# More on Deterministic and Nondeterministic Finite Cover Automata[☆],[☆☆]

Hermann Gruber[a], Markus Holzer[b], Sebastian Jakobi[b]

[a]*knowledgepark AG, Leonrodstr. 68,*
*80636 München, Germany*
[b]*Institut für Informatik, Universität Giessen,*
*Arndtstr. 2, 35392 Giessen, Germany*

## Abstract

Finite languages are an important sub-regular language family, which were intensively studied during the last two decades in particular from a descriptional complexity perspective. An important contribution to the theory of finite languages are the deterministic and the recently introduced nondeterministic finite *cover* automata (DFCAs and NFCAs, respectively) as an alternative representation of finite languages by ordinary finite automata. We compare these two types of cover automata from a descriptional complexity point of view, showing that these devices have a lot in common with ordinary finite automata. In particular, we study how to adapt lower bound techniques for nondeterministic finite automata to NFCAs such as, e.g., the biclique edge cover technique, solving an open problem from the literature. Moreover, the trade-off of conversions between DFCAs and NFCAs as well as between finite cover automata and ordinary finite automata are investigated. Finally, we present some results on the average size of finite cover automata.

---

## 1. Introduction

If one tries to describe formal objects such as, e.g., Boolean functions, graphs, trees, languages, as compact as possible we are faced with the question, which representation to use. This quest for compact representations of formal objects dates back to the early beginnings of theoretical computer science. For instance, one can prove by a *simple* counting argument that most Boolean functions have exponential circuit complexity [1]. For other representations of Boolean functions than circuits, such as formulas, ordered binary decision diagrams, etc. a similar result applies. This incompressibility is inherent in almost all possible representations of formal objects.

When considering formal languages, automata are the preferred choice of representation. In particular, for regular languages and subfamilies one may use deterministic (DFAs) or nondeterministic finite automata (NFAs) or variants thereof to describe these languages. It is well known that these two formalisms are equivalent. The obvious way to obtain a DFA from a given NFA is by applying the *subset* or *power-set construction* [2]. This construction allows to show an upper bound of $2^n$ states in the DFA obtained from an $n$-state NFA, and this bound is known to be tight. For finite languages a slightly smaller bound on the determinization problem is given in [3]. Here the tight bound depends on the alphabet size $k$ and reads as $\Theta(k^{\frac{n}{1+\log_2 k}})$. Thus, for a two-letter input alphabet $\Theta(2^{\frac{n}{2}})$ states are sufficient and necessary in the worst case for a DFA to accept a language specified by an $n$-state NFA. There are a lot of other results known for finite automata accepting finite languages such as, e.g., the maximal number of states of the minimal DFA accepting a subset of $\Sigma^\ell$ or $\Sigma^{\leq \ell}$ [4, 5], or the average case size of DFAs and NFAs w.r.t. the number of states and transitions accepting a subset of $\Sigma^\ell$ or $\Sigma^{\leq \ell}$ [6].

Since regular languages and finite automata are widely used in applications, and most of them use actually finite languages only, it is worth considering further representations for finite languages that may be more compact, but still bare nice handling in applications. Such a representation is based on finite automata and is known as finite cover automata. The idea is quite simple, namely a finite cover automaton $A$ of a finite language $L \subseteq \Sigma^*$ is a fi-

nite automaton that accepts all words in $L$ and possibly other words that are longer than any word in $L$. Formally, this reads as $L = L(A) \cap \Sigma^{\leq \ell}$, where $\ell$ is the length of the longest word(s) in $L$; then we say that $A$ *covers* the finite language $L$. Originally deterministic finite cover automata (DFCAs) were introduced in [7], where an efficient minimization algorithm for these devices was given. Further results on important aspects of DFCAs can be found in, e.g., [7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17]. Recently, DFCAs were generalized to nondeterministic finite cover automata (NFCAs) in [18] and it was shown that they can even give a more compact representation of finite languages than both NFAs and DFCAs. To our knowledge this was the first systematic study on this subject, although it has been suggested already earlier in a survey paper on cover automata [19].

We further develop the theory of finite cover automata in this paper. At first we introduce the necessary definitions in the next section. Then we briefly recall what is known on lower bound techniques for both types of finite cover automata. In particular, we first reconsider the fooling set techniques known for nondeterministic finite automata (NFAs) and secondly we show how to alter the biclique edge cover technique from [20] to make it applicable for NFCAs, too. This positively answers a question stated in [18], whether the biclique edge cover technique can be used at all to prove lower bounds for NFCAs. As a byproduct we develop a lower bound method for $E$-equivalent NFAs. This concept was recently introduced in [21]. Two languages are $E$-equivalent if their symmetric difference lies in the so called *error language $E$*. Thus, $E$-equivalence is a generalization of ordinary equivalence and also of cover-automata. In particular, setting $E = \Sigma^{>\ell}$, thus not taking care of words that are too long, we are back to covering languages and cover automata. Section 4 is devoted to conversions between finite automata and finite cover automata. First we provide a large family of languages where cover state complexity meets ordinary state complexity (up to one state for deterministic devices). Hence, for the conversions from finite automata to finite cover automata not much state savings are possible. For the opposite direction we show that an $n$-state finite cover automaton for a language of order $\ell$ can be converted to an equivalent finite automaton with about $n \cdot \ell$ states; the exact bounds are shown to be tight for all $n$ and $\ell$. In particular, this shows that roughly speaking the number of states of a finite cover automaton is *at least* an $\ell$th fraction of the state size of the equivalent finite automaton. Then we take a closer look on determinizing NFCAs by the well known power-set construction. We show that here the state blow-up heavily

depends on the order $\ell$ of the finite language represented by the NFCA. When the order is large enough, we get a tight exponential blow-up of $2^n$, just as in the case of ordinary finite automata. We give a range of conditions that imply sub-exponential, polynomial, and even linear determinization blow-ups. These results are presented in Section 5. In the penultimate section, we perform average case comparisons of the descriptional complexity of finite cover automata. For ordinary finite automata this was already done in, e.g., [6], where it was shown that almost all DFAs accepting finite languages of order $\ell$ over a binary input alphabet have state complexity $\Theta(2^\ell/\ell)$, while NFAs are shown to perform better, namely the nondeterministic state complexity is in $\Theta(\sqrt{2^\ell})$. Interestingly, in both cases the aforementioned bounds are asymptotically like in the worst case. For finite cover automata exactly the same picture as for ordinary finite automata emerges. Finally, we summarize our results in the conclusions section and state some open problems for future research.

## 2. Preliminaries

We recall some definitions on finite automata as contained in [22]. A *nondeterministic finite automaton* (NFA) is a quintuple $A = (Q, \Sigma, \delta, q_0, F)$, where $Q$ is the finite set of *states*, $\Sigma$ is the finite set of *input symbols*, $q_0 \in Q$ is the *initial state*, $F \subseteq Q$ is the set of *accepting states*, and $\delta\colon Q \times \Sigma \to 2^Q$ is the *transition function*. The *language accepted* by the NFA $A$ is defined as

$$L(A) = \{\, w \in \Sigma^* \mid \delta(q_0, w) \cap F \neq \emptyset \,\},$$

where the transition function is recursively extended to $\delta\colon Q \times \Sigma^* \to 2^Q$. An NFA is *deterministic* (DFA), if and only if $|\delta(q, a)| = 1$, for every $q \in Q$ and $a \in \Sigma$. In this case we simply write $\delta(q, a) = p$ instead of $\delta(q, a) = \{p\}$, assuming that the transition function $\delta\colon Q \times \Sigma \to Q$ is a *total* mapping. Two automata $A$ and $B$ are *equivalent* if they accept the same language, that is, $L(A) = L(B)$. An NFA (DFA, respectively) $A$ is *minimal* if any equivalent NFA (DFA, respectively) needs at least as many states as $A$. It is a well known fact that minimal DFAs are unique up to isomorphism, while minimal NFAs are *not* necessarily unique in general. Let $\mathsf{nsc}(L)$ ($\mathsf{sc}(L)$, respectively) refer to the number of states a minimal NFA (DFA, respectively) needs to accept the language $L$. By definition and the seminal result in [2] we have $\mathsf{nsc}(L) \leq \mathsf{sc}(L) \leq 2^{\mathsf{nsc}(L)}$, if $L$ is a language accepted by a finite automaton.

Proving lower bounds for $\mathsf{nsc}(L)$ can be done by applying, e.g., the extended fooling set technique, which reads as follows [23]:

**Theorem 1.** *Let $L \subseteq \Sigma^*$ be a regular language and suppose there exists a set of pairs $S = \{ (x_i, y_i) \mid 1 \leq i \leq n \}$ such that*

1. *$x_i y_i \in L$, for $1 \leq i \leq n$, and*
2. *$i \neq j$ implies $x_i y_j \notin L$ or $x_j y_i \notin L$, for $1 \leq i, j \leq n$.*

*Then any nondeterministic finite automaton for $L$ has at least $n$ states, i.e., $n \leq \mathsf{nsc}(L)$. Here $S$ is called an* extended fooling set *for $L$.*

A non-empty finite language $L \subseteq \Sigma^*$ is said to be of *order* $\ell$, if $\ell$ is the length of the longest word(s) in the set $L$, i.e., $L \subseteq \Sigma^{\leq \ell}$, where $\Sigma^{\leq \ell}$ refers to the set $\{ w \in \Sigma^* \mid |w| \leq \ell \}$, where $|w|$ denotes the length of the word $w$. In particular, the length of the empty word $\lambda$ is zero.

