



# Alternating complexity of counting first-order logic for the subword order

Dietrich Kuske<sup>1</sup> · Christian Schwarz<sup>1</sup>

Received: 1 June 2021 / Accepted: 1 May 2022 / Published online: 26 June 2022  
© The Author(s) 2022

## 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 twofold exponential alternating time with linearly many alternations (and therefore in particular in twofold exponential space as announced in the conference version (Kuske and Schwarz, in: MFCS'20, Leibniz International Proceedings in Informatics (LIPIcs) vol. 170, pp 56:1–56:13. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020) of this paper) 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 in the space bound. 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 generalises 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.

## 1 Introduction

The subword relation is one of the simplest nontrivial examples of a well-quasi-ordering [5] and can be used in the verification of infinite state systems [4]. 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. [7, 8, 12, 13, 15, 23, 24].

Many of these papers study logical aspects of the embeddability relation. Regarding the subword relation, the 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 [14] and the  $\exists^*\forall^*$ -theory is undecidable [9]. For the subword order together with regular predicates, the two-variable theory is decidable [9] (this holds even for

---

✉ Dietrich Kuske  
dietrich.kuske@tu-ilmenau.de

<sup>1</sup> Technische Universität Ilmenau, Ilmenau, Germany

the two-variable fragment of the logic  $C+MOD$ , i.e., the extension of first-order logic by threshold- and modulo-counting quantifiers [16]) and the three-variable theory [9] as well as the  $\exists^*$ -theory are undecidable [6] (these two undecidabilities already hold if we only consider singleton predicates, i.e., constants). Recently, Baumann et al. [1] strengthened the last undecidability result by showing that all semi-decidable languages can be defined by an existential formula using constants (even more, a language belongs to the  $n^{th}$  existential level of the arithmetical hierarchy if, and only if, it can be defined by a  $\Sigma_n$ -formula).

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 [9]. 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 [9] 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 deterministic finite automaton (henceforth dfa) accepting this set. The only known upper bound for the size of 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 a nondeterministic finite automaton, i.e., nfa) and then to determinise and minimise this nfa. The crucial insight from the follow-up paper [10] 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 threefold exponential space<sup>1</sup> which is also an upper bound for the said theory (the best lower bound we know so far is PSPACE [9]). This bound on the size of the minimal dfas is possible since all defined languages are piecewise testable [20]. A useful complexity measure for piecewise testable languages is their height. The new and innovative contribution of the proof from [10] 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).<sup>2</sup>

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

1. We prove an upper bound of twofold exponential alternating time with linearly many alternations (which implies an upper bound of twofold exponential space, i.e., the result we announced in the conference version of this paper [11]).
2. We allow piecewise testable predicates given by so-called pt-nfas [17, 18] (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 the two-variable fragment of first-order logic by threshold counting quantifiers  $\exists^{\geq t}$  (from [16], we know that this theory is decidable, even with regular predicates).

<sup>1</sup> The claim of threefold exponential time from [10, Thm. 7.5] is not supported by the proof idea [19].

<sup>2</sup> This view indicates that the result from [10] can be improved by allowing, instead of singleton predicates, piecewise testable predicates given by minimal dfas. Also then, the algorithm from [9] should run in threefold exponential space.

Following and extending the ideas from [10], 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 at least  $t$  subwords in  $L$  (and similarly for  $L\downarrow_{\geq t}$  and  $L\parallel_{\geq t}$ ). These considerations can be found in Sect. 3.

From these results, it follows that a 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 4.3).

**Remark** Consequently,  $L$  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 4.5).

So far, this parallels the development in [10] where the corresponding result was shown for first-order logic. But at this point, instead of building automata (as done in [10]), we follow another path of argument, that is an adaptation of Ferrante and Rackoff’s method [3].

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 predicates (Corollary 4.4). From this, we derive that quantification in formulas can be restricted to words of doubly exponential length. This implies that the two-variable fragment of the threshold counting extension of first-order logic becomes decidable in twofold exponential alternating time with linearly many alternations (allowing piecewise testable predicates in the formula given by pt-nfas).

