

Complexity of Counting First-Order Logic for the Subword Order

Dietrich Kuske

Technische Universität Ilmenau, Germany

Christian Schwarz

Technische Universität Ilmenau, Germany

Abstract

This paper considers the structure consisting of the set of all words over a given alphabet together with the subword relation, regular predicates, and constants for every word. We are interested in the counting extension of first-order logic by threshold counting quantifiers. The main result shows that the two-variable fragment of this logic can be decided in two-fold exponential space provided the regular predicates are restricted to piecewise testable ones. This result improves prior insights by Karandikar and Schnoebelen by extending the logic and saving one exponent. Its proof consists of two main parts: First, we provide a quantifier elimination procedure that results in a formula with constants of bounded length (this generalizes the procedure by Karandikar and Schnoebelen for first-order logic). From this, it follows that quantification in formulas can be restricted to words of bounded length, i.e., the second part of the proof is an adaptation of the method by Ferrante and Rackoff to counting logic and deviates significantly from the path of reasoning by Karandikar and Schnoebelen.

2012 ACM Subject Classification Theory of computation → Logic and verification; Theory of computation → Regular languages

Keywords and phrases Counting logic, piecewise testable languages

Digital Object Identifier 10.4230/LIPIcs.MFCS.2020.61

1 Introduction

The subword relation is one of the simplest nontrivial examples of a well-quasi ordering [4] and can be used in the verification of infinite state systems [2]. It can be understood as embeddability of one word into another. This embeddability relation has been considered for other classes of structures like trees, posets, semilattices, lattices, graphs, Mazurkiewicz traces etc. [8, 10, 9, 5, 20, 21, 12].

Many of these papers study logical aspects of the embeddability relation. Regarding the subword relation, literature provides a rather sharp description of the border between decidable and undecidable fragments of first-order logic: For the subword order alone, the \exists^* -theory is decidable [11] and the $\exists^*\forall^*$ -theory is undecidable [6]. For the subword order together with regular predicates, the two-variable theory is decidable [6] (this holds even for the two-variable fragment of the logic $C+MOD$, i.e., the extension of first-order logic by threshold- and modulo-counting quantifiers [13]) and the three-variable theory [6] as well as the \exists^* -theory are undecidable [3] (these two undecidabilities already hold if we only consider singleton predicates, i.e., constants). If one restricts the universe from all words to a particular language, an even more diverse picture appears [13].

We next sketch the decision procedure for the 2-variable fragment of the first-order theory of the subword relation together with regular predicates from [6]. Let $\varphi(x)$ be a formula with a single free variable. It may contain regular predicates that are given in any familiar formalism. Then the crucial insight from [6] is that the set of words satisfying $\varphi(x)$ can be obtained from the regular predicates by a fixed set of rational transductions and Boolean operations. Hence, one can inductively build the minimal dfa accepting this set. The only



© Dietrich Kuske and Christian Schwarz;
licensed under Creative Commons License CC-BY

45th International Symposium on Mathematical Foundations of Computer Science (MFCS 2020).

Editors: Javier Esparza and Daniel Král'; Article No. 61; pp. 61:1–61:12

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

known upper bound for this minimal dfa is non-elementary since any quantification requires to apply one of the rational transductions to the language of a minimal dfa (which leads to an nfa) and then to determinize and minimize this nfa. The crucial insight from the follow-up paper [7] by the same authors is that the size of these minimal dfas is at most triply exponential if, instead of regular predicates, one allows constants, only (alternatively: singleton predicates). Since determinisation and minimisation of an nfa can be done in space polynomial in the resulting minimal dfa (and logarithmic in the nfa), the above construction can be carried out in three-fold exponential space¹ which is also an upper bound for the said theory (the best lower bound we know so far is PSPACE [6]). This bound on the size of the minimal dfas is possible since all defined languages are piecewise testable [18]. A useful complexity measure for piecewise testable languages is their height. The new and innovative contribution of the proof from [7] are bounds for the height of the upwards closure $L\uparrow$, the downwards closure $L\downarrow$, and the incomparability set $L\parallel$ of a piecewise testable language L ; these new bounds are polynomial in the height of L (assuming a fixed alphabet).²

We improve this 3EXPSPACE upper bound for the theory in three aspects:

1. We prove an upper bound of two-fold exponential space.
 2. We allow piecewise testable predicates given by so-called pt-nfas [15, 16] (which are more succinct than minimal dfas). Further, the upper bound is measured in the *depth* of these pt-nfas as opposed to their *size*.
- **Remark.** Any piecewise testable predicate can be defined in the one-variable fragment of first-order logic. Consequently, these predicates do not increase the expressive power. Since a pt-nfa of depth k accepts a piecewise testable language of height k , the naive translation of a pt-nfa into a formula yields a formula of size exponential in the depth of the pt-nfa. As to whether this size increase is necessary seems not to be known.
3. We extend first-order logic by threshold counting quantifiers $\exists^{\geq t}$ (from [13], we know that this theory is decidable, even with regular predicates).

Following and extending the ideas from [7], we first prove new results on the height of piecewise testable languages. Namely, we extend the above mentioned results about $L\uparrow$, $L\downarrow$, and $L\parallel$ to, e.g., $L\uparrow_{\geq t}$, the set of words that have t subwords in L (and similarly for $L\downarrow_{\geq t}$ and $L\parallel_{\geq t}$). These considerations can be found in Section 3.

From these results, it follows that the language L defined by a formula (that uses threshold counting quantifiers and piecewise testable predicates given by pt-nfas) is piecewise testable of height at most doubly exponential in the size of the formula (Theorem 19).

► **Remark.** Consequently, the language L can be defined by a quantifier-free first-order formula. It follows that also the addition of counting quantifiers $\exists^{\geq t}$ does not increase the expressive power of the logic. But the use of counting quantifiers allows to write exponentially more succinct formulas (Theorem 21).

So far, this parallels the development in [7] where the corresponding result was shown for first-order logic. But at this point, instead of building automata (as done in [7]), we follow another path of argument, that is an adaptation of Ferrante and Rackoff's method [1].

