation of n corresponds to a solution of β .

SMT^σ is the decision procedure for a theory over some sort σ , say integers, such as the theory of integer linear arithmetic. This algebra can be implemented through an interface to an SMT solver. Ψ contains in this case the set of all formulas $\varphi(x)$ in that theory with one fixed free integer variable x . For example, a formula $(x \bmod k) = 0$, say div_k , denotes the set of all numbers divisible by k . Then $div_2 \wedge div_3$ denotes the set of numbers divisible by six.

We can now define symbolic finite automata. Intuitively, a symbolic finite automaton is a finite automaton over a symbolic alphabet, where edge labels are replaced by predicates. In order to preserve the classical closure operations (intersection, complement, etc.), the predicates must form an effective Boolean algebra.

DEFINITION 1. A *symbolic finite automaton (SFA)* M is a tuple $(\mathcal{A}, Q, q^0, F, \Delta)$ where \mathcal{A} is an effective Boolean algebra, called the *alphabet*, Q is a finite set of *states*, $q^0 \in Q$ is the *initial state*, $F \subseteq Q$ is the set of *final states*, and $\Delta \subseteq Q \times \Psi_{\mathcal{A}} \times Q$ is a finite set of *moves* or *transitions*. \square

Elements of $\mathcal{Q}_{\mathcal{A}}$ are called *characters* and finite sequences of characters, elements of $\mathcal{Q}_{\mathcal{A}}^*$, are called *words*; ϵ denotes the empty word. A move $\rho = (p, \varphi, q) \in \Delta$ is also denoted by $p \xrightarrow{\varphi}_M q$ (or $p \xrightarrow{\varphi} q$ when M is clear), where p is the *source state*, denoted $Src(\rho)$, q is the *target state*, denoted $Tgt(\rho)$, and φ is the *guard* or *predicate* of the move, denoted $Grd(\rho)$. A move is *feasible* if its guard is satisfiable. Given a character $a \in \mathcal{Q}_{\mathcal{A}}$, an *a-move* of M is a move $p \xrightarrow{\varphi} q$ such that $a \in \llbracket \varphi \rrbracket$, also denoted $p \xrightarrow{a}_M q$ (or $p \xrightarrow{a} q$ when M is clear). In the following let $M = (\mathcal{A}, Q, q^0, F, \Delta)$ be an SFA.

DEFINITION 2. A word $w = a_1 a_2 \dots a_k \in \mathcal{Q}_{\mathcal{A}}^*$, is *accepted at state* p of M , denoted $w \in \mathcal{L}_p(M)$, if there exist $p_{i-1} \xrightarrow{a_i}_M p_i$ for $1 \leq i \leq k$, such that $p_0 = p$, and $p_k \in F$. The *language accepted by* M is $\mathcal{L}(M) \stackrel{\text{def}}{=} \mathcal{L}_{q^0}(M)$. \square

¹The variable order of the BDD is the reverse bit order of the binary representation of a number, in particular, the most significant bit has the lowest ordinal.

Given a state $q \in Q$, we use the following definitions for the set transitions from and to q :

$$\overrightarrow{\Delta}(q) \stackrel{\text{def}}{=} \{\rho \in \Delta \mid \text{Src}(\rho) = q\}, \quad \overleftarrow{\Delta}(q) \stackrel{\text{def}}{=} \{\rho \in \Delta \mid \text{Tgt}(\rho) = q\}.$$

The definitions are lifted to sets in the usual manner. The following terminology is used to characterize various key properties of M . A state p of M is called *partial* if there exists a character a for which there exist no a -move with source state p .

- M is *deterministic*: for all $p \xrightarrow{\varphi} q, p \xrightarrow{\varphi'} q' \in \Delta$, if $\text{IsSat}(\varphi \wedge \varphi')$ then $q = q'$.
- M is *complete*: there are no partial states.
- M is *clean*: for all $p \xrightarrow{\varphi} q \in \Delta$, p is reachable from q^0 and $\text{IsSat}(\varphi)$,
- M is *normalized*: for all $p, q \in Q$, there is at most one move from p to q .
- M is *minimal*: M is deterministic, complete, clean, normalized, and for all $p, q \in Q$, $p = q$ if and only if $\mathcal{L}_p(M) = \mathcal{L}_q(M)$.²

For the special case in which M is *deterministic and complete*, we denote the transition function using the function $\delta_M : \mathfrak{D}_A \times Q \rightarrow Q$, such that for all $a \in \mathfrak{D}_A$ and $p \in Q$, $\delta_M(a, p) \stackrel{\text{def}}{=} q$, where q is the state such that $p \xrightarrow{a} q$. Observe that, because of determinism, if $p \xrightarrow{\varphi_1} q_1, p \xrightarrow{\varphi_2} q_2$, and $a \in \llbracket \varphi_1 \wedge \varphi_2 \rrbracket$, then $q_1 = q_2$. Moreover, due to completeness, there exists some q and φ such that $p \xrightarrow{\varphi} q$ and $a \in \llbracket \varphi \rrbracket$.

Determinization of SFAs is always possible and is studied in [42]. Completion is straightforward: if M is not complete then add a new state q_0 and the self-loop $q_0 \xrightarrow{\top} q_0$, and for each partial state q add the move $(q, \bigwedge_{\rho \in \overrightarrow{\Delta}(q)} \neg \text{Grd}(\rho), q_0)$. Observe that completion requires negation of predicates.

Normalization is obvious: if there exist states p and q and two distinct transitions $p \xrightarrow{\varphi} q$ and $p \xrightarrow{\psi} q$, then replace these transitions with the single transition $p \xrightarrow{\varphi \vee \psi} q$. This does not affect $\mathcal{L}_p(M)$ for any p .

We always assume that M is clean. Cleaning amounts to running a standard forward reachability algorithm that keeps only reachable states, and eliminates infeasible moves. Observe that infeasible moves $p \xrightarrow{\perp} q$ do not add expressiveness and might cause unnecessary state space explosion.

It is important to show that minimality of SFAs is in fact well-defined in the sense that minimal SFAs are unique up to renaming of states and equivalence of predicates. To do so we use the following construction.

Assume $M = (\mathcal{A}, Q, q^0, F, \Delta)$ to be deterministic and complete. Let $\Sigma_{\mathcal{A}}$ denote the *first-order language* that contains the *unary relation symbol* \bar{F} and the *unary function symbol* \bar{a} for each $a \in \mathfrak{D}_A$. We define the $\Sigma_{\mathcal{A}}$ -*structure* of M , denoted by \mathfrak{M} , to have the universe Q , and the interpretation function:

$$\bar{F}^{\mathfrak{M}} \stackrel{\text{def}}{=} F, \quad \forall \bar{a} \in \Sigma_{\mathcal{A}}, p \in Q \quad (\bar{a}^{\mathfrak{M}}(p) \stackrel{\text{def}}{=} \delta_M(a, p)).$$

Also, let

$$M(\epsilon) \stackrel{\text{def}}{=} q^0,$$

$$M(w \cdot a) \stackrel{\text{def}}{=} \delta_M(a, M(w)) \quad \text{for } a \in \mathfrak{D}_A \text{ and } w \in \mathfrak{D}_A^*.$$

In other words, $M(w)$ is the state reached in M for the word $w \in \mathfrak{D}_A^*$. Recall that two Σ -structures are *isomorphic* if there

² It is sometimes convenient to define minimality over incomplete SFAs, in which case the *dead-end* state q ($q \neq q^0$ and $\mathcal{L}_q(M) = \emptyset$) is eliminated if it is present.

exists a bijective mapping between their universes that preserves the meaning of all symbols in Σ .

THEOREM 1. *If M and N are minimal SFAs over the same alphabet \mathcal{A} such that $\mathcal{L}(M) = \mathcal{L}(N)$, then \mathfrak{M} and \mathfrak{N} are isomorphic $\Sigma_{\mathcal{A}}$ -structures.*

Proof: Assume M and N to be minimal SFAs over \mathcal{A} such that $\mathcal{L}(M) = \mathcal{L}(N)$. We define $\iota : \mathfrak{M} \cong \mathfrak{N}$ as follows:

$$\forall w \in \mathfrak{D}_A^* \quad (\iota(M(w)) \stackrel{\text{def}}{=} N(w)).$$

To show that ι is well-defined as a function, observe first that all states of M correspond to some w because M is clean. Second, we prove that for all words v and w , if $M(v) = M(w)$ then $N(v) = N(w)$. Fix $v, w \in \mathfrak{D}_A^*$ such that $M(v) = M(w)$. Then, for all $u \in \mathfrak{D}_A^*$,

$$\begin{aligned} u \in \mathcal{L}_{N(v)}(N) & \Leftrightarrow v \cdot u \in \mathcal{L}(N) \\ & \stackrel{(\text{by } \mathcal{L}(M) = \mathcal{L}(N))}{\Leftrightarrow} v \cdot u \in \mathcal{L}(M) \\ & \Leftrightarrow u \in \mathcal{L}_{M(v)}(M) \\ & \stackrel{(\text{by } M(w) = M(v))}{\Leftrightarrow} u \in \mathcal{L}_{M(w)}(M) \\ & \Leftrightarrow w \cdot u \in \mathcal{L}(M) \\ & \stackrel{(\text{by } \mathcal{L}(M) = \mathcal{L}(N))}{\Leftrightarrow} w \cdot u \in \mathcal{L}(N) \\ & \Leftrightarrow u \in \mathcal{L}_{N(w)}(N) \end{aligned}$$

So, $\mathcal{L}_{N(v)}(N) = \mathcal{L}_{N(w)}(N)$, and thus $N(v) = N(w)$ by minimality of N . So, ι is well-defined as a function. By switching the roles of M and N we also get the opposite direction. Thus,

$$(*) \quad \forall v, w \in \mathfrak{D}_A^* \quad (M(v) = M(w) \Leftrightarrow N(v) = N(w))$$

Next, we show that ι is an isomorphism.

First, we show that ι is bijective: ι is onto because N is clean, i.e., each state of N corresponds to $N(w)$ for some word w ; ι is into because if $M(v) \neq M(w)$ then, by (*), $\iota(M(v)) = N(v) \neq N(w) = \iota(M(w))$.

Finally, we show that ι is an embedding of \mathfrak{M} into \mathfrak{N} , i.e., that ι preserves all the functions and the relations: for all $p \in Q_M$, $p \in \bar{F}^{\mathfrak{M}} \Leftrightarrow \iota(p) \in \bar{F}^{\mathfrak{N}}$, and for all $p \in Q_M$ and $a \in \mathfrak{D}_A$, $\iota(\bar{a}^{\mathfrak{M}}(p)) = \bar{a}^{\mathfrak{N}}(\iota(p))$.

Let $p \in Q_M$. Let w be any word such that $p = M(w)$. Then

$$\begin{aligned} p \in \bar{F}^{\mathfrak{M}} & \Leftrightarrow M(w) \in \bar{F}^{\mathfrak{M}} \Leftrightarrow w \in \mathcal{L}(M) \Leftrightarrow w \in \mathcal{L}(N) \\ & \Leftrightarrow N(w) \in \bar{F}^{\mathfrak{N}} \Leftrightarrow \iota(M(w)) \in \bar{F}^{\mathfrak{N}} \Leftrightarrow \iota(p) \in \bar{F}^{\mathfrak{N}}. \end{aligned}$$

and, for any $a \in \mathfrak{D}_A$,

$$\begin{aligned} \iota(\bar{a}^{\mathfrak{M}}(p)) & = \iota(\bar{a}^{\mathfrak{M}}(M(w))) = \iota(\delta_M(a, M(w))) \\ & = \iota(M(w \cdot a)) = N(w \cdot a) = \delta_N(a, N(w)) \\ & = \bar{a}^{\mathfrak{N}}(N(w)) = \bar{a}^{\mathfrak{N}}(\iota(M(w))) = \bar{a}^{\mathfrak{N}}(\iota(p)) \end{aligned}$$

Thus \mathfrak{M} and \mathfrak{N} are isomorphic. \square

The theorem implies that minimal SFAs are unique up to renaming of states and up to equivalence of predicates due to normalization.