## 2 Definitions and main results

Throughout this paper, we fix an alphabet  $\Sigma$ . We denote by  $\Sigma^*$  the set of (finite) words over  $\Sigma$ . A word  $u \in \Sigma^*$  is a *subword* of  $v \in \Sigma^*$  if  $u = u_1u_2 \dots u_n$  and  $v = v_0u_1v_1u_2v_2 \dots u_nv_n$  for some  $n \in \mathbb{N}$  and  $u_i, v_i \in \Sigma^*$ . We write  $u \sqsubseteq v$  for this fact and alternatively say that  $v$  is a *superword* of  $u$ . Finally, we write  $u \parallel v$  if neither  $u$  is a subword of  $v$  nor *vice versa*; we say that  $u$  and  $v$  are *incomparable*. Note that for any two distinct words  $u$  and  $v$ , we have precisely one of the three relations  $u \sqsubseteq v$ ,  $u \supseteq v$ , or  $u \parallel v$ .

Let  $L \subseteq \Sigma^*$  be a language. Its *upwards closure* is the language  $L\uparrow = \{v \in \Sigma^* \mid \exists u \in L : u \sqsubseteq v\}$  of all words  $v$  that have some subword  $u$  in  $L$ . Dually, the *downwards closure* of  $L$  is the language  $L\downarrow = \{u \in \Sigma^* \mid \exists v \in L : u \sqsubseteq v\}$  of all words  $u$  that have some superword  $v$  in  $L$ . Finally, the *incomparability set* of  $L$  is the language  $L\parallel = \{u \in \Sigma^* \mid \exists v \in L : u \parallel v\}$  of all words  $u$  that have some incomparable word  $v$  in  $L$ .

Note that, for any language  $L$ , we have  $L \subseteq L\uparrow \cap L\downarrow$ , i.e., these two sets need not be disjoint. For, e.g.,  $L = \{aa, bb\}^*$ , we even get  $L\uparrow = L\downarrow = L\parallel = \Sigma^*$  provided  $\Sigma = \{a, b\}$ .

### 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** 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  $u \in \Sigma^*$  (which is equivalent to saying that  $L$  is closed under  $\sim_n$ ). 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  $\text{PT}(n) \subseteq \text{PT}(n + 1)$ , and that both  $\emptyset$  and  $\Sigma^*$  are of height 0. Since the set of equivalence classes  $[u]_n$  forms a partition of  $\Sigma^*$ , the class  $\text{PT}(n)$  is closed under Boolean operations. 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  $\text{PT}(n)$ .

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  $\Sigma$ ) [10]. 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

$$L\uparrow_{\geq t} = \{v \in \Sigma^* \mid \exists u_1, \dots, u_t \in L \text{ pairwise distinct: } u_i \sqsubseteq v \text{ for all } 1 \leq i \leq t\}$$

denotes the set of words  $v$  that have  $\geq t$  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 language  $L\uparrow_{\geq t}$  is *upwards closed* (i.e., satisfies  $(L\uparrow_{\geq t})\uparrow = L\uparrow_{\geq t}$ ) and therefore piecewise testable.

Dually, the set

$$L\downarrow_{\geq t} = \{u \in \Sigma^* \mid \exists v_1, \dots, v_t \in L \text{ pairwise distinct: } u \sqsupseteq v_i \text{ for all } 1 \leq i \leq t\}$$

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

Let

$$L\parallel_{\geq t} = \{u \in \Sigma^* \mid \exists v_1, \dots, v_t \in L \text{ pairwise distinct: } u \parallel v_i \text{ for all } 1 \leq i \leq t\}$$

contain all words  $u$  that are incomparable with  $\geq t$  words from  $L$ .

We will also write, e.g.,  $L\parallel^{<t}$  for the complement of  $L\parallel_{\geq t}$ , i.e., for the set of words that are incomparable with at most  $t - 1$  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 [10, Thm. 3.7 & 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 Sect. 3): it generalises [10, Theorems 4.4, 5.5, and 6.1] from  $t = 1$  to general thresholds.

