

# Finite Automata for the Sub- and Superword Closure of CFLs: Descriptive and Computational Complexity\*

Georg Bachmeier, Michael Luttenberger, and Maximilian Schlund

Technische Universität München, {bachmeie,luttenbe,schlund}@in.tum.de

**Abstract.** We answer two open questions by (Gruber, Holzer, Kutrib, 2009) on the state-complexity of representing sub- or superword closures of context-free grammars (CFGs): (1) We prove a (tight) upper bound of  $2^{\mathcal{O}(n)}$  on the size of nondeterministic finite automata (NFAs) representing the subword closure of a CFG of size  $n$ . (2) We present a family of CFGs for which the minimal deterministic finite automata representing their subword closure matches the upper-bound of  $2^{2^{\mathcal{O}(n)}}$  following from (1). Furthermore, we prove that the inequivalence problem for NFAs representing sub- or superword-closed languages is only NP-complete as opposed to PSPACE-complete for general NFAs. Finally, we extend our results into an approximation method to attack inequivalence problems for CFGs.

## 1 Introduction

Given a (finite) word  $w = w_1w_2 \dots w_n$  over some alphabet  $\Sigma$ , we say that  $u$  is a (*scattered*) *subword* or *subsequence* of  $w$  if  $u$  can be obtained from  $w$  by erasing some letters of  $w$ . We denote the fact that  $u$  is a subword of  $w$  by  $u \preceq w$ , and alternatively say that  $w$  is a *superword* of  $u$ . As shown by Higman [11] in 1952  $\preceq$  is a well-quasi-order on  $\Sigma^*$ , implying that *every* language  $L \subseteq \Sigma^*$  has a finite set of  $\preceq$ -minimal elements. This proves that both the subword (also: downward) closure  $\nabla L := \{u \in \Sigma^* \mid \exists w \in L: u \preceq w\}$  and the superword (also: upward) closure  $\Delta L := \{w \in \Sigma^* \mid \exists u \in L \mid u \preceq w\}$  are regular for *any* language  $L$ . While in general, we cannot effectively construct a finite automaton accepting  $\nabla L$  resp.  $\Delta L$ , for specific classes of languages effective constructions are known.

It is well-known that this is the case when  $L$  is given as a context-free grammar (CFG). This was first shown by van Leeuwen [13] in 1978. Later, Courcelle gave an alternative proof of this result in [6]. Section 3 builds up on these results by Courcelle. We also mention that for Petri-net languages an effective construction is known thanks to Habermehl, Meyer, and Wimmel [10].

These results can be used to tackle undecidable questions regarding the ambiguity, inclusion, equivalence, universality or emptiness of languages by overapproximating one or both languages by suitable regular languages [15,14,8,10]:

---

\* This work was partially funded by the DFG project “Polynomial Systems on Semirings: Foundations, Algorithms, Applications”.

For instance, consider the scenario where we are given a procedural program whose runs can be described as a pushdown automaton resp. a CFG  $G_1$  and a context-free specification  $G_2$  of all safe executions, and we want to check whether all runs of the system conform to the safety specification  $\mathcal{L}(G_1) \subseteq \mathcal{L}(G_2)$ . As  $\mathcal{L}(G_1) \cap \overline{\mathcal{L}(G_2)} \neq \emptyset \Rightarrow \mathcal{L}(G_1) \not\subseteq \mathcal{L}(G_2)$ , we can obtain at least a partial answer to the otherwise undecidable question. Of course, in the case  $\mathcal{L}(G_1) \subseteq \mathcal{L}(G_2)$  no information is gained, and one needs to refine the problem e.g. by using some sort of counter-example guided abstraction refinement as done e.g. in [14].

*Contributions and Outline* Our first results (Sections 3 and 4) concern the blow-up incurred when constructing a (non-)deterministic finite automaton (NFA resp. DFA) for the subword closure of a language given by a context-free grammar  $G$  where we improve the results of [9]: For a CFG  $G$  of size  $n$ , [9] shows that an NFA recognizing  $\nabla\mathcal{L}(G)$  has at most  $2^{2^{\mathcal{O}(n)}}$  states, and there are CFGs requiring at least  $2^{\Omega(n)}$  states. (For linear CFGs the upper and lower bounds are both single exponential.) The upper bound of [9] is established by analyzing the inductive construction of [13]. We improve this result in Section 3 to  $2^{\mathcal{O}(n)}$  by slightly adapting Courcelle’s construction [6] (we also briefly discuss that naively applying Courcelle’s construction cannot do better than  $2^{\Omega(n \log n)}$  in general). This result of course yields immediately an upper bound of  $2^{2^{\mathcal{O}(n)}}$  on the size of minimal DFA representing  $\nabla\mathcal{L}(G)$ . In Section 4 we show this bound is tight already over a binary alphabet. To the best of our knowledge, so far only examples were known which showcase the single-exponential blow-up when constructing an NFA accepting the subword closure of a context-free grammar [9] resp. a DFA accepting the subword closure of a DFA or NFA [17]. We then study in Section 5 the equivalence problem for NFAs recognizing subword- resp. supword-closed languages. While for general NFAs this problem is PSPACE-complete, we show that it becomes coNP-complete under this restriction. We combine these results in Section 6 to derive a conceptual simple semi-decision procedure for checking language-inequivalence of two CFGs  $G_1, G_2$ : we first construct NFAs for  $\nabla\mathcal{L}(G_1)$  and  $\nabla\mathcal{L}(G_2)$ , and check language-inequivalence of these NFAs; if the NFAs are inequivalent, we construct a witness of the language-inequivalence of  $G_1$  and  $G_2$ ; otherwise we refine the grammars, and repeat the test on the so obtained new grammars. This approach is motivated by the abstraction-refinement approach of [14] for checking if the intersection of two context-free languages is empty. We experimentally evaluate our approach by comparing it to *cfg-analyzer* of [2] which uses incremental SAT-solving to tackle the language-inequivalence problem. Missing proofs can be found in the extended version of the paper [3].