A *deterministic finite cover automaton* (DFCA) for a language $L \subseteq \Sigma^*$ of order $\ell$ is a DFA $A$ such that $L(A) \cap \Sigma^{\leq \ell} = L$; these devices were introduced in [7]. This definition naturally carries over to NFAs, hence leading to *nondeterministic finite cover automata* (NFCA), which were recently introduced in [18]. Two cover automata $A$ and $B$ are *equivalent* if they cover the same finite language $L \subseteq \Sigma^*$, that is, $L(A) \cap \Sigma^{\leq \ell} = L(B) \cap \Sigma^{\leq \ell}$, where $\ell$ is the order of $L$. A DFCA (NFCA, respectively) $A$ for a finite language $L$ is *minimal* if any equivalent automaton of same type needs at least as many states as $A$. Let $\mathsf{ncsc}(L)$ ($\mathsf{csc}(L)$, respectively) refer to the number of states a minimal NFCA (DFCA, respectively) needs to accept the finite language $L$. By definition we have $\mathsf{ncsc}(L) \leq \mathsf{csc}(L)$, if $L$ is a finite language. Moreover, since any cover automaton can be at most as large as an ordinary finite automaton of the same type for a finite language $L$, we have $\mathsf{csc}(L) \leq \mathsf{sc}(L)$ as well as $\mathsf{ncsc}(L) \leq \mathsf{nsc}(L)$. A useful tool for the study of minimal DFCAs is the notion the similarity relation, which plays a similar role as the Myhill-Nerode relation[1] in case of DFAs. For a finite language $L \subseteq \Sigma^*$ of order $\ell$ the similarity relation $\approx_L$ on words is defined as follows: for $u, v \in \Sigma^*$ let $u \approx_L v$ if and only if we have $uw \in L \iff vw \in L$, for all $w \in \Sigma^*$, whenever $|uw| \leq \ell$ and $|vw| \leq \ell$. Observe, that $\approx_L$ is not a equivalence relation in general. The relation $\approx_L$ can also be defined for

---

[1]For a language $L \subseteq \Sigma^*$ define the Myhill-Nerode relation $\equiv_L$ on words as follows: for $u, v \in \Sigma^*$ let $u \equiv_L v$ if and only if $uw \in L \iff vw \in L$, for all $w \in \Sigma^*$.

states of a DFCA $A = (Q, \Sigma, \delta, q_0, F)$. Two states $p$ and $q$ are *similar*, denoted by $p \approx_L q$, if $\delta(p, w) \in F \iff \delta(q, w) \in F$ holds for all $w \in \Sigma^{\leq \ell - m}$, with $m = \max(lev_A(p), lev_A(q))$—here $lev_A(p) = \min\{\, |u| \mid \delta(q_0, u) = p \,\}$. If $p \not\approx_L q$ then $p$ and $q$ are *dissimilar*. It is known [7] that a DFCA is minimal if all its states are pairwise dissimilar.

## 3. Lower Bound Techniques For Cover Automata

It is well known that the minimal DFA for a language $L$ is isomorphic to the DFA induced by the Myhill-Nerode equivalence relation for $L$. Hence, the number of states of the minimal DFA accepting the language $L \subseteq \Sigma^*$ equals the index, i.e., the cardinality of the set of equivalence classes of the Myhill-Nerode equivalence relation. On the other hand, the problem to estimate the necessary number of states of a minimal NFA accepting a given regular language is complicated, and stated as open in [24] and [25]. Several authors have introduced methods for proving lower bounds, using communication complexity methods for proving such lower bounds, see, for example, [23, 26, 27]. The results of [27] have been generalized by the advent of so-called multi-party nondeterministic message complexity [28]. The most widely used lower bound techniques for NFAs are the so-called *fooling set* techniques—the fooling set technique [26] and the extended fooling set method [23]—and the biclique edge cover technique [20].

A similar situation applies for finite cover automata. Although DFCA are not unique in general, it was shown in [7] that one can still estimate the exact number of states of a minimal DFCA by determining the size of the canonical maximal dissimilar sequence for the finite language $L$, which is based on the similarity relation for $L$. Recently, in [18] both fooling set methods were adapted to work for NFCAs as well. In fact, it was shown that the fooling set lower bound techniques for state complexity of NFAs mentioned above, can be modified to prove lower bounds for minimal NFCAs when the order of the language is explicitly considered in the lower bound statement. Here we first reconsider the fooling set techniques and then show how to modify yet another lower bound method, the biclique edge cover technique of [20], to work with NFCAs. Whether this latter technique can be generalized to NFCAs was stated as an open problem in [18].

In [18] it was argued that there is no doubt that any fooling set type technique used to prove a lower bound for NFCAs must explicitly consider the order of the language under consideration. In this vein, both fooling

set techniques were adapted. In fact, we show that the original fooling set technique of [26] (not the extended version of [23]) already gives a lower bound for NFCAs without modifying the technique to explicitly deal with the order of the language under consideration.

**Theorem 2.** *Let $L \subseteq \Sigma^*$ be a finite language and suppose there exists a set of pairs $S = \{ (x_i, y_i) \mid 1 \leq i \leq n \}$ such that*

1. *$x_i y_i \in L$, for $1 \leq i \leq n$, and*
2. *$x_i y_j \notin L$, for $1 \leq i, j \leq n$, and $i \neq j$.*

*Then any nondeterministic finite* cover *automaton for $L$ has at least $n$ states, i.e., $n \leq \mathsf{ncsc}(L)$. Here $S$ is called a fooling set for $L$.*

PROOF. Let $A = (Q, \Sigma, \delta, q_0, F)$ be any NFCA covering the finite language $L$ of order $\ell$. Since $x_i y_i \in L$, there is a state $q_i$ in $Q$ such that $q_i \in \delta(q_0, x_i)$ and $\delta(q_i, y_i) \cap F \neq \emptyset$. Assume that a fixed choice of $q_i$ has been made for any $i$ with $1 \leq i \leq n$. We prove that $q_i \neq q_j$ for $i \neq j$. For the sake of a contradiction assume that $q_i = q_j$ for some $i \neq j$. We consider two cases: (i) At least one of $|x_i| + |y_j|$ and $|x_j| + |y_i|$ is at most $\ell$. Without loss of generality we assume that $|x_i| + |y_j| \leq \ell$. Then clearly the automaton $A$ accepts the word $x_i y_j$, which has length at most $\ell$, but is not in $L$, contradicting the assumption that $A$ covers $L$. (ii) Both of $|x_i| + |y_j|$ and $|x_j| + |y_i|$ exceed $\ell$. We show that this case cannot apply. Since the words $x_i y_i$ and $x_j y_j$ belong to $L$, we have $|x_i y_i| \leq \ell$ and $|x_j y_j| \leq \ell$, and therefore $|x_i y_i x_j y_j| \leq 2\ell$. Now if both $|x_i y_j| > \ell$ and $|x_j y_i| > \ell$, then we get $|x_i y_i x_j y_j| > 2\ell$, a contradiction.□

In contrast the more powerful extended fooling set technique presented in [23] does not work as a lower bound technique for NFCAs as the following example illustrates, and therefore the modification of this technique presented in [18] is the right generalization.

**Example 1.** Consider the unary finite language $L = \{a\}^{\leq \ell}$, for $\ell \geq 1$. Clearly, this language can be covered by an NFCA with a single state. However, the set $S = \{ (a^i, a^{\ell - i}) \mid 0 \leq i \leq \ell \}$ is an extended fooling set for $L$, proving a lower bound of $\ell + 1$ on the nondeterministic state complexity of $L$. □

In the remainder of this subsection we turn our attention to the biclique edge cover technique from [20].

A central role in this technique plays the notion of the *bipartite dimension* $\dim(G)$ of a bipartite graph $G$, which is the minimum number of bicliques in $G$ needed to cover all edges of $G$. For a trivial upper bound, observe that each bipartite graph $G = (X, Y, E)$ can be covered by stars[2], so $d(G) \leq \min(|X|, |Y|)$. But note that each edge may be covered by one or more bicliques according to the definition of bipartite dimension. A nontrivial example is the crown graph $K_{n,n}^-$, which is obtained from a complete bipartite graph $K_{n,n}$ by removing a perfect matching. Using a computer program, one can find that $d(K_{5,5}^-) = 4$. The bipartite dimension $d(K_{n,n}^-)$ in general was determined in [29, Corollaries 1,2] to be equal to $\sigma(n)$, where $\sigma(n) = \min\left\{ k \mid n \leq \binom{k}{\lfloor k/2 \rfloor} \right\}$. Computing the bipartite dimension of a bipartite graph is an **NP**-complete problem [30].

Now the biclique edge cover technique asserts that for a regular language $L$, the bipartite dimension of the graph $G = (X, Y, E)$ with $X = Y = L$ and $E = \{ (x, y) \in X \times Y \mid xy \in L \}$ is a lower bound for the number of states of every NFA accepting $L$; see [20] for a proof.

The following example shows that the biclique edge cover technique cannot be applied to NFCAs without modification.

