

Rational index of languages with bounded dimension of parse trees^{*}

Ekaterina Shemetova^{1,2,3}, Alexander Okhotin¹, and Semyon Grigorev^{1,3}

¹ Department of Mathematics and Computer Science, St. Petersburg State University, 7/9 Universitetskaya nab., Saint Petersburg 199034, Russia

² St. Petersburg Academic University, ul. Khlopina, 8, Saint Petersburg 194021, Russia

³ JetBrains Research, Primorskiy prospekt 68-70, Building 1, St. Petersburg, 197374, Russia

Abstract. The rational index ρ_L of a language L is an integer function, where $\rho_L(n)$ is the maximum length of the shortest string in $L \cap R$, over all regular languages R recognized by n -state nondeterministic finite automata (NFA). This paper investigates the rational index of languages defined by (context-free) grammars with bounded tree dimension, and shows that it is of polynomial in n . More precisely, it is proved that for a grammar with tree dimension bounded by d , its rational index is $O(n^{2d})$, and that this estimation is precise, as there exists a grammar with rational index $\Theta(n^{2d})$.

Keywords. Dimension of a parse tree; rational index; CFL-reachability; parallel complexity; context-free languages; Datalog programs.

1 Introduction

The notion of a rational index was introduced by Boasson, Courcelle and Nivat [2] as a complexity measure for context-free languages. The rational index ρ_L of a language L is an integer function, where $\rho_L(n)$ is the maximum length of the shortest string in a language of the form $L \cap R$, where R is a regular language recognized by n -state nondeterministic finite automata (NFA), and the maximum is taken over all such languages R . The rational index plays an important role in determining the parallel complexity of such practical problems as the CFL-reachability problem and the more general Datalog query evaluation.

The rational index of a context-free language L is an integer function ρ_L , where $\rho_L(n)$ is the maximum length of the shortest string in the intersection of L with a regular language R recognized by a finite automaton with n states,

The *CFL-reachability problem* is stated as follows: given a context-free grammar G and an NFA A over the same alphabet, determine whether $L(G) \cap L(A)$ is non-empty. With A is regarded as a labelled graph, this is a kind of graph reachability problem with path constraints given by context-free languages. This is an important problem used in static code analysis [13] and graph database query evaluation [15].

^{*} Research supported by the Russian Science Foundation, project 18-11-00100.

The CFL-reachability problem is P-complete already for a fixed context-free grammar [6]. The question on the parallel complexity of this problem was investigated by Ullman and Van Gelder [14] in a much more general case, with a rich logic for database queries instead of grammars, and it was proved that under an assumption called the *polynomial fringe property* the problem is decidable in NC [14]. In the special case of grammars, the result of Ullman and Van Gelder [14] gives an NC^2 algorithm for the CFL-reachability problem, under the assumption that the grammar's rational index is polynomial.

Theoretical properties of the rational index have received some attention in the literature. Pierre and Farinone [12] proved that for every algebraic number γ , a language with the rational index in $\Theta(n^\gamma)$ exists. An upper bound on the rational index, shown by Pierre [11], is $2^{\Theta(n^2/\ln n)}$, and this bound is reached on the Dyck language on two pairs of parentheses. For several important sub-families of grammars, such as the linear and the one-counter languages, there are polynomial upper bounds on the rational index, which imply that the CFL-reachability problem is in NC^2 ; they can be proved to lie in NL by direct methods not involving the rational index [8, 9].

In this paper we investigate the rational index of a generalization of linear languages: the *languages of bounded tree dimension*, that is, those defined by grammars with a certain limit on branching in the parse trees; see Lutzenberger and Schlund [10] for a survey of tree dimension. Linear languages are languages of tree dimension 1, their rational index is known to be $O(n^2)$ [2]. The result of this paper is that languages of tree dimension bounded by d have rational index $O(n^{2^d})$, and furthermore, for every d , there is a language of tree dimension bounded by d with rational index $\Theta(n^{2^d})$.

2 Definitions

A (*context-free*) *grammar* is a quadruple $G = (\Sigma, N, R, S)$, where Σ is an alphabet; N is a set of nonterminal symbols; R is a set of rules, each of the form $A \rightarrow \alpha$, with $A \in N$ and $\alpha \in (\Sigma \cup N)^*$; and $S \in N$ is the start symbol. A parse tree is a tree, in which every leaf is labelled with a symbol from Σ , while every internal node is labelled with a nonterminal symbol $A \in N$ and has an associated rule $A \rightarrow X_1 \dots X_\ell \in R$, so that the node has ℓ ordered children labelled with X_1, \dots, X_ℓ . The language defined by each nonterminal symbol $A \in N$, denoted by $L_G(A)$, is the set of all strings $w \in \Sigma^*$, for which there exists a parse tree, with A as a root and with the leaves forming the string w . The language defined by the grammar is $L(G) = L_G(S)$.