## 2 Preliminaries

By  $\Sigma$  we denote a finite alphabet. For every natural number  $n$ , let  $\Sigma^{\leq n}$  denote the words of length at most  $n$  over  $\Sigma$ . The empty word is denoted by  $\varepsilon$ ; the set of all finite words by  $\Sigma^*$ .

We measure the *size*  $|G|$  of a CFG  $G$  as the total number of symbols on the right hand sides of all productions. The size of an NFA is simply measured as the number of states (this is an adequate measure for a constant alphabet, since the number of transitions is at most quadratic in the number of states).

Throughout the paper we will always assume that all CFGs are reduced, i.e. do not contain any unproductive or unreachable nonterminals (any CFG can be reduced in polynomial time). Let  $X$  be a nonterminal in a CFG  $G$ . We define  $\mathcal{L}(X)$  as the set of all words  $w \in \Sigma^*$  derivable from  $X$ . If  $S$  is the start symbol of  $G$ , then  $\mathcal{L}(G) := \mathcal{L}(S)$ . Moreover,  $\Sigma_X \subseteq \Sigma$  denotes the set of all terminals reachable from  $X$ . Overloading notation we sometimes write  $\nabla X$  for  $\nabla \mathcal{L}(X)$ .

The dependency graph of a CFG  $G$  is the finite graph with nodes the nonterminals of  $G$  where there is an edge from  $X$  to  $Y$  if there is a production  $X \rightarrow \alpha Y \beta$  in  $G$ . We say that  $X$  *depends directly on*  $Y$  (written as  $X \triangleright Y$ ) if  $X \neq Y$  and there is an edge from  $X$  to  $Y$ . The reflexive and transitive closure of  $\triangleright$  is denoted by  $\triangleright^*$ . We write  $X \equiv Y$  if  $X \triangleright^* Y \wedge Y \triangleright^* X$ , i.e. if  $X$  and  $Y$  are located in a common strongly-connected component of the dependency graph. We say that  $G$  is strongly connected if the dependency graph is strongly connected.

From [6] we recall some useful facts concerning the subword closure:

**Lemma 1.** *For any nonterminals  $X, Y, Z$  in a CFG  $G$  it holds that:*

1.  $\nabla(\mathcal{L}(X) \cup \mathcal{L}(Y)) = \nabla \mathcal{L}(X) \cup \nabla \mathcal{L}(Y)$
2.  $\nabla(\mathcal{L}(X) \cdot \mathcal{L}(Y)) = \nabla \mathcal{L}(X) \cdot \nabla \mathcal{L}(Y)$
3.  $X \equiv Y \Rightarrow \nabla X = \nabla Y$
4. If  $X \rightarrow^* \alpha Y \beta Z \gamma$  for  $Y, Z \equiv X$  then  $\nabla X = \Sigma_X^*$

### 3 Computing the Subword Closure of CFGs

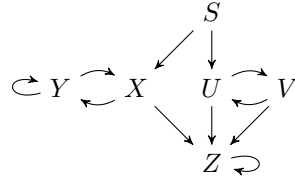
In this section we describe an optimized version of the construction in [6] to compute an NFA for the subword closure of a CFG  $G$  of size  $2^{\mathcal{O}(|G|)}$ , which is asymptotically optimal. We first illustrate the construction by a simple example.

As explained at the end of the next section, a naive implementation of the construction of [6] leads to an automaton of size  $2^{\Omega(n)} n! = 2^{\Omega(n \log n)}$  whereas our approach achieves the (optimal) bound of  $2^{\mathcal{O}(n)}$ .

#### 3.1 Construction by Example

Consider the grammar  $G$  with start symbol  $S$  defined by the productions:

$$\begin{array}{ll}
 S \rightarrow XaU \mid UaU \mid X & X \rightarrow ZbY \mid \varepsilon \\
 Y \rightarrow XYa \mid b & U \rightarrow VZ \mid acb \\
 V \rightarrow ZU \mid \varepsilon & Z \rightarrow cZ \mid bc
 \end{array}$$



On the right-hand side, the dependency graph is shown where an edge  $x \rightarrow y$  stands for  $x \triangleright y$ . To simplify the construction, we first transform the grammar

$G$  into a certain normal form  $G'$  (with  $\nabla\mathcal{L}(G) = \nabla\mathcal{L}(G')$ ) and then construct an NFA from  $G'$ .