DEFINITION 3. Two states $p, q \in Q$ are *M-equivalent*, $p \equiv_M q$, when $\mathcal{L}_p(M) = \mathcal{L}_q(M)$. \square

We have that \equiv_M is an equivalence relation. If \equiv is an equivalence relation over Q , then for $q \in Q$, q/\equiv denotes the equivalence class containing q , for $X \subseteq Q$, X/\equiv denotes $\{q/\equiv \mid q \in X\}$, and M/\equiv denotes the SFA $M/\equiv \stackrel{\text{def}}{=} (\mathcal{A}, Q/\equiv, q^0/\equiv, F/\equiv, \Delta/\equiv)$ where:

$$\Delta/\equiv \stackrel{\text{def}}{=} \{(p/\equiv, \bigvee_{(p, \varphi, q) \in \Delta} \varphi, q/\equiv) \mid p, q \in Q, \exists \varphi((p, \varphi, q) \in \Delta)\}$$

Observe that $M_{/\equiv}$ is normalized by construction. We need the following theorem that shows that minimization of SFAs preserves their intended semantics.

THEOREM 2. *Let M be a clean, complete and deterministic SFA. Then $M_{/\equiv_M}$ is minimal and $\mathcal{L}(M) = \mathcal{L}(M_{/\equiv_M})$.*

Proof: Let \equiv be \equiv_M . Clearly, $M_{/\equiv}$ is clean and complete because M is clean and complete. To show determinism, let $\mathbf{p} \xrightarrow{a}_{M_{/\equiv}} \mathbf{q}_1$, and $\mathbf{p} \xrightarrow{a}_{M_{/\equiv}} \mathbf{q}_2$. Take $p_1, p_2 \in \mathbf{p}$, $q_1 \in \mathbf{q}_1$ and $q_2 \in \mathbf{q}_2$ such that $p_1 \xrightarrow{a}_M q_1$ and $p_2 \xrightarrow{a}_M q_2$. Since $\mathcal{L}_{p_1}(M) = \mathcal{L}_{p_2}(M)$, and M is deterministic, it follows that $\mathcal{L}_{q_1}(M) = \mathcal{L}_{q_2}(M)$ i.e., $\mathbf{q}_1 = \mathbf{q}_2$. Thus,

$$(*) \quad \forall a \in \mathcal{Q}_A, p \in Q_M (\delta_{M_{/\equiv}}(a, p_{/\equiv}) = \delta_M(a, p)_{/\equiv}).$$

Minimality of $M_{/\equiv}$ follows from the definition. Next, we show by induction over the length of w that

$$(\star) \quad \forall w \in \mathcal{Q}_A^* (M(w)_{/\equiv} = M_{/\equiv}(w))$$

For $w = \epsilon$ we have that $M(\epsilon) = q_M^0$, $M_{/\equiv}(\epsilon) = q_{M_{/\equiv}}^0$, and $q_{M_{/\equiv}}^0 = (q_M^0)_{/\equiv}$.

For $w = v \cdot a$, where $a \in \mathcal{Q}_A$ and $v \in \mathcal{Q}_A^*$, we have that

$$\begin{aligned} M(v \cdot a)_{/\equiv} &= \delta_M(a, M(v))_{/\equiv} \stackrel{(\text{by } *)}{=} \delta_{M_{/\equiv}}(a, M(v)_{/\equiv}) \\ &\stackrel{(\text{by IH})}{=} \delta_{M_{/\equiv}}(a, M_{/\equiv}(v)) = M_{/\equiv}(v \cdot a). \end{aligned}$$

It follows that, for all words $w \in \mathcal{Q}_A^*$, $w \in \mathcal{L}(M)$ iff $M(w) \in F_M$ iff $M(w)_{/\equiv} \in F_{M_{/\equiv}}$ iff (by (\star)) $M_{/\equiv}(w) \in F_{M_{/\equiv}}$ iff $w \in \mathcal{L}(M_{/\equiv})$. Thus $\mathcal{L}(M) = \mathcal{L}(M_{/\equiv})$. \square

Theorems 1 and 2 are classical theorems lifted to arbitrary (possibly infinite) alphabets. Theorem 2 implies that SFAs have equivalent minimal forms that, by Theorem 1, are unique up to relabeling of states, and modulo equivalence of predicates in \mathcal{A} . In particular, since for all Q and all equivalence relations E over Q , $|Q_{/E}| \leq |Q|$, each equivalent form has minimal number of states.

3. Moore's algorithm over symbolic alphabets

Moore's minimization algorithm [34] of DFAs (also due to Huffman [29]) is commonly known as the *standard* algorithm. Even though the classical version of Moore's algorithm depends on the alphabet being finite, the general idea can be lifted to SFAs as follows. Given an SFA M , initially, let D be the binary relation $(F \times F^c) \cup (F^c \times F)$, where F^c is $Q \setminus F$. Compute the fixpoint of D as follows: if $\langle p, q \rangle \in D$ and there exist moves (p', φ, p) and (q', ψ, q) where $\varphi \wedge \psi$ is *satisfiable*³ then add $\langle p', q' \rangle$ to D . This process clearly terminates. Upon termination, $E = (Q \times Q) \setminus D$ is the equivalence relation \equiv_M , and the SFA $M_{/E}$ is therefore minimal. We refer to this algorithm as Min_{SFA}^M . Observe that Min_{SFA}^M checks only *satisfiability of conjunctions* of conditions and does not depend on the full power of the alphabet algebra, in particular it does not require the ability to complement predicates (assuming the initial completion of M is viewed as a separate preprocessing step). This is in contrast with the generalization of Hopcroft's algorithm discussed in the next section.

Complexity. In the finite alphabet case Moore's algorithm can be implemented in $\mathcal{O}(kn \log n)$ (using the approach described in [10]), where n is the number of states of the DFA, and k the number of characters in the input alphabet. However, such implementation relies on the alphabet being finite.

³In a concrete implementation a dictionary can be used to maintain D similar to the case of DFAs [28, Section 3.4].

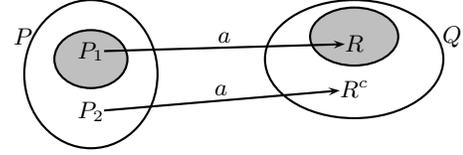


Figure 1. The idea behind an (a, R) -split of a part P .

Our implementation of the version of Moore's algorithm presented above has the following complexity. Given an SFA A , let n be the number of states of A , m the number of moves of A , and ℓ the size of the largest guard appearing in a transition of A (the biggest predicate). If the complexity of the alphabet theory for instances of size l is $f(l)$, then Min_{SFA}^M has complexity $\mathcal{O}(m^2 f(\ell))$, which is $\mathcal{O}(n^4 f(\ell))$ assuming normalized A where m is $\mathcal{O}(n^2)$.

4. Hopcroft's algorithm over symbolic alphabets

Hopcroft's algorithm [27] for minimizing DFAs is based on a technique called *partition refinement* of states. The symbolic version of Hopcroft's algorithm is given in Figure 3. Initially the set of states is partitioned into two sets: *final* states and *non-final* states, i.e., $\mathcal{P} = \{F, Q \setminus F\}$. Here we assume that the SFA M is complete and nontrivial, so that both F and $Q \setminus F$ are nonempty. The partition \mathcal{P} induces the equivalence relation $\equiv_{\mathcal{P}}$ (or \equiv when \mathcal{P} is clear from the context) over Q , such that $Q_{/\equiv} = \mathcal{P}$. At a high level, \mathcal{P} is refined as follows: suppose there exist parts $P, R \in \mathcal{P}$ and a character $a \in \mathcal{Q}_A$ such that some a -move from P goes into R , and some a -move from P goes into R^c ($Q \setminus R$). Then the (a, R) -split of P is the set $\{P_1, P_2\}$, where P_1 (resp. P_2) is the set of all $p \in P$ from which there is an a -move into R (resp. R^c). P is then replaced in \mathcal{P} by P_1 and P_2 . See Figure 1. The (a, R) -split of P is well-defined by determinism and completeness of M . The following invariant is maintained by splitting.

LEMMA 1. *After each refinement step of \mathcal{P} , for all $p, q \in Q$, if $p \not\equiv q$ then $\mathcal{L}_p(M) \neq \mathcal{L}_q(M)$.*

Proof: Initially, for $p_1 \in F$: $\epsilon \in \mathcal{L}_{p_1}(M)$, and $p_2 \in Q \setminus F$: $\epsilon \notin \mathcal{L}_{p_2}(M)$. So, the statement holds for step $i = 0$. To see that the statement holds after step $i + 1$, fix $P, R \in \mathcal{P}_i$ and $a \in \mathcal{Q}_A$ such that (P_1, P_2) is the (a, R) -split of P . It is enough to show that, for all $p_1 \in P_1, p_2 \in P_2$, $\mathcal{L}_{p_1}(M) \neq \mathcal{L}_{p_2}(M)$. Consider $p_1 \xrightarrow{a} q_1$ and $p_2 \xrightarrow{a} q_2$. Since q_1 and q_2 belong to distinct parts of \mathcal{P}_i , by IH, there exists a word w such that $w \in \mathcal{L}_{q_1}(M)$ iff $w \notin \mathcal{L}_{q_2}(M)$. Moreover, by determinism of M , we have $a \cdot w \in \mathcal{L}_{p_j}(M)$ iff $w \in \mathcal{L}_{q_j}(M)$ for $j \in \{1, 2\}$. So, $a \cdot w \in \mathcal{L}_{p_1}(M)$ iff $a \cdot w \notin \mathcal{L}_{p_2}(M)$, and thus $\mathcal{L}_{p_1}(M) \neq \mathcal{L}_{p_2}(M)$. \square

Assuming a simple iterative refinement loop of \mathcal{P} , splitting is repeated until no further splits are possible, i.e., until \mathcal{P} is *finest*, at which point the following property holds.

LEMMA 2. *Upon termination of refinement of \mathcal{P} , for all $p, q \in Q$, if $p \equiv q$ then $\mathcal{L}_p(M) = \mathcal{L}_q(M)$.*

Proof: By way of contradiction. Suppose

- (*) there exists a word x and states $p_1 \equiv p_2$, such that $x \in \mathcal{L}_{p_1}(M) \Leftrightarrow x \notin \mathcal{L}_{p_2}(M)$.