**Example 2.** Let $\ell \geq 1$ and consider the finite language $L = \{a\}^{\leq \ell}$. Clearly the single-state DFA accepting for the language $\{a\}^*$ is a cover automaton for $L$, hence we have $\mathsf{ncsc}(L) = 1$. However, the bipartite dimension of the graph $G = (X, Y, E)$, with $X = Y = L$ and $E = \{ (x, y) \in X \times Y \mid xy \in L \}$, is $\ell + 1 > 1$. This can be seen as follows. Notice that $(a^i, a^j) \in E$ if and only if $i + j \leq \ell$. In particular, for $0 \leq i \leq \ell$, the edge $e_i = (a^i, a^{\ell-i})$ belongs to $E$. Therefore, every such $e_i$ has to be covered by some biclique $H_i = (X_i, Y_i, E_i)$ with $a^i \in X_i$, $a^{\ell-i} \in Y_i$, and $E_i = X_i \times Y_i$. Now we see that distinct edges $e_i$ and $e_j$ must be covered by distinct bicliques, that is, $H_i \neq H_j$, for $1 \leq i, j \leq \ell$, with $i \neq j$: if $H_i = H_j$ then we have $a^i, a^j \in X_i$ and $a^{\ell-i}, a^{\ell-j} \in Y_i$, and since $H_i$ is a biclique, its set of edges $E_i$ contains both $(a^i, a^{\ell-j})$ and $(a^j, a^{\ell-i})$. But since $i \neq j$, either $i + \ell - j > \ell$ or $j + \ell - i > \ell$, which means that one of the two edges does not belong to $E$—a contradiction to $H_0, H_1, \ldots H_\ell$ being a biclique edge cover. This shows that the bipartite dimension of $G$ is at least $\ell + 1$. Equality is witnessed by the bicliques $H_i = (X_i, Y_i, E_i)$ with

---

[2]A *star* is a bipartite graph where all other vertices are adjacent to a single vertex, and all other vertices are mutually not adjacent.

$X_i = \{a^i\}$, $Y_i = \{a\}^{\leq \ell - i}$, and $E_i = X_i \times Y_i$, for $0 \leq i \leq \ell$. $\qquad\square$

In the following we want to generalize the biclique edge cover technique so that it can also be used to prove lower bounds for the size of NFCAs. In fact, we present a generalization that can be used even for the more general notion of $E$-equivalent automata, which was recently introduced in [21]. In order to avoid confusion with the set of edges of a graph, we use here the term $D$-equivalence instead of $E$-equivalence. Let $D \subseteq \Sigma^*$ be some language, the so called *error language*. Two languages $L$ and $L'$ over the alphabet $\Sigma$ are called $D$-*equivalent* if they differ only on elements from the error language $D$, that is, if

$$(L \setminus L') \cup (L' \setminus L) \subseteq D.$$

In this case we write $L \sim_D L'$. Similarly, two automata $A$ and $B$ are $D$-equivalent, if $L(A) \sim_D L(B)$. The connection between $D$-equivalence and cover automata is as follows. Assume $L \subseteq \Sigma^{\leq \ell}$ is some finite language of order $\ell$. Then a language $L' \subseteq \Sigma^*$ is a cover language for $L$ if and only if $L \sim_D L'$, for the error language $D = \Sigma^{> \ell}$. In other words, any two cover languages $L'$ and $L''$ for a finite language of order $\ell$ are $D$-equivalent, for $D = \Sigma^{> \ell}$.

We now come to our generalization of the biclique edge cover technique. In the original technique we have to find bicliques $H_i = (X_i, Y_i, E_i)$ with $1 \leq i \leq k$, for some $k$, of a bipartite graph $G = (X, Y, E)$, such that $E = \bigcup_{i=1}^{k} E_i$. In our generalization, we use two sets of edges in the bipartite graph $G$, namely a set $\underline{E}$ of edges that *must* be covered, and a set $\overline{E}$, with $\underline{E} \subseteq \overline{E}$, of edges that *may* be covered by bicliques. We use the notation $G = (X, Y, \underline{E}, \overline{E})$ to denote such a bipartite graph. Now an $(\underline{E}, \overline{E})$-*approximation* of $G$ is a collection of bicliques $H_i = (X_i, Y_i, E_i)$ of $G$, with $1 \leq i \leq k$ for some $k$, such that

$$\underline{E} \subseteq \bigcup_{i=1}^{k} E_i \subseteq \overline{E}.$$

The $(\underline{E}, \overline{E})$-*dimension* of $G$, denoted by $\dim^*(G)$, is defined as the minimal number of bicliques that constitute an $(\underline{E}, \overline{E})$-*approximation* of $G$.

Now we are ready to present our lower bound technique for $D$-equivalent automata. Notice that the sets $\underline{E}$ and $\overline{E}$ of edges of graph $G$ in the following theorem depend on the given language $L$ and error set $D$ by definition.

**Theorem 3.** *Let $L$ and $D$ be languages over some alphabet $\Sigma$. Moreover, let $X, Y \subseteq \Sigma^*$ and $G = (X, Y, \underline{E}, \overline{E})$, with $\underline{E} = \{ (x, y) \in X \times Y \mid xy \in L \setminus D \}$ and $\overline{E} = \{ (x, y) \in X \times Y \mid xy \in L \cup D \}$. Then the number of states of any nondeterministic finite automaton $A$, with $L(A) \sim_D L$, is at least $\dim^*(G)$.*

PROOF. We use a similar argumentation as in [20]. Let $A = (Q, \Sigma, \delta, q_0, F)$ be a nondeterministic finite automaton with $L(A) \sim_D L$. It suffices to show that there exists an $(\underline{E}, \overline{E})$-approximation of $G$ by $n = |Q|$ bicliques. For each state $q \in Q$ we define the biclique $H_q = (X_q, Y_q, E_q)$ with

$$X_q = \{ x \in X \mid q \in \delta(q_0, x) \}, \qquad Y_q = \{ y \in Y \mid \delta(q, y) \cap F \neq \emptyset \},$$

and $E_q = X_q \times Y_q$. Clearly $H_q$ is a biclique of $G$. To show that the bicliques form an $(\underline{E}, \overline{E})$-*approximation*, we have to prove that $\bigcup_{q \in Q} E_q$ contains $\underline{E}$, and is itself contained in $\overline{E}$.

Let $(x, y) \in \underline{E}$, which means that the word $xy$ belongs to the language $L$ but not to the error language $D$. Because $L(A) \sim_D L$, the word $xy$ must also belong to $L(A)$. Hence, there must be some state $q \in Q$ such that $q \in \delta(q_0, x)$ and $\delta(q, y) \cap F \neq \emptyset$, which implies $(x, y) \in E_q$. This proves $\underline{E} \subseteq \bigcup_{q \in Q} E_q$.

Finally let $(x, y) \in E_q$, for some state $q \in Q$. This means that $q \in \delta(q_0, x)$ and $\delta(q, y) \cap F \neq \emptyset$, so $xy \in L(A)$. Since $L(A) \sim_D L$, we obtain $xy \in L \cup D$, which in turn implies $(x, y) \in \overline{E}$. This concludes our proof. □

Notice that Theorem 3 yields the original biclique edge cover technique when choosing the error language $D = \emptyset$, that is, when considering the special case of classical language equivalence. Moreover, with the error language $D = \Sigma^{>\ell}$ we obtain the following technique for proving lower bounds on the state complexity of nondeterministic cover automata for finite languages of order $\ell$.

**Corollary 4.** *Let $L \subseteq \Sigma^*$ be some finite language of order $\ell$. Moreover, let $X, Y \subseteq \Sigma^*$ and $G = (X, Y, \underline{E}, \overline{E})$, with $\underline{E} = \{ (x, y) \in X \times Y \mid xy \in L \}$ and $\overline{E} = \{ (x, y) \in X \times Y \mid xy \in L \cup \Sigma^{>\ell}, \}$. Then the number of states of any nondeterministic finite cover automaton for $L$ is at least $\dim^*(G)$, that is, $\dim^*(G) \leq \mathsf{ncsc}(L)$.* □

## 4. Conversions Between Finite Automata and Cover Automata

In this section we compare the descriptional complexity of finite automata and cover automata, by studying the cost of conversions between these models. We consider nondeterministic as well as deterministic automata.

*4.1. From Finite Automata to Cover Automata*

Clearly, a finite automaton for a finite language $L$ is also a cover automaton for that language. So the bounds $\mathsf{ncsc}(L) \leq \mathsf{nsc}(L)$ and $\mathsf{csc}(L) \leq \mathsf{sc}(L)$ are obvious. However, the question is whether these bounds are tight in the following sense: does there exist, for every integer $n \geq 1$, a regular language $L_n$ that is accepted by a DFA (NFA, respectively) with $n$ states such that the minimal DFCA (NFCA, respectively) needs $n$ states, too? The next result answers this question in the affirmative for nondeterministic automata, while for deterministic devices the bound is off by one.

In the proof of the result we use the following notion: a state $q$ of an NFCA $A = (Q, \Sigma, \delta, q_0, F)$ for a finite language $L \subseteq \Sigma^{\leq \ell}$ is *productive* if there are words $u, v \in \Sigma^*$, with $|uv| \leq \ell$, such that $q \in \delta(q_0, u)$ and $\delta(q, v) \cap F \neq \emptyset$. Of course, states that are not productive can be safely removed without changing the accepted language.

**Theorem 5.** *If $L$ is a finite language with all words having the same length $\ell$, then $\mathsf{ncsc}(L) = \mathsf{nsc}(L)$ and $\mathsf{csc}(L) = \mathsf{sc}(L) - 1$.*

PROOF. Let $A = (Q, \Sigma, \delta, q_0, F)$ be a minimal NFCA covering $L$. Let $[\ell]$ denote the set $\{0, 1, \ldots, \ell\}$ and define the NFA

$$B = (Q \times [\ell], \Sigma, \delta', (q_0, 0), F \times [\ell]),$$

where the transition function satisfies $\delta'((q, i), a) = \{ (p, i+1) \mid p \in \delta(q, a) \}$, for $q \in Q$ and $i \in [\ell - 1]$. Clearly $L(B) = L$. Recall that $\mathsf{ncsc}(L) \leq \mathsf{nsc}(L)$, so in order to prove the statement $\mathsf{ncsc}(L) = \mathsf{nsc}(L)$ it is sufficient to show that the number of productive states in the NFA $B$ is at most $\mathsf{ncsc}(L) = |Q|$.