In the first step we compute the strongly connected components (SCCs) of  $G$ , here  $\{X, Y\}$  and  $\{U, V\}$ . Since  $Y \rightarrow XYa$  (with  $Y \equiv X$  and  $X \equiv X$ ), we know that  $\nabla Y = \nabla X = \Sigma_X^* = \{a, b, c\}^*$ . We therefore can replace any occurrence of  $Y$  by  $X$  (thereby removing  $Y$  from the grammar) and redefine the rules for  $X$  to  $X \rightarrow aX \mid bX \mid cX \mid \varepsilon$ . In case of the SCC  $\{U, V\}$  the grammar is linear w.r.t.  $U$  and  $V$ , i.e. starting from either of the two we can never produce sentential forms in which the total number of occurrences of  $U$  and  $V$  exceeds one. Hence, we can identify  $U$  and  $V$  without changing the subword closure. Finally, we introduce unique nonterminals for each terminal symbol and restrict the right-hand side of each production to at most two symbols by introducing auxiliary nonterminals  $W$  and  $T$ :

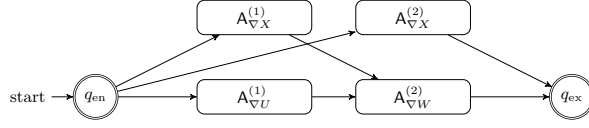
$$\begin{array}{ll}
S \rightarrow XW \mid UW \mid X & W \rightarrow A_a U \\
X \rightarrow A_a X \mid A_b X \mid A_c X \mid A_\varepsilon & U \rightarrow UZ \mid ZU \mid A_a T \mid A_\varepsilon \\
T \rightarrow A_c A_b & Z \rightarrow A_c Z \mid A_b A_c \\
A_a \rightarrow a & A_b \rightarrow b \\
A_c \rightarrow c & A_\varepsilon \rightarrow \varepsilon
\end{array}$$

Note that the dependency graph of this transformed grammar is now acyclic apart from self-loops. Because of this, we can directly transform the grammar into an *acyclic* equation system (or straight-line program, or algebraic circuit) whose solution is a regular expression for  $\nabla S$ :

$$\begin{array}{ll}
\nabla A_a = (a + \varepsilon) & \nabla A_b = (b + \varepsilon) \\
\nabla A_c = (c + \varepsilon) & \nabla A_\varepsilon = \varepsilon \\
\nabla Z = c^*(\nabla A_b \nabla A_c) & \nabla T = \nabla A_c \nabla A_b \\
\nabla U = \Sigma_Z^*(\nabla A_a \nabla T) \Sigma_Z^* & \nabla W = \nabla A_a \nabla U \\
\nabla X = \Sigma_X^* & \nabla S = \nabla X \nabla W + \nabla U \nabla W + \nabla X
\end{array}$$

In order to obtain an NFA for  $\nabla S$ , we evaluate this equation system from bottom to top while re-using as many of the already constructed automata as possible. For instance, consider the equation:  $\nabla S = \nabla X \nabla W + \nabla U \nabla W + \varepsilon \cdot \nabla X$ . Because of acyclicity of the equation system, we may assume inductively that we have already constructed NFAs  $A_{\nabla X}$ ,  $A_{\nabla W}$ , and  $A_{\nabla U}$  for  $\nabla X$ ,  $\nabla W$ , and  $\nabla U$ , respectively. To construct the NFA for  $\nabla S$ , we first make two copies  $A^{(1)}$ ,  $A^{(2)}$  of each of these automata. Automata with superscript (1) will be used exclusively for variable occurrences to the left of the concatenation operator, while automata with superscript (2) will be used for the remaining occurrences. We then read quadratic monomials, like  $\nabla X \nabla W$ , as an  $\varepsilon$ -transition connecting  $A_{\nabla X}^{(1)}$  with  $A_{\nabla W}^{(2)}$  as shown in Figure 1 where all edges represent  $\varepsilon$ -transitions.

We do not claim that this construction yields the smallest NFA, but it is easy to describe and yields an NFA of sufficiently small size in order to deduce in the following subsections an asymptotically tight upper bound on the number of states. We recall that using a CFG of size  $3n + 2$  to succinctly represent the singleton language  $\{a^{2^n}\}$ , the bound of  $2^{\Theta(n)}$  follows [9].



**Fig. 1.** Efficient re-use of re-occurring NFAs in Courcelle’s construction.

In [1] it is remarked that a straight-forward implementation of Courcelle’s construction yields an NFA “single exponential” size w.r.t.  $|G|$ . However, no detailed complexity analysis is given. Consider the CFG with start-symbol  $A_n$  and consisting of the rules  $A_0 \rightarrow a$  and for all  $1 \leq k \leq n$  :  $A_k \rightarrow A_i A_j \quad \forall 0 \leq i, j \leq (k-1)$ . If we compute an NFA for  $\nabla A_n$  via the straight-forward bottom-up construction it will have size  $a_n := |\mathbf{A}_{\nabla A_n}|$  with  $a_n = 2 + \sum_{0 \leq i, j \leq (n-1)} (a_i + a_j)$ . It is easy to show that  $a_n \geq 2^n n! \in 2^{\Omega(n \log n)}$ . Hence, the crucial part to achieve the optimal bound of  $2^{\mathcal{O}(n)}$  is to reuse already computed automata. We just remark that one can also achieve similar savings by factoring out common terms in the right hand side of the acyclic equations. A subsequent bottom-up construction leads to an NFA of size  $2^{\mathcal{O}(n)}$  as well but the constant hidden in the  $\mathcal{O}$  is larger and the analysis is more involved. Note that this also shows that we can construct a regular expression of size  $2^{\mathcal{O}(n)}$  representing the subword closure.