**Theorem 2.1** *Let  $\Sigma$  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 to a consequence in logic, we shortly recall some results on the relation of nondeterministic finite automata (abbreviated nfa) and piecewise testable languages.

There are different characterisations of piecewise testable languages using nfas; we only rely on one by Masopust and Thomazo [17, 18] (see following remark for missing definition). They define a class of nondeterministic finite automata, called pt-nfa and prove the following:

- A language is piecewise testable iff it is accepted by some pt-nfa [18, Thm. 25].
- Further, the depth  $||\mathcal{A}||$  of a pt-nfa (i.e., the maximal length of a simple path) bounds the height of the accepted language [17, Thm. 8].

**Remark** The concrete definition of a pt-nfa is of no importance for this paper; we only recall it for the convenience of the interested reader.

An nfa is a tuple  $\mathcal{A} = (Q, I, T, F)$  such that  $Q$  is a finite set,  $I, F \subseteq Q$ , and  $T \subseteq Q \times \Sigma \times Q$ . For  $p, q \in Q$  and  $\Gamma \subseteq \Sigma$ , we write  $p \xrightarrow{\Gamma^*} q$  whenever there exists a word over  $\Gamma$  that labels some path from  $p$  to  $q$ . The *depth* of the nfa  $\mathcal{A}$  is the maximal length of a simple path. The *language*  $L(\mathcal{A})$  of the nfa  $\mathcal{A}$  is the set of words over  $\Sigma$  that label some path from some element of  $I$  to some element of  $F$ .

Let  $\mathcal{A} = (Q, I, T, F)$  be an nfa. For  $r \in Q$ , we write  $\Sigma_r$  for the set of letters  $a \in \Sigma$  with  $(r, a, r) \in T$ . The nfa  $\mathcal{A}$  is a *pt-nfa* [17, Def. 3] if the following hold:

- The reachability relation is a partial order (i.e.,  $p \xrightarrow{\Sigma^*} q \xrightarrow{\Sigma^*} p$  implies  $p = q$ ). (An nfa satisfying this property is called *acyclic*.)
- For all  $p, q \in Q$ ,  $p \xrightarrow{\Sigma_p^*} q$  implies  $p = q$ .
- For all  $p, q, r \in Q$ ,  $p \xrightarrow{\Sigma_r^*} q, r$  implies  $q \xrightarrow{\Sigma_r^*} r$ .

## 2.2 The logic $C^2$ and the main result for logicians

Let  $NFA$  be the set of all nfes over the alphabet  $\Sigma$  (to make this a set as opposed to a class, we require that states of these nfes belong to  $\mathbb{N}$ ). Consider the structure

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

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  $\mathcal{A} \in NFA$  and a constant for every word over  $\Sigma$ .

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

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

where  $c, d \in \{x, y\} \cup \Sigma^*$  are variables from  $\{x, y\}$  or words over  $\Sigma$ ,  $\mathcal{A} \in A$  is some nfa over  $\Sigma$ ,  $t \in \mathbb{N}$ , 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  $\varphi$  true. Consequently,  $\exists^{\geq 1}$  is the usual existential quantifier and  $\exists^{\geq 0} x \varphi$  is always true. Let  $FO^2_A$  denote the subset of  $C^2_A$  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  $C^2_{ptNFA}$  and  $FO^2_{\emptyset}$  are equally expressive by Corollary 4.4, but  $C^2_{ptNFA}$  is exponentially more succinct than  $FO^2_{\emptyset}$  by Theorem 4.5.

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

**Theorem 2.2** *Let  $A \subseteq NFA$ . For every formula  $\varphi \in C^2_A$ , there exists an equivalent formula  $\psi \in C^2_A$  that uses constants of length  $\leq 2$ , only. The same applies to the logic  $FO^2_A$ .*

**Proof** We show that, for every word  $w \in \Sigma^*$ , there exists a formula  $\lambda_w(x) \in \text{FO}_{\emptyset}^2$  using at most constants of length  $\leq 2$  such that  $w$  is the only word satisfying  $\lambda_w(x)$ .