Assume that the number of productive states in $Q \times [\ell]$ is greater than $|Q|$. Then there are integers $i, j \in [\ell]$, with $i > j$, and a state $q \in Q$ of $A$ such that $(q, i)$ and $(q, j)$ are productive in $B$. This means that there are words $u_i$ and $v_i$, with $|u_i v_i| \leq \ell$, such that $(q, i) \in \delta'((q_0, 0), u_i)$ and $\delta'((q, i), v_i)$ contains a state $(f_i, k_i) \in F \times [\ell]$; and similarly there are words $u_j$ and $v_j$ with analogous conditions. In fact, from the construction of $B$ and the fact that all words in $L$ have length $\ell$, we conclude $|u_i| = i$ and $|v_i| = \ell - i$, and similarly $|u_j| = j$ and $|v_j| = \ell - j$. But then $B$ also accepts the word $u_j v_i$ with $|u_j v_i| < \ell$ since

$$\delta'((q_0, 0), u_j v_i) \supseteq \delta'((q, j), v_i) \ni (f_i, j + \ell - i).$$

This is a contradiction to $L \subseteq \Sigma^\ell$, so the number of productive states in $Q \times [\ell]$ cannot be greater than $|Q|$.

Now let us prove the statement $\mathsf{csc}(L) = \mathsf{sc}(L) - 1$. From a minimal DFCA $A$ we can construct a DFA $B$ similar as described above, but since the transition function of a deterministic machines is total, the constructed DFA $B$ needs an additional non-accepting sink state which is reached after reading more than $\ell$ input symbols. As before we can show that the number of productive states in $B$ is at most $|Q|$, so $B$ needs most $|Q| + 1$ states. This shows $\mathsf{csc}(L) \geq \mathsf{sc}(L) - 1$. To see that also $\mathsf{csc}(L) \leq \mathsf{sc}(L) - 1$ we argue as follows. Notice that every DFA for a language $L \subseteq \Sigma^\ell$ has a non-accepting sink state. We can obtain an equivalent DFCA with one state less by deleting this sink state, and re-routing the incoming transitions to the initial state of the DFA. Surely this allows the constructed automaton to accept new words, but the length of such words is at least $\ell + 1$ because the shortest word leading from the initial state to an accepting state has length $\ell$. $\qquad \square$

From Theorem 5 and the obvious upper bound $\mathsf{ncsc}(L) \leq \mathsf{nsc}(L)$ we obtain the following result. In fact, Theorem 5 provides the lower bound already by *unary* witness languages.

**Corollary 6.** *Let $n \geq 1$ and $L$ be a finite language accepted by a nondeterministic finite automaton with $n$ states. Then $n$ states are sufficient and necessary in the worst case for a nondeterministic finite cover automaton to accept $L$. This bound is tight already for a unary alphabet.* $\qquad \square$

Next we want to close the gap between the lower and upper bound for the conversion from DFAs to DFCAs.

**Theorem 7.** *Let $L$ be a finite language accepted by a deterministic finite automaton with $n$ states. If $n = 1$ or $n \geq 4$ then $n$ states are sufficient and necessary in the worst case for a deterministic finite cover automaton to accept $L$. These bounds are tight already for binary alphabets. If $n \in \{2, 3\}$, or if $n \geq 2$ and $L$ is a unary language, then $n - 1$ states are sufficient and necessary in the worst case.*

PROOF. The upper bound of $n$ states is clear. Moreover, notice that the only *finite* language that is accepted by a two-state DFA is the language $\{\lambda\}$, which is covered by the single-state automaton for the language $\Sigma^*$ with cover length 0. Next assume that $L \subseteq \Sigma^*$ is a finite language that is accepted by a

minimal DFA with three states. Then the order of $L$ must be 1, which means that we have $L = \Sigma_1$ or $L = \{\lambda\} \cup \Sigma_1$ for some non-empty subset $\Sigma_1 \subseteq \Sigma$. In both cases the language $L$ can be covered by a two-state automaton with cover length $\ell$: if $L = \Sigma_1$, we choose the cover language $L' = (\Sigma \setminus \Sigma_1)^* \cdot \Sigma_1 \cdot \Sigma^*$, and for the case $L = \{\lambda\} \cup \Sigma_1$ we use $L' = \Sigma_1^*$. This proves the statement of the theorem for the cases $n = 2$ and $n = 3$. The case $n = 1$ is clear.

Next let us briefly discuss the unary case. If $L$ is a unary language of order $\ell$, then the minimal DFA for $L$ is a chain of exactly $n = \ell + 2$ states, with a non-accepting sink state at the end. Such a DFA can be transformed into a DFCA with one state less by deleting the sink state and re-routing its incoming transitions to an arbitrary state. Hence we have a DFCA with $n-1$ states. Tightness of this bound is provided by Theorem 5.

The lower bound of $n$ states for the case $n = 4$ can be seen by considering the finite language $L_4 = \{\lambda, a, b, ab\}$ over alphabet $\Sigma = \{a, b\}$, which can be accepted by a DFA with four states. To see that any DFCA $B = (Q, \Sigma, \delta, q_0, F)$ with $L(B) \cap \Sigma^{\leq 2} = L_4$ has at least four states, we show that the four states $q_0$, $q_1 = \delta(q_0, a)$, $q_2 = \delta(q_0, b)$, and $q_3 = \delta(q_0, aa)$ must be pairwise distinct. Since $\lambda, a, b \in L_4$ and $aa \in \Sigma^{\leq 2} \setminus L_4$, we know that state $q_3$ is non-accepting and the other three states must be accepting. Moreover, it cannot be $q_0 = q_1$ nor $q_0 = q_2$ because this would implies $aa \in L_4$ or $bb \in L_4$, respectively. Finally states $q_1$ and $q_2$ must also be distinct since otherwise we would get $ab \in L_4$ if and only if $bb \in L_4$. Hence, every DFCA for the language $L_4$ needs four states.

It remains to consider the case $n \geq 5$. Here we use the witness language $L_n = \{a^{n-2}\} \cup \{a^i b \mid 0 \leq i \leq n - 3\}$ over alphabet $\Sigma = \{a, b\}$, which can be accepted by an $n$-state DFA. Notice that the order of $L_n$ is $\ell = n - 2$. Let $B = (Q', \Sigma, \delta', q_0', F')$ be some DFCA for $L_n$ and consider the sequence of states $q_i' = \delta'(q_0', a^i)$, for $0 \leq i \leq n - 3$. Clearly all these states must be non-accepting, and each of them leads to an accepting state on input symbol $b$. Moreover, state $q_{n-3}$ also leads to an accepting state on input $a$, and therefore all the states $q_i'$ with $0 \leq i \leq n - 3$ must be pairwise distinct—otherwise reading $a^{n-2}$ cannot lead to an accepting state. So far we have shown that $B$ has at least the $n - 2$ non-accepting states $q_i'$, with $0 \leq i \leq n - 3$ and at least one accepting state. In fact, if $B$ has *more* than one accepting state, then it has at least $n$ states in total. Therefore assume that $B$ has only one accepting state $q_{n-2}'$. This state is the target of every $b$-transition from states $q_i'$ for $0 \leq i \leq n - 3$. Now consider the state $q_{n-1}' = \delta'(q_{n-2}', b)$. This state must be non-accepting because it is reached from the initial state $q_0'$ by

13

reading the word $bb$, which is shorter than $\ell = n - 2 \geq 3$ symbols, and which does not belong to the language $L_n$. Moreover, state $q'_{n-1}$ must be different from all the other non-accepting states $q'_i$, which is seen as follows. Assume for the sake of contradiction that $q'_{n-1} = q'_i$ for some $i$ with $0 \leq i \leq n - 3$. This means that $B$ allows the computation

$$\delta'(q'_0, bbb) = \delta'(q'_{n-2}, bb) = \delta'(q'_i, b) = q_{n-2},$$

which means that the word $bbb$ is accepted by $B$. Since $bbb \notin L_n$ but $bbb \in \Sigma^{\leq \ell}$ (recall that $\ell = n - 2 \geq 3$), we obtain a contradiction to $L(B) \cap \Sigma^{\leq \ell} = L_n$. Therefore, automaton $B$ has at least $n$ states. $\qquad\square$

We also note that the conversion from NFAs to DFCAs was investigated already in [31], where binary languages $L_n$ were presented that can be accepted by an $n$-state NFA, while $2^{n-t} - 2^{t-2} + 2^t - 1$ states are necessary, with $t = \lfloor \frac{n}{2} \rfloor$, for a deterministic finite cover automaton to accept $L_n$. Then they generalize their examples to larger alphabets. The lower bound is known to be tight if $n$ is even, but the tight bound for odd $n$ remains to be determined. Asymptotically, their lower and upper bound read as $\Theta(k^{\frac{n}{1+\log_2 k}})$, which is (up to a possible constant factor) the same behavior in the worst case as for the conversion from NFAs to DFAs, which was determined in [3].

### 4.2. From Cover Automata to Finite Automata

In the previous subsection we have seen that there are finite languages where the description size cannot be reduced when changing the descriptional model from finite automata to cover automata. In this section we now consider the inverse conversion: given a cover automaton for a finite language, how large can a minimal finite automaton for that language become? In this setting we will see that the number of states of a cover automaton alone is not a fair size measure. In fact, we propose that a reasonable size measure for cover automata must also take the cover length into account: for every integer $\ell \geq 0$ the finite language $\{a\}^{\leq \ell}$ can be covered by a single-state cover automaton, but a NFA for this language has at least $\ell + 1$ states. Therefore, if we start with a cover automaton with $n$ states that describes a finite language of order $\ell$, then the number of states of an equivalent finite automaton should be a function in $n$ *and* $\ell$.