### 3.2 Normal Form for Computing the Subword Closure

To simplify our construction, we will assume that our grammar has a special form which is similar to CNF but with unary rules allowed. Any CFG can be transformed into this form with at most linear blowup in size preserving its subword closure (but not its language).

**Definition 2.** A CFG  $G$  is in quadratic normal form (QNF) if for every terminal  $x \in \Sigma \cup \{\varepsilon\}$  there is a unique nonterminal  $A_x$  with the only production  $A_x \rightarrow x$  and every other production is in one of the following forms:

- $X \rightarrow YX$  or  $X \rightarrow XY$  (with  $Y \neq X$ )
- $X \rightarrow Y$  or  $X \rightarrow YZ$  (with  $Y, Z \neq X$ )

A grammar in QNF is called simple if

- for all  $X \rightarrow YX$  or  $X \rightarrow XY$ , we have  $X \triangleright Y$
- for all  $X \rightarrow Y$  or  $X \rightarrow YZ$ , we have  $X \triangleright Y, Z$ .

Note that the dependency graph associated with a grammar in simple QNF is acyclic with the exception of self-loops.

First, we need a small lemma that allows us to eliminate all linear productions “within” some SCC, i.e. productions of the form  $X \rightarrow \alpha Y \beta$  such that  $X \neq Y$  but  $Y \succeq^* X$ .

**Lemma 3.** *Let  $G$  be a strongly connected linear CFG with nonterminals  $\mathcal{X} = \{X_1, \dots, X_n\}$  so that every production is either of the form  $X \rightarrow \alpha Y \beta$  or  $X \rightarrow \alpha$  for  $\alpha, \beta \in \Sigma^*$ . Consider the grammar  $G'$  which we obtain from  $G$  by replacing in every production of  $G$  every occurrence of a nonterminal  $X_i$  by  $Z$ . We then have that  $\nabla \mathcal{L}(Z) = \nabla \mathcal{L}(X_i)$  for all  $i \in [n]$ .*

Using the preceding lemma, we can show that it suffices to consider only CFG in simple QNF in the following.

**Theorem 4.** *Every CFG  $G$  can be transformed into a CFG  $G'$  in simple QNF such that  $\nabla \mathcal{L}(G) = \nabla \mathcal{L}(G')$  and  $|G'| \in \mathcal{O}(|G|)$ .*

*Proof (sketch).* First, we use Lemma 1 to simplify all productions involving an  $X$  with  $X \Rightarrow^* \alpha X \beta X \gamma$ . Then we apply Lemma 3 to contract SCCs to a single non-terminal. Finally, we introduce auxiliary variables for the terminals and we binarize the grammar (keeping unary rules like [12]).

**Theorem 5.** *For any CFG  $G$  in simple QNF with  $n$  nonterminals there is an NFA  $A$  with at most  $2 \cdot 3^{n-1}$  states which recognizes the subword closure of  $G$ , i.e.  $\nabla \mathcal{L}(G) = \mathcal{L}(A)$ .*

*Proof (sketch).* Since the dependency graph of a grammar in simple QNF is a DAG (if we ignore self-loops), we can order the nonterminals according to a topological ordering of this graph. We proceed bottom-up to inductively build an NFA for  $\nabla \mathcal{L}(G) = \nabla S$  as in section 3.1. Since our grammar is in QNF, at each stage we only have to produce at most two copies of every automaton representing the subword-closure of a “lower” nonterminal  $Y$ . Inductively, for each of these  $Y$  we can build an NFA with at most  $2 \cdot 3^i$  many states where  $i$  is  $Y$ ’s position in the topological ordering. Using the “bipartite wiring” sketched in Figure 1 the size of the automaton for  $X$  can then be estimated as

$$|A_S| \leq 2 + \sum_{Y: S \triangleright Y} 2 \cdot |A_Y| \leq 2 + 4 \cdot \sum_{i=0}^{n-2} 3^i = 2 \cdot 3^{n-1}.$$

**Corollary 6.** *For every CFG  $G$  of size  $n$  there is an NFA  $A$  of size  $2^{\mathcal{O}(n)}$  and a DFA  $D$  of size  $2^{2^{\mathcal{O}(n)}}$  with  $\nabla \mathcal{L}(G) = \mathcal{L}(A) = \mathcal{L}(D)$ .*

## 4 CFG $\rightarrow$ DFA: Double-exponential Blowup

As seen in the preceding section, moving from a context-free grammar  $G$  representing a subword-closed language to a language-equivalent NFA  $\mathcal{A}$ , the size of the automaton is bounded from above by  $2^{\mathcal{O}(|G|)}$ . For superword closures [9] prove the same upper bound for the size of the NFA. From both results we immediately obtain the upper bound  $2^{2^{\mathcal{O}(|G|)}}$  on the size of the minimal language-equivalent DFA recognizing the sub- or superword closure of a CFG  $G$ . This bound is essentially tight as witnessed by the family of finite languages

$$L_k = \bigcup_{j=1}^k \{0, 1\}^{j-1} \{0\} \{0, 1\}^k \{0\} \{0, 1\}^{k-j}.$$

$L_k$  contains exactly all those words  $w \in \{0, 1\}^{2^{k+1}}$  which contain two 0s which are separated by exactly  $k$  letters. Using the idea of iterated squaring in order to succinctly encode the language  $\{a^{2^n}\}$  as a context-free grammar (resp. straight-line program) of size  $\mathcal{O}(n)$ , the language  $L_{2^n}$  can be represented by a context-free grammar of size  $\mathcal{O}(n)$  as well. One then easily shows that the Myhill-Nerode relation w.r.t.  $L_{2^n}$ ,  $\nabla L_{2^n}$ , and  $\Delta L_{2^n}$ , respectively, has at least  $2^{2^n}$  equivalence classes:

**Theorem 7.** *There exists a family of CFGs  $G_n$  of size  $\mathcal{O}(n)$  (generating finite languages) such that the minimal DFAs accepting either  $L(G_n)$ , or  $\nabla L(G_n)$ , or  $\Delta L(G_n)$ , have at least  $2^{2^n}$  states.*

## 5 Equivalence of NFAs modulo Sub-/Superword Closure

As hinted at in the introduction, one application of the sub- resp. superword closure is (in-)equivalence checking of CFGs by regular over-approximation. For this, we must solve the equivalence problems for NFAs representing sub/superword closed languages. Naturally, the question arises how hard this is.

Let  $A$  and  $B$  denote NFAs over the common alphabet  $\Sigma$ , having  $n_A$  and  $n_B$  many states, respectively. Recall that the universality problem for NFAs, i.e.  $\mathcal{L}(A) \stackrel{?}{=} \Sigma^*$ , and hence also the equivalence problem  $\mathcal{L}(A) \stackrel{?}{=} \mathcal{L}(B)$  are PSPACE-complete. Only recently, it was shown in [18] that these problems *stay* PSPACE-complete even when restricted to NFAs representing languages which are closed w.r.t. either prefixes or suffixes or factors. However, in [18] it was also shown that for subword-closed NFAs (i.e.  $\nabla \mathcal{L}(A) = \mathcal{L}(A)$ ), universality is decidable in linear time as  $\mathcal{L}(A) = \Sigma^*$  holds if and only if there is an SCC in  $A$  whose labels cover all of  $\Sigma$ . It is easily shown that a similar result also holds for superword-closed NFAs (i.e.  $\Delta \mathcal{L}(A) = \mathcal{L}(A)$ ): We have  $\mathcal{L}(A) = \Sigma^*$  if and only if  $\varepsilon \in \mathcal{L}(A)$ .

In this section we show that both equivalence problems, i.e.  $\nabla \mathcal{L}(A) \stackrel{?}{=} \nabla \mathcal{L}(B)$  and  $\Delta \mathcal{L}(A) \stackrel{?}{=} \Delta \mathcal{L}(B)$ , are coNP-complete, hence are easier than in the general case (unless  $\text{NP} = \text{PSPACE}$ ). In the following, we write more succinctly  $A \stackrel{?}{\equiv}_{\nabla} B$  and  $A \stackrel{?}{\equiv}_{\Delta} B$  for these two problems. The following lemma is easy to prove:

**Lemma 8.** *Let  $A$  be an NFA. Define  $A^{\nabla}$  as the NFA we obtain from  $A$  by adding for every transition  $q \xrightarrow{a} q'$  of  $A$  the  $\varepsilon$ -transition  $q \xrightarrow{\varepsilon} q'$ . Similarly, define  $A^{\Delta}$  to be the NFA we obtain by adding the loops  $q \xrightarrow{a} q$  for every state  $q$  and every terminal  $a \in \Sigma$  to  $A$ . Then  $\nabla \mathcal{L}(A) = \mathcal{L}(A^{\nabla})$  and  $\Delta \mathcal{L}(A) = \mathcal{L}(A^{\Delta})$ .*

To prove that both  $A \stackrel{?}{\equiv}_{\Delta} B$  and  $A \stackrel{?}{\equiv}_{\nabla} B$  are coNP-complete we will give a polynomial bound on the length of a *separating word*, i.e. a word  $w$  in the symmetric difference of  $\mathcal{L}(A^{\nabla})$  and  $\mathcal{L}(B^{\nabla})$  resp. of  $\mathcal{L}(A^{\Delta})$  and  $\mathcal{L}(B^{\Delta})$ .

We first show that the DFA obtained from  $A^{\nabla}$  resp.  $A^{\Delta}$  using the powerset construction has a particular simple structure (this was also observed in [17]).

**Lemma 9.** *Let  $A$  be an NFA. Let  $D_A^\nabla$  (resp.  $D_A^\Delta$ ) be the DFA we obtain from  $A^\nabla$  (resp.  $A^\Delta$ ) by means of the powerset construction. For any transition  $S \xrightarrow{a} T$  of  $D_A^\nabla$  ( $D_A^\Delta$ ) it holds that  $S \supseteq T$  (resp.  $S \subseteq T$ ).*

Thus, the transition relation of  $D_A^\nabla$  (disregarding self-loops) can be “embedded” into the lattice of subsets of the states of  $A$ , which has height  $n_A - 1$ .

**Corollary 10.** *With the assumptions of the preceding lemma: The length of the longest simple path in  $D_A^\nabla$  (resp.  $D_A^\Delta$ ) is at most  $n_A - 1$ .*

It now immediately follows that a shortest separating word for sub- resp. supword closed NFAs – if one exists – has at most length linear in the size of the two NFAs.