The language-theoretic considerations imply that any formula is equivalent to a quantifier-free formula that uses constants of doubly exponential length and no piecewise testable

¹ The claim of three-fold exponential time from [7, Theorem 7.5] is not supported by the proof idea [17].

² This view indicates that the result from [7] can be improved by allowing, instead of singleton predicates, piecewise testable predicates given by minimal dfas. Also then, the algorithm from [6] should run in three-fold exponential space.

predicates (Corollary 20). From this, we derive that quantification in formulas can be restricted to words of doubly exponential length. This implies that the 2-variable fragment of the threshold counting extension of first-order logic becomes decidable in two-fold exponential space (allowing piecewise testable predicates in the formula given by pt-nfas).

2 Definitions and Main Results

Throughout this paper, we fix an alphabet Σ . We denote by Σ^* the set of (finite) words over Σ . A word $u \in \Sigma^*$ is a *subword* of $v \in \Sigma^*$ if $u = u_1 u_2 \dots u_n$ and $v = v_0 u_1 v_1 u_2 v_2 \dots u_n v_n$ for some $n \in \mathbb{N}$ and $u_i, v_i \in \Sigma^*$. We write $u \sqsubseteq v$ for this fact.

2.1 Piecewise Testable Languages and the Main Result for Language Theorists

The length of a word $u \in \Sigma^*$ is denoted $|u|$, $\Sigma^{\leq n}$ denotes the set of words of length $\leq n$. We next define Simon's congruences \sim_n that play an important role in our considerations.

► **Definition 1.** *Let $u, v \in \Sigma^*$ and $n \in \mathbb{N}$. Then u and v are n -equivalent (denoted $u \sim_n v$) if they have the same subwords of length $\leq n$. We denote by $[u]_n$ the equivalence class containing the word u wrt. the equivalence relation \sim_n .*

A language $L \subseteq \Sigma^*$ is *piecewise testable* if there exists $n \in \mathbb{N}$ such that L is a union of languages $[u]_n$ for some words $u \in \Sigma^*$. The minimal such n is called the *height* of L . We write $\text{PT}(n)$ for the class of piecewise testable languages of height $\leq n$. Note that $\text{PT}(n) \subseteq \text{PT}(n+1)$.

Let $L \subseteq \Sigma^*$ be piecewise testable. Then the upwards closure $L\uparrow$, the downwards closure $L\downarrow$ and the incomparability set $L\parallel$ are all piecewise testable of height polynomial in that of L (the degree of the polynomial is the size of the alphabet Σ) [7]. We will extend these results to the following more general operations.

Let $L \subseteq \Sigma^*$ be some language and $t \in \mathbb{N}$ some threshold. Then

$$L\uparrow_{\geq t} = \{v \in \Sigma^* \mid \exists u_1, \dots, u_t \in L: \begin{array}{l} u_i \sqsubseteq v \text{ for all } 1 \leq i \leq t \text{ and} \\ u_i \neq u_j \text{ for all } 1 \leq i < j \leq t \end{array}\}$$

denotes the set of words v that have t subwords in L . In particular, $L\uparrow_{\geq 0} = \Sigma^*$ and $L\uparrow_{\geq 1}$ is the usual upwards closure $L\uparrow$ of L . Note that any set $L\uparrow_{\geq t}$ is *upwards closed* and therefore piecewise testable.

The set

$$L\downarrow_{\geq t} = \{u \in \Sigma^* \mid \exists v_1, \dots, v_t \in L: \begin{array}{l} u \sqsubseteq v_i \text{ for all } 1 \leq i \leq t \text{ and} \\ v_i \neq v_j \text{ for all } 1 \leq i < j \leq t \end{array}\}$$

consists of all words u that have t superwords in L ; the above remarks on $L\uparrow_{\geq t}$ apply *mutatis mutandis*.

For two words u and v , we write $u\parallel v$ if neither u is a subword of v nor *vice versa*; we say that u and v are *incomparable*. Then let

$$L\parallel_{\geq t} = \{v \in \Sigma^* \mid \exists u_1, \dots, u_t \in L: \begin{array}{l} u_i \parallel v \text{ for all } 1 \leq i \leq t \text{ and} \\ u_i \neq u_j \text{ for all } 1 \leq i < j \leq t \end{array}\}$$

contain all words v that are incomparable with t words from L .

The function $g_{|\Sigma|}$ that will bound the height of the resulting languages $L \uparrow_{\geq t}$ etc. is defined as follows: Let $n \in \mathbb{N}$. Then \sim_n has only finitely many equivalence classes. Let $g_{|\Sigma|}(n)$ be minimal such that every equivalence class $[x]_n$ contains some word of length $\leq g_{|\Sigma|}(n)$. Then $n \leq g_{|\Sigma|}(n) \leq g_{|\Sigma|}(n+1)$ for all $n \in \mathbb{N}$. From [7, Theorem 3.7 and Eq. (3.12)], we know that $g_{|\Sigma|}(n) \leq (n+2)^{|\Sigma|}$.

The main result for language theorists now reads as follows (for the proof, see Section 3), it generalizes [7, Theorems 4.4, 5.5, and 6.1] from $t = 1$ to general thresholds.

► **Theorem 2.** *Let Σ be some alphabet, $n, t \in \mathbb{N}$, and $L \subseteq \Sigma^*$ be a piecewise testable language of height $\leq n$. Then the following hold:*

1. $L \uparrow_{\geq t}$ is piecewise testable of height $\leq g_{|\Sigma|}(n) + t - 1$.
2. $L \downarrow_{\geq t}$ is piecewise testable of height $\leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1)$ (note that this upper bound does not depend on t).
3. $L \parallel_{\geq t}$ is piecewise testable of height $\leq g_{|\Sigma|}(n) + t$.

Before we turn attention to a consequence in logic, we shortly recall some results on the relation of nfas and piecewise testable languages.

There are different characterisations of piecewise testable languages using nfas, we only rely on one by Masopust and Thomazo [15, 16]: a *pt-nfa* is a partially ordered and complete nfa that satisfies the UMS-property (the details are of no importance for our considerations) [15, Definition 3]. A language is piecewise testable iff it is accepted by some pt-nfa [16, Theorem 25]. Further, the depth $\|\mathcal{A}\|$ of a pt-nfa (i.e., the maximal length of a simple path in the nfa) bounds the height of the accepted language [15, Theorem 8].