Choose w to be a *shortest* such x . Since w cannot be ϵ (because $\{p_1, p_2\} \subseteq F$ or $\{p_1, p_2\} \subseteq Q \setminus F$), there are a and v such that $a \cdot v = w$, and moves $p_1 \xrightarrow{\varphi_1} q_1$ and $p_2 \xrightarrow{\varphi_2} q_2$ such that $a \in \llbracket \varphi_1 \rrbracket$ and $a \in \llbracket \varphi_2 \rrbracket$. The choice of q_1 and q_2 is unique for given a by determinism, thus $w \in \mathcal{L}_{p_i}(M) \Leftrightarrow v \in \mathcal{L}_{q_i}(M)$ for $i \in \{1, 2\}$. So, by (*), $v \in \mathcal{L}_{q_1}(M) \Leftrightarrow v \notin \mathcal{L}_{q_2}(M)$.

```

1  $Minterms_{\mathcal{A}}(\bar{\psi}) \stackrel{\text{def}}{=}$ 
2  $tree := \text{new Tree}(\top_{\mathcal{A}}, \text{null}, \text{null});$ 
3 foreach ( $\psi$  in  $\bar{\psi}$ )  $tree.Refine(\psi);$ 
4 return  $Leaves(tree);$  //return the set of all the leaf predicates
5 class  $Tree$ 
6  $Predicate \varphi; Tree \text{left}; Tree \text{right};$ 
7  $Refine(\psi) \stackrel{\text{def}}{=}$ 
8 if ( $IsSat_{\mathcal{A}}(\varphi \wedge_{\mathcal{A}} \psi)$  and  $IsSat_{\mathcal{A}}(\varphi \wedge_{\mathcal{A}} \neg_{\mathcal{A}} \psi)$ )
9 if ( $\text{left} = \text{null}$ ) //if the tree is a leaf then split  $\varphi$  into two parts
10  $\text{left} := \text{new Tree}(\varphi \wedge_{\mathcal{A}} \psi, \text{null}, \text{null});$  // $[\varphi] \cap [\psi]$ 
11  $\text{right} := \text{new Tree}(\varphi \wedge_{\mathcal{A}} \neg_{\mathcal{A}} \psi, \text{null}, \text{null});$  // $[\varphi] \setminus [\psi]$ 
12 else  $\text{left.Refine}(\psi); \text{right.Refine}(\psi);$  //refine subtrees recursively

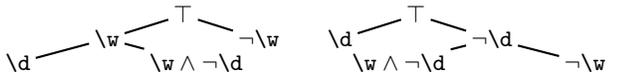
```

Figure 2. Minterm generation for $\bar{\psi} \subseteq \Psi_{\mathcal{A}}$ modulo \mathcal{A} .

We now have the following contradiction. There are two cases: 1) if $q_1 \neq q_2$, say $q_1 \in R$ and $q_2 \in R^c$ for some part R , then \mathcal{P} is not finest because we can split the part $P \in \mathcal{P}$ such that $p_1, p_2 \in P$ using the (a, R) -split of P . 2) if $q_1 \equiv q_2$ then, because $|v| < |w|$ and $(*)$ holds for $x = v$ and q_1, q_2 in place of p_1, p_2 , it follows that w is not the shortest x such that $(*)$ holds \downarrow . \square

In addition to the state partition refinement, in the symbolic case, we also use *predicate refinement*. Predicate refinement builds a set of *minterms*. A minterm is a minimal satisfiable Boolean combinations of all guards that occur in the SFA. The algorithm is shown in Figure 2. It uses a binary tree whose leaves define the partition. Initially the tree is the leaf \top . Each time a predicate ψ is used to refine the tree it may cause splitting of its leaves into finer predicates. The following example illustrates such minterm generation.

EXAMPLE 1. Consider the alphabet algebra $2^{\text{BV}7}$ (ASCII characters). We use standard regex notation for character classes. Suppose that the following two guards occur in the given SFA: $\backslash w$ ($[\backslash w] = [[a-zA-Z0-9_]]$), and $\backslash d$ ($[\backslash d] = [[0-9]]$). Then, the value of $tree$ in $Minterms_{2^{\text{BV}7}}(\{\backslash w, \backslash d\})$ in line 4 in Figure 2 is either the first of the two trees below if $\backslash w$ is selected first in the loop of line 3, or else the second tree (note that $[\backslash d] \subseteq [\backslash w]$ so $[\backslash d \wedge \backslash w] = [\backslash d]$):



The minterms are the leaves $\backslash d$, $\backslash w \wedge \neg \backslash d$ and $\neg \backslash w$. \square

Next, we analyze a property of the set of minterms. Given a predicate $\psi \in \Psi_{\mathcal{A}}$ and a state $p \in Q$, define:

$$\begin{aligned} \delta(\psi, p) &\stackrel{\text{def}}{=} \{Tgt(\rho) \mid \rho \in \overrightarrow{\Delta}(p), IsSat(Grd(\rho) \wedge \psi)\} \\ \delta^{-1}(\psi, p) &\stackrel{\text{def}}{=} \{Src(\rho) \mid \rho \in \overleftarrow{\Delta}(p), IsSat(Grd(\rho) \wedge \psi)\} \\ \delta^{-1}(\psi, P) &\stackrel{\text{def}}{=} \bigcup_{p \in P} \delta^{-1}(\psi, p) \quad (\text{for } P \subseteq Q) \end{aligned}$$

Let $Minterms(M) \stackrel{\text{def}}{=} Minterms_{\mathcal{A}}(\bigcup_{\rho \in \Delta} Grd(\rho))$. The following proposition implies that all characters that occur in one minterm are indistinguishable.

PROPOSITION 1. *Let M be deterministic and complete. For all $\psi \in Minterms(M)$ and $p \in Q$, $|\delta(\psi, p)| = 1$.*

We can therefore treat δ as a function from $Minterms \times Q$ to Q and reduce minimization of the SFA to minimization of the DFA with alphabet $Minterms$ and transition function δ . In particular, we can

```

1  $Min_{SFA}^H(M = (\mathcal{A}, Q, q^0, F, \Delta)) \stackrel{\text{def}}{=}$ 
2  $\mathcal{P} := \{F, Q \setminus F\};$  //initial partition
3  $W := \{\text{if } (|F| \leq |Q \setminus F|) \text{ then } F \text{ else } Q \setminus F\};$ 
4  $\Psi := Minterms_{\mathcal{A}}(\{Grd(\rho) \mid \rho \in \Delta\});$  //compute the minterms
5 while ( $W \neq \emptyset$ ) //iterate over unvisited parts
6  $R := \text{choose}(W); W := W \setminus \{R\};$ 
7 foreach ( $\psi$  in  $\Psi$ ) //iterate over all minterms
8  $S := \delta^{-1}(\psi, R);$  //all states leading into  $R$  for given minterm
9 while ( $\text{exists } (P \text{ in } \mathcal{P}) \text{ where } P \cap S \neq \emptyset \text{ and } P \setminus S \neq \emptyset$ )
10  $\langle \mathcal{P}, W \rangle := Split_{\mathcal{P}, W}(P, P \cap S, P \setminus S);$  //split  $P$ 
11 return  $M /_{\equiv_{\mathcal{P}}}$ ;
12  $Split_{\mathcal{P}, W}(P, P_1, P_2) \stackrel{\text{def}}{=} \langle \mathcal{P}', W' \rangle$  where
13  $\mathcal{P}' = (\mathcal{P} \setminus \{P\}) \cup \{P_1, P_2\}$  //refine  $\mathcal{P}$ 
14  $W' = \text{if } (P \in W) \text{ then } (W \setminus \{P\}) \cup \{P_1, P_2\}$  //both parts
15 else  $W \cup \{\text{if } (|P_1| \leq |P_2|) \text{ then } P_1 \text{ else } P_2\}$  //smaller part

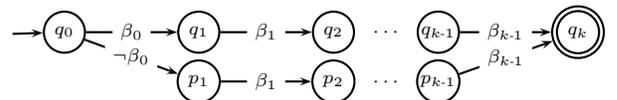
```

Figure 3. Hopcroft's minimization algorithm lifted to deterministic SFAs. M is assumed to be clean, complete, and nontrivial ($F \neq \emptyset$ and $Q \setminus F \neq \emptyset$).

use Hopcroft's algorithm. We refer to the resulting algorithm by Min_{SFA}^H . Such an algorithm is shown in Figure 3.

Lifting the transition relation to be over *minterms* is a powerful hammer that can be used to adapt most classical automata algorithms to the symbolic setting. However, one drawback of minterm generation is that, in the worst case, the number of minterms of an SFA is exponential in the number of guards occurring in the SFA. The following example illustrates a worst case scenario in which, due to such a problem, Min_{SFA}^H runs in exponential time.

EXAMPLE 2. Let the character domain be nonnegative integers $< 2^k$. Suppose $\beta_i(x)$ is a predicate that is true for x iff the i 'th bit of the binary representation of x is 1, e.g. $\beta_3(8)$ is true and $\beta_3(7)$ is false. Predicate β_3 can be defined as $\neg((x \& 8) = 0)$, provided that, besides equality, the bitwise-and operator $\&$ is a built-in function symbol of $\Psi_{\mathcal{A}}$ (e.g., consider the bit-vector theory of an SMT solver). Similarly, we may also use the algebra $2^{\text{BV}k}$, where the size of the concrete BDD representation for β_i is linear in k , it has one node that is labeled by i and whose left child (case bit is 0) is false and whose right child (case bit is 1) is true. The point is that predicates are small (essentially constant) in size. Consider the following SFA M_k with such an alphabet \mathcal{A} .



Then $Minterms(M_k) = Minterms_{\mathcal{A}}(\{\neg \beta_i, \beta_i\}_{i < k}) = \{\hat{n}\}_{n < 2^k}$ has 2^k elements, where $[\hat{n}] = \{n\}$. For example, suppose $k = 3$, then $[\beta_2 \wedge \neg \beta_1 \wedge \beta_0] = \{5\}$. The minimal automaton is



The dead-end state q_0 necessary for completion is implicit. \square

Complexity. In the finite alphabet case, Hopcroft's algorithm has complexity $\mathcal{O}(kn \log n)$, where n is the number of states of the DFA and k is the number of characters in the input alphabet [27], assuming k is treated as a constant. It is shown in [31] that if k is $\mathcal{O}(n)$ then the complexity of the algorithm presented in [27] is

```

1  $Min_{SFA}^N(M = (\mathcal{A}, Q, q^0, F, \Delta)) \stackrel{\text{def}}{=}$ 
2  $\mathcal{P} := \{F, Q \setminus F\};$  //initial partition
3  $W := \{\text{if } (|F| \leq |Q \setminus F|) \text{ then } F \text{ else } Q \setminus F\};$ 
4 while ( $W \neq \emptyset$ ) //main loop
5    $R := \text{choose}(W); W := W \setminus \{R\};$ 
6    $S := \{Src(\tau) \mid \tau \in \overleftarrow{\Delta}(R)\};$  //all states leading into R
   // $\Gamma(p)$  denotes the set of all characters leading from  $p \in S$  into R
7    $\Gamma := \{p \mapsto \bigvee_{\tau \in \overleftarrow{\Delta}(R), Src(\tau)=p} Grd(\tau) \mid p \in S\};$ 
8   while (exists ( $P$  in  $\mathcal{P}$ ) where  $P \cap S \neq \emptyset$  and  $P \setminus S \neq \emptyset$ )
9      $\langle \mathcal{P}, W \rangle := Split_{\mathcal{P}, W}(P, P \cap S, P \setminus S);$  //( $\_, R$ )-split
10    while (exists ( $P$  in  $\mathcal{P}$ ) where  $P \cap S \neq \emptyset$  and
11      exists ( $p_1, p_2$  in  $P$ ) where  $IsSat(\neg(\Gamma(p_1) \Leftrightarrow \Gamma(p_2)))$ )
12       $a := \text{choose}(\llbracket \neg(\Gamma(p_1) \Leftrightarrow \Gamma(p_2)) \rrbracket);$ 
13       $P_1 := \{p \in P \mid a \in \llbracket \Gamma(p) \rrbracket\};$ 
14       $\langle \mathcal{P}, W \rangle := Split_{\mathcal{P}, W}(P, P_1, P \setminus P_1);$  //( $a, R$ )-split
15  return  $M_{/ \equiv_{\mathcal{P}}}$ ;

```

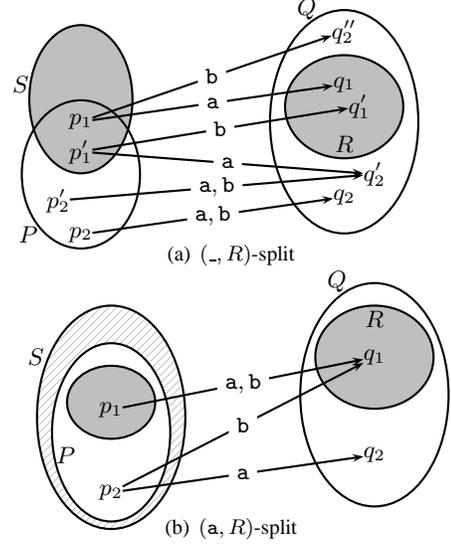


Figure 5. Split cases of P in Min_{SFA}^N . Suppose $\mathcal{D}_A = \{a, b\}$.