A grammar G is said to be in the *Chomsky normal form*, if all rules of R are of the form $A \rightarrow BC$, with $B, C \in N$, or of the form $A \rightarrow a$, with $a \in \Sigma$.

A *nondeterministic finite automaton* (NFA) is a quintuple $\mathcal{A} = (\Sigma, Q, Q_0, \delta, F)$, where Q is a finite set of states, Σ is a finite set of input symbols, $Q_0 \subseteq Q$ is the set of initial states, $\delta: Q \times \Sigma \rightarrow 2^Q$ is the transition function, $F \subseteq Q$ is the set of accepting states. It accepts a string $w = a_1 \dots a_n$ if there is

a sequence of states $q_0, \dots, q_n \in Q$ with $q_0 \in Q_0$, $q_i \in \delta(q_{i-1}, a_i)$ for all i , and $q_n \in F$. The language of all strings accepted by \mathcal{A} is denoted by $L(\mathcal{A})$.

For a language L over an alphabet Σ , its rational index ρ_L is a function defined as follows:

$$\rho_L(n) = \max_{\substack{\mathcal{A}: \text{NFA with } n \text{ states} \\ L \cap L(\mathcal{A}) \neq \emptyset}} \min_{w \in L \cap L(\mathcal{A})} |w|$$

Tree dimension. For each node v in a parse tree t , its *dimension* $\dim v$ is an integer representing the amount of branching in its subtree. It is defined inductively: a leaf v has dimension 0. For an internal node v , if one of its children v_1, v_2, \dots, v_k , with $k \geq 1$, has a greater dimension than all the others, then v has the same dimension, and if there are multiple children of maximum dimension, then the dimension of v is greater by one.

$$\dim v = \begin{cases} \max_{i \in \{1, \dots, k\}} \dim v_i & \text{if there is a unique maximum} \\ \max_{i \in \{1, \dots, k\}} \dim v_i + 1 & \text{otherwise} \end{cases}$$

The dimension of a parse tree t , denoted by $\dim t$, is the dimension of its root.

Definition 1 (Grammars of bounded tree dimension). A grammar G is of d -bounded tree dimension if every parse tree t of G has $\dim t \leq d$, where d is some constant. This constant is called the dimension of G , denoted by $\dim G = d$.

Classical transformation to the Chomsky normal form preserves the class of grammars of d -bounded tree dimensions. Languages defined by such grammars are called *languages of d -bounded tree dimension*.

3 Upper bound on the rational index

Lemma 1. Let $G = (\Sigma, N, R, S)$ be a grammar in the Chomsky normal form, and let \mathcal{A} be an NFA with n states. Let $A \in N$ and $p, q \in Q$. Let w be the shortest string in $L_G(A)$, such that \mathcal{A} can read w starting from state p and ending in state q . Let t be a parse tree of w as A , and let d be the dimension of t . Then, $|w| \leq |N|^d n^{2d}$.

Proof. The proof is by induction on the dimension of t .

Basis: $\dim t = 0$. Then, t should use a rule $A \rightarrow a$, and $w = a$. The desired inequality holds as an equality $|w| = 1 = |N|^0 n^0$.

Inductive step: dimension $d - 1 \rightarrow d$. Let t have dimension d , and consider the path from the root of t passing through nodes of dimension d . If both children of the root of t have dimension $d - 1$, then this path consists of a single node; otherwise, one of the children of the root of t has dimension d , and the other has dimension less than d , and the path continues to the child of dimension d , etc. Let the path contain h edges, and let A_0, A_1, \dots, A_h be the nonterminals in the labels of nodes on this path, with $A_0 = A$.

At each node of this path except the last one, a subtree of dimension less than d spawns off to the left or to the right, whereas the last node on this path has two children of dimension $d-1$. Let s_1, \dots, s_k be the subtrees spawned off to the left, let $B_i \in N$ be the root of each s_i , and let u_i be the substring corresponding to s_i . Similarly, let t_1, \dots, t_ℓ be the subtrees spawned off to the right, with each t_i having a root C_i and corresponding to a substring v_i . Then $k + \ell = h$. Let the last node on the path have two children labelled with B_{k+1} and $C_{\ell+1}$, with the corresponding substrings u_{k+1} and $v_{\ell+1}$. Then, $w = u_1 \dots u_k u_{k+1} v_{\ell+1} v_\ell \dots v_1$, as illustrated in Figure 1.