2.2 The Logic C^2 and the Main Result for Logicians

Let NFA be the set of all nfas over the alphabet Σ . Consider the structure

$$\mathcal{S} = (\Sigma^*, \sqsubseteq, (L(\mathcal{A}))_{\mathcal{A} \in \text{NFA}}, (w)_{w \in \Sigma^*})$$

whose universe is the set of words, whose only binary relation is the subword relation, that has a unary relation $L(\mathcal{A})$ for each nfa $\mathcal{A} \in \text{NFA}$ and a constant for every word over Σ .

We can make statements about this structure using some variant of classical first-order logic. To control the use of nfas in these formulas, let $A \subseteq \text{NFA}$ be a set of nfas (e.g., $A = \text{NFA}$, $A = \emptyset$, or $A = \text{ptNFA} \subseteq \text{NFA}$ which is the set of pt-nfas). Then formulas from C_A^2 are defined by the following syntax:

$$\varphi := c \sqsubseteq d \mid c = d \mid c \in L(\mathcal{A}) \mid \varphi \vee \psi \mid \neg \varphi \mid \exists^{\geq t} z \varphi$$

where c, d are variables from $\{x, y\}$ or words from Σ^* , $\mathcal{A} \in A$ is some nfa over Σ , $t \in \mathbb{N}$, and $z \in \{x, y\}$ is a variable. Note that we allow only the variables x and y . The semantics of these formulas is defined in the obvious way with the understanding that $\exists^{\geq t} x \varphi$ holds if there are t mutually distinct words that all make the formula φ true. Consequently $\exists^{\geq 1}$ is the usual existential quantifier and $\exists^{\geq 0} x \varphi$ is always true. Let FO_A^2 denote the subset of C_A^2 that only uses the quantifier $\exists^{\geq 1}$, i.e., the classical first-order quantifier.

For arbitrary structures, the introduction of threshold counting quantifiers $\exists^{\geq t}$ in conjunction with the restriction to two variables extends the expressive power. Later, we will see that in our context, the logics C_{ptNFA}^2 and FO_{\emptyset}^2 are equally expressive (Corollary 20), but C_{ptNFA}^2 is exponentially more succinct than FO_{\emptyset}^2 by Theorem 21.

As a side remark, we prove that constants of length ≤ 2 suffice for the whole expressive power.

► **Theorem 3.** *Let $A \subseteq \text{NFA}$. For every formula $\varphi \in C_A^2$, there exists an equivalent formula $\psi \in C_A^2$ that uses constants of length ≤ 2 , only. The same applies to the logic FO_A^2 .*

Proof. It suffices to produce, for every word $w \in \Sigma^*$, a formula $\lambda_w(x) \in \text{FO}_{\emptyset}^2$ using at most constants of length ≤ 2 such that w is the only word satisfying $\lambda_w(x)$.

For $n \in \mathbb{N}$, there are formulas $\alpha_n(x)$ expressing $|x| \geq n$.

For $|w| \leq 2$, the formula $\lambda_w(x)$ is simply $x = w$. Now let $|w| > 2$. Set $m = \lfloor n/2 \rfloor + 1 < n$ and $S_w = \{u \in \Sigma^{\leq m} \mid u \sqsubseteq w\}$. By [14, Theorem 6.2.16], w is the only word of length $|w|$ such that S_w is the set of subwords of w of length $\leq m$. Since this can, using the formulas $\alpha_n(x)$ and induction, be expressed by a formula from FO_{\emptyset}^2 , the claim follows. ◀

The size of a formula is defined with the understanding that the size $|\mathcal{A}|$ of an nfa \mathcal{A} is its number of states, the size of a variable is 1, the size of a word is its length, and the size of the quantifier $\exists^{\geq t}$ is the length $|\text{bin}(t)|$ of the binary encoding of t .

Besides the size, we also define the *norm* $\|\varphi\|$ of a formula φ from C_{ptNFA}^2 :

$$\begin{aligned} \|c \sqsubseteq d\| &= \|c = d\| = \max(|c|, |d|), & \|c \in L(\mathcal{A})\| &= \max(|c|, |\mathcal{A}|), \\ \|\alpha \vee \beta\| &= \max(\|\alpha\|, \|\beta\|), & \|\neg\beta\| &= \|\beta\|, \text{ and} \\ \|\exists^{\geq t} x \varphi\| &= |\text{bin}(t)| + \|\varphi\|. \end{aligned}$$

This norm $\|\varphi\|$ is similar to the quantifier depth. Note that, in particular, $\|\varphi\|$ bounds the length of constants and the depth of pt-nfas occurring in φ . Note further that $\|\varphi\| \leq |\varphi|$ holds for any formula φ .

From Theorem 2, we infer in Section 4 that all definable languages are piecewise testable of bounded height (Theorem 19). This allows to derive a quantifier elimination result that reads as follows:

► **Corollary 20.** *Let $c = 2 \cdot |\Sigma|$. Every C_{ptNFA}^2 -formula φ is equivalent to some quantifier- and automata-free formula $\psi \in \text{FO}_{\emptyset}^2$ with $\|\psi\| < 2^{c^{2\|\varphi\|}}$.*

Karandikar and Schnoebelen [7] showed that any non-empty piecewise testable language of height n has elements of length polynomial in n . Based on Corollary 20, we can therefore restrict quantification in a formula φ to “short words” implying our main result for logicians.

► **Theorem 24.** *The C_{ptNFA}^2 -theory of S is decidable in doubly exponential space.*

3 Closure of the Class of Piecewise Testable Languages

The purpose of this section is to indicate how Theorem 2 can be proved.

3.1 Notions and Results Used in the Proof

A set of words L is *convex* if $u, w \in L$ and $u \sqsubseteq v \sqsubseteq w$ imply $v \in L$.