Figure 4. Minimization of deterministic SFAs. M is assumed to be clean, complete and nontrivial. $Split_{\mathcal{P}, W}$ is defined in Figure 3.

$O(n^3)$, but it is also shown that $O(kn \log n)$ complexity can be maintained with a more careful implementation.

Given an SFA, let n be the number of states, let m be the number of moves, and let ℓ be the size of the largest label (guard) of a move. In a normalized SFA there are at most n^2 moves. Let μ be the number of minterms of the SFA; μ is bounded by 2^m . Each minterm is of the form $\phi_1 \wedge \dots \wedge \phi_m$, where each ϕ_i is a guard or a negated guard, and thus minterms have size $O(m\ell)$. The computation of minterms is $O(2^m f(m\ell))$ where f is the complexity to decide satisfiability of formulas of given size. The rest of the algorithm (after line 4) can be seen as a standard implementation of Hopcroft’s algorithm, where the SFA is first transformed into a DFA with an alphabet of minterm identifiers where each SFA label has been replaced by identifiers of all relevant minterms. Then Ψ is viewed as a concrete alphabet (of such minterm identifiers). In the final result, the identifiers are mapped back to minterms and the resulting SFA is normalized. The overall complexity is then $O(2^m f(m\ell) + 2^m n \log n)$.

5. Minimization without minterm generation

When inspecting the worst case complexity of Min_{SFA}^H , the factor that stands out the most is the exponential blowup caused by the predicate refinement (minterm generation). In this section we investigate a new technique which is based on the symbolic representation of the alphabet, and that avoids minterm generation. Figure 4 shows Min_{SFA}^N , a new minimization algorithm, that does not require predicate refinement.

The intuition behind Min_{SFA}^N is the following: when splitting a partition’s part, it is not necessary to provide the exact witness that defines the split, instead it is enough to check if some witness (or witness set) exists. The main steps are the two inner while loops in Figure 4, they both refine the partition \mathcal{P} with respect to R . In the first loop (lines 8–9) a part P is split into $P \cap S$ and $P \setminus S$ without using a fixed witness (see Figure 5(a) where $\{\{p_1, p_1'\}, \{p_2, p_2'\}\}$ is a $(_, R)$ -split of P , but it is neither an (a, R) -split nor a (b, R) -split of P). The second loop (lines 10–14) splits P if there exists some a that produces an (a, R) -split of P . Such an element a

must somehow distinguish at least two states p_1 and p_2 of P , (see Figure 5(b)), so that the split is indeed proper and guarantees progress.

The first loop (lines 8–9) is an optimization that can be omitted without sacrificing correctness, provided that for $p \notin S$ we let $\Gamma(p) \stackrel{\text{def}}{=} \perp$. The conditions $P \setminus S \neq \emptyset$ and $P \cap S \neq \emptyset$ together imply that there exist $p_1 \in P \cap S$ and $p_2 \in P \setminus S$ and thus $\llbracket \Gamma(p_1) \rrbracket \neq \emptyset$ but $\llbracket \Gamma(p_2) \rrbracket = \emptyset$, so $\neg(\Gamma(p_1) \Leftrightarrow \Gamma(p_2))$ is satisfiable.

The concrete implementation is shown in Figure 6. It differs from the abstract algorithm in that it computes local minterms, and does not compute any concrete witnesses (the element a is not computed) in the second loop. In the concrete implementation it is important to keep the first loop for the following reasons. One is efficiency, the first loop is cheaper. Second is simplicity, it is useful to work with Γ as a dictionary or array whose index set is S , and it is practical to assume that the invariant $P \subseteq S$ holds during the second loop.

Next, we formally prove the correctness of Min_{SFA}^N . The proof provides more intuition on how the algorithm works. We then provide more details on how the concrete implementation works.

THEOREM 3. $Min_{SFA}^N(M)$ is minimal, and $\mathcal{L}(Min_{SFA}^N(M)) = \mathcal{L}(M)$.

Proof: We show first that the invariant of Lemma 1 holds. The invariant clearly holds initially. We show that it is preserved by each split (lines 8–9, and lines 10–14).

First, consider the first splitting loop in Figure 4. Fix $R \in \mathcal{P}$ and let $S = \{Src(\tau) \mid \tau \in \overleftarrow{\Delta}(R)\}$, and choose $P \in \mathcal{P}$ such that $P_1 = P \cap S \neq \emptyset$ and $P_2 = P \setminus S \neq \emptyset$. Fix $p_1 \in P_1$ and $p_2 \in P_2$. Then, there is a move $p_1 \xrightarrow{\varphi_1} q_1$ for some $q_1 \in R$. Let $a \in \llbracket \varphi_1 \rrbracket$. Since $p_2 \notin S$ and M is complete, for some $q_2 \in R^c$ there exist a move $p_2 \xrightarrow{\varphi_2} q_2$ such that $a \in \llbracket \varphi_2 \rrbracket$. The situation is illustrated in Figure 5(a). By using the invariant, $\mathcal{L}_{q_1}(M) \neq \mathcal{L}_{q_2}(M)$, there is a word w such that $w \in \mathcal{L}_{q_1}(M) \Leftrightarrow w \notin \mathcal{L}_{q_2}(M)$. Thus, by using the fact that M is deterministic, $a \cdot w \in \mathcal{L}_{p_1}(M) \Leftrightarrow a \cdot w \notin \mathcal{L}_{p_2}(M)$. Therefore, $\mathcal{L}_{p_1}(M) \neq \mathcal{L}_{p_2}(M)$.

Second, consider the second splitting loop. Fix P and $p_1, p_2 \in P$ that satisfy the loop condition. All parts in \mathcal{P} that intersect with S must be subsets of S due to the first splitting loop, so, since M is clean, $P \subseteq S$, $IsSat(\Gamma(p_1))$ and $IsSat(\Gamma(p_2))$. The condi-

```

MinSFA(Automaton<S> fa)
{
  var fB = new Block(fa.GetFinalStates());
  var nfB = new Block(fa.GetNonFinalStates());
  var blocks = new Dictionary<int, Block>();
  foreach (var q in fa.GetFinalStates()) blocks[q] = fB;
  foreach (var q in fa.GetNonFinalStates()) blocks[q] = nfB;
  var W = new BlockStack();
  if (nfB.Count < fB.Count) W.Push(nfB); else W.Push(fB);

  while (!W.IsEmpty) {
    var R = W.Pop();
    var G = ... //Γ in Figure 4
    var S = G.Keys;
    var relevant = ... //blocks intersecting with S
    foreach (var P in relevant){ //lines 8-9 in Figure 4
      var P1 = ... //P ∩ S
      if (P1.Count < P.Count) { //(L,R)-split of P
        foreach (var p in P1) { P.Remove(p); blocks[p] = P1;
          if (W.Contains(P)) W.Push(P1);
          else if (P.Count <= P1.Count) W.Push(P);
          else W.Push(P1);
        }
      }
    }
    bool iterate = true;
    while (iterate) { //lines 10-14 in Figure 4
      iterate = false;
      relevant = ... //blocks intersecting with S
      foreach (var P in relevant) {
        var P1 = new Block();
        var psi = G[P.Current]; //start with some element of P
        bool splitterFound = false;
        P1.Add(P.Current);

        while (P.MoveNext()) {
          var q = P.Current;
          var phi = G[q];
          if (splitterFound) {
            if (IsSat(psi & phi)) { P1.Add(q); psi = psi & phi; }
          } else {
            if (IsSat(psi & !phi)) {
              psi = psi & !phi; //refine the local minterm
              splitterFound = true;
            } else { //psi implies phi
              if (IsSat(phi & !psi)) {
                P1.Clear(); P1.Add(q); //set P1 to {q}
                psi = phi & !psi; //swap the local minterm
                splitterFound = true;
              } else P1.Add(q); //psi is equivalent to phi
            }
          }
        }
      }
      if (P1.Count < P.Count) { //(a,R)-split of P for some a
        iterate = (iterate || (P.Count > 2));
        foreach (var p in P1) { P.Remove(p); blocks[p] = P1; }
        if (W.Contains(P)) W.Push(P1);
        else if (P.Count <= P1.Count) W.Push(P);
        else W.Push(P1);
      }
    }
  }
  ... //construct the result using blocks and normalize it
}

```

Figure 6. Concrete implementation of Min_{SFA}^N .

tion $IsSat(\neg(\Gamma(p_1) \Leftrightarrow \Gamma(p_2)))$ means that either $IsSat(\Gamma(p_1) \wedge \neg\Gamma(p_2))$ or $IsSat(\Gamma(p_2) \wedge \neg\Gamma(p_1))$. Assume the former case and choose $a \in \llbracket \Gamma(p_1) \wedge \neg\Gamma(p_2) \rrbracket$. By definition of Γ , we know that there is a move $p_1 \xrightarrow{\varphi_1} q_1$ where $q_1 \in R$ such that $a \in \llbracket \varphi_1 \rrbracket$. Moreover, since $a \notin \llbracket \Gamma(p_2) \rrbracket$, and $\Gamma(p_2)$ covers all the characters that lead from p_2 to R , there must be (by completeness and determinism of M) a move $p_2 \xrightarrow{\varphi_2} q_2$ where $q_2 \in R^c$ and $a \in \llbracket \varphi_2 \rrbracket$. See Figure 5(b). It follows as above, by using $\mathcal{L}_{q_1}(M) \neq \mathcal{L}_{q_2}(M)$, that $\mathcal{L}_{p_1}(M) \neq \mathcal{L}_{p_2}(M)$. Since each step properly refines the partition, Lemma 1 follows.

We now show that Lemma 2 holds. The proof is by way of contradiction.

(*) Assume there exists a word x , a part $P \in \mathcal{P}$, and two states $p_1, p_2 \in P$ such that $x \in \mathcal{L}_{p_1}(M) \Leftrightarrow x \notin \mathcal{L}_{p_2}(M)$.