Consider the computation of \mathcal{A} on w , which starts in p and ends in q , and consider the following intermediate states in this computation: let p_i be the state after reading each u_i , and let q_i be the state before reach each v_i . Also let $p_0 = p$ and $q_0 = q$ for uniformity. Then each string u_i is a string in $L_G(B_i)$, which \mathcal{A} can read starting in p_{i-1} and finishing in p_i ; this must be the shortest string with this property, for otherwise w could be shortened by replacing u_i with a shorter one. By the same reasoning, each v_i is the shortest string in $L_G(C_i)$, on which \mathcal{A} has a computation that begins in q_i and ends in q_{i-1} . Then, by the induction hypothesis, $|u_i|, |v_i| \leq |N|^{d-1} n^{2(d-1)}$ for all applicable i . This upper bound holds for $k + \ell + 2 = h + 2$ subtrees shown in Figure 1.

Next, it is claimed that the number $h + 2$ is bounded by $|N| \cdot n^2$. Suppose, for the sake of a contradiction, that $h + 2 > |N| \cdot n^2$. Consider the nodes A_0, A_1, \dots, A_h on the main path, followed by the left child B_{k+1} . Each of them has a corresponding substring of the form $u_{i+1} \dots u_{k+1} v_{k+1} \dots v_{j+1}$ (or u_{k+1} in the case of B_{k+1}), and the computation of \mathcal{A} on w begins reading this substring in the state p_i and ends reading it in the state q_j . Accordingly, each of these $h + 2$ nodes has an associated triple that consists of a nonterminal symbol and two states. In total, there are $|N| \cdot n^2$ distinct triples of this form. Since $h + 2 > |N| \cdot n^2$ by the assumption, some triple must occur on this path multiple times. Then the segment of the path between these two nodes can be contracted, resulting in a valid parse tree of a string shorter than w as A , while the computation of \mathcal{A} on w is accordingly contracted to remove two fragments, resulting in a computation of \mathcal{A} on a shorter string. This contradicts the assumption that w is the shortest string with this property, and thus confirms that $h + 2 \leq |N| \cdot n^2$.

Then the length of w is bounded as follows.

$$|w| = \sum_i |u_i| + \sum_i |v_i| \leq (h+2) \cdot |N|^{d-1} n^{2(d-1)} \leq |N| \cdot n^2 \cdot |N|^{d-1} n^{2(d-1)} = |N|^d n^{2d}$$

□

Theorem 1. *Let G be a grammar of d -bounded tree dimension, and let \mathcal{A} be an NFA with n states. Then the length of the shortest string in $L(G) \cap L(\mathcal{A})$ is at most $|G|^d n^{2d}$, where $|G|$ is sum of lengths of all rules of the grammar.*

Proof. A given grammar G is first transformed to the Chomsky normal form, resulting in a grammar with the same bound on the tree dimension and with at most $|G|$ nonterminal symbols.

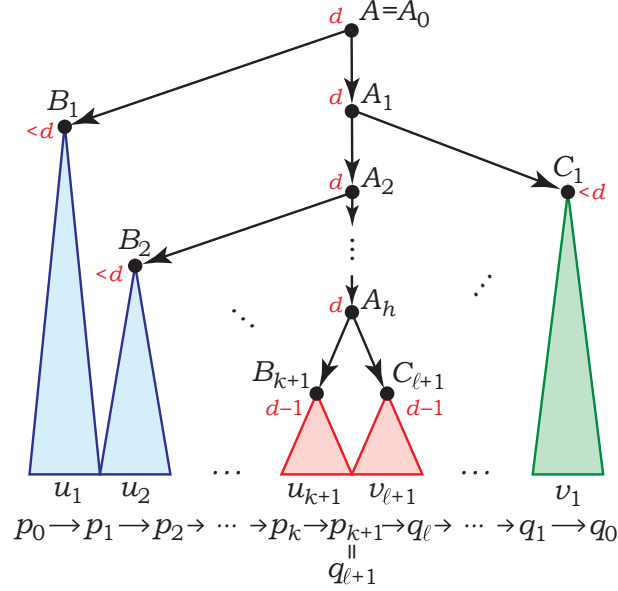


Fig. 1. A parse tree, with a path from the root passing through nodes of dimension d .