**Lemma 11.** *Let  $A$  and  $B$  be two NFAs. If  $A \not\equiv_\nabla B$  (resp.  $A \not\equiv_\Delta B$ ), then there exists a separating word of length at most  $n_A + n_B - 2$ .*

**Theorem 12.** *The decision problems  $A \stackrel{?}{\equiv}_\nabla B$  and  $A \stackrel{?}{\equiv}_\Delta B$  are in coNP.*

To show coNP-hardness, recall the proof that the equivalence problem for star-free regular expressions is coNP-hard by reduction from TAUT: Given a formula  $\phi$  in propositional calculus, we build a regular expression  $\rho$  (without Kleene stars) over  $\Sigma = \{0, 1\}$  that enumerates exactly the satisfying assignments of  $\phi$ . Hence,  $\phi \in \text{TAUT}$  iff  $\mathcal{L}(\rho) = \Sigma^n$  iff  $\nabla\mathcal{L}(\rho) = \Sigma^{\leq n}$ , since the subword closure can only add new words of length less than  $n$  (analogously for  $\Delta$ ).

**Theorem 13.** *The decision problems  $A \stackrel{?}{\equiv}_\nabla B$  and  $A \stackrel{?}{\equiv}_\Delta B$  are coNP-hard.*

## 6 Application to Grammar Problems

We apply our results to devise an approximation approach for the well-known undecidable problem whether  $\mathcal{L}(G_1) = \mathcal{L}(G_2)$  for two CFGs  $G_1, G_2$ . Possible attacks on this problem include exhaustive search for a word in the symmetric difference  $w \in (L_1 \oplus L_2) \cap \Sigma^{\leq n}$  w.r.t. some increasing bound  $n$  e.g. by using incremental SAT-solving [2]. Unfortunately, this quickly becomes infeasible for large problems. Previous work has successfully applied regular approximation for ambiguity detection [19,5] or intersection non-emptiness of CFGs [14].

A high-level description of our approach to (in-)equivalence-checking is given in Figure 2. Of course the procedure will not terminate if  $\mathcal{L}(G_1) = \mathcal{L}(G_2)$ , so in practice a timeout will be used after which the algorithm will terminate itself and output “Maybe equal”. Steps (1) and (2) might take time (at most) double exponential in the size of the grammars  $G_1$  and  $G_2$ : Recall that the construction of Section 3 yields in the worst-case an NFA  $A_i$  whose number of states is exponential in the size of the given CFG  $G_i$ . To check if  $\nabla\mathcal{L}(G_1) = \nabla\mathcal{L}(G_2)$ , an on-the-fly construction of the power-set automaton for  $A_1 \times A_2$  can be used which terminates as soon as a set of states is reached which contains at least one accepting state of, say,  $A_1$  but no accepting state of  $A_2$ . Using Lemma 11, we



1. Compute NFAs  $A_1$  and  $A_2$  for the subword closures of  $G_1$  and  $G_2$ , respectively.
2. Check, if  $\mathcal{L}(A_1) = \mathcal{L}(A_2)$ .
  - (a) Case “Not equal”: Generate a witness  $w \in \mathcal{L}(G_1) \oplus \mathcal{L}(G_2)$ .
  - (b) Case “Equal”: Refine the grammars and restart at **1**.

**Fig. 2.** Equivalence checking via subword closure approximation.

can safely terminate the exploration of simple paths if their length exceeds the bound stated in Lemma 11. In the worst case this might take time exponential in the size of  $A_1$  and  $A_2$ , so at most double exponential in the size of  $G_1$  and  $G_2$ .

In the following, we describe in greater detail how we generate a separating word  $w'$  in  $\mathcal{L}(G_1)$  or  $\mathcal{L}(G_2)$  if we find a separating word  $w \in \nabla\mathcal{L}(G_1) \oplus \nabla\mathcal{L}(G_2)$ , resp. how we refine  $G_1$  and  $G_2$  if  $\nabla\mathcal{L}(G_1) = \nabla\mathcal{L}(G_2)$ .

### 6.1 Witness Generation for $\mathcal{L}(G_1) \neq \mathcal{L}(G_2)$

If our check in step (2) returns “Not equal” we know that  $\nabla\mathcal{L}(G_1) \neq \nabla\mathcal{L}(G_2)$  and we obtain a word  $w \in \nabla\mathcal{L}(G_1) \oplus \nabla\mathcal{L}(G_2)$ , w.l.o.g. assume in the following  $w \in \nabla\mathcal{L}(G_1) \setminus \nabla\mathcal{L}(G_2)$ . This word has length linear in  $|A_1|$  and  $|A_2|$ , i.e. at most exponential w.r.t.  $|G_1|$  and  $|G_2|$ .

To obtain a (direct) certificate for the fact that  $\mathcal{L}(G_1) \neq \mathcal{L}(G_2)$ , we construct a superword  $w' \succ w$  with  $w' \in \mathcal{L}(G_1)$  – such a  $w'$  is guaranteed to exist as it is the reason for  $w \in \nabla\mathcal{L}(G_1)$ . Straight-forward induction on  $w$  shows:

**Lemma 14.** *For  $w \in \Sigma^*$  a DFA recognizing  $\nabla\mathcal{L}(\{w\})$  resp.  $\Delta\mathcal{L}(\{w\})$  and having at most  $|w| + 2$  states can be constructed in time polynomial in  $|w|$ .*

We can therefore intersect  $G_1$  with a DFA accepting  $\Delta\mathcal{L}(\{w\})$ , to obtain a new CFG  $G'_1$  whose size is at most cubic in  $|w|$  [4,16], i.e. exponential in the size of  $G_1$ . From this grammar, we can obtain in time linear in  $|G'_1|$  a shortest word  $w'$  in  $\mathcal{L}(G'_1) = \mathcal{L}(G_1) \cap \Delta\mathcal{L}(\{w\})$ . The length of  $w'$  is at most exponential in  $|G'_1|$ , i.e. at most double exponential in  $|G_1|$ .