Let w be *shortest* such x . Since w cannot be ϵ , there exist a and v such that $w = a \cdot v$. So, there are, by determinism and completeness of M , unique q_1 and q_2 such that $p_1 \xrightarrow{a} q_1$ and $p_2 \xrightarrow{a} q_2$. It follows that $v \in \mathcal{L}_{q_1}(M) \Leftrightarrow v \notin \mathcal{L}_{q_2}(M)$. So $q_1 \neq q_2$ or else v satisfies (*) and v is shorter than w .

Consider any fixed computation of Min_{SFA}^N . It follows from the definition of $Split_{\mathcal{P}, W}$, that W is always a subset of \mathcal{P} and (due to the first condition of the update to W), if W ever contains a part containing a state q , then W will keep containing a part that contains q until such a part is removed from W in line 5.

Next, we show that the following W -invariant must hold at all times: for all $R \in W$, $q_1 \in R \Leftrightarrow q_2 \in R$. Let $\{i, \bar{i}\} = \{1, 2\}$. Suppose, by way of contradiction, that at some point in line 5 a part R is chosen from W such that $q_i \in R$, and $q_{\bar{i}} \notin R$. So, $p_i \in S$, with S as in line 6. We then have two following cases.

1. If $p_{\bar{i}} \notin S$, then p_i and $p_{\bar{i}}$ are split apart in the first splitting loop. This contradicts the fact that $p_1 \equiv p_2$.
2. Assume $p_i \in S$ and consider the second splitting loop. By choice of the character a above, we know that there exist moves $p_i \xrightarrow{\varphi_i} q_i$ and $p_{\bar{i}} \xrightarrow{\varphi_{\bar{i}}} q_{\bar{i}}$ where $a \in \llbracket \varphi_i \rrbracket$ and $a \in \llbracket \varphi_{\bar{i}} \rrbracket$. It follows that $a \notin \llbracket \Gamma(p_{\bar{i}}) \rrbracket$ (because M is deterministic and $q_i \notin R$) while $a \in \llbracket \Gamma(p_i) \rrbracket$. So $a \in \llbracket \Gamma(p_i) \rrbracket \setminus \llbracket \Gamma(p_{\bar{i}}) \rrbracket$, or, in other words, $a \in \llbracket \Gamma(p_i) \wedge \neg\Gamma(p_{\bar{i}}) \rrbracket$, and therefore $IsSat(\neg(\Gamma(p_i) \Leftrightarrow \Gamma(p_{\bar{i}})))$ holds. Consequently, p_i and $p_{\bar{i}}$ end up in distinct parts upon termination of the second splitting loop. This again contradicts the fact that $p_1 \equiv p_2$.

So, initially $q_1 \in F \Leftrightarrow q_2 \in F$, or else the initial part of W violates the invariant. But now consider the point when the part containing both q_1 and q_2 is split into two parts containing q_1 and q_2 respectively. But at this point, at least one of those parts will be added to W by definition of $Split_{\mathcal{P}, W}$. Thus, we have reached the desired contradiction, because the W -invariant is violated at that point.

We have shown that, upon termination of $Min_{SFA}^N(M)$, $\equiv_{\mathcal{P}}$ coincides with \equiv_M . It follows from Theorem 2 that $M_{/\equiv_{\mathcal{P}}}$ is minimal and accepts $\mathcal{L}(M)$. \square

Implementation. A simplified version of our concrete C# implementation of Min_{SFA}^N is shown in Figure 6. Parts of a state partition are represented by mutable sets called *blocks*, i.e., objects of type `Block`, and states are represented by integers. Each block contains a `HashSet` of states. The search frontier W is maintained as a stack of blocks, and the partition \mathcal{P} is an array of blocks, called `blocks`, that is indexed by states.

The first inner while loop (lines 8–9 in Figure 4) is implemented by iterating over all blocks P that intersect with S . The content of block $P1$ becomes $P \cap S$, while the content of block P is updated to $P \setminus S$. Observe that blocks are objects (pointers), thus if W contains P , after the (L, R) -split it will still contain P as well as the new block $P1$.

The second inner while loop (lines 10–14 in Figure 4) is implemented by an efficient encoding of the search for p_1 and p_2 in line 11 of Figure 4. Moreover, no concrete witness a is computed. Instead, a “local minterm” (called `psi` or ψ in Figure 6), is computed using Γ . This avoids the use of *model generation* that is more expensive than satisfiability checking (i.e. checking if there *exists* a model).⁴ Thus, the implementation does not rely on model generation. Observe also that the second inner while loop relies on the fact that all remaining (`relevant`) blocks that intersect with S must be *contained* in S due to the first inner loop. The split of P

⁴Although $IsSat_{\mathcal{A}}(\varphi)$ is formally defined as $\llbracket \varphi \rrbracket_{\mathcal{A}} \neq \emptyset$, satisfiability checking is a more lightweight operation than actual model generation, in particular, in the context of SMT solvers.

then happens with respect to ψ , which, by construction, is a member of $Min_{terms}(\{\Gamma(p) \mid p \in P\})$, so some member a of $\llbracket \psi \rrbracket$ would have produced the same split.

Relation to classical techniques. Implementation of Hopcroft’s algorithm is discussed in detail in [7]. In the classical setting, the notions of blocks and (a, R) -splits are standard, the pair (R, a) is called a *splitter* in classical algorithms, where it plays a key role. The idea of splitting a block into two and keeping only the smaller block in the waiting set W is a core classical feature. In the case of *partial* DFAs, the algorithms are similar except that W must be initialized with both F and F^c [40].

In classical implementations the waiting set W consists of *splitters* rather than blocks, where characters have been *preselected*. The same is true for partial DFAs except that only those characters that are relevant for a given block are being used, which is beneficial for DFAs with sparse transition graphs. In the symbolic setting the character selection is indirect and achieved only through the satisfiability checks using local minterms during the second loop. The alphabet algebra may be infinite in the case of SFAs, this is not possible in the classical setting. As far as we know, the idea behind the first inner loop of Min_{SFA}^N and the notion of $(-, R)$ -splits (recall Figure 5), say *free splitters* $(R, -)$, have not been studied in the classical setting.

Complexity. The complexity of Min_{SFA}^N depends on several factors, most importantly on the representation of predicates and the concrete representation of the partition refinement data structure, that is explained above.

First of all, observe that each $\Gamma(p)$ has size at most $\mathcal{O}(n\ell)$, and the total size of Γ is $\mathcal{O}(n^2\ell)$. Since the split operator always adds to W only the smallest partition, the outer loop is run at most $\log n$ times. The two internal loops have different running times:

- the first loop is run at most n times, with an internal complexity of $\mathcal{O}(n)$, due to the split operation, and
- the second loop is run at most n times, and if the complexity of the label theory for instances of size l is $f(l)$, the complexity of each iteration is at most $\mathcal{O}(n \cdot f(n\ell))$. This is due to the n internal iterations over the current part P . We also notice that each iteration calls the solver on a predicate of size at most $\mathcal{O}(n\ell)$.

We can conclude that Min_{SFA}^N has complexity $\mathcal{O}(n^2 \log n \cdot f(n\ell))$. As we will observe in the next section, the quadratic behavior is not observed in practice.

6. Evaluation

We evaluate the performance of Min_{SFA}^N , Min_{SFA}^H , and Min_{SFA}^M with the following experiments:

1. We minimize the randomly generated DFAs from the benchmark presented in [3]. This experiment analyzes the performance in the presence of small finite alphabets and allows us to validate our implementation of Min_{SFA}^H against the results shown in [3];
2. We minimize the SFAs generated by common regular expressions taken from the web. This experiment measures performance in the case of typical character alphabets (Unicode);
3. We minimize the SFAs M_k from Example 2 to show the worst case exponential behavior of Min_{SFA}^H . In this experiment we also compare against the `brics` library, a state-of-the-art automata library that uses character ranges as a symbolic representations of finite alphabets;

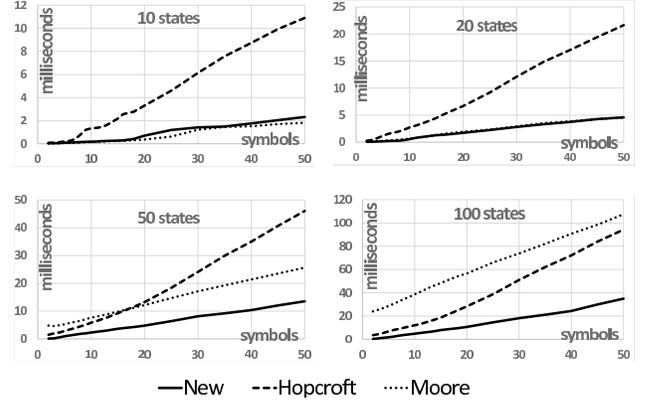


Figure 7. Running times on benchmark of randomly generated DFAs. The corresponding SFAs are over the theory of bitvectors and each set of input symbols is represented as a predicate using a binary decision diagram (BDD).

4. We minimize randomly generated SFAs over the theory of pairs of integers and strings. Here, we analyze the performance in the case in which the alphabet is infinite; and
5. We implemented the classical procedure for transforming Monadic Second Order (MSO) logic formulae into DFAs [38] and measure how the running time of such transformation is affected by different minimization algorithms. This last test aims at understanding the performance in the case of very large alphabets. We also compare our implementation against the tool `Mona` [22], the state of the art library for deciding MSO formulae.

All the experiments are run on a 64 bits Intel(R) Xeon(R) 3.60GHz processor with 8 GB of RAM memory.

6.1 Small randomly generated DFAs

In [3] the performance of several minimization algorithms are compared against a set of randomly generated DFAs. Such benchmark contains 5 million DFAs with number of states ranging between 5 and 1000, and alphabet sizes ranging between 2 and 50. The set of DFAs in [3] is uniformly generated at random, and therefore it offers good statistical coverage of the set of all possible DFAs with such number of states and alphabet sizes.

We run Min_{SFA}^N , Min_{SFA}^M and Min_{SFA}^H on such set of DFAs. Figure 7 shows the results. For simplicity we only plot the results for DFAs with 10, 20, 50, and 100 states. Each plot contains the running time for each algorithm, where the x axis represents the number of symbols in the alphabet. In [3], Moore’s algorithm is not considered and we weren’t therefore able to validate the accuracy of our implementation of Min_{SFA}^M . However, we were able to replicate the behavior of Hopcroft’s algorithm shown [3] using Min_{SFA}^H .

Figure 7 shows how the complexity of Min_{SFA}^M highly depends on the number of states, while the complexity of Min_{SFA}^H is mainly affected by the number of input alphabets. We can indeed see that, already for 100 states, Min_{SFA}^M performs worse than Min_{SFA}^H . In this experiment, for most of the input DFAs the number of minterms was the same as the number of input symbols.

It is not a surprise that Min_{SFA}^N is much faster than both Min_{SFA}^M and Min_{SFA}^H . Indeed, Min_{SFA}^N ’s performance seems to be resistant to both bigger state space and bigger alphabets. Moreover, for no single input Min_{SFA}^N is slower than Min_{SFA}^M or Min_{SFA}^H .