► **Lemma 4** (compiled in [7]). *Let $u, v, v' \in \Sigma^*$, $a, b \in \Sigma$, and $n \in \mathbb{N}$.*

1. *The equivalence class $[u]_n$ is convex.*
2. *If $u \sim_n v$, then there exists $w \in [u]_n$ with $u, v \sqsubseteq w$.*
3. *If $uv \sim_n uav$, then $uv \sim_n ua^\ell v$ for all $\ell \in \mathbb{N}$.*
4. *The equivalence class $[u]_n$ is infinite or a singleton.*

Proof (cited from [7]). (1) is by combining the definition of \sim_n with the observation $\{u\} \downarrow \subseteq \{v\} \downarrow$ provided $u \sqsubseteq v$. (2) is [19, Lemma 6] (cf. [14, Thm. 6.2.6] for an alternative proof). (3) is in the proof of [14, Corollary 6.2.8]. Finally, (4) follows from (1), (2), and (3). ◀

An example of a singleton equivalence class is $[u]_{|u|+1}$ for any $u \in \Sigma^*$; if u contains two distinct letters, then even $[u]_{|u|} = \{u\}$ (but $[aa]_2 = aaa^*$).

Since $\Sigma^{\leq n}$ is finite, there are only finitely many equivalence classes of \sim_n . Hence, for any $n \in \mathbb{N}$, there are only finitely many languages $L \subseteq \Sigma^*$ in $\text{PT}(n)$ and this class is closed under Boolean operations.

For a set $L \subseteq \Sigma^*$ of words, let $\min(L)$ denote the set of words $v \in L$ that have no proper subword in L . Since the subword relation is well-founded, any word from L is a superword of some word from $\min(L)$, i.e., $L \subseteq \min(L) \uparrow$.

Imre Simon found a description of the set of minimal elements of an equivalence class $[u]_n$ that uses the following concept. For a set $B \subseteq \Sigma$, let $\text{Perm}(B) \subseteq \Sigma^*$ denote the set of permutations of B seen as words (e.g., $\text{Perm}(\{a, b\}) = \{ab, ba\}$ and $\text{Perm}(\emptyset) = \{\varepsilon\}$). For sets $B_i \subseteq \Sigma$, define $\text{Perm}(B_1, B_2, \dots, B_k) = \text{Perm}(B_1) \text{Perm}(B_2) \cdots \text{Perm}(B_k)$. For instance, $\text{Perm}(\{a\}, \{b\}, \{c\}) = \{abc\}$ while $\text{Perm}(\{a, b\}, \{c\}) = \{abc, bac\}$ for all letters $a, b, c \in \Sigma$.

► **Theorem 5** ([18], cf. [14, Theorem 6.2.9]). *Let $n \in \mathbb{N}$ and $u \in \Sigma^*$. Then there exist $k \in \mathbb{N}$ and $B_1, B_2, \dots, B_k \subseteq \Sigma$ with $\min([u]_n) = \text{Perm}(B_1, B_2, \dots, B_k)$.*

Deleting all empty sets from the tuple (B_1, B_2, \dots, B_k) makes the above presentation of $\min([u]_n)$ unique. By Theorem 5, all words from $\min([u]_n)$ have the same length $\sum_{1 \leq i \leq k} |B_i|$ which is $\leq g_{|\Sigma|}(n)$ (by the very definition of that function) and therefore $\leq (n+2)^{|\Sigma|}$ (by [7, Theorem 3.7 and Eq. (3.12)]).

► **Theorem 6.** *Let Σ be an alphabet, $w \in \Sigma^*$, and $n \in \mathbb{N}$. Then there exists a word $v \sim_n w$ with $|v| \leq g_{|\Sigma|}(n)$ and $v \sqsubseteq w$.*

Recall that $g_{|\Sigma|}(n) \leq (n+2)^{|\Sigma|}$ by [7, Theorem 3.7 and Eq. (3.12)].

Proof. Since the subword order is well-founded, there exist $u, v \in \min([w]_n)$ with $|u| \leq g_{|\Sigma|}(n)$ and $v \sqsubseteq w$. Theorem 5 implies $|v| = |u| \leq g_{|\Sigma|}(n)$. ◀

3.2 Upward Closures

The following result verifies Theorem 2(1).

► **Proposition 7.** *Let $L \in \text{PT}(n)$ be a language over Σ and $t \in \mathbb{N}$. Then the language $L \uparrow_{\geq t}$ belongs to $\text{PT}(g_{|\Sigma|}(n) + t - 1)$.*

Proof. Let $z \in L \uparrow_{\geq t}$ and $z' \sim_{g_{|\Sigma|}(n)+t-1} z$. Then there exists a t -elements set $Y \subseteq L$ with $y \sqsubseteq z$ for all $y \in Y$. Choosing the elements of Y as short as possible, we can assume $Y \downarrow \cap L = Y$. Using the definition of $g_{|\Sigma|}$, Lemma 4(1), and Theorem 5, one can show that all words from Y are of length $\leq g_{|\Sigma|}(n) + t - 1$. Consequently, they are subwords of z' implying $z' \in L \uparrow_{\geq t}$. ◀

3.3 Downward Closures

The following result verifies Theorem 2(2).

► **Proposition 8.** *Let $L \in \text{PT}(n)$ be a language over Σ and $t \in \mathbb{N}$. Then the language $L \downarrow_{\geq t}$ belongs to $\text{PT}((|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1))$.*

Proof. Let $F \subseteq L$ denote the set of elements $x \in L$ with $[x]_n$ finite (i.e., a singleton by Lemma 4(4)) and $I = L \setminus F$. From Lemma 4(2) and (3), it follows that I has no maximal element implying $L \downarrow^{\geq t} = F \downarrow^{\geq t} \cup I \downarrow$.

By [7, Theorem 5.5], the height of $I \downarrow$ is $\leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1)$. By definition of $g_{|\Sigma|}(n)$, all words in F have length $\leq g_{|\Sigma|}(n)$. Hence the height of $F \downarrow^{\geq t}$ is $\leq g_{|\Sigma|}(n) + 1 \leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1)$. \blacktriangleleft

3.4 Incomparability Set

There are three types of equivalence classes $[x]_n$: the singletons, the chains, and the infinite ones which are no chains. Propositions 9, 11, and 15, respectively, bound the heights of $[x]_n \downarrow^{\geq t}$ for these three types of equivalence classes and collectively verify Theorem 2(3).