In practice, shorter witnesses are preferable, so we construct the shortest word in  $\overline{\mathcal{L}(A_2)} \cap \mathcal{L}(G_1)$ . In theory this might incur a triple exponential blow-up resulting from complementing  $A_2$ , but this way we can find a separating word  $w'$  which is *not* a superword of  $w$  and hence is usually shorter.

### 6.2 Refinement

In case that the test in step (2) returns “Equal”, we refine both grammars such that subsequent subword-approximations may find a counterexample to equality. Assume that our equivalence check yields  $\nabla\mathcal{L}(G_1) = \nabla\mathcal{L}(G_2)$ . A possible refinement strategy is to cover  $L := \nabla\mathcal{L}(G_1)$  using a finite number of regular languages  $L \subseteq L' := L_0 \cup L_1 \cup \dots \cup L_k$  and then to repeat the equivalence check for all pairs of refined languages  $\mathcal{L}(G_1) \cap L_i$  and  $\mathcal{L}(G_2) \cap L_i$  for all  $i$ . The requirement  $L' \supseteq L$  protects the refinement from cutting off potential witnesses.

A simple method is covering using prefixes: Here we generate all prefixes  $p_1, \dots, p_k$  of words in  $L$  of increasing length (up to some small bound  $d$  called the *refinement depth*) and set  $L_i := p_i \Sigma^*$  and  $L_0 = \nabla\{p_i \mid i \in [k]\}$ . Since  $\bigcup_i L_i \supseteq L$  this strategy preserves potential witnesses and since any counterexample eventually appears as a prefix, this yields a semi-decision procedure for grammar inequivalence. In our experiments we disregard the finite language  $L_0$  (which can also be checked by enumeration) and only check refinement using the infinite sets  $p_i \Sigma^*$  with the goal of quickly finding *some* (not the shortest) distinguishing word. This strategy is often able to tell apart different CFLs after few iterations as shown in the following.

### 6.3 Implementation and Experiments

We implemented the inequivalence check in an extension<sup>1</sup> of the FPSOLVE tool [7]. The additional code comprises roughly 1800 lines of C++ and uses libfa<sup>2</sup> to handle finite automata.

Our worst-case descriptonal complexity results for the subword closure of CFGs (exponential sized NFA, double-exponential sized DFA) and our remarks on the length of possible counterexamples might suggest that our inequivalence checking procedure is merely of academic interest. Here we briefly show that this is not the case, and that overapproximation via subword closures is actually quite fast in practice.

The paper [2] presents cfg-analyzer, a tool that uses SAT-solving to attack several undecidable grammar problems by exhaustive enumeration. We demonstrate the feasibility of our approximation approach on several slightly altered grammars (cf. [20]) for the PASCAL programming language<sup>3</sup>. The altered grammars were obtained by adding, deleting, or mutating a single rule from the original grammar [20]. We used FPSOLVE and cfg-analyzer to check equivalence of the altered grammar with the original. Both tools were given a timeout of 30 seconds. We want to stress that we do not strive to replace enumeration-based tools like cfg-analyzer, but rather envision a combined approach: Use overapproximations like the subword closure (with small refinement depth) as a quick check and resort to more computationally demanding techniques like SAT-solving for a thorough test. Also note that it is not too hard to find examples where enumeration-based tools cannot detect inequivalence anymore, e.g. by considering grammars with large alphabet (like C# or Java) for which the shortest word in the language is already longer than 20 tokens. Here we just showcase an example where both approaches can be fruitfully combined.

Table 1 demonstrates that even if our tool uses the very simple prefix-refinement (which is the main bottleneck in terms of speed), we can successfully solve 100 cases where cfg-analyzer has to give up after 30 seconds and even in cases where both tools find a difference, FPSOLVE does so much faster.

<sup>1</sup> The fork is available from <https://github.com/regularApproximation/newton>.

<sup>2</sup> <http://augeas.net/libfa/>

<sup>3</sup> Available from <https://github.com/nvasudevan/experiment/tree/master/grammars/mutlang/acc>.

scenario	# instances	# CA	$t_{CA}$	#FP	$t_{FP}$	$\#(CF \wedge FP)$	$t_{CA}^\wedge$	$t_{FP}^\wedge$
add	700	190	17.9	18	2.43	8	10.7	4.97
delete	284	61	17.8	34	0.424	10	14.4	0.464
empty	69	32	18.7	1	1.35	1	5.62	1.35
mutate	700	167	19.1	100	1.3	36	15.8	2.87
switchadj	187	16	20.5	2	5.46	1	9.68	0.34
switchany	328	35	18	9	3.72	8	9.09	2.84
$\Sigma$	2268	501	–	164	–	64	–	–

**Table 1.** Numbers of solved instances for different scenarios and respective average times: #CA: solved by cfg-analyzer, #FP: solved by FPSOLVE,  $\#(CA \wedge FP)$ : solved by both tools,  $t_{tool}^\wedge$ : time needed by *tool* on instances from  $(CA \wedge FP)$ .