Before we start the construction of  $\lambda_w(x)$ , consider the following inductively defined formula  $\alpha_n(z)$  (where  $z$  is any variable from  $\{x, y\}$  and  $z'$  is the other variable):

$$\alpha_n(z) = \begin{cases} z = z & \text{if } n = 0 \\ \exists z': z' \sqsubseteq z \wedge z' \neq z \wedge \alpha_{n-1}(z') & \text{otherwise.} \end{cases}$$

Then,  $\mathcal{S} \models \alpha_n(u)$  iff  $|u| \geq n$ .

We now come to the construction of  $\lambda_w(x)$  by induction on the length of  $w$ . If  $|w| \leq 2$ , we simply set  $\lambda_w(x) = w$ .

Now let  $n = |w| > 2$  and define  $m = \lfloor n/2 \rfloor + 1 < n$ . Define  $S_w = \{u \in \Sigma^{\leq m} \mid u \sqsubseteq w\}$ . Then,  $\lambda_w(x)$  is the following formula:

$$\begin{aligned} & \alpha_n(x) \wedge \neg \alpha_{n+1}(x) \\ & \wedge \bigwedge_{u \in S_w} \exists y: \lambda_u(y) \wedge y \sqsubseteq x \\ & \wedge \bigwedge_{u \in \Sigma^{\leq m} \setminus S_w} \exists y: \lambda_u(y) \wedge \neg y \sqsubseteq x. \end{aligned}$$

The first two conjuncts express  $|x| = n$ , i.e., the length of  $x$  equals that of  $w$ . By the induction hypothesis,  $\lambda_u(y)$  expresses  $y = u$ . Consequently, the latter two conjuncts are equivalent to  $x \sim_m w$ .

In other words,  $\mathcal{S} \models \lambda_w(v)$  iff  $|v| = |w|$  and  $v \sim_m w$ . But this is equivalent to  $v = w$  [22, Thm. 6.2.16]. □

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  $\varphi$  from  $\text{C}_{\text{ptNFA}}^2$  (recall that  $\|\mathcal{A}\|$  denotes the depth of the pt-nfa  $\mathcal{A}$ ):

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

Note that this norm  $\|\varphi\|$  forms a mixture between the size of a formula and its quantifier depth: It depends on the maximal size of constants and simple paths in automata appearing in  $\varphi$  as well as on the quantifier depth (where the quantifier  $\exists^{\geq t}$ , that intuitively corresponds to a sequence of  $t$  quantifiers, contributes only  $\lceil \log(t) \rceil$  to the norm). In particular,  $\|\varphi\|$  bounds the length of constants and the depth of pt-nfas occurring in  $\varphi$ . Note further that the norm  $\|\varphi\|$  of any formula  $\varphi$  is at most its size  $|\varphi|$ , i.e.,  $\|\varphi\| \leq |\varphi|$ .

From Theorem 2.1, we infer in Sect. 4 that all definable languages are piecewise testable of bounded height (Theorem 4.3). This allows to derive a quantifier elimination result that reads as follows:

**Corollary 4.4.** *Let  $c = 2 \cdot |\Sigma|$ . Every  $\text{C}_{\text{ptNFA}}^2$ -formula  $\varphi$  is equivalent to some quantifier- and automata-free formula  $\psi \in \text{FO}_{\emptyset}^2$  with  $\|\psi\| < 2^{c^{2\|\varphi\|}}$ .*

Karandikar and Schnoebelen [10] showed that any non-empty piecewise testable language of height  $n$  has elements of length polynomial in  $n$ . Based on Corollary 4.4, we can therefore restrict quantification in a formula  $\varphi$  to words of bounded length, implying our main result for logicians.