► **Proposition 9.** *Let $n, t \in \mathbb{N}$ and $x \in \Sigma^*$ such that $L = [x]_n$ is a singleton. Then $L \downarrow^{\geq t} \in \text{PT}(g_{|\Sigma|}(n))$.*

Proof. If $t \neq 1$, then $L \downarrow^{\geq t} \in \{\emptyset, \Sigma^*\}$, hence of height $0 \leq g_{|\Sigma|}(n)$. If $t = 1$, then $L \downarrow^{\geq t}$ is the complement of $\{x\} \uparrow \cup \{x\} \downarrow \setminus \{x\}$, two languages of height $\leq |x|$. Since $[x]_n$ is a singleton, the definition of $g_{|\Sigma|}$ implies $|x| \leq g_{|\Sigma|}(n)$. \blacktriangleleft

Next, we consider the case that $[x]_n$ is a chain and bound the height of $[x]_n \downarrow^{\geq t}$. The following lemma provides the central argument that will also be used later.

► **Lemma 10.** *Let $t \geq 1$ and let C be the convex chain $x_0 \sqsubset x_1 \sqsubset \dots$. Then $C \downarrow^{< t} = C \cup \{x_{t-1}\} \downarrow$.*

Proof. Since subwords of x_{t-1} are, at most, incomparable with x_0, x_1, \dots, x_{t-2} , the inclusion “ \supseteq ” is obvious. Now let $x \notin C \cup \{x_{t-1}\} \downarrow$. Since $x_0 \sqsubseteq x_{t-1}$, this implies $x \sqsupseteq x_0$ or $x \parallel x_0, x_{t-1}$. In the first case, convexity of C and $x \notin C$ imply that x is incomparable with infinitely many elements of C . In the latter case, x is incomparable with all x_i for $0 \leq i \leq t-1$. \blacktriangleleft

► **Proposition 11.** *Let $n, t \in \mathbb{N}$ and $x \in \Sigma^*$ such that $C = [x]_n$ is a chain. Then $C \downarrow^{\geq t} \in \text{PT}(g_{|\Sigma|}(n) + t)$.*

Proof. Since $C \downarrow^{\geq 0} = \Sigma^* \in \text{PT}(0)$, it remains to consider the case $t > 0$.

List the elements of C in increasing order: $x_0 \sqsubset x_1 \sqsubset x_2 \dots$. Since C is convex by Lemma 4(1), we obtain $|x_i| = |x_0| + i$ for all $i \geq 0$. By Lemma 10, $C \downarrow^{\geq t}$ is the complement of $C \cup \{x_{t-1}\} \downarrow$. But C is of height n . Furthermore $|x_{t-1}| = |x_0| + t - 1 < g_{|\Sigma|}(n) + t$ since $x_0 \in \min(C)$ implying that the height of $\{x_{t-1}\} \downarrow$ is $\leq g_{|\Sigma|}(n) + t$. \blacktriangleleft

It remains to prove a similar statement for infinite equivalence classes $[x]_n$ that are not a chain. The proof of the case $t = 1$ from [7] first shows that $[x]_n$ contains two elements of every length $> |x|$. Consequently, every word of length $> |x|$ is incomparable with some word from $[x]_n$, i.e., $[x]_n \downarrow^{\geq 1}$ is cofinite and therefore piecewise testable.

Our proof for $t > 1$ shows that the set of pairs of words of equal length can be grouped into two convex chains, i.e., the equivalence class $[x]_n$ contains two convex chains that intersect, at most, in $\min([x]_n)$ (Lemma 14). Then we apply Lemma 10. But first, we need some insight into convex chains which is the topic of the following considerations.

► **Lemma 12.** *Let $x, y \in \Sigma^*$ and $a \in \Sigma$. Then xa^*y is a convex chain.*

Proof. The language xa^*y is clearly a chain and one shows that $xa^k y \sqsubseteq z \sqsubseteq xa^m y$ implies the existence of ℓ with $z = xa^\ell y$. \blacktriangleleft

The third item of the following lemma implies, together with Theorem 5, that the maximal a -prefixes of two words from $\min([x]_n)$ differ in length by at most one.

► **Lemma 13.** *Let $B_1, B_2, \dots, B_k \subseteq \Sigma$ be non-empty, $a \in \Sigma$ and $u, v \in \Sigma^*$.*

(1) *If $au \in \text{Perm}(B_1, \dots, B_k)$, then $a \in B_1$ and $u \in \text{Perm}(B_1 \setminus \{a\}, B_2, \dots, B_k)$.*

(2) *If $aa \in \text{Perm}(B_1, \dots, B_k)$, then $B_1 = \{a\}$.*

(3) *If $u, v \notin a\Sigma^*$ and $m, n \in \mathbb{N}$ with $a^m u, a^n v \in \text{Perm}(B_1, \dots, B_k)$, then $|m - n| \leq 1$.*

► **Lemma 14.** *Let $u \in \Sigma^*$ such that $[u]_n$ is infinite but not a single chain. Then $[u]_n$ contains two infinite convex chains C_1 and C_2 with $C_1 \cap C_2 \subseteq \min([u]_n)$ and $C_i \cap \min([u]_n) \neq \emptyset$ for $i \in \{1, 2\}$.*

Proof. By [7, Lemmas 6.2 and 6.3], there are words $x_1, x_2, y_1, y_2 \in \Sigma^*$ and letters $a, b \in \Sigma$ such that $x_1 x_2, y_1 y_2 \in \min([u]_n)$, $x_1 a x_2, y_1 b y_2 \in [u]_n$, and $x_1 a x_2 \neq y_1 b y_2$. Then $x_1 a^* x_2$ and $y_1 b^* y_2$ form infinite convex chains in $[u]_n$ (Lemmas 4(3) and 12) and it remains to be shown that their intersection contains words from $\min([u]_n)$, only.