Let w be the shortest string in $L(G) \cap L(\mathcal{A})$, and let $q_0 \in Q_0$ and $r \in F$ be the first and the last states in the computation of \mathcal{A} on w . Then, w is the shortest string in $L_G(S)$, which \mathcal{A} can read starting in the state q_0 and ending in the state r . Then, by Lemma 1, the length of w is at most $|G|^d n^{2d}$. \square

4 Lower bound on the rational index

The upper bound $O(n^{2d})$ on the rational index of a language defined by grammar with tree dimension bounded by d has a matching lower bound $\Omega(n^{2d})$. It is first established for a convenient infinite set of values of n , to be extended to arbitrary n in the following.

Lemma 2. *For every $d \geq 1$, there is a grammar G of bounded tree dimension d , such that for every $n \geq 2^{d+1}$ divisible by 2^d there is an n -state NFA \mathcal{A} , such that the shortest string w in $L(G) \cap L(\mathcal{A})$ is of length at least $\frac{1}{2^{d^2+3d-3}} n^{2d}$.*

Proof. The grammar and the automaton are constructed inductively on d , for every d and only for n divisible by 2^d . Each constructed NFA shall have a unique initial state, which is also the unique accepting state.

Basis. $\dim(G) = 1$. The family of languages having dimension $d = 1$ coincides with the family of linear languages. Let G be a linear grammar with the rules $S \rightarrow aSb \mid ab$, which defines the language $L(G) = \{a^i b^i \mid i \geq 1\}$.

For every $n \geq 4$ divisible by $2^d = 2$, let $\ell = \frac{n}{2}$, $m = \frac{n}{2} + 1$. Then ℓ and m are coprime integers. Define an NFA \mathcal{A} over the alphabet $\{a, b\}$, which consists

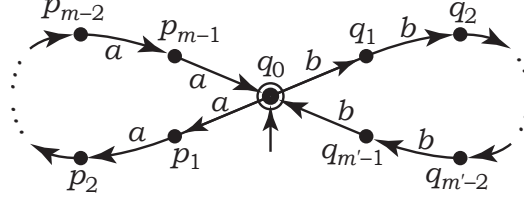


Fig. 2. Worst-case NFA \mathcal{A} for grammar G of bounded dimension $d = 1$

of two cycles sharing one node, q_0 , which is both the initial and the unique accepting state. The cycle of length ℓ has all transitions by a , and the other by b , as shown in Figure 2. The automaton has $\ell + m - 1 = n$ states.

Every string in $L(G) \cap L(\mathcal{A})$ is of the form $a^i b^i$, with $i \geq 1$. For the automaton to accept it, i must be divisible both by ℓ and by m . Since the cycle lengths are relatively prime, the shortest string w with this property has $i = \ell m$, and is accordingly of length $2\ell m$. Its growth with n is estimated as follows.

$$|w| = 2\ell m = 2 \frac{n}{2} \cdot \left(\frac{n}{2} + 1 \right) = \frac{1}{2} n^2 + n$$

This example is well-known to the community [7, 15].

Inductive step. $\dim(G) = d$.

By the induction hypothesis, there is a grammar $\widehat{G} = (\widehat{\Sigma}, \widehat{N}, \widehat{R}, \widehat{S})$ of bounded dimension $\dim(\widehat{G}) = d - 1$, which satisfies the statement of the lemma. The new grammar $G = (\Sigma, N, R, S)$ of dimension d is defined over the alphabet $\Sigma = \widehat{\Sigma} \cup \{a, b, c\}$, where $a, b, c \notin \widehat{\Sigma}$ are new symbols. It uses nonterminal symbols $N = \widehat{N} \cup \{S, A\}$, adding two new nonterminals $A, S \notin \widehat{N}$ to those in \widehat{G} , where S is the new initial symbol. Its set of rules includes all rules from \widehat{G} and the following new rules.

$$\begin{aligned} S &\rightarrow ASc \mid Ac \\ A &\rightarrow aAb \mid a\widehat{S}b \end{aligned}$$

To see that the dimension of the new grammar is greater by 1 than the dimension of \widehat{G} , first consider the dimension of any parse tree t with the root labeled by the nonterminal A , shown in Figure 3(right). The dimension of the \widehat{S} -subtree at the bottom is at most $d - 1$ by the properties of \widehat{G} . This dimension is inherited by all A -nodes in the tree, because their remaining children are leaves.