Since the language $L$ described by a cover automaton $A$ with cover length $\ell$ satisfies $L = L(A) \cap \Sigma^{\leq \ell}$, a finite automaton for $L$ can be obtained

14

by applying a cross product construction on $A$ and an automaton for $\Sigma^{\leq \ell}$, similar as in the proof of Theorem 5. The states of the constructed automaton are pairs $(q, i)$, where $q$ is a state of $A$, and $i$ is a counter for the word length. An NFA for the language $\Sigma^{\leq \ell}$ has $\ell + 1$ states, while a DFA has $\ell + 2$ states. Therefore, if $A$ is a nondeterministic cover automaton with $n$ states, then one can construct an equivalent NFA with at most $n \cdot (\ell + 1)$ states, and if $A$ is deterministic, then one constructs a DFA with $n \cdot (\ell + 2)$ states. This yields the upper bounds $\mathsf{nsc}(L) \leq \mathsf{ncsc}(L) \cdot (\ell + 1)$ and $\mathsf{sc}(L) \leq \mathsf{csc}(L) \cdot (\ell + 2)$ for finite languages $L$ of order $\ell$. In the upcoming lemma we show that these bounds can be slightly reduced. In the following we do not consider languages of order $\ell = 0$, because the only such language is $\{\lambda\}$, which is accepted by a single-state NFA and a two-state DFA. Moreover, the case where $\mathsf{ncsc}(L) = 1$ is also omitted—here it is easy to see that the upper bounds $\mathsf{nsc}(L) \leq \ell + 1$ and $\mathsf{sc}(L) \leq \ell + 2$ apply, and optimality is witnessed by the language $L = \Sigma^{\leq \ell}$.

**Lemma 8.** *Let $n \geq 2$ and $A$ be an $n$-state nondeterministic cover automaton for a finite language $L$ of order $\ell \geq 1$. Then one can construct a nondeterministic finite automaton for $L$ that has at most $n \cdot (\ell - 1) + 2$ states. If $A$ is deterministic, then one can construct a deterministic finite automaton for $L$ with $n \cdot (\ell - 1) + 3$ states.*

PROOF. Let $A = (Q, \Sigma, \delta, q_0, F)$ be an NFCA for the language $L \subseteq \Sigma^{\leq \ell}$, that is, with $L(A) \cap \Sigma^{\leq \ell} = L$. We construct an NFA $A' = (Q \times [\ell], \Sigma, \delta', (q_0, 0), F \times [\ell])$ where the transition function is defined by $\delta'((p, i), a) = \{ (q, i + 1) \mid q \in \delta(p, a) \}$, for $q \in Q$, $0 \leq i \leq \ell - 1$, and $a \in \Sigma$. In states $(q, \ell)$ no transitions are defined, that is, $\delta((q, \ell), a) = \emptyset$, for $q \in Q$ and $a \in \Sigma$. Clearly this automaton accepts the language $L(A) \cap \Sigma^{\leq \ell} = L$ and has $n \cdot (\ell + 1)$ states. We now show that some states of $A'$ can be omitted. First notice that the only reachable state with $0$ in its second component is $(q_0, 0)$. Moreover, all states of the form $(q, \ell)$, with $q \in Q \setminus F$, can be removed because no transitions are defined in these states. Finally, all states $(q_f, \ell)$ with $q_f \in F$ can be merged to a single accepting state $(\bullet, \ell)$. Hence the state set of $A'$ can be restricted to $\{(q_0, 0), (\bullet, \ell)\} \cup Q \times \{1, 2, \ldots, \ell - 1\}$ which gives a total number of $n \cdot (\ell - 1) + 2$ states.

In case $A$ is a *deterministic* cover automaton, a similar construction can be applied to obtain an equivalent DFA. Again, states $(q_f, \ell)$, with $q_f \in F$, are merged into a single accepting state $(\bullet, \ell)$, and states $(q, \ell)$, with $q \in Q \setminus F$,

are removed. To deal with the resulting undefined transitions, we add an additional non-accepting sink state $(\bullet, >\ell)$, which is also the target state for transitions from state $(\bullet, \ell)$. The obtained automaton has $n \cdot (\ell - 1) + 3$ states. $\qquad \square$

Next we show that the constructions from Lemma 8 cannot be improved in general, by providing a matching lower bound. Observe that the following lemma even provides a lower bound for the conversion from *deterministic* cover automata to *nondeterministic* finite automata.

**Lemma 9.** *For all integers $n \geq 2$ and $\ell \geq 1$ there exists a finite language $L$ of order $\ell$ that is described by a deterministic $n$-state cover automaton, such that any nondeterministic finite automaton for $L$ needs $n \cdot (\ell - 1) + 2$ states, and any deterministic finite automaton for $L$ needs $n \cdot (\ell - 1) + 3$ states.*

PROOF. Let $n \geq 2$ and $\ell \geq 1$, and define the DFCA $A = (Q, \Sigma, \delta, q_0, F)$ with input alphabet $\Sigma = \{\, a_i, b_i \mid 1 \leq i \leq n - 1 \,\}$, state set $Q = \{q_0, q_1, \ldots, q_{n-1}\}$, final states $F = Q \setminus \{q_0\}$. The transition function $\delta$ is defined for $1 \leq i, j, k \leq n - 1$ with $i \neq j$ as follows:

$$\delta(q_0, a_i) = q_i, \quad \delta(q_0, b_i) = q_0, \quad \delta(q_i, a_k) = q_0, \quad \delta(q_i, b_i) = q_i, \quad \delta(q_i, b_j) = q_0.$$

Let $L = L(A) \cap \Sigma^{\leq \ell}$. We show that every NFA for $L$ has at least $n \cdot (\ell-1) + 2$ states, by proving that the set

$$S = \{\, (b_1^k, b_1^{\ell-1-k} a_1) \mid 0 \leq k \leq \ell - 1 \,\} \cup \{(b_1^{\ell-1} a_1, \lambda)\}$$
$$\cup \{\, (a_i b_i^{j-1}, b_i^{\ell-j}) \mid 1 \leq i \leq n - 1, \ 1 \leq j \leq \ell - 1 \,\}$$

is an extended fooling set for $L$. First notice that for each pair $(x, y) \in S$ we have $xy \in L$, because $b_1^{\ell-1} a_1$ leads to the accepting state $q_1$, and $a_i b_i^{\ell-1}$ to the accepting state $q_i$. Now let $(x, y)$ and $(x', y')$ be two distinct pairs from $S$. We have to show that $xy' \notin L$ or $x'y \notin L$:

1. First assume $(x, y) = (b_1^{\ell-1} a_1, \lambda)$. Notice that $(b_1^{\ell-1} a_1, \lambda)$ is the only pair in $S$ with $\lambda$ as second component. Hence the word $xy'$ cannot belong to $L$ because $|xy'| > \ell$. The case where $(x', y') = (b_1^{\ell-1} a_1, \lambda)$ is symmetric.
2. Similarly, when combining strings from two pairs $(x, y) = (b_1^k, b_1^{\ell-1-k} a_1)$ and $(x', y') = (b_1^{k'}, b_1^{\ell-1-k'} a_1)$, with $k \neq k'$, one obtains a word of length greater than $\ell$, which cannot belong to the language $L$.

16

3. Next consider the pairs $(x, y) = (b_1^k, b_1^{\ell-1-k}a_1)$ and $(x', y') = (a_i b_i^{j-1}, b_i^{\ell-j})$. Here we find $b_1^k b_i^{\ell-j} \notin L$ because this word leads to the non-accepting state $q_0$ in $A$. The case where $(x, y)$ and $(x', y')$ are interchanged is symmetric.

4. Finally it remains to consider the case where $(x, y) = (a_i b_i^{j-1}, b_i^{\ell-j})$ and $(x', y') = (a_{i'} b_{i'}^{j'-1}, b_{i'}^{\ell-j'})$, with $(i, j) \neq (i', j')$. If $j \neq j'$ then we have $|xy'| > \ell$ or $|x'y| > \ell$, hence one of those words does not belong to $L$. In the case $j = j'$ we have $i \neq i'$, and then we have $xy' = a_i b_i^{j-1} b_{i'}^{\ell-j'} \notin L$, because reading the prefix $a_i b_i^{j-1}$ takes $A$ to state $q_i$ and from there, the suffix $b_{i'}^{\ell-j'}$ takes $A$ back to the non-accepting state $q_0$.

We have shown that a minimal NFA for the language $L$ has at least $n \cdot (\ell - 1) + 2$ states. Now one readily sees that a minimal DFA for $L$ has at least $n \cdot (\ell - 1) + 3$ states: because $L$ is a finite language, a minimal DFA for $L$ needs a non-accepting sink-state, which is of course not present in a minimal NFA. $\qquad\square$

From Lemmata 8 and 9 we obtain the following result.