Let $\ell, m \in \mathbb{N}$ such that $x_1 a^\ell x_2 = y_1 b^m y_2$. Since $x_1 x_2, y_1 y_2 \in \min([u]_n)$, they have the same Parikh image by Theorem 5 implying $a = b$ and $\ell = m$, i.e., $x_1 a^\ell x_2 = y_1 a^\ell y_2$.

One first shows, w.l.o.g., $|x_1| < |y_1|$ and then distinguishes the cases $|y_1| \leq |x_1 a^\ell|$ and $|x_1 a^\ell| < |y_1|$.

In the first case there exists $k \in \mathbb{N}$ with $0 < k \leq \ell$ and $x_1 a^k = y_1$ and therefore $a^k y_2 = x_2$. This leads to $x_1 a x_2 = y_1 b y_2$, contradicting our assumption.

It remains to consider the case $|x_1 a^\ell| < |y_1|$. Then there exists a word $h \neq \varepsilon$ with $x_1 a^\ell h = y_1$ and therefore $x_2 = h a^\ell y_2$.

By Theorem 5, there is a tuple \overline{B} of nonempty subsets of Σ with $\min([u]_n) = \text{Perm}(\overline{B}) = \min([u]_n) \ni x_1 a^\ell h y_2, x_1 h a^\ell y_2$. Applying Lemma 13(1) and its dual, we obtain a tuple \overline{C} of nonempty subsets of Σ with $a^\ell h, h a^\ell \in \text{Perm}(\overline{C})$. Assuming $h \in a^* \Sigma^*$ leads to $x_1 a x_2 = y_1 b y_2$ which contradicts our assumption. Hence we can write h as $a^k v$ with $v \in \Sigma^+ \setminus a\Sigma^*$. Then $a^{\ell+k} v = a^\ell h$ and $a^k v a^\ell = h a^\ell$ both belong to $\text{Perm}(\overline{C})$. Hence Lemma 13(3) implies $\ell \leq 1$. But $\ell = 1$ is impossible since $x_1 a x_2 \neq y_1 a y_2$. Thus, we obtain $\ell = 0$, i.e., the two chains $x_1 a^* x_2$ and $y_1 b^* y_2$ intersect, at most, in $\min([u]_n)$. ◀

Now we can handle the remaining equivalence classes, i.e., bound the height of $[x]_n \|^{\geq t}$ provided $[x]_n$ is infinite but not a chain.

► **Proposition 15.** *Let $n, t \in \mathbb{N}$ and $x \in \Sigma^*$ such that $L = [x]_n$ is infinite but not a chain. Then $L \|^{\geq t} \in \text{PT}(g_{|\Sigma|}(n) + t)$.*

Proof. Since $L \|^{\geq 0} = \Sigma^* \in \text{PT}(0)$, it remains to consider the case $t > 0$.

By Lemma 14, there exist two infinite convex chains $C_1, C_2 \subseteq L$ such that $C_1 \cap C_2 \subseteq \min(L)$ and $C_i \cap \min(L) \neq \emptyset$ for $i \in \{1, 2\}$. By Lemma 10, there are words $x_i \in C_i$ of length $< g_{|\Sigma|}(n) + t$ such that

$$L \|^{\leq t} \subseteq C_1 \|^{\leq t} \cap C_2 \|^{\leq t} = (C_1 \cup \{x_1\} \downarrow) \cap (C_2 \cup \{x_2\} \downarrow) \subseteq \Sigma^{< g_{|\Sigma|}(n) + t}. \quad \blacktriangleleft$$

We can now put these three propositions together to verify Theorem 2(3).

► **Proposition 16.** *Let $L \in \text{PT}(n)$ be a language over Σ and $t \in \mathbb{N}$. Then $L \|^{\geq t} \in \text{PT}(g_{|\Sigma|}(n) + t)$.*

Proof. The language L is a finite disjoint union of equivalence classes $[x]_n$. Hence $L \|^{\geq t}$ can be written as a Boolean combination of languages of the form $[x]_n \|^{\geq s}$ for $s \in \{0, 1, \dots, t\}$. But all these languages have height $\leq g_{|\Sigma|}(n) + t$. ◀

4 Expressive Power and Quantifier Elimination

In this section, we show that every language definable in C_{ptNFA}^2 is piecewise testable of height bounded in terms of the norm of the defining formula. But first a simple result on the expressive power of quantifier-free formulas.

► **Lemma 17.** *Let $n \in \mathbb{N}$.*

- (1) *Any language $L \in \text{PT}(n)$ is defined by some quantifier- and automata-free formula $\varphi(x) \in \text{FO}_{\emptyset}^2$ with $\|\varphi(x)\| \leq n$.*
- (2) *If $\varphi(x) \in \text{FO}_{\text{ptNFA}}^2$ is a quantifier-free formula with $\|\varphi\| \leq n$, then it defines a language from $\text{PT}(n+1)$.*

Proof. By the very definition of \sim_n , the first claim holds for all equivalence classes $[x]_n$ and therefore for any language from $\text{PT}(n)$.

For the second claim, we use that the depth of any pt-nfa in φ is bounded by n and the same holds for the length of constants in φ . ◀

► **Example 18.** The language $\{aaa\}$ belongs to $\text{PT}(4) \setminus \text{PT}(3)$, and it can be defined by the formula $x = aaa$ of norm 3. Hence, in the first statement of the above lemma, the converse implication does not hold.

Regarding the second statement, the language $\{aaa\}a^*$ belongs to $\text{PT}(3)$, but cannot be defined by a formula of norm ≤ 2 . Hence, also that implication cannot be inverted.

► **Theorem 19.** *Let $c = 2 \cdot |\Sigma|$ and $\varphi(x) \in C_{\text{ptNFA}}^2$. Then the language $L_\varphi = \{u \in \Sigma^* \mid \mathcal{S} \models \varphi(u)\}$ is piecewise testable of height $< 2^{c^{2\|\varphi\|}}$.*

Proof. The claim is shown by induction on the construction of the formula φ . The cases that $\varphi(x)$ is quantifier-free or a Boolean combination are straightforward. So let $\varphi(x) = \exists^{\geq t} y: \varphi'(x, y)$.