Now consider the dimension of a complete parse tree t with the start symbol S in the root, as in Figure 3(left). All A -subtrees in this tree have dimension at most $d - 1$. Then the bottom S -subtree, which uses the rule $S \rightarrow Ac$, also has dimension at most $d - 1$. Every S -subtree higher up in the tree uses a rule $S \rightarrow ASc$, and its dimension is at most d , because getting a higher dimension would require two subtrees of dimension d , which is never the case.

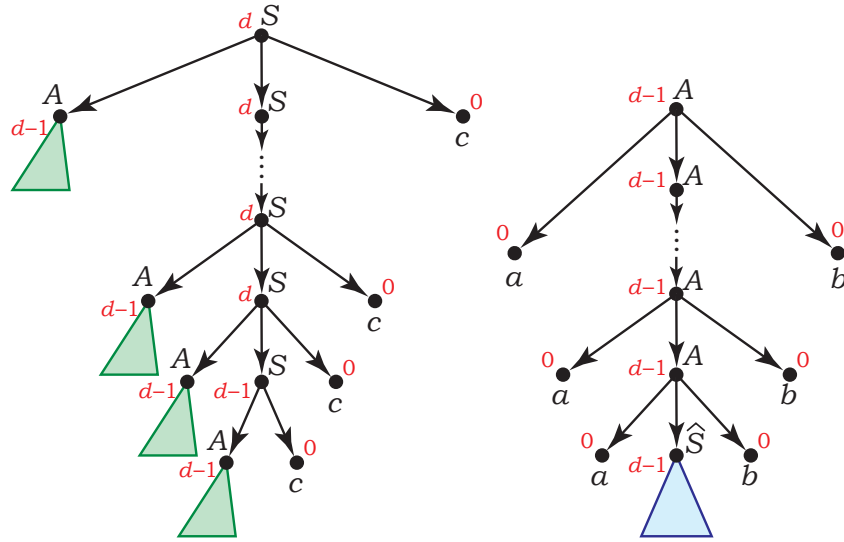


Fig. 3. Parse trees for S and for A , annotated with dimensions of their vertices

Now, for every $n \geq 2^{d+1}$ divisible by 2^d , the goal is to construct an n -state NFA over the alphabet Σ , so that the shortest string w in $L(G) \cap L(\mathcal{A})$. Since the number $\frac{n}{2}$ is at least 2^d and is divisible by 2^{d-1} , the induction hypothesis for the grammar \hat{G} asserts that there is an NFA $\hat{\mathcal{A}} = (\hat{Q}, \hat{\Sigma}, \hat{\delta}, \hat{q}_0, \{\hat{q}_0\})$, with $\frac{n}{2}$ states, with the shortest string \hat{w} in $L(\hat{G}) \cap L(\hat{\mathcal{A}})$ of length $\frac{1}{2^{(d-1)^2+3(d-1)-3}} (\frac{n}{2})^{2(d-1)}$.

The desired n -state NFA $\mathcal{A} = (\Sigma, Q, q_0, \delta, \{q_0\})$ is constructed as follows. Let $\ell = \frac{n}{4}$ and $m = \frac{n}{4} + 1$, these are two coprime integers. The set of states of \mathcal{A} contains all $\frac{n}{2}$ states from \hat{Q} , in which \mathcal{A} it operates as $\hat{\mathcal{A}}$, and $m + \ell - 1 = \frac{n}{2}$ new states forming a cycle of length ℓ and a chain of length m , which share a state.

$$Q = \hat{Q} \cup \{p_1, \dots, p_{\ell-1}, q_0, \dots, q_{m-1}\}$$

The new initial state q_0 has a transition by a leading to the initial state of $\hat{\mathcal{A}}$, from where one can return to q_1 by b .

$$\delta(q_0, a) = \{\hat{q}_0\}$$

$$\delta(\hat{q}_0, b) = \{q_1\}$$

There is a chain of transitions by a from q_{m-1} to q_0 , and another chain b in the opposite direction, from q_1 to q_{m-1} and back to q_0 .