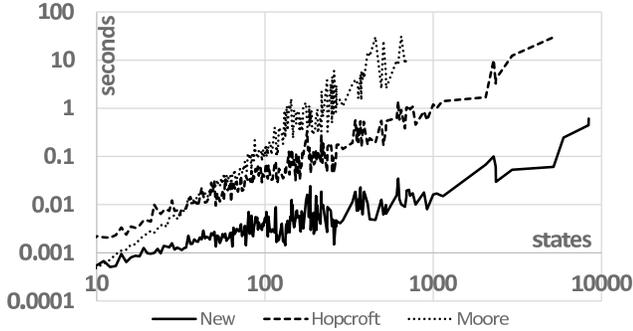


Figure 8. Running times on regexes from regexlib. Both axes are in log-scale.

6.2 Unicode regular expressions

Next, we compared the performance of different minimization algorithms over a sample set of 1850 SFAs over the alphabet 2^{BV16} (of Unicode characters) constructed from typical regexes (taken from a public web-site⁵ of popular regexes). Figure 8 shows the running times. The figure clearly shows how Min_{SFA}^N is faster than both Min_{SFA}^H and Min_{SFA}^M for every input instance. Moreover, it can be appreciated how for bigger state sizes (70-80 states) Min_{SFA}^H starts outperforming Min_{SFA}^M .

In all cases, the number of minterms turned out to be *smaller* (by a factor between 2 and 3) than the total number of predicates, and the exponential blowup of minterms never occurred. The size of the generated SFAs ranged from 2 states to 15800 states with an average of 100 states. After the minimization each SFA was 32% smaller (number of states) than the original SFA.

The following is a typical regex from the sample set:

```
[NS]\d{1,}\(\:[0-5]\d\){2}\. {0,1}\d{0,}, [EW]\d{1,}\(\:[0-5]\d\){2}
```

The generated SFA uses 16 predicates (such as $\backslash d$)⁶ while there are only 7 minterms (such as $[:]$). For this regex, the determinized SFA has 47 states and the minimized SFA has 23 states.

Since the number of minterms does not blowup, is there a performance incentive in using Min_{SFA}^N in this context? The answer is yes since the total time used to minimize all 1700 SFAs was 20 seconds when using Min_{SFA}^H , and only 0.8 seconds when using Min_{SFA}^N , thus showing a 24x speedup.

We also measured the performance of the minimization algorithm implemented in the `brics.automaton` library (version 1.11-8) that uses *symbolic integer ranges* to represent Unicode characters, but does not implement them as a Boolean algebra (since ranges are not closed under complement and union). We observed that the `brics` implementation of Hopcroft’s algorithm is comparable to our implementation of Min_{SFA}^H , and for 95% of the regexes the running times of `brics` Hopcroft’s algorithm and Min_{SFA}^H were at most 5% apart. The trend of Min_{SFA}^H in Figure 8 is very similar to that of the Hopcroft’s implementation in the `brics` library. We decided not to include the latter in the plot for the following reasons: 1) `brics` is implemented in Java, while our implementations are all in C#, and 2) in `brics`, the code corresponding to the minterm computation is not part of the minimization algorithm, therefore the comparison would not be completely fair (especially for big instances).

⁵ <http://regexlib.com/>

⁶ We use standard regex character class notation for character predicates.

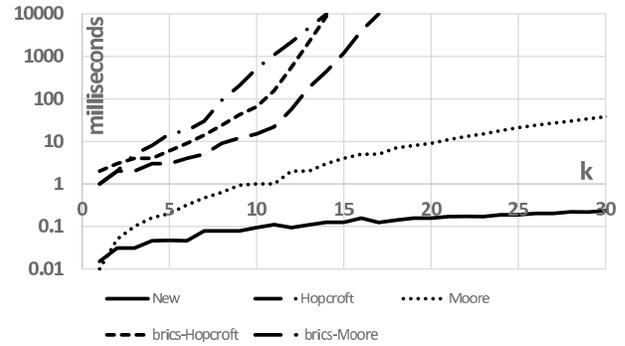


Figure 9. Running times on the M_k SFAs of Example 2. The y -axis is in log-scale.

6.3 Alphabet corner cases

The next experiment shows the importance of avoiding the minterm generation used by Min_{SFA}^H . We consider the SFAs M_k from Example 2, for values of k ranging between 1 and 31, over the alphabet 2^{BV32} of 32-bit bit-vectors. Figure 9 shows the running times for Min_{SFA}^N , Min_{SFA}^H , and Min_{SFA}^M .

As expected, the performance of Min_{SFA}^H degrades exponentially. Already for $k = 11$, minimizing the SFA using Min_{SFA}^H took 1 second due to minterm generation, while both Min_{SFA}^M and Min_{SFA}^N took less than 1ms. For Min_{SFA}^M the time increased from 1ms to .1sec with $k = 31$, while for Min_{SFA}^N the time remained below 1ms for all the values of k .

Similarly to the previous experiment, we also compared the running time of the sample set against Hopcroft’s and Moore’s minimization algorithm in the `brics.automaton` library (version 1.11-8). For this experiment, we decided to show the `brics` performance in order to appreciate how this particular example causes the number of ranges to grow exponentially, causing Hopcroft’s algorithm to be very slow. It is interesting to notice that the `brics` implementation of Moore’s algorithm is also slow. This is again due to the alphabet’s blowup caused by the ranges algebra.

6.4 Complex theories

We compare the performance of Min_{SFA}^N , Min_{SFA}^H and Min_{SFA}^M over a sample set of 220 randomly generated SFAs over the theory of pairs over $string \times int$. The guards of each SFAs are conjunctions of simple predicates over strings and integers. In the theory of strings we only allow comparison against a constant string, while for integers we generate unary predicates containing $+$, $-$, $<$, $=$ and integer constants. The set of generated SFAs is created as follows. We first generated a set S of 10 SFAs with at most 10 states. For each SFA $a \in S$ we also compute the complement, and for every pair $a, b \in S$ we compute the union, intersection, and difference of a and b .

Figure 10 shows the running time of each algorithm for different numbers of states. We first observe how the performance of Min_{SFA}^H quickly degrades when increasing the number of states. This is mainly due to the large number of minterms. Since the predicates are randomly generated, many overlaps are possible and the number of minterms grows quickly. Next, we can see how the performance of Min_{SFA}^M is affected by the increasing number of states. Finally, Min_{SFA}^N can quickly minimize SFAs with up to 96 states in less than 1.5 seconds. This experiment shows how Min_{SFA}^N is not affected by complex theories, while Min_{SFA}^H and Min_{SFA}^M are both impractical in this setting.

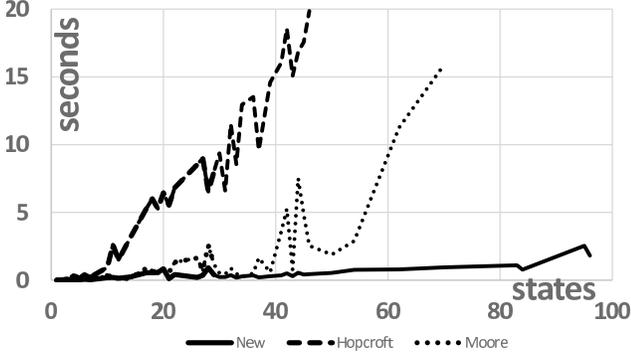


Figure 10. Running times for SFAs over the theory $String \times Int$ with 20 seconds timeout.

6.5 From Monadic Second-Order logic predicates to DFAs

The relation between regular languages and logic has been extensively investigated in the past [38]. Particularly, every regular language can be expressed as a formula in Monadic Second-Order logic (MSO) over strings, and vice versa. In the rest of the section we assume the alphabet to be a finite set Σ , say $\{a, b\}$. The following is an example of an MSO formula: $\phi = \exists x. \hat{a}(x)$ where \hat{a} is a unary relation symbol. A string $s \in \Sigma^*$ is a model of ϕ , iff there exists a *position* x in the string with label a . Transforming MSO formulas to automata gives us an algorithm for deciding satisfiability of MSO formulas.

The procedure for converting an MSO formula into a DFA inductively transforms each sub-formula into the corresponding DFA and then combines such DFAs using classical regular language operations. The complexity arises from the presence of free variables, that range over positions. For example, in the formula $\hat{a}(x)$, x occurs free. In order to represent such a language, the alphabet Σ is extended with one extra bit to $\Sigma \times \{0, 1\}$. The string $\langle a, 0 \rangle \langle b, 0 \rangle \langle a, 1 \rangle$ over the extended alphabet will then represent that the third position of the string aba is assigned to variable x . Following this intuition, every sub-formula ψ is compiled into a DFA over the alphabet $\Sigma \times \{0, 1\}^n$ where n is the number of quantified variables around ψ . For each existential quantification a projection is performed which leads to a non-deterministic automaton that must be determinized when a negation on it is performed. This means that each quantifier alternation might therefore lead to an exponential blowup, and in general the procedure has non-elementary complexity.

Despite the non-elementary complexity, practical algorithms that can translate non-trivial MSO formulas are presented in [22]. The tool implementing such algorithms is called Mona. The two key-features of such algorithms are:

1. determinizing and minimizing the intermediate DFA at every step in the transformation, and
2. using BDDs for representing the lifted bit-vector alphabets.

We implemented the same transformation using the BDD solver in our library, and compared the performance using different minimization algorithms.

Figure 11 shows the performance of the transformation for different MSO formulas. The running time for Mona are also shown. The four sub-figures depict the running time for the MSO to DFA transformation for the following formulas:

- a) $\exists x_1, \dots, x_k. x_1 < \dots < x_k$, for k between 2 and 40,
- b) $\exists x_1, \dots, x_k. x_1 < \dots < x_k \wedge a(x_1) \wedge \dots \wedge a(x_k)$, for k between 2 and 40,

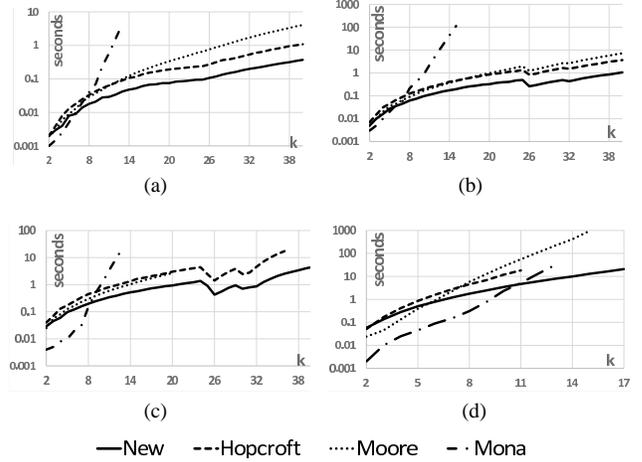


Figure 11. MSO to DFA running times. The y -axes are in log-scale.

- c) $(\exists x_1, \dots, x_k. x_1 < \dots < x_k \wedge a(x_1) \wedge \dots \wedge a(x_k)) \vee \exists y. c(y)$, for k between 2 and 40,
- d) $\exists x_1, \dots, x_k. f_k$ where $f_2 = (x_1 < x_2 \wedge a(x_1)) \vee c(x_1)$, and for $n > 2$, $f_n = f_{n-1} \wedge ((x_{n-1} < x_n \wedge a(x_{n-1})) \vee c(x_{n-1}))$, for k between 2 and 17.