## 7 Discussion and Future Work

Motivated by the language-equivalence problem for context-free languages, we have studied the problems of the space requirements of representing the subword closure of CFGs by NFAs and DFAs, and the computational complexity of the equivalence problem of subword-closed NFAs. We have shown how to construct from a context-free grammar  $G$  an NFA accepting  $\nabla\mathcal{L}(G)$  consisting of at most  $2^{\mathcal{O}(|G|)}$  states – a small gap between the lower bound of  $\Omega(2^{|G|})$  and our upper bound of  $\mathcal{O}(3^{|G|})$  for grammars in QNF remains for future work. A further question is if this bound can be improved in the case of languages given as deterministic pushdown automata. We have further shown that the upper bound on the size of a DFA accepting  $\nabla\mathcal{L}(G)$  of  $2^{\mathcal{O}(|G|)}$  is tight. Interestingly, a binary alphabet suffices for the presented language family  $L_k$ : for instance the worst-case example of [17], which showcases the exponential blow-up suffered when constructing an DFA for the subword closure of a language given as DFA or NFA, requires an unbounded alphabet. We note that a unary context-free language cannot lead to this double exponential blow-up – this follows from the proof of Theorem 3.14 in [9] (see also Lemma 14 here). Regarding the language-equivalence problem, we have shown that it becomes coNP-complete when restricted to subresp. superword-closed NFAs. This is somewhat surprising given the fact that it stays PSPACE-complete for many related families (e.g. for prefix-, suffix-, or factor-closed languages). Finally, we have briefly described an approach to tackle the equivalence problem for CFGs using the presented results, though much work remains to turn our current implementation into a mature tool: In particular, since the intersection of two regular overapproximations is again a regular overapproximation, it could be fruitful to combine the subword closure (or variants like [14]) with other regular approximation techniques like [15]. We also need to improve the refinement of the approximations when scaling the problem size.

## References

1. Atig, M.F., Bouajjani, A., Touili, T.: On the Reachability Analysis of Acyclic Networks of Pushdown Systems. In: CONCUR 2008. pp. 356–371 (2008)
2. Axelsson, R., Heljanko, K., Lange, M.: Analyzing Context-Free Grammars Using an Incremental SAT Solver. In: ICALP (2). pp. 410–422 (2008)
3. Bachmeier, G., Luttenberger, M., Schlund, M.: Finite Automata for the Sub- and Superword Closure of CFLs: Descriptive and Computational Complexity (extended version). arXiv:1410.2737
4. Bar-Hillel, Y., Perles, M., Shamir, E.: On Formal Properties of Simple Phrase Structure Grammars. *Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung* 14, 143–172 (1961)
5. Brabrand, C., Giegerich, R., Møller, A.: Analyzing Ambiguity of Context-Free Grammars. *Sci. Comput. Program.* 75(3), 176–191 (2010)
6. Courcelle, B.: On Constructing Obstruction Sets of Words. *Bulletin of the EATCS* 44, 178–186 (1991)
7. Esparza, J., Luttenberger, M., Schlund, M.: FPSolve: A Generic Solver for Fixpoint Equations over Semirings. In: CIAA, LNCS, vol. 8587, pp. 1–15 (2014)
8. Ganty, P., Majumdar, R., Monmege, B.: Bounded underapproximations. *Formal Methods in System Design* 40(2), 206–231 (2012)
9. Gruber, H., Holzer, M., Kutrib, M.: More on the Size of Higman-Haines Sets: Effective Constructions. *Fundam. Inf.* 91(1), 105–121 (Jan 2009)
10. Habermehl, P., Meyer, R., Wimmel, H.: The Downward-Closure of Petri Net Languages. In: ICALP (2). pp. 466–477 (2010)
11. Higman, G.: Ordering by Divisibility in Abstract Algebras. *Proc. London Math. Soc.* s3-2(1), 326–336 (Jan 1952)
12. Lange, M., Leiß, H.: To CNF or not to CNF? An Efficient Yet Presentable Version of the CYK Algorithm. *Informatica Didactica* 8 (2009)
13. van Leeuwen, J.: Effective constructions in well-partially-ordered free monoids. *Discrete Mathematics* 21(3), 237 – 252 (1978)
14. Long, Z., Calin, G., Majumdar, R., Meyer, R.: Language-Theoretic Abstraction Refinement. In: FASE. pp. 362–376 (2012)
15. Mohri, M., Nederhof, M.J.: Regular Approximation of Context-Free Grammars through Transformation. In: *Robustness in Language and Speech Technology, Text, Speech and Language Technology*, vol. 17, pp. 153–163 (2001)
16. Nederhof, M.J., Satta, G.: Probabilistic Parsing. In: *New Developments in Formal Languages and Applications*, pp. 229–258 (2008)
17. Okhotin, A.: On the State Complexity of Scattered Substrings and Superstrings. *Fundam. Inform.* 99(3), 325–338 (2010)
18. Rampersad, N., Shallit, J., Xu, Z.: The Computational Complexity of Universality Problems for Prefixes, Suffixes, Factors, and Subwords of Regular Languages. *Fundam. Inform.* 116(1-4), 223–236 (2012)
19. Schmitz, S.: Conservative Ambiguity Detection in Context-Free Grammars. In: ICALP. pp. 692–703 (2007)
20. Vasudevan, N., Tratt, L.: Detecting Ambiguity in Programming Language Grammars. In: *Software Language Engineering, LNCS*, vol. 8225, pp. 157–176 (2013)