$$\delta(q_i, a) = \{q_{i-1}\}, \quad \text{with } 1 \leq i \leq m-1$$

$$\delta(q_i, b) = \{q_{i+1}\}, \quad \text{with } 1 \leq i \leq m-2$$

$$\delta(q_{m-1}, b) = \{q_0\}$$

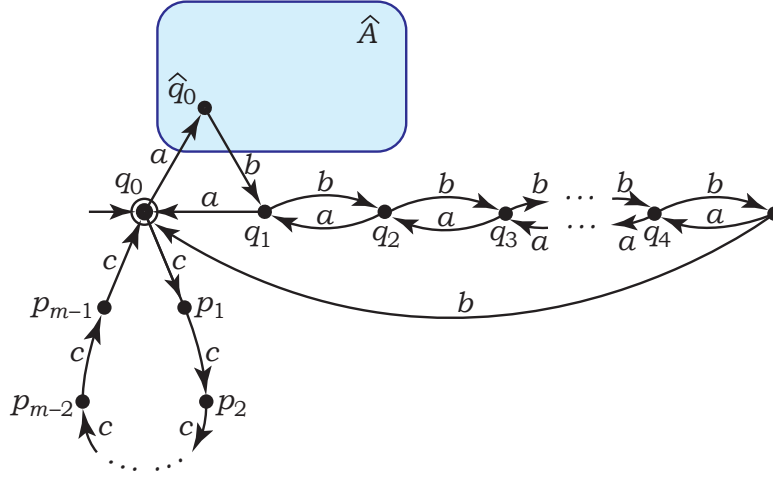


Fig. 4. NFA \mathcal{A} for grammar G of bounded dimension d

There is a cycle by c in the states $q_0, p_1, \dots, p_{\ell-1}$; for uniformity, denote $p_0 = q_0$.

$$\delta(p_i, c) = \{p_{i+1 \bmod \ell}\}, \quad \text{with } 0 \leq i \leq \ell - 1$$

The general form of \mathcal{A} is shown in Figure 4.

Let w be the shortest string in $L(G) \cap L(\mathcal{A})$. Consider how w is formed. Start state is q_0 . According to the grammar rule $S \rightarrow ASc \mid Ac$, the string w should start with a substring u in $L_G(A)$. There is the only one outgoing edge labeled with a , so the next state is \hat{q}_0 . The next part of w should be a symbol a or a string v in $L(\hat{G})$. As there is no outgoing edge labeled with a , the string v is the shortest string in $L(\hat{G}) \cap L(\hat{\mathcal{A}})$, and, hence, $v = \hat{w}$. Now the first part of w is $a\hat{w}$. To complete a substring derived by the nonterminal A , there is only one possible transition, which is an edge from \hat{q}_0 to q_1 labeled with b . The next substring should be symbol c (the rule $S \rightarrow Ac$) or a string derived by A . The only suitable transition here is from q_1 to q_0 by a , so a substring in $L(A)$ is started. Again, to complete the string generated by A , one goes to the state q_2 , and w now starts with $a\hat{w}baa\hat{w}bb$. By the construction of NFA \mathcal{A} , this process continues until one comes to the state q_0 without starting a substring derived by the nonterminal A (notice that such substrings are the shortest possible). Clearly, it happens after m iterations. Then it is left to read m symbols c by going from q_0 to q_0 . But m and ℓ are coprime, so to balance the number of substrings derived by the nonterminal A and the number of symbols c , one needs to repeat the first cycle ℓ times and the second cycle m times.

Accordingly, the shortest string w has the following structure. Let w_i be the shortest string such that there exists computation $q_{i-1} \xrightarrow{w_i} q_i$ ($q_{m-1} \xrightarrow{w_m} q_0$ for w_m) for $1 \leq i \leq m$ in \mathcal{A} and $w_i \in L(A)$. Notice that $w_i = aw_{i-1}b$ and $w_0 = \hat{w}$, and there exists computation $q_0 \xrightarrow{w_1} q_1 \xrightarrow{w_2} q_2 \xrightarrow{w_3} \dots \xrightarrow{w_{m-1}} q_m \xrightarrow{w_m} q_0$ in \mathcal{A} .

Considering the above and the rules $S \rightarrow ASc \mid Ac$ of the grammar G , the string w is of the following form:

$$w = \left(\prod_{i=1}^m w_i \right)^\ell c^{m\ell}$$

Then the length of w can be bounded as follows.

$$\begin{aligned} |w| &= \left(\sum_{i=1}^m |w_i| \right) \ell + \ell m = \left(\sum_{i=1}^m (|\hat{w}| + 2i) \right) \ell + \ell m = \left(\sum_{i=1}^m |\hat{w}| + \sum_{i=1}^m 2i \right) \ell + \ell m = \\ &= |\hat{w}| m \ell + (m+1) \ell m + \ell m \geq \ell m |\hat{w}| \end{aligned}$$