All the values missing from the graphs are due to the algorithms running out of memory. The figure shows the following behaviors:

- for all of the four classes of formulas, the transformation based on Min_{SFA}^N is able to create the DFAs for higher values of k (number of nested variables) than those supported by Mona;
- for small instances Mona is faster than our implementation. However for higher values of k our implementation is faster, even when using Min_{SFA}^H or Min_{SFA}^M ;
- while Mona immediately shows an exponential behavior (very steep slope), we couldn't observe such trend in our implementation; and
- Min_{SFA}^N is faster than Min_{SFA}^H and Min_{SFA}^M , and it is also less memory intensive (it runs out of memory for higher k).

Even though such formulas are only a small benchmark, they are quite representative of the kind of inputs that cause Mona's transformation to be impractical. We believe that the performance improvement of our implementation with respect to Mona is primarily due to a better use of the BDD solver, rather than due to the minimization algorithm.

7. Applications

The development of the theory of symbolic automata, including use of minimization, is motivated by several concrete practical problems. Here we discuss three such applications. In each case we illustrate what kind of character theory we are working with, the role of minimization, and focus on the benefits of the symbolic representation.

7.1 Regex processing

Practical applications of regular expressions or *regexes* is ubiquitous. Practical regexes differ from schoolbook regular expressions in the following aspects:

- they support, besides non-regular features that go beyond capabilities of finite state automata representations, constructs such

as *bounded quantifiers* and *character classes*, that make them more appealing (more succinct) than their classical counterparts, and

- the size of the alphabet is 2^{16} due to the widely adopted UTF16 standard of Unicode characters.

As a somewhat unusual example, the regex $\hat{[\uFF10-\uFF19]} \$$ matches the set of digits in the so-called Wide Latin range of Unicode. We let the alphabet algebra be 2^{BV16} . Let the BDD β_w^7 represent all ASCII word characters as the set of character codes $\{‘0’, \dots, ‘9’, ‘A’, \dots, ‘Z’, ‘_’, ‘a’, \dots, ‘z’\}$. We write ‘0’ for the code 48, ‘a’ for the code 97, etc. Let also β_d^7 represent the set of all decimal digits $\{‘0’, \dots, ‘9’\}$ and let β_- represent underscore $\{‘_’\}$. By using the Boolean operations, e.g., $\beta_w^7 \wedge \neg(\beta_d^7 \vee \beta_-)$ represents the set of all upper- and lower-case ASCII letters. As a regex character class it is expressible as $[\d_-\x7F-\uFFFF]$.

Regexes are used in many different contexts. A common use of regexes is as a constraint language over strings for *checking* presence or absence of different patterns, e.g., for *security validation* of packet headers in network protocols. Another application, is the use of regexes for *generating* strings that match certain criteria, e.g., for *fuzz testing* applications. As another application consider the *password generation* problem based on constraints given in form of regexes. Here is a concrete set of constraints adopted in password generation:⁷

1. Length is k and characters are in visible ASCII range:
 $\hat{[\x21-\x7E]} \{k\} \$$
2. There are at least two letters: $[a-zA-Z] . * [a-zA-Z]$
3. There is at least one digit: $\backslash d$
4. There is at least one non-word character: $\backslash w$

The smallest feasible value for k is obviously 4. Consider SFAs A_1, A_2, A_3 , and A_4 for each case and let A be their product. For $k = 4$ the minimized version of A has 12 states, and minimizing A with any of Min_{SFA}^N , Min_{SFA}^H , or Min_{SFA}^M takes a few ms. When k is increased to 40,⁸ the number of states increases to 444 and minimizing using Min_{SFA}^N or Min_{SFA}^H takes respectively 15 and 90 ms, while Min_{SFA}^M becomes impractical, taking more than 43 sec (450x slower).

The minimal (canonical) form of A plays an important role here. Together with finiteness of the language accepted by A , minimality is used to guarantee a *uniform distribution* of the passwords generated from A . Uniformity is achieved by using the canonical form of the state graph together with the canonical form of the guards. While uniform sampling of $a \in [\varphi]$ is not possible for arbitrary alphabet theories, it is a well-known feature of BDDs.

7.2 Sanitizer analysis

Sanitizers are string transformation routines (special purpose encoders) that are extensively used in web applications, in particular as the first line of defense against cross site scripting (XSS) attacks. There are at least three, semantically very different, string sanitizers involved in a single web page: CssEncoder, UrlEncoder, and HtmlEncoder. A large class of sanitizers (including all the ones mentioned above) can be described and analyzed using *symbolic finite state transducers* (SFTs) [25]. SFAs are used in that context for certain operations over SFTs, for example for checking domain equivalence of SFTs [44]. Since several basic operations are quadratic

⁷ Recall the standard convention: a regex without the start-anchor $\hat{}$ matches any prefix and a regex without the end-anchor $\$$ matches any suffix.

⁸ Relevant values of k in the above scenario depend on the concrete context, but range between 4 and 128

in the number of states (product, union), minimization comes in handy as a tool for reducing the state space size,

In this setting we choose a different character theory. We use integer linear modular arithmetic (or bit-vector arithmetic of an SMT solver, in our case Z3 [19]) in place of BDDs. The main advantage of this choice is that it makes it possible to seamlessly combine the guards over characters with expressions over *yields*, i.e. the *symbolic outputs* of SFT moves. A concrete example of a yield is the following transformation that takes a character and encodes it as a sequence of five other characters:

$$f : \lambda x. [‘\&’, ‘\#’, ((x \div 10) \bmod 10) + 48), ((x \bmod 10) + 48), ‘;’, ‘;’]$$

In general, a yield denotes a function from \mathcal{Q}_A to \mathcal{Q}_A^* (or to \mathcal{Q}_B when the output alphabet \mathcal{B} is different from the input alphabet \mathcal{A}). In the yield above, for example, $f(‘a’)$ is the sequence $[‘\&’, ‘\#’, ‘9’, ‘7’, ‘;’, ‘;’]$ (or the string "&\#97;"). Given f as above, a typical SFT move ρ looks like:

$$\rho : q \xrightarrow{(\lambda x. 0 < x < 32) / \lambda x. f(x)} q$$

This specific rule happens to be an HtmlEncoder rule for encoding control characters in state q and remaining in that state. What is the connection with SFAs and the theory mentioned above? For analyzing the idempotence of an encoder with such rules, the encoder is sequentially composed with itself. As a result, this leads to an SFT with more complex guards and outputs (SFTs are closed under composition). When composing the move ρ with itself (i.e., roughly speaking, feeding the five output characters as its inputs again five times in a row), the guard of the composed rule will have sub-conditions such as $0 < (((x \div 10) \bmod 10) + 48) < 32$, which may involve nontrivial arithmetic operations (in this particular case the guard of the composed move will be infeasible). One of the operations performed during *idempotence* analysis is checking whether the original SFT and the composed one have the same *domain*. This reduces to *language equivalence* of SFAs for which the guards involve arithmetic operations of the above kind, that are not readily expressible using the earlier BDD based approach. Domain equivalence of two SFTs checks that both the SFTs accept/reject the same input sequences. Maintaining the SFAs minimal speeds up the equivalence check, and in general provides a better understanding of the structure of the domain language.

In general, the alphabet theory may be a Boolean combination of other decidable theories that are available for example in state-of-the-art SMT solvers. In the context of sanitizers, encoders, and decoders, the alphabet theory is a combination of lists, tuples, bit-vectors and integer linear arithmetic. Lists are used for example to represent composite characters or characters that represent “look-ahead” [17]. We demonstrated in the evaluation section that, when the alphabet theory is complex the algorithm Min_{SFA}^N outperforms both Min_{SFA}^H and Min_{SFA}^M by several orders of magnitude enabling analysis of bigger sanitizers/encoders.

7.3 Solving Monadic Second Order logic

We already anticipated in Section 6.5 how Monadic Second Order (MSO) logic predicates can be transformed into equivalent DFAs using a non-elementary algorithm. Such algorithm also provides a decision procedure for MSO.

We already discussed how several techniques have been introduced by the tool Mona [22] in order to make such transformation practical. Keeping the DFA minimal at any step in the transformation is one of the key techniques. We gave experimental evidence of the fact that the new algorithm presented in this paper, and in general the use of symbolic automata, can further move the barrier of solvable MSO formulas, in terms of both formula’s size (number of nested quantifiers) and running times.

8. Related work

DFA Minimization: The classical algorithms for DFA minimization have been studied and analyzed extensively in several different aspects. In particular, Moore’s algorithm is studied in [10] where it is shown that, by a clever change of data structures, its complexity can be reduced from $O(kn^2)$ to $O(kn \log n)$. The bound $O(kn \log n)$ has been shown to be tight for Hopcroft’s algorithm [6, 8]. Brzozowski [13] observed that a DFA can be minimized by reversing its transitions, then determinizing, then reversing its transitions again, and finally determinizing again. However, due to the determinization steps, this procedure is exponential and we decided not to consider it in this paper. Linear time algorithms have been studied for minimizing acyclic automata [32, 36] and automata containing only simple cycles [2]. The book chapter [7] provides an in-depth study of the state-of-the-art techniques for automata minimization, including the approaches mentioned above and several other ones. Watson [45] also provides an elegant classification of the classical minimization algorithms based on their structural properties.

In the case of DFAs, it matters whether the DFA to be minimized is partial (incomplete), because it may be useful to avoid completion of DFAs with sparse transition graphs. Minimization of partial DFAs is studied in [5, 39, 40]. In this setting the complexity of the algorithm depends on the number of transitions rather than on the size of the alphabet, however, in the case of DFAs these two quantities are generally related. In contrast, in a normalized SFA the number of transitions is independent of the alphabet size, and it is at most n^2 where n is the number of states. One concrete difference between the minimization algorithms of complete DFAs versus partial DFAs is that the initial value of the waiting set W (see Figure 3) for the partial case must contain both the sets F and F^c [40]. We believe that similar modifications may be applied to Min_{SFA}^N , even though we expect that in the case of SFAs the benefit of using partiality might not be as visible. In fact, a complete SFA has at most n transitions more than a partial one. Moreover, in the case of SFAs, there are different ways for completing an SFA, e.g., one can effectively restrict the alphabet to only those characters that are mentioned in the trimmed SFA (SFA without dead-end states) prior to completion.

Different Notions of Minimization: A notion of incremental minimization is investigated in [47]. An incremental minimization algorithm can be halted at any point and produce a partially minimal DFA which is equivalent to the starting DFA, and has less or equal number of states. If the algorithm runs till completion, a minimal DFA is computed. In this paper we did not address incremental computation, however it would be interesting to identify variants of the presented algorithms with such a property. A similar idea, based on intermediate computation results, is used for a modular minimization algorithm in [14]. An algorithm for building a minimal DFA accepting a given finite set of strings is presented in [15]. The same paper investigates the problem of “maintaining” a minimal DFA when applying modification such as node, edge, or string deletion and insertion to an already minimal DFA. This class of problems is called dynamic minimization.

A parallel version of Moore’s algorithm is presented in [37]. We are not aware of a parallel version of Hopcroft’s algorithm. We leave as an open problem identifying parallel algorithms corresponding to Min_{SFA}^H and Min_{SFA}^N .

A variant of minimization called hyper-minimization is investigated in [24]. Given an input DFA A , a DFA A' is hyper minimal with respect to A if it is the smallest DFA that differs from A only on a finite number of strings. In the case of symbolic automata, this definition doesn’t extend naturally due to potentially infinite alphabet. A less restrictive notion called k -minimization is studied

in [21]. Given an input DFA A , a DFA A' is k -minimal with respect to A if it is the smallest DFA that differs from A only on strings of length smaller or equal than k . This second restriction naturally extends to symbolic automata. A more general notion is that of minimization up to E -equivalence, where given a regular language E , the DFA A' is allowed to differ from A on a set of strings $L \subseteq E$ [23]. Extending the results of [21, 23] to symbolic automata is an interesting open research direction. Relationships between minimality and selection of final states are studied in [35] where *uniform* minimality is defined as minimality for any choice of final states.