In a first step, one expresses $\varphi(x)$ as a Boolean combination of formulas $\exists^{\geq s} y: (x\theta y \wedge \delta(x, y))$ where $s \leq t$, $\theta \in \{\sqsubseteq, \supseteq, =, \|\}$, and $\delta(x, y)$ is a conjunction of possibly negated atomic and existential (i.e., starting with $\exists^{\geq s}$) subformulas of $\varphi'(x, y)$.

In any such formula, $\delta(x, y)$ can be written as $\alpha(x) \wedge \beta(x, y) \wedge \gamma(y)$ with $\|\alpha\|, \|\gamma\| \leq \|\varphi'\|$ and $\beta(x, y)$ a conjunction of formulas of the form $x \sqsubseteq y$, $x \supseteq y$, and their negations. Depending on whether $x\theta y$ is consistent with $\beta(x, y)$ or not, the formula $\exists^{\geq s} y: x\theta y \wedge \delta(x, y)$ is equivalent to \perp or to

$$\alpha(x) \wedge \exists^{\geq s} y: (x\theta y \wedge \gamma(y)).$$

Since $\|\alpha\|, \|\gamma\| \leq \|\varphi'\| < \|\varphi\|$, we can apply the induction hypothesis, i.e., the languages defined by $\alpha(x)$ and by $\gamma(y)$ are of bounded height. Then Theorem 2 allows to also bound the height of the language defined by $\exists^{\geq s} y: (x\theta y \wedge \gamma(y))$ (this requires tedious and non-illuminating calculations that use, in particular, the estimate $(|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + m) < (m + n + 2)^c$). ◀

Since piecewise testable languages of bounded height can be defined by quantifier-free formulas from FO_{\emptyset}^2 , we obtain the following quantifier-elimination result (that does not only hold for formulas with a single free variable since only atomic formulas make “proper” use of more than one variable).

► **Corollary 20.** *Let $c = 2 \cdot |\Sigma|$. Every C_{ptNFA}^2 -formula φ is equivalent to some quantifier- and automata-free formula $\psi \in \text{FO}_{\emptyset}^2$ with $\|\psi\| < 2^{c^{2\|\varphi\|}}$.*

61:10 Complexity of Counting First-Order Logic for the Subword Order

For first-order formulas φ , this result can be found in [7, Cor. 7.4 and Thm. 7.5].

Note that the above corollary implies in particular that the logics C_{ptNFA}^2 and FO_{\emptyset}^2 are equally expressive (a description of this expressive power in terms of subword-piecewise testable relations can be found in [7, Theorem 7.2(ii)]). But we have the following result on the succinctness.

► **Theorem 21.** *The logic C_{\emptyset}^2 is exponentially more succinct than FO_{\emptyset}^2 .*

Proof. Let $\Sigma = \{a\}$ and $n \in \mathbb{N}$. Then $w_n = a^{2^n - 1}$ is the only word satisfying the formula

$$\varphi_n(x) = \exists^{\geq 2^n} y: y \sqsubseteq x \wedge \neg \exists^{\geq 2^n + 1} y: y \sqsubseteq x.$$

The size of this formula is in $O(n)$ since the thresholds are encoded in binary.

Now let ψ_n be an equivalent first-order formula. Since Σ is a singleton, one can eliminate all constants from ψ_n (incurring a linear increase in size). Since $\psi_n(x)$ distinguishes w_n from $w_n a$, its quantifier depth is $\geq 2^n - 1$ which implies the claim. ◀

5 Complexity of the C_{ptNFA}^2 -Theory

We now adapt the technique by Ferrante and Rackoff from first-order logic to its extension by threshold counting quantifiers to derive our upper complexity bound from Corollary 20.³ Central to this proof is the following lemma expressing that quantification in formulas can be restricted to words of bounded length. This property is the core of the method by Ferrante and Rackoff [1].

► **Lemma 22.** *Let $\varphi(x) = \exists^{\geq t} y: \psi(x, y)$ be a formula from C_{ptNFA}^2 . Let $c = 2 \cdot |\Sigma|$, $N \in \mathbb{N}$ with $2^{c^{2^{|\varphi|}}} \leq N$, and $u \in \Sigma^*$ with $|u| < N$. Then $\mathcal{S} \models \varphi(u)$ iff there are t words v of length $< N^{2^c}$ such that $\mathcal{S} \models \psi(u, v)$.*

Proof. For the non-trivial implication assume there are t words in the language $L := \{v \in \Sigma^* \mid \mathcal{S} \models \psi(u, v)\}$.

Corollary 20 yields a formula $\psi'(x, y) \in \text{FO}_{\emptyset}^2$ equivalent to $\psi(x, y)$ such that $\|\psi'\| < N$. Substituting u for x does not increase the norm since ψ' is quantifier-free and $|u| < N$. Since L is defined by this formula $\psi'(u, y)$, we get $L \in \text{PT}(N)$. Since $|L| \geq t$, the definition of $g_{|\Sigma|}$ and Lemma 4(1) allow to find t words in L of length $< N^{2^c}$. ◀

► **Proposition 23.** *There is an algorithm that, on input of a formula $\varphi(x, y) \in C_{\text{ptNFA}}^2$ and words u and v , decides whether $\mathcal{S} \models \varphi(u, v)$. This algorithm uses working space polylogarithmic in $|\varphi|$ and doubly exponential in $\|\varphi(u, v)\|$.*

Proof. We use a recursive procedure `check` whose parameters are

- a C_{ptNFA}^2 -formula $\alpha(x, y)$,
- two words w_x and w_y , and
- a natural number N .

It evaluates the formula $\alpha(w_x, w_y)$ recursively. When encountering a formula $\alpha(x, y) = \exists^{\geq t} y: \psi(x, y)$, it calls `check`(ψ, w_x, w'_y, N^{2^c}) for all words w'_y of length $< N^{2^c}$ and counts those that return true.

This procedure returns, if started with the correct value for N , the correct value due to Lemma 22. ◀

³ For first-order logic, the use of Corollary 20 can be replaced by the corresponding statements from [7, Cor. 7.4 and Thm. 7.5].