**Theorem 5.3.** *The  $C^2_{\text{ptNFA}}$ -theory of  $\mathcal{S}$  belongs to  $\text{STA}(*, 2^{2^{\text{poly}(n)}}, O(n))$ , i.e., can be decided in doubly exponential alternating time with linearly many alternations.*

Recall that, by [2],  $\text{STA}(s, t, a)$  is the class of all languages, for which membership can be decided by an alternating Turing machine whose space, time, and alternations are bounded by the functions  $s, t$ , and  $a$ , respectively. Typically,  $*$  is used to denote that no restriction is placed on a specific resource. Thus,  $\text{STA}$  is a combined complexity measure that is particularly useful when describing the complexity of logical theories (see, e.g., [2, 3]).

### 3 Closure of the class of piecewise testable languages

The purpose of this section is to prove Theorem 2.1, i.e., our main result for language theorists.

#### 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  $v \in L$ . It is a *chain* if it is linearly ordered by the subword order and if it is infinite. Since the subword order is well-founded, any chain is isomorphic to  $(\mathbb{N}, \leq)$ .

**Lemma 3.1** (compiled in [10]) *Let  $u, v \in \Sigma^*$ ,  $a \in \Sigma$ , 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 [10]) (1) is by combining the definition of  $\sim_n$  with the observation  $\{u\} \downarrow \subseteq \{v\} \downarrow$  provided  $u \sqsubseteq v$ . (2) is [21, Lemma 6] (cf. [22, Thm. 6.2.6] for an alternative proof). (3) is in the proof of [22, Cor. 6.2.8]. Finally, (4) follows from (1), (2), and (3).  $\square$

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  $[aa]_2 = aaa^*$ ).

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, i.e.,  $\text{Perm}(\emptyset) = \{\varepsilon\}$  and  $\text{Perm}(B) = \bigcup_{b \in B} b \text{Perm}(B \setminus \{b\})$  for  $B \neq \emptyset$ . For sets  $B_i \subseteq \Sigma$ , define  $\text{Perm}(B_1, B_2, \dots, B_k) = \text{Perm}(B_1) \text{Perm}(B_2) \dots \text{Perm}(B_k)$ . For instance,  $\text{Perm}(\{a\}, \{b\}, \{c\}) = \{abc\}$  while  $\text{Perm}(\{a, b\}, \{c\}) = \{abc, bac\}$  for all letters  $a, b, c \in \Sigma$ . For  $k = 0$ , we set  $\text{Perm}(\ ) = \{\varepsilon\}$ .

**Theorem 3.2** ([20], cf. [22, Thm. 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. The theorem implies in particular that all words from  $\min([u]_n)$  have the same Parikh image. Further, they all 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 [10, Thm. 3.7 and Eq. (3.12)]).

**Theorem 3.3** *Let  $\Sigma$  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  $v \sqsubseteq w$ .*

**Proof** The definition of the function  $g_{|\Sigma|}$  implies the existence of some word  $u' \sim_n w$  with  $|u'| \leq g_{|\Sigma|}(n)$ . Since the subword order is well-founded, there exist words  $u, v \in \min([w]_n)$  with  $u \sqsubseteq u'$  and  $v \sqsubseteq w$ . Now Theorem 3.2 implies  $|v| = |u| \leq |u'| \leq g_{|\Sigma|}(n)$ .  $\square$

### 3.2 Upward closures

The following result verifies the first claim of Theorem 2.1.<sup>3</sup>

**Proposition 3.4** *Let  $L \in \text{PT}(n)$  be a piecewise testable language of height  $\leq n$  and  $t \in \mathbb{N}$ . Then, the language  $L \uparrow_{\geq t}$  is piecewise testable of height  $\leq g_{|\Sigma|}(n) + t - 1$ .*

**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  $Y \downarrow \cap L = Y$ .

Now let  $y \in Y$ . By Theorem 3.3, there is a word  $x$  with  $x \sim_n y$ ,  $x \sqsubseteq y$ , and  $|x| \leq g_{|\Sigma|}(n)$ .