Automata with Predicates: The concept of automata with predicates instead of concrete symbols was first mentioned in [46] and was first discussed in [41] in the context of natural language processing. A symbolic generalization of Moore’s algorithm was first discussed in [43]. To the best of our knowledge, no other minimization algorithms have been studied for SFAs. The MONA implementation [22] provides decision procedures for monadic second-order logic, and it relies on a highly-optimized BDD-based representation for automata which has seen extensive engineering effort [30]. Therefore, the use of BDDs in the context of automata is not new, but is used here as an example of a Boolean algebra that seems particularly well suited for working with Unicode alphabets.

Minimization of Other Automata: The problem of automata minimization has been investigated in several other settings. A new approach for minimizing *nondeterministic* automata and nondeterministic Büchi automata has been studied in [33]. The problem of minimizing weighted automata (automata computing functions from strings to a semi-ring) is studied in [20]. Classical minimization algorithms are used in this setting, and we hope that the results shown in this paper can extend to such domain. The minimization problem has also been studied for timed automata [11] and register automata [16]. These models are able to represent infinite domains, but not arbitrary theories. An orthogonal direction is to extend the techniques presented in this paper to minimization of symbolic *tree* automata [18].

Other Applications: The problem of learning of symbolic transducers has recently been studied in [12]. Classical automata learning is based on Angluin’s algorithm [4], which fundamentally relies on DFA minimality. An almost unexplored topic is whether such learning techniques can be extended to SFAs.

9. Conclusions

We presented three algorithms for minimizing Symbolic Finite Automata (SFAs). We first extended Moore’s algorithm and Hopcroft’s algorithm to the symbolic setting. We then show that, while in the classical setting Hopcroft’s algorithm is the fastest known minimization algorithm, in the presence of symbolic alphabets, it might incur in an exponential blowup. To address this issue, we introduced a new minimization algorithm that fully benefits from the symbolic representation of the alphabet and does not suffer from the exponential blowup. The new algorithm is a refinement of Hopcroft’s one in which splits are computed locally without having to consider the entire input alphabet. The new algorithm can also be adopted in the classical setting. Finally, we implemented all the algorithms and provided experimental evidence that the new minimization algorithm is the most efficient one in practice.

Acknowledgments. We thank Marco Almeida for helping us accessing his benchmark of randomly generated DFAs. Loris D’Antoni was supported by the NSF Expeditions in Computing award CCF 1138996, and this work was done as part of an internship at Microsoft Research. We also thank the anonymous reviewers for their insightful comments.

References

- [1] BRICS finite state automata utilities. <http://www.brics.dk/automaton/>.
- [2] J. Almeida and M. Zeitoun. Description and analysis of a bottom-up DFA minimization algorithm. *Information Processing Letters*, 107(2):52–59, 2008.
- [3] M. Almeida, N. Moreira, and R. Reis. On the performance of automata minimization algorithms. Technical Report DCC-2007-03, University of Porto, 2007.
- [4] D. Angluin. Learning regular sets from queries and counterexamples. *Inf. Comput.*, 75(2):87–106, 1987.
- [5] M.-P. Béal and M. Crochemore. Minimizing incomplete automata. In *Finite-State Methods and Natural Language Processing, 7th International Workshop*, pages 9–16, 2008.
- [6] J. Berstel, L. Boasson, and O. Carton. Hopcroft’s automaton minimization algorithm and Sturmian words. In *DMTCS’2008*, pages 355–366, 2008.
- [7] J. Berstel, L. Boasson, O. Carton, and I. Fagnot. Minimization of automata. To appear in *Handbook of Automata*, 2011.
- [8] J. Berstel and O. Carton. On the complexity of Hopcroft’s state minimization algorithm. In *CIAA’2004*, volume 3317, pages 35–44, 2004.
- [9] N. Bjørner, V. Ganesh, R. Michel, and M. Veanes. An SMT-LIB format for sequences and regular expressions. In P. Fontaine and A. Goel, editors, *SMT’12*, pages 76–86, 2012.
- [10] N. Blum. An $O(n \log n)$ implementation of the standard method for minimizing n -state finite automata. *Information Processing Letters*, 57:65–69, 1996.
- [11] M. Bojaczek and S. Lasota. A machine-independent characterization of timed languages. In A. Czumaj, K. Mehlhorn, A. Pitts, and R. Wattenhofer, editors, *Automata, Languages, and Programming*, volume 7392 of *LNCS*, pages 92–103. Springer Berlin Heidelberg, 2012.
- [12] M. Botinčan and D. Babić. Sigma*: symbolic learning of input-output specifications. In *POPL ’13*, pages 443–456, New York, NY, USA, 2013. ACM.
- [13] J. A. Brzozowski. Canonical regular expressions and minimal state graphs for definite events. In *Proc. Sympos. Math. Theory of Automata*, pages 529–561, New York, 1963.
- [14] D. Bustan. Modular minimization of deterministic finite-state machines. In *In 6th International Workshop on Formal Methods for Industrial Critical Systems*, pages 163–178, 2001.
- [15] R. C. Carrasco and M. L. Forcada. Incremental construction and maintenance of minimal finite-state automata. *Comput. Linguist.*, 28(2):207–216, June 2002.
- [16] S. Cassel, B. Jonsson, F. Howar, and B. Steffen. A succinct canonical register automaton model for data domains with binary relations. In S. Chakraborty and M. Mukund, editors, *ATVA 2012*, volume 7561 of *LNCS*, pages 57–71. Springer, 2012.
- [17] L. D’Antoni and M. Veanes. Equivalence of extended symbolic finite transducers. In N. Sharygina and H. Veith, editors, *CAV 2013*, volume 8044 of *LNCS*, pages 624–639. Springer, 2013.
- [18] L. D’Antoni, M. Veanes, B. Livshits, and D. Molnar. Fast: A transducer-based language for tree manipulation. Technical Report MSR-TR-2012-123, Microsoft Research, November 2012.
- [19] L. de Moura and N. Bjørner. Z3: An Efficient SMT Solver. In *TACAS’08*, LNCS. Springer, 2008.
- [20] M. Droste, W. Kuich, and H. Vogler. *Handbook of Weighted Automata*. Springer Publishing Company, Incorporated, 1st edition, 2009.
- [21] P. Gawrychowski, A. Jež, and A. Maletti. On minimising automata with errors. In *MFCS’11*, pages 327–338, Berlin, Heidelberg, 2011. Springer.
- [22] J. Henriksen, J. Jensen, M. Jørgensen, N. Klarlund, B. Paige, T. Rauhe, and A. Sandholm. Mona: Monadic second-order logic in practice. In *TACAS ’95*, volume 1019 of *LNCS*. Springer, 1995.
- [23] M. Holzer and S. Jakobi. From equivalence to almost-equivalence, and beyond – minimizing automata with errors. In H.-C. Yen and O. Ibarra, editors, *Developments in Language Theory*, volume 7410 of *LNCS*, pages 190–201. Springer, 2012.
- [24] M. Holzer and A. Maletti. An $n \log n$ algorithm for hyper-minimizing a (minimized) deterministic automaton. *Theor. Comput. Sci.*, 411(38–39):3404–3413, 2010.
- [25] P. Hooimeijer, B. Livshits, D. Molnar, P. Saxena, and M. Veanes. Fast and precise sanitizer analysis with Bek. In *USENIX Security*, August 2011.
- [26] P. Hooimeijer and M. Veanes. An evaluation of automata algorithms for string analysis. In *VMCAI’11*, volume 6538 of *LNCS*, pages 248–262. Springer, 2011.
- [27] J. Hopcroft. An $n \log n$ algorithm for minimizing states in a finite automaton. In Z. Kohavi, editor, *Theory of machines and computations, Proc. Internat. Sympos., Technion, Haifa, 1971*, pages 189–196, New York, 1971. Academic Press.
- [28] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison Wesley, 1979.
- [29] D. Huffman. The synthesis of sequential switching circuits. *Journal of the Franklin Institute*, 257(3–4):161–190, 275–303, 1954.
- [30] N. Klarlund, A. Møller, and M. I. Schwartzbach. MONA implementation secrets. *International Journal of Foundations of Computer Science*, 13(4):571–586, 2002.
- [31] T. Knuutila. Re-describing an algorithm by Hopcroft. *Theor. Comput. Sci.*, 250(1-2):333–363, 2001.
- [32] S. L. Krivol. Algorithms for minimization of finite acyclic automata and pattern matching in terms. *Cybernetics*, 27:324–331, 1991.
- [33] R. Mayr and L. Clemente. Advanced automata minimization. In *POPL’13*, pages 63–74, 2013.
- [34] E. F. Moore. Gedanken-experiments on sequential machines. *Automata studies, Annals of mathematics studies*, (34):129–153, 1956.
- [35] A. Restivo and R. Vaglica. Some remarks on automata minimality. In G. Mauri and A. Leporati, editors, *DLT’2011*, volume 6795 of *LNCS*, pages 15–27. Springer, 2011.
- [36] D. Revuz. Minimisation of acyclic deterministic automata in linear time. *Theoret. Comput. Sci.*, 92:181–189, 1992.
- [37] A. Tewari, U. Srivastava, and P. Gupta. A parallel DFA minimization algorithm. In *HiPC 2002*, volume 2552 of *LNCS*, pages 34–40. Springer, 2002.
- [38] W. Thomas. Languages, automata, and logic. In *Handbook of Formal Languages*, pages 389–455. Springer, 1996.
- [39] A. Valmari. Fast brief practical DFA minimization. *Information Processing Letters*, 112:213–217, 2012.
- [40] A. Valmari and P. Lehtinen. Efficient minimization of DFAs with partial transition functions. In S. Albers and P. Weil, editors, *STACS 2008*, pages 645–656, Dagstuhl, 2008.
- [41] G. van Noord and D. Gerdemann. Finite state transducers with predicates and identities. *Grammars*, 4(3):263–286, 2001.
- [42] M. Veanes, N. Bjørner, and L. de Moura. Symbolic automata constraint solving. In C. Fermüller and A. Voronkov, editors, *LPAR-17*, volume 6397 of *LNCS/ARCoSS*, pages 640–654. Springer, 2010.
- [43] M. Veanes, P. de Halleux, and N. Tillmann. Rex: Symbolic regular expression explorer. In *ICST’10*, pages 498–507. IEEE, 2010.
- [44] M. Veanes, P. Hooimeijer, B. Livshits, D. Molnar, and N. Bjørner. Symbolic finite state transducers: Algorithms and applications. In *POPL’12*, pages 137–150, 2012.
- [45] B. W. Watson. A taxonomy of finite automata minimization algorithms. Computing Science Report 93/44, Eindhoven University of Technology, January 1995.
- [46] B. W. Watson. Implementing and using finite automata toolkits. In *Extended finite state models of language*, pages 19–36, New York, NY, USA, 1999. Cambridge University Press.
- [47] B. W. Watson and J. Daciuk. An efficient incremental DFA minimization algorithm. *Nat. Lang. Eng.*, 9(1):49–64, Mar. 2003.