**Theorem 10.** *Let $L$ be a finite language of order $\ell \geq 1$ that is described by a nondeterministic cover automaton $A$ with $n \geq 2$ states. Then $n \cdot (\ell - 1) + 2$ states are sufficient and necessary in the worst case for a nondeterministic finite automaton to accept $L$. Moreover, if $A$ is a deterministic cover automaton for $L$, then $n \cdot (\ell - 1) + 3$ states are sufficient and necessary in the worst case for a deterministic finite automaton to accept $L$.* $\qquad\square$

Observe that the proof for the lower bound from Lemma 9 uses $2n - 2$ alphabet symbols. In fact, one can also show that the bounds $\mathsf{nsc}(L) \leq \mathsf{ncsc}(L) \cdot (\ell - 1) + 2$ and $\mathsf{sc}(L) \leq \mathsf{csc}(L) \cdot (\ell - 1) + 3$ for the conversions from cover automata to finite automata are *not* tight for languages over an alphabet of constant size. For the deterministic case, this is easy to see: assuming a $k$-letter alphabet $\Sigma$, at most $k$ different states of the form $(q, 1)$ are reachable from the initial state $(q_0, 0)$ in the DFA constructed from a DFCA as shown in the proof of Lemma 8.

Although this argumentation does not hold for nondeterministic automata, where every state of the given NFCA could be reachable in one step from the initial state, the number of states of an equivalent minimal NFA still depends on the number of alphabet symbols: when using the construction from

Lemma 8 to obtain an NFA $A'$ for the language $L \subseteq \Sigma^{\leq \ell}$, the automaton $A'$ has a distinguished "last" accepting state $(\bullet, \ell)$, which has no outgoing transitions. This state is only reachable from states of the form $(q, \ell - 1)$, and from such states no other state is reachable. Assume that two such states $(p, \ell - 1)$ and $(q, \ell - 1)$ go to state $(\bullet, \ell)$ on the same set of input letters. If additionally $p$ and $q$ are of same acceptance value, then clearly they can be merged into a single state. Since a $k$-letter alphabet $\Sigma$ has $2^k - 1$ non-empty subsets, the number of accepting states of the form $(q, \ell - 1)$ can always be reduced to $2^k - 1$, and similarly for the non-accepting states. So in total there are at most $2 \cdot (2^k - 1)$ states of the form $(q, \ell - 1)$, which may be large compared to $k$, but it is still a constant.

The search for exact bounds depending on the size of the input alphabet is left for further research.

## 5. Determinization of Finite Cover Automata

In this section we continue our descriptional complexity studies of cover automata: we investigate the cost of determinization, that is, the conversion from a nondeterministic to a deterministic cover automaton. A classical result in the theory of finite automata is that every $n$-state NFA can be converted by the so-called power-set construction to an equivalent DFA with at most $2^n$ states [2]. Moreover, it is known that this bound is tight in the sense that for every $n \geq 1$ there exists a language accepted by a minimal $n$-state NFA, and for which the minimal DFA needs exactly $2^n$ states [32, 33, 34]. Now the question is to which extent these results carry over to cover automata. Clearly, since the power-set construction for finite automata preserves the accepted language, it can be used to convert an NFCA into an equivalent DFCA. Thus, the following is immediate.

**Lemma 11.** *Let $L$ be a finite language described by a nondeterministic cover automaton with $n \geq 1$ states. Then one can construct a deterministic cover automaton for $L$ that has at most $2^n$ states.* $\qquad \square$

Our next goal is to prove a matching lower bound of $2^n$ states for the determinization of $n$-state NFCAs. The next fact we present is useful to show that a number of worst case results known for the state complexity of deterministic finite automata carry over to the setting of cover automata.

**Theorem 12.** *Assume $L$ is a regular language over $\Sigma$ with $\mathsf{sc}(L) = n$, and let $L' = L \cap \Sigma^{\leq n + 2^n}$. Then $\mathsf{csc}(L') = n$.*

PROOF. Assume $A = (Q, \Sigma, \delta, q_0, F)$ is a minimal DFA accepting $L^R$, and let $\mathcal{P}(A^R) = (Q', \Sigma, \delta', q_0', F')$ be the DFA obtained by first reverting this automaton, thus obtaining a nondeterministic finite automaton with multiple initial states, and then applying the lazy power-set construction to this automaton. Here, lazy means that only the subsets reachable from the start state are constructed. Formally we have $Q' \subseteq 2^Q$, $q_0' = F$, $F' = \{ P \in Q' \mid q_0 \in P \}$ and $\delta'(P, a) = \{ q \mid \delta(q, a) \in P \}$, for $P \in Q'$ and $a \in \Sigma$. Then Brzozowski's theorem [35] asserts that this automaton is a minimal DFA for the reversed language $L$, hence it has $n$ states.

Our goal is to show that every two distinct states $R, S \in Q'$ of $\mathcal{P}(A^R)$ are dissimilar with respect to $L'$. Assume $q \in S \setminus R$ (if $S \subset R$, exchange the roles of $R$ and $S$). Let $v_q$ be a word of minimal length such that the automaton $A$ reaches state $q$ on reading $v_q$ backwards, that is $\delta(q_0, v_q^R) = q$. Then in the reversed automaton $\mathcal{P}(A^R)$ we have $q_0 \in \delta'(S, v_q)$, but $q_0 \notin \delta'(R, v_q)$, since $A$ is deterministic. By definition of $F'$ this means $\delta'(S, v_q) \in F'$ and $\delta'(R, v_q) \notin F'$. Next we want to estimate the length of $v_q$. By symmetry, $A$ can be also obtained by reverting $\mathcal{P}(A^R)$ and then applying the lazy power-set construction. Thus the number of states in $A$ is at most $2^n$, and we can conclude that $|v_q| \leq 2^n$. We now choose for every state $P$ in the automaton $\mathcal{P}(A^R)$ a word $u_P$ such that $|u_P| \leq n$ and $\delta'(F, u_P) = P$. Then we have $|u_R v_q|, |u_S v_q| \leq n + 2^n$, and

$$u_R v_q \notin L \text{ but } u_S v_q \in L.$$

As the languages $L$ and $L'$ agree on all words of length at most $n + 2^n$, this establishes that $R$ and $S$ are dissimilar with respect to $L'$. $\qquad\square$

Theorem 12 implies that if the order of the language is large compared to the size of the NFA, then determinization of cover automata is as expensive as for usual finite automata. In particular, classical examples for finite automata [32, 33, 34] show that the full blow-up from $n$ states to $2^n$ states may be necessary for converting an NFCA into an equivalent DFCA. Together with Lemma 11 we obtain the following result.

**Corollary 13.** *Let $L$ be a finite language that is described by a nondeterministic cover automaton with $n \geq 1$ states. Then $2^n$ states are sufficient and*

*necessary in the worst case for a deterministic cover automaton to accept L.*
□

A natural question is now whether the full blow-up can be reached if the order of the described language is small compared to the number of states in the given NFCA. First, recall that every finite language $L$ of order $\ell$ over a $k$-letter alphabet satisfies $\mathsf{sc}(L) \leq (1+\mathrm{o}(1))\frac{k^{\ell+2}}{d_k \ell}$ with $d_k = (k-1)^2 \log k$; see [4]. This shows that the full blow-up cannot be reached if $\ell$ is too small compared to $n$. From that result and the fact that $\mathsf{csc}(L) \leq \mathsf{sc}(L)$, the following bounds for the size of a deterministic cover automaton can be derived. In fact, since the proof of the next result only uses the above bound on $\mathsf{sc}(L)$, the statements also hold for the determinization of finite automata.

**Theorem 14.** *Let $L$ be a finite language of order $\ell$ over a $k$-letter alphabet $\Sigma$ and assume $L$ is described by a nondeterministic finite cover automaton with $n$ states.*

1. *If $(\ell + 2) \cdot \log k - \log \ell + 1 < n$, then $\mathsf{csc}(L) < 2^n$, for large enough $n$.*
2. *if $\ell \in \mathrm{o}(n)$, then $\mathsf{csc}(L) \in 2^{\mathrm{o}(n)}$,*
3. *if $\ell \in \mathrm{O}(\log n)$, then $\mathsf{csc}(L) \in n^{\mathrm{O}(1)}$,*
4. *if $(\ell + 2) \cdot \log k - \log \ell + 1 < \log n$, then $\mathsf{csc}(L) < n$, for large enough $n$.*

PROOF. We solve the inequality $(1 + \mathrm{o}(1))\frac{k^{\ell+2}}{\ell(k-1)^2 \log k} < 2^n$, where the left-hand side is the maximum state complexity among all languages $L \subseteq \Sigma^{\leq \ell}$, see [4]. By taking logarithms we obtain

$$\log(1 + \mathrm{o}(1)) + (\ell + 2) \cdot \log k - \log \ell - (2 \cdot \log(k-1) + \log \log k) < n$$

By observing that $\log_2(1 + \mathrm{o}(1)) < 1$ for large enough $n$, and omitting the negative term $-(2 \cdot \log(k-1) + \log \log k)$, we get the first statement. Clearly the fourth statement can be shown in the same way when starting with the inequality $(1+\mathrm{o}(1))\frac{k^{\ell+2}}{\ell(k-1)^2 \log k} < n$. Similarly the second and third statement can be derived by rewriting $k$ as $2^{\log_2 k}$. □

The fourth statement in the above theorem is of particular practical relevance: in this case, the given $n$-state NFCA is not minimal, and determinization followed by minimization yields a smaller cover automaton. Thus, for a *minimal* $n$-state NFCA we always have $(\ell + 2) \cdot \log k - \log \ell + 1 \geq \log n$.