Since  $x$  is a subword of  $y$ , there are more than  $|y| - |x|$  words  $x'$  with  $x \sqsubseteq x' \sqsubseteq y$ . Since the equivalence class  $[y]_n$  is convex, any such word  $x'$  satisfies  $x' \sim_n y$  and therefore  $x' \in L$ . Consequently,

$$t = |Y| = |Y \downarrow \cap L| \geq |y \downarrow \cap L| > |y| - |x| \geq |y| - g_{|\Sigma|}(n).$$

But this implies  $|y| \leq g_{|\Sigma|}(n) + t - 1$ .

So far, we proved that all words from  $Y$  have length at most  $g_{|\Sigma|}(n) + t - 1$ . Since they all are subwords of  $z \sim_{g_{|\Sigma|}(n)+t-1} z'$ , we obtain  $y \sqsubseteq z'$  for all  $y \in Y$ . From  $|Y| = t$  and  $Y \subseteq L$ , we derive  $z' \in L \uparrow_{\geq t}$ , i.e.,  $L \uparrow_{\geq t}$  is closed under  $\sim_{g_{|\Sigma|}(n)+t-1}$ .  $\square$

### 3.3 Downward closures

To verify the second claim of Theorem 2.1, we first prove that only singleton equivalence classes  $[x]_n$  have maximal elements. We will use this lemma in the following proof when showing that  $[x]_n \downarrow^{\geq t} = [x]_n \downarrow$  if the equivalence class  $[x]_n$  is not a singleton.

**Lemma 3.5** *Let  $n \in \mathbb{N}$  and  $x, y \in \Sigma^*$  be distinct with  $x \sim_n y$ . Then, there exists  $z \in \Sigma^*$  with  $y \sim_n z$ ,  $y \sqsubseteq z$ , and  $y \neq z$ .*

**Proof** Since  $[x]_n = [y]_n$  is not a singleton, it is infinite by Lemma 3.1(4), thus contains in particular a word  $w$  of length  $|w| > |y|$ . By Lemma 3.1(2), there exists a  $z \in [y]_n$  with  $y, w \sqsubseteq z$ , implying  $|z| \geq |w|$ , and therefore  $z \neq y$ .  $\square$

<sup>3</sup> The operation  $\uparrow_{\geq t}$  for  $t \geq 2$  is not idempotent and therefore not a closure operator. For convenience, we stick to this name.



**Proposition 3.6** *Let  $L \in \text{PT}(n)$  be a language over  $\Sigma$  and  $t \in \mathbb{N}$ . Then, the language  $L \downarrow^{\geq t}$  belongs to  $\text{PT}((|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1))$ .*

**Proof** Since  $L \in \text{PT}(n)$  and since  $\sim_n$  has finite index, there are finitely many words  $x_1, \dots, x_m$  with  $L = \bigcup_{1 \leq i \leq m} [x_i]_n$  and  $x_i \not\sim_n x_j$  for all  $1 \leq i < j \leq m$ . By the definition of the function  $g_{|\Sigma|}$ , we can assume  $|x_i| \leq g_{|\Sigma|}(n)$  for all  $1 \leq i \leq m$ .

Set

$$F = \bigcup_{\substack{1 \leq j \leq m \\ [x_j]_n \text{ finite}}} [x_j]_n \text{ and } I = \bigcup_{\substack{1 \leq i \leq m \\ [x_i]_n \text{ infinite}}} [x_i]_n$$

such that in particular  $L = F \cup I$ .

We first show

$$L \downarrow^{\geq t} = F \downarrow^{\geq t} \cup I \downarrow.$$

For the inclusion “ $\subseteq$ ”, let  $x \in L \downarrow^{\geq t} \setminus F \downarrow^{\geq t}$ . Then,  $x$  has  $t$  superwords in  $L = F \cup I$ , but at most  $t - 1$  many in  $F$ . Hence, it has at least one superword in  $I$ , i.e.,  $x \in I \downarrow$ . For the converse inclusion, note that  $F \downarrow^{\geq t} \subseteq L \downarrow^{\geq t}$  is trivial since  $F \subseteq L$ . So let  $x \in I \downarrow$ . Then, there exists  $y \in I$  with  $x \sqsupseteq y$ . Since  $y \in I$ , the equivalence class  $[y]_n \subseteq I$  is infinite and therefore contains no maximal element by Lemma 3.5. Hence, there are infinitely many (and therefore in particular  $\geq t$ ) superwords of  $y \sqsupseteq x$  in  $I \subseteq L$ . Consequently,  $x \in L \downarrow^{\geq t}$ .