Using the lower bound on the length of \hat{w} , the desired lower bound on the length of w is obtained.

$$\begin{aligned} \ell m |\hat{w}| &\geq \frac{n}{4} \cdot \frac{n}{4} \cdot \frac{1}{2^{(d-1)^2+3(d-1)-3}} \left(\frac{n}{2} \right)^{2(d-1)} = \\ &= \frac{n^2}{16} \cdot \frac{1}{2^{d^2+d-5}} \cdot \frac{n^{2d-2}}{2^{2d-2}} = \frac{1}{2^{d^2+3d-3}} n^{2d} \end{aligned}$$

□

Theorem 2. *For every $d \geq 1$, there is a grammar G of bounded tree dimension d , such that for every $n \geq 2^{d+1}$ there is an n -state NFA \mathcal{A} , such that the shortest string w in $L(G) \cap L(\mathcal{A})$ is of length at least $\frac{1}{2^{d^2+d-3} 3^{2d}} n^{2d}$.*

Proof. Let G be the grammar given for d by Lemma 2. Let $2^d r \leq n < 2^d(r+1)$, for some integer r . Then $r \geq 2$ (for otherwise $n < 2^{d+1}$), and $2^d r \geq 2^{d+1}$.

Since $2^d r$ is divisible by 2^d , by Lemma 2, there is an NFA \mathcal{A} with $2^d r \leq n$ states, such that the length of the shortest string w in $L(G) \cap L(\mathcal{A})$ is at least $\frac{1}{2^{d^2+3d-3}} (2^d r)^{2d}$. This is the desired n -state NFA.

The inequality $n < 2^d(r+1)$ implies that $n < 2^d \frac{3r}{2}$, because $r+1$ is at most $\frac{3r}{2}$ for $r \geq 2$. Then $2^d r > \frac{2}{3} n$, and the lower bound on the length of w is expressed as a function of n as follows.

$$|w| \geq \frac{1}{2^{d^2+3d-3}} (2^d r)^{2d} \geq \frac{1}{2^{d^2+3d-3}} \left(\frac{2}{3} n \right)^{2d} = \frac{1}{2^{d^2+d-3} 3^{2d}} n^{2d}$$

□

Accordingly, the rational index of grammars with tree dimension bounded by d is $\Theta(n^{2^d})$ in the worst case.

5 Rational indices for some language families

Superlinear languages. A grammar $G = (\Sigma, N, R, S)$ is *superlinear* (Brzozowski [3]) if its nonterminal symbols split into two classes, $N = N_{lin} \cup N_{nonlin}$,

where rules for each nonterminal $A \in N_{lin}$ are of the form $A \rightarrow uBv$ or $A \rightarrow w$, with $B \in N_{lin}$, $u, v, w \in \Sigma^*$, while rules for a nonterminal $A \in N_{nonlin}$ are of the form $A \rightarrow \alpha B \beta$, with $B \in N$ and $\alpha, \beta \in (\Sigma \cup N_{lin})^*$. A language is *superlinear* if it is generated by some superlinear grammar.

Corollary 1. *For every superlinear grammar G , the rational index $\rho_{L(G)}$ is at most $|G|^2 \cdot n^4$.*

Proof. Parse trees in a superlinear grammar G have dimension at most 2. Then, by Theorem 1, the rational index $\rho_{L(G)}$ is bounded by $|G|^2 \cdot n^4$.

Turning to a lower bound, note that the grammar constructed in Theorem 2 for $d = 2$ is actually superlinear.

Corollary 2. *There exists a superlinear grammar G with rational index $\rho_{L(G)}(n) \geq \frac{1}{648} n^4$.*

Bounded-oscillation languages Bounded-oscillation languages were introduced by Ganty and Valput [5] as a generalization of the linear languages. They introduce the notion of harmonics of well-nested sequence of brackets, and give two equivalent definitions of *languages with oscillation bounded by k* : one in terms for runs of nondeterministic pushdown automata (NPDA), and the other using grammars and their parse trees.

Among other results Ganty and Valput [5] prove that the oscillation of a parse tree is closely related with its dimension.

Lemma 3 (Ganty and Valput [5]). *Let $G = (\Sigma, N, R, S)$ be a grammar in the Chomsky normal form, and let t be a parse tree in G . Then, $\text{osc } t - 1 \leq \dim t \leq 2 \text{ osc } t$.*

Thus, k -bounded-oscillation grammars have dimension of parse trees bounded by $2k$, and Theorem 1 gives the following upper bound on the rational index of these languages.