Since $||\varphi|| \leq |\varphi|$, we immediately obtain

► **Theorem 24.** *The C_{ptNFA}^2 -theory of \mathcal{S} is decidable in doubly exponential space.*

6 Summary and Open Question

We considered the extension of first-order logic by threshold-counting quantifiers over the subword order with piecewise testable predicates and constants. We showed that the 2-variable fragment of this theory is decidable in two-fold exponential space. This extends a result from [7] in two aspects: first, we add threshold counting quantifiers and piecewise testable predicates to first-order logic and, secondly, we improve their upper bound by one exponent. Our proof relies on two independent aspects: the consideration of the height of definable languages (which is a direct continuation from [7]) and an adaptation of Ferrante and Rackoff’s method [1].

The work done in this paper can be continued in the following directions:

- Addition of further binary relations: Let \mathcal{C} be some collection of binary relations on Σ^* such that Boolean combinations of relations from $\mathcal{C} \cup \{\sqsubseteq\}$ are effectively rational. This holds, e.g., if \mathcal{C} consists of the prefix relation, the relation “have equal length”, the cover relation as well as powers thereof (e.g., the relation “ $u \sqsubseteq v$ and $|v| - |u| = k$ ” for fixed $k \in \mathbb{N}$). Then the proof of [6, Theorem 5.5] can be extended to show the following result: The FO_{NFA}^2 -theory of the extension of the structure \mathcal{S} with the binary relations from \mathcal{C} is decidable. If the Boolean combinations are even effectively unambiguous rational, then the C_{NFA}^2 -theory becomes decidable using the arguments from [13] (where the result is demonstrated in case \mathcal{C} contains the cover relation, only).
It is not clear for which sets \mathcal{C} the C_{ptNFA}^2 -theory becomes decidable in elementary space (which is the case for $\mathcal{C} = \emptyset$ as demonstrated in this paper). The same question applies already for the FO_{\emptyset}^2 -theory.
- Addition of regular predicates: By [13], the C_{NFA}^2 -theory is decidable, but the only known algorithm is non-elementary. On the other hand, the C_{ptNFA}^2 -theory is decidable using elementary space. It is not clear whether there are other classes of nfAs $A \subseteq \text{NFA}$ such that the C_A^2 - or FO_A^2 -theory are decidable in elementary space.

References

- 1 J. Ferrante and Ch. Rackoff. *The computational complexity of logical theories*. Lecture Notes in Mathematics, vol. 718. Springer, 1979.
- 2 A. Finkel and Ph. Schnoebelen. Well-structured transition systems everywhere! *Theoretical Computer Science*, 256:63–92, 2001.
- 3 S. Halfon, Ph. Schnoebelen, and G. Zetsche. Decidability, complexity, and expressiveness of first-order logic over the subword ordering. In *LICS’17*, pages 1–12. IEEE Computer Society, 2017.
- 4 G. Higman. Ordering by divisibility in abstract algebras. *Proc. London Math. Soc.*, 2:326–336, 1952.
- 5 J. Ježek and R. McKenzie. Definability in substructure orderings. I: Finite semilattices. *Algebra Univers.*, 61(1):59–75, 2009.
- 6 P. Karandikar and Ph. Schnoebelen. Decidability in the logic of subsequences and super-sequences. In *FSTTCS’15*, Leibniz International Proceedings in Informatics (LIPIcs) vol. 45, pages 84–97. Leibniz-Zentrum für Informatik, 2015.
- 7 P. Karandikar and Ph. Schnoebelen. The height of piecewise-testable languages and the complexity of the logic of subwords. *Logical Methods in Computer Science*, 15(2), 2019.

61:12 Complexity of Counting First-Order Logic for the Subword Order

- 8 O. V. Kudinov and V. L. Selivanov. Undecidability in the homomorphic quasiorder of finite labelled forests. *J. Log. Comput.*, 17(6):1135–1151, 2007.
- 9 O. V. Kudinov, V. L. Selivanov, and L. V. Yartseva. Definability in the subword order. In *CiE'10*, Lecture Notes in Comp. Science vol. 6158, pages 246–255. Springer, 2010.
- 10 O. V. Kudinov, V. L. Selivanov, and A. V. Zhukov. Definability in the h-quasiorder of labeled forests. *Ann. Pure Appl. Logic*, 159(3):318–332, 2009.
- 11 D. Kuske. Theories of orders on the set of words. *Theoretical Informatics and Applications*, 40:53–74, 2006.
- 12 D. Kuske. The subtrace order and counting first-order logic. In *CSR'20*, Lecture Notes in Comp. Science vol. 12159, pages 289–302. Springer, 2020.
- 13 D. Kuske and G. Zetsche. Languages ordered by the subword order. In *FoSSaCS'19*, Lecture Notes in Comp. Science vol. 11425, pages 348–364. Springer, 2019.
- 14 M. Lothaire. *Combinatorics on Words*. Addison-Wesley, 1983.
- 15 T. Masopust. Piecewise testable languages and nondeterministic automata. In *MFCS'16*, Leibniz International Proceedings in Informatics (LIPIcs) vol. 58, pages 67:1–67:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2016.
- 16 T. Masopust and M. Thomazo. On boolean combinations forming piecewise testable languages. *Theoretical Computer Science*, 682:165–179, 2017.
- 17 Ph. Schnoebelen. personal communication, February 2020.
- 18 I. Simon. *Hierarchies of events with dot-depth one*. PhD thesis, University of Waterloo, 1972.
- 19 I. Simon. Piecewise testable events. In *Automata Theory and Formal Languages*, Lecture Notes in Comp. Science vol. 33, pages 214–222. Springer, 1975.
- 20 Ramanathan S. Thinniyam. Definability of recursive predicates in the induced subgraph order. In *7th Indian Conference on Logic and Its Applications (ICLA'17)*, Lecture Notes in Comp. Science vol. 10119, pages 211–223. Springer, 2017.
- 21 Ramanathan S. Thinniyam. Defining recursive predicates in graph orders. *Logical Methods in Computer Science*, 14(3), 2018.