Note that the height of  $I$  is  $\leq n$  since it is a union of equivalence classes of  $\sim_n$ . Consequently, the height of  $I \downarrow$  is  $\leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1)$  by [10, Thm. 5.5].

Since every finite equivalence class  $[x_j]_n$  is a singleton, we obtain

$$F = \{x_j \mid [x_j]_n \text{ finite}\},$$

implying that all words from  $F$  and therefore from  $F \downarrow^{\geq t}$  have length at most  $g_{|\Sigma|}(n)$ . Hence,  $F \downarrow^{\geq t}$  is finite and thus of height  $\leq g_{|\Sigma|}(n) + 1 \leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + 1)$ .

We showed that both,  $F \downarrow^{\geq t}$  and  $I \downarrow$ , are closed under  $\sim_{(|\Sigma|+1) \cdot (g_{|\Sigma|}(n)+1)}$ ; hence, the same holds for their union  $L \downarrow^{\geq t}$ . □

### 3.4 Incomparability set

There are three types of equivalence classes  $[x]_n$ : the singletons, the chains (i.e., infinite languages ordered linearly by the subword order), and the infinite ones which are no chains. Note that by Lemma 3.1(4) this is a complete characterization of the equivalence classes. Propositions 3.7, 3.9, and 3.14, respectively, bound the heights of  $[x]_n \parallel^{\geq t}$  for these three types of equivalence classes and collectively verify Theorem 2.1(3).

**Proposition 3.7** *Let  $n, t \in \mathbb{N}$  and  $x \in \Sigma^*$  such that  $L = [x]_n$  is a singleton. Then,  $L \parallel^{\geq t} \in \text{PT}(g_{|\Sigma|}(n))$ .*

**Proof** If  $t \geq 2$ , then  $L \parallel^{\geq t} = \emptyset$  since  $L$  is a singleton. If  $t = 0$ , then  $L \parallel^{\geq t} = \Sigma^*$ . Note that both these languages belong to  $\text{PT}(0) \subseteq \text{PT}(g_{|\Sigma|}(n))$ .

Finally, consider the case  $t = 1$ . Then,  $L \parallel^{\geq t} = \Sigma^* \setminus (L \uparrow \cup L \downarrow)$  since  $L$  is a singleton.

Note that  $L \uparrow \cup L \downarrow = L \uparrow \cup (L \downarrow \setminus \{x\})$  since  $x \in L \uparrow$ . The height of the former language is  $\leq |x|$ . The latter is finite, and all its elements have length  $< |x|$ ; hence, the height of that language is  $\leq |x|$  as well. Thus, the height of  $L \uparrow \cup L \downarrow$  is  $\leq |x|$  and the same applies to its

complement  $L \parallel^{\geq t}$ . Since  $L = [x]_n$  is a singleton, the definition of the function  $g_{|\Sigma|}$  implies  $|x| \leq g_{|\Sigma|}(n)$ . □

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

**Lemma 3.8** *Let  $t \geq 1$  and let  $C$  be the convex chain  $x_0 \sqsubsetneq x_1 \sqsubsetneq \dots$ . Then,  $C \parallel^{< t} = C \cup \{x_{t-1}\} \downarrow$ .*

Note that, provided  $x_0 \neq \varepsilon$ , the chain  $C$  is not maximal since it can be extended to the left.

**Proof** We first demonstrate the inclusion “ $\supseteq$ ”. Since any two elements of  $C$  are comparable, we clearly have  $C \parallel^{< t} \supseteq C$ . Further, any subword of  $x_{t-1}$  is a subword of all words  $x_{t-1+i}$  for  $i \in \mathbb{N}$  and therefore at most incomparable with the  $t - 1$  words  $x_0, x_1, \dots, x_{t-2}$  from  $C$ .