20

In contrast to languages of order less than $n$, where the blow-up of $2^n$ states cannot be achieved, there are quite natural examples reaching the full blow-up already for order linear in the number of states of the NFCA. The example used in the following proof is essentially due to [36, Lemma 2]:

**Theorem 15.** *Let* $L_n = (a + (a \cdot b^*)^{n-1} \cdot a)^* \cap \Sigma^{\leq 5n-2}$. *Then* $L_n$ *can be covered by an $n$-state nondeterministic cover automaton, but the smallest deterministic cover automaton for* $L_n$ *has at least $2^n$ states.*

PROOF. Let $K_n$ be the language $(a + (a \cdot b^*)^{n-1} \cdot a)^*$ as introduced in [36], and define $L_n = K_n \cap \Sigma^{\leq 5n-2}$. The language $K_n$ is accepted by an $n$-state NFA, thus it is covered by an NFCA with $n$ states. Along the lines of [36, Lemma 2], it can be proved that when applying the power-set construction, all $2^n$ subsets are reachable by a string of length at most $3n$, and each pair of subsets can be distinguished by a string of length at most $2n - 2$. Hence all $2^n$ states of the resulting DFA are pairwise dissimilar, so the minimal DFCA has $2^n$ states. □

## 6. Average Size Comparisons of Finite Cover Automata

This section is devoted to the average case state complexity of DFCAs and NFCAs, when choosing a finite language of a certain "size" $\ell$ uniformly at random from all finite languages of that particular size. Here size means that all words of the language are either of the same length $\ell$, or of length at most $\ell$. This model was used in [6] to compare the number of states or transitions of ordinary finite automata on average. There it is shown that almost all DFAs accepting finite languages over a binary input alphabet have state complexity $\Theta(2^\ell/\ell)$, while NFAs are shown to perform better, namely the nondeterministic state complexity is in $\Theta(\sqrt{2^\ell})$. Interestingly, in both cases the aforementioned bounds are asymptotically like in the worst case. As we will see, a similar situation emerges for finite cover automata as well. The first theorem gives us the expected number of states a DFCA has on average, if we assume that all finite languages from $\mathfrak{P}(\Sigma^{\leq \ell})$, that is, the power-set of $\Sigma^{\leq \ell}$, are equiprobable.

**Theorem 16.** *Let $\Sigma$ be an alphabet of size $k$ and $c_k = (k-1) \log k$. Then* $\mathbb{E}[\mathsf{csc}(L)] \geq (1 - \mathrm{o}(1))\frac{k^\ell}{c_k \ell}$, *if $L$ is a language drawn uniformly at random from the power-set of $\Sigma^{\leq \ell}$.*

PROOF. The argument is similar as for an analogous result on finite automata, in place of cover automata, from [6]. It suffices to show that the number of languages acceptable by DFCAs with at most $(1-\delta)\frac{k^\ell}{c_k\ell}$ states, for $0 < \delta < 1$, is in $o(|\mathfrak{P}(\Sigma^{\leq\ell})|)$; the result then follows using Markov's inequality.

Let $g_k(m)$ be the function counting the number of languages over $\Sigma$ acceptable by DFAs with at most $m$ states over a $k$-letter alphabet. In [37, Theorem 9] it was shown that $g_k(m) \leq m2^m\frac{m^{km}}{m!}$. For finite languages of order at most $\ell$, each $m$-state minimal DFA gives rise to at most $\ell+1$ different minimal DFCAs, so the number of subsets $\Sigma^{\leq\ell}$ acceptable by DFCAs with at most $m$ states is at most $(\ell+1)\cdot g_k(m)$. Using

$$\log m! > \int_1^m \log x\,dx = m\log m - \frac{1}{\ln 2}(m-1),$$

and the fact that $\frac{1}{\ln 2} < \frac{3}{2}$, we obtain $\log\left((\ell+1)\cdot g_k(m)\right) = \log(g_k(m)) + \log(\ell+1) < (k-1)m\log m + \frac{5}{2}m + \log m + \log(\ell+1)$. Thus for every constant $\delta$ with $0 < \delta < 1$,

$$\log\left((\ell+1)\cdot g_k\left((1-\delta)\frac{k^\ell}{c_k\ell}\right)\right) < \log(\ell+1) + (1-\delta)\left(1+\frac{5}{2c_k\ell}\right)k^\ell + \ell\log k$$
$$= (1-\delta)k^\ell + o(k^\ell),$$

and for $\ell$ large enough, this is much smaller than $k^\ell < \log|\mathfrak{P}\left(\Sigma^{\leq\ell}\right)|$, that is

$$\log(\ell+1) + \log g_k\left((1-\delta)\frac{k^\ell}{c_k\ell}\right) - \log|\mathfrak{P}\left(\Sigma^{\leq\ell}\right)|$$

tends to $-\infty$. We can deduce that $\lim_{\ell\to\infty}(\ell+1)\cdot g_k((1-\delta)\frac{k^\ell}{c_k\ell})/|\mathfrak{P}\left(\Sigma^{\leq\ell}\right)| = 0$, for every such $\delta$. □

Regarding an upper bound, it is known from [4] that $\mathsf{sc}(L) \leq (1 + o(1))\frac{k^{\ell+2}}{d_k\ell}$, as $\ell$ tends to infinity, with $d_k = (k-1)^2\log_2 k$, for languages $L \subseteq \Sigma^{\leq\ell}$ and alphabet size $k$. This generalized a previous result of [5]. Recall that the size of a minimal DFA for a finite language is an upper bound for the size of a minimum DFCA; and the state complexity in the worst case is of course an upper bound for the average state complexity. So the above average case result is tight up to a factor of at most $(1+o(1))\frac{k^2}{(k-1)}$. Next we turn our attention to the average state complexity of NFCAs.

22

**Theorem 17.** *Let $\Sigma$ be an alphabet of size $k$. Then for large enough $\ell$ we have $\mathbb{E}[\mathsf{ncsc}(L)] > k^{\frac{\ell}{2}-1}$, if $L$ is a language drawn uniformly at random from the power-set of $\Sigma^{\leq \ell}$.*

PROOF. We generalize a result from [38] regarding nondeterministic finite automata over binary alphabets.

Let $G_k(m)$ denote the number of distinct regular languages accepted by NFAs with $m$ states over a $k$-letter alphabet. Then $G_k(m) \leq 2m2^{km^2}$, as shown in [37]. If we consider finite languages of order at most $\ell$, every NFA counted by $G_k(m)$ gives rise to at most $\ell+1$ different finite languages covered by NFCAs with at most $m$ states. Now we have

$$(\ell + 1) \cdot G_k(k^{\frac{\ell}{2}-1}) \leq (\ell + 1) \cdot 2 \cdot k^{\frac{\ell}{2}-1} 2^{k \cdot (k^{\frac{\ell}{2}-1})^2}$$
$$= \underbrace{(\ell + 1) \cdot 2 \cdot k^{\frac{\ell}{2}-1}}_{=o\left(\sqrt{2^{k^\ell}}\right)} \sqrt[k]{2^{k^\ell}} = o\left(2^{k^\ell}\right);$$

and thus the fraction of "nice" subsets of $\Sigma^{\leq \ell}$, which can be covered by NFCAs having at most $k^{\frac{\ell}{2}-1}$ states, tends to zero as $\ell$ grows large. $\qquad\square$

A worst case upper bound for the nondeterministic state complexity of subsets of $\Sigma^{\leq \ell}$ is given in [6] for binary alphabets. Generalizing this result to cover automata and larger alphabets, the bound reads as follows:

**Theorem 18.** *Let $\Sigma$ be an alphabet of size $k$. Then*

$$\mathsf{ncsc}(L) \leq \mathsf{nsc}(L) < \frac{3}{k-1}\sqrt{k^\ell},$$

*if $L$ is any subset of $\Sigma^{\leq \ell}$, i.e., $L \subseteq \Sigma^{\leq \ell}$.*

PROOF. There was a minor mistake in the computation of the final estimation of the corresponding result in [6]. So, for convenience, we include a proof of the corrected bound here, which also covers the case of non-binary alphabet size.

Let $\mu = \lfloor (\ell - 1)/2 \rfloor$ and $\nu = \lceil (\ell - 1)/2 \rceil$. We construct an NFA $A = (Q, \{0, 1\}, \delta, p_\lambda, F)$, where the state set $Q = P_1 \cup P_2$ (the union is disjoint) with

$$P_1 = \{\, p_w \mid w \in \{0, 1\}^* \text{ and } |w| \leq \mu \,\}$$

and

23

$$P_2 = \{\, q_w \mid w \in \{0,1\}^* \text{ and } |w| \leq \nu \,\},$$

the set $F = \{q_\lambda\} \cup \{\, p_\lambda \mid \lambda \in L \,\}$, and the transition function is specified as follows:

1. For all $p_w \in P_1$ and $a \in \{0,1\}$, the set $\delta(p_w, a)$ contains the element $p_{wa}$.
2. For all $w \in L \setminus \{\lambda\}$, if $w = xay$ is the unique decomposition, where $|x| = \lfloor (|w| - 1)/2 \rfloor$, $a$ is a single letter, and $|y| = \lceil (|w| - 1)/2 \rceil$, then let $\delta(p_x, a)$ contain the element $q_y$.
3. For all $q_w \in P_2 \setminus \{q_\lambda\}$ and $a \in \{0,1\}$, the set $\delta(p_{aw}, a)$ contains the element $q_w$.

This completes the construction of the NFA (which can also be interpreted as a NFCA). It is easy to see that for the number of states in $A$, we have

$$|P_1| + |P_2| = \frac{k^{\mu+1} - 1}{k - 1} + \frac{k^{\nu+1} - 1}{k - 1} < \frac{3}{k - 1}\sqrt{k^\ell}.$$

It remains to show that $L(A) = L$. Note that every state $p_w$ in $P_1$ is only reachable by the word $w$ from the initial state $p_\lambda$, and that for every state $q_w$ in $P_2$ there is only one path leading to the final state $q_\lambda$. So every transition leading from $P_1$ to $P_2$ leads to the acceptance of exactly one word in $L$. This proves the stated claim. $\square$

## 7. Conclusions

We completed the picture of lower bound techniques for nondeterministic finite cover automata, and solved the problems left open in [18]. Then we determined the precise best-case and worst-case bounds for conversions between DFCAs and DFAs, as well as between NFCAs and NFAs. In [31], almost tight bounds for the conversion between NFAs and DFCAs were given. Determining the precise bound in this case remains an open problem.