Corollary 3. *Let L be a k -bounded-oscillation language. Then $\rho_L(n) = O(n^{4k})$.*

6 Conclusion and open problems

We have proved that the languages of bounded tree dimension have polynomial rational index. This implies, in particular, that the CFL-reachability problem and Datalog query evaluation for these languages is in NC.

There is a family of languages which has polynomial rational index, the *one-counter languages*. Their rational index is known to be $O(n^2)$ [4]. Could this class be generalized in the same manner as linear languages, preserving the polynomial order of the rational index? One can consider the Polynomial Stack Lemma by Afrati et al. [1], where some restriction on the PDA stack contents are given, or investigate the properties of the substitution closure of the one-counter languages, which is known to have polynomial rational index [2].

References

1. Afrati, F., Papadimitriou, C.: The parallel complexity of simple chain queries. In: Proceedings of the Sixth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems. pp. 210–213. PODS '87, ACM, New York, NY, USA (1987). <https://doi.org/10.1145/28659.28682>
2. Boasson, L., Courcelle, B., Nivat, M.: The rational index: a complexity measure for languages. *SIAM Journal on Computing* **10**(2), 284–296 (1981)
3. Brzozowski, J.A.: Regular-like expressions for some irregular languages. In: 9th Annual Symposium on Switching and Automata Theory (swat 1968). pp. 278–286 (Oct 1968). <https://doi.org/10.1109/SWAT.1968.24>
4. Chistikov, D., Czerwinski, W., Hofman, P., Pilipczuk, M., Wehar, M.: Shortest paths in one-counter systems. *Log. Methods Comput. Sci.* **15**(1) (2019). [https://doi.org/10.23638/LMCS-15\(1:19\)2019](https://doi.org/10.23638/LMCS-15(1:19)2019), [https://doi.org/10.23638/LMCS-15\(1:19\)2019](https://doi.org/10.23638/LMCS-15(1:19)2019)
5. Ganty, P., Valput, D.: Bounded-oscillation pushdown automata. *Electronic Proceedings in Theoretical Computer Science* **226**, 178–197 (Sep 2016). <https://doi.org/10.4204/eptcs.226.13>
6. Greenlaw, R., Hoover, H.J., Ruzzo, W.L.: Limits to Parallel Computation: P-completeness Theory. Oxford University Press, Inc., New York, NY, USA (1995)
7. Hellings, J.: Path results for context-free grammar queries on graphs. *CoRR abs/1502.02242* (2015)
8. Holzer, M., Kutrib, M., Leiter, U.: Nodes connected by path languages. In: Mauri, G., Leporati, A. (eds.) *Developments in Language Theory*. pp. 276–287. Springer Berlin Heidelberg, Berlin, Heidelberg (2011)
9. Komarath, B., Sarma, J., Sunil, K.S.: On the complexity of l-reachability. In: Jürgensen, H., Karhumäki, J., Okhotin, A. (eds.) *Descriptional Complexity of Formal Systems*. pp. 258–269. Springer International Publishing, Cham (2014)
10. Luttenberger, M., Schlund, M.: Convergence of newton's method over commutative semirings. *Inf. Comput.* **246**, 43–61 (2016). <https://doi.org/10.1016/j.ic.2015.11.008>
11. Pierre, L.: Rational indexes of generators of the cone of context-free languages. *Theoretical Computer Science* **95**(2), 279–305 (1992). [https://doi.org/10.1016/0304-3975\(92\)90269-L](https://doi.org/10.1016/0304-3975(92)90269-L)
12. Pierre, L., Farinone, J.M.: Context-free languages with rational index in $\theta(n^\gamma)$ for algebraic numbers γ . *RAIRO - Theoretical Informatics and Applications - Informatique Théorique et Applications* **24**(3), 275–322 (1990)
13. Reps, T.W.: Program analysis via graph reachability. *Information & Software Technology* **40**, 701–726 (1997)
14. Ullman, J.D., Van Gelder, A.: Parallel complexity of logical query programs. In: 27th Annual Symposium on Foundations of Computer Science (sfcs 1986). pp. 438–454 (Oct 1986). <https://doi.org/10.1109/SFCS.1986.40>
15. Yannakakis, M.: Graph-theoretic methods in database theory. In: Rosenkrantz, D.J., Sagiv, Y. (eds.) *Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, April 2-4, 1990, Nashville, Tennessee, USA. pp. 230–242. ACM Press (1990). <https://doi.org/10.1145/298514.298576>