When the length $\ell$ of the longest word is much smaller than the number $n$ of states in a minimal cover automaton, then the succinctness gain offered by finite cover automata over finite automata is very modest, even in the best case. We note that this is the case in the area of natural language processing: in [39] they construct a minimum 29317-state DFA accepting 81142 English words. Of course, almost all common English words have $\ell < 20$. Similarly, in [40] they construct an NFA accepting roughly 230000 Greek words, whose number of states is of order $10^5$.

24

Finally, a recent experimental study [41] showed that for binary finite languages, the expected reduction in the number of states provided by DFCAs is negligible. Our analysis of the average case provides a theoretical underpinning of their observations. One may study further random models of finite languages, e.g., a Bernoulli-type model [6], and one based on the sum of word lengths [42].

## References

[1] C. E. Shannon, The synthesis of two-terminal switching circuits, Bell Systems Technical Journal 28 (1) (1949) 59–98.

[2] M. O. Rabin, D. Scott, Finite automata and their decision problems, IBM J. Res. Dev. 3 (1959) 114–125.

[3] K. Salomaa, S. Yu, NFA to DFA transformation for finite language over arbitrary alphabets, J. Autom., Lang. Comb. 2 (3) (1997) 177–186.

[4] C. Câmpeanu, W. H. Ho, The maximum state complexity for finite languages, J. Autom., Lang. Comb. 9 (2–3) (2004) 189–202.

[5] J.-M. Champarnaud, J.-E. Pin, A maxmin problem on finite automata, Discrete Appl. Math. 23 (1989) 91–96.

[6] H. Gruber, M. Holzer, Results on the average state and transition complexity of finite automata accepting finite languages, Theoret. Comput. Sci. 387 (2) (2007) 155–166.

[7] C. Câmpeanu, N. Sântean, S. Yu, Minimal cover-automata for finite languages, Theoret. Comput. Sci. 267 (1–2) (2001) 3–16.

[8] C. Câmpeanu, A. Păun, The number of similarity relations and the number of minimal deterministic finite cover automata, in: J.-M. Champarnaud, D. Maurel (Eds.), Proceedings of the 7th International Conference Implementation and Application of Automata, no. 2608 in LNCS, Springer, Tours, France, 2003, pp. 67–76.

[9] C. Câmpeanu, A. Păun, J. R. Smith, Incremental construction of minimal deterministic finite cover automata, Theoret. Comput. Sci. 363 (2) (2006) 135–148.

[10] C. Câmpeanu, A. Păun, J. R. Smith, Tight bounds for the state complexity of deterministic cover automata, in: H. Leung, G. Pighizzini (Eds.), Proceedings of the 8th Workshop on Descriptional Complexity of Formal Systems, Las Cruces, New Mexico, USA, 2006, pp. 223–231, Computer Science Technical Report NMSU-CS-2006-001.

[11] C. Câmpeanu, A. Păun, S. Yu, An efficient algorithm for constructing minimal cover automata for finite languages, Internat. J. Found. Comput. Sci. 13 (1) (2002) 83–97.

[12] J.-M. Champarnaud, F. Guingne, G. Hansel, Similarity relations and cover automata, RAIRO–Informatique théorique et Applications / Theoretical Informatics and Applications 39 (1) (2005) 115–123.

[13] P. García, J. Ruiz, A note on cover automata for finite languages, Bull. EATCS 83 (2004) 193–199.

[14] A. Jéz, A. Maletti, Computing all $\ell$-cover automata fast, in: B. Bouchou-Markhoff, P. Caron, J.-M. Champarnaud, D. Maurel (Eds.), Proceedings of the 16th Conference on Implementation and Application of Automata, no. 6807 in LNCS, Springer, Blois, France, 2011, pp. 203–214.

[15] H. Körner, On minimizing cover automata for finite languages in $O(n \log n)$ time, in: J.-M. Champarnaud, D. Maurel (Eds.), Proceedings of the 7th International Conference Implementation and Application of Automata, no. 2608 in LNCS, Springer, Tours, France, 2003, pp. 117–127.

[16] H. Körner, A time and space efficient algorithm for minimizing cover automata for finite languages, Internat. J. Found. Comput. Sci. 14 (6) (2003) 1071–1086.

[17] A. Păun, N. Sântean, S. Yu, An $O(n^2)$ algorithm for constructing minimal cover automata for finite languages, in: A. Păun, S. Yu (Eds.), Proceedings of the 5th International Conference Implementation and Application of Automata, no. 2088 in LNCS, Springer, London, Ontario, Canada, 2001, pp. 243–251.

[18] C. Câmpeanu, Non-deterministic finite cover automata, Scientific Annals of Computer Science 25 (1) (2015) 3–28.

[19] S. Yu, Cover automata for finite language, Bull. EATCS 92 (2007) 65–74.

[20] H. Gruber, M. Holzer, Finding lower bounds for nondeterministic state complexity is hard (extended abstract), in: O. H. Ibarra, Z. Dang (Eds.), Proceedings of the 10th International Conference on Developments in Language Theory, no. 4036 in LNCS, Springer, Santa Barbara, California, USA, 2006, pp. 363–374.

[21] M. Holzer, S. Jakobi, From equivalence to almost-equivalence, and beyond: Minimizing automata with errors, Internat. J. Found. Comput. Sci. 24 (7) (2013) 1083–1134.

[22] M. A. Harrison, Introduction to Formal Language Theory, Addison-Wesley, 1978.

[23] J.-C. Birget, Intersection and union of regular languages and state complexity, Inform. Process. Lett. 43 (1992) 185–190.

[24] H. N. Adorna, Some descriptional complexity problems in finite automata theory, in: R. P. Saldaña, C. Chua (Eds.), Proceedings of the 5th Philippine Computing Science Congress, Computing Society of the Philippines, Cebu City, Philippines, 2005, pp. 27–32.

[25] J. Hromkovič, Descriptional complexity of finite automata: Concepts and open problems, J. Autom., Lang. Comb. 7 (4) (2002) 519–531.

[26] I. Glaister, J. Shallit, A lower bound technique for the size of nondeterministic finite automata, Inform. Process. Lett. 59 (1996) 75–77.

[27] J. Hromkovič, Communication Complexity and Parallel Computing, Springer, 1997.

[28] H. N. Adorna, 3-party message complexity is better than 2-party ones for proving lower bounds on the size of minimal nondeterministic finite automata, J. Autom., Lang. Comb. 7 (4) (2002) 419–432.

[29] D. de Caen, D. A. Gregory, N. J. Pullman, The Boolean rank of zero-one matrices, in: C. C. Cadogan (Ed.), 3rd Caribbean Conference on Combinatorics and Computing, Department of Mathematics, University of the West Indies, 1981, pp. 169–173.

[30] J. Orlin, Contentment in graph theory: Covering graphs with cliques, Indigationes Mathematicae 80 (1977) 406–424.

[31] C. Câmpeanu, L. Kari, A. Păun, Results on transforming NFA into DFCA, Fund. Inform. 64 (1-4) (2005) 53–63.

[32] O. B. Lupanov, Über den Vergleich zweier Typen endlicher Quellen, Probleme der Kybernetik 6 (1966) 328–335.

[33] A. R. Meyer, M. J. Fischer, Economy of description by automata, grammars, and formal systems, in: Proceedings of the 12th Annual Symposium on Switching and Automata Theory, IEEE Computer Society Press, 1971, pp. 188–191.

[34] F. R. Moore, On the bounds for state-set size in the proofs of equivalence between deterministic, nondeterministic, and two-way finite automata, IEEE Transaction on Computing C-20 (1971) 1211–1219.

[35] J. A. Brzozowski, Canonical regular expressions and minimal state graphs for definite events, in: Mathematical Theory of Automata, Vol. 12 of MRI Symposia Series, Polytechnic Press, NY, 1962, pp. 529–561.

[36] H. Leung, Separating exponentially ambiguous finite automata from polynomially ambiguous finite automata, SIAM J. Comput. 27 (4) (1998) 1073–1082.

[37] M. Domaratzki, D. Kisman, J. Shallit, On the number of distinct languages accepted by finite automata with $n$ states, J. Autom., Lang. Comb. 7 (4) (2002) 469–486.

[38] G. Gramlich, G. Schnitger, Minimizing nfa's and regular expressions, J. Comput. System Sci. 73 (6) (2007) 908–923.

[39] C. L. Lucchesi, T. Kowaltowski, Applications of finite automata representing large vocabularies, Software—Practice and Experience 23 (1) (1993) 15–30.

[40] K. N. Sgarbas, N. D. Fakotakis, G. K. Kokkinakis, Incremental construction of compact acyclic NFAs, in: Proceedings of the 39th Annual Meeting on Association for Computer Linguistic and 10th Conference

of the European Chapter, Morgan Kaufmann, Toulouse, France, 2001, pp. 474–481.

[41] C. Câmpeanu, N. Moreira, R. Reis, Expected compression ratio for DFCA: experimental average case analysis, Tech. rep., Departamento de Ciência de Computadores, Universidade do Porto (2011).

[42] F. Bassino, L. Giambruno, C. Nicaud, The average state complexity of rational operations on finite languages, Internat. J. Found. Comput. Sci. 21 (4) (2010) 495–516.