For the converse inclusion, let  $y \in C \parallel^{< t}$ . Then,  $y$  is comparable with infinitely many words from  $C$ . Since  $C$  has only finitely many words that are shorter than  $y$ , there is  $\ell \in \mathbb{N}$  with  $y \sqsubseteq x_\ell$ . Let  $\ell \in \mathbb{N}$  be minimal with this property. We distinguish three cases of the relation between  $x_0$  and  $y$ :

- If  $y \sqsubseteq x_0$ , then  $y \in \{x_{t-1}\} \downarrow$ .
- If  $y \parallel x_0$ , then  $x_0, x_1, \dots, x_{\ell-1} \parallel y$  since  $\ell$  was chosen minimal with  $y \sqsubseteq x_\ell$  and  $x_0 \sqsubseteq x_1 \sqsubseteq \dots \sqsubseteq x_{t-1}$  is a chain. From  $y \in C \parallel^{< t}$ , we obtain  $\ell < t$  and therefore  $y \sqsubseteq x_\ell \sqsubseteq x_{t-1}$ .
- If  $y \supseteq x_0$ , we have  $x_0 \sqsubseteq y \sqsubseteq x_\ell$  and therefore  $y \in C$  since  $C$  is convex.

Thus, in any case,  $y \in C \cup \{x_{t-1}\} \downarrow$ . □

**Proposition 3.9** *Let  $n, t \in \mathbb{N}$  and  $x \in \Sigma^*$  such that  $C = [x]_n$  is a chain. Then,  $C \parallel^{\geq t} \in \text{PT}(g_{|\Sigma|}(n) + t)$ .*

**Proof** If  $t = 0$ , then  $C \parallel^{\geq t} = \Sigma^* \in \text{PT}(0)$ . It therefore suffices to consider the case  $t > 0$ .

We list the elements of the chain  $C$  in increasing order:

$$x_0 \sqsubsetneq x_1 \sqsubsetneq x_2 \dots$$

Since  $C = [x]_n$  is a chain, it is a convex chain by Lemma 3.1(1) such that  $|x_i| = |x_0| + i$  holds for all  $i \geq 0$ . From Lemma 3.8, we obtain

$$C \parallel^{< t} = C \cup \{x_{t-1}\} \downarrow.$$

Since  $\{x_{t-1}\} \downarrow$  is finite, its height is  $\leq |x_{t-1}| + 1 = |x_0| + t - 1 + 1 \leq g_{|\Sigma|}(n) + t$ . The height of  $C$  is  $\leq n \leq g_{|\Sigma|}(n)$  by assumption; thus, the height of  $C \parallel^{< t}$  is  $\leq g_{|\Sigma|}(n) + t$ . But then the same bound applies to the height of  $C \parallel^{\geq t} = \Sigma^* \setminus C \parallel^{< t}$ . □

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 [10] first shows that  $[x]_n$  contains at least two elements of every length  $> |x|$ . Consequently, every word of length  $> |x|$  is incomparable with some word from  $[x]_n$ , i.e.,  $[x]_n \parallel^{\geq 1}$  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 3.13). Then, we apply Lemma 3.8. But first, we need some insight into convex chains which is the topic of the following considerations.

**Lemma 3.10** *Let  $x, y \in \Sigma^*$  and  $a \in \Sigma$ . Then,  $xa^*y$  is a convex chain.*

**Proof** Let  $x'$  be the longest prefix of  $x$  not ending with  $a$  (i.e.,  $x' \in \Sigma^* \setminus \Sigma^*a$  and  $x \in x'a^*$ ) and  $y'$  the longest suffix of  $y$  not beginning with  $a$ . Then,  $xa^*y \subseteq x'a^*y'$  is convex in  $(x'a^*y', \sqsubseteq)$  and we prove the stronger claim that the latter set is a convex chain.

To simplify notation, suppose  $x \in \Sigma^* \setminus \Sigma^*a$  and  $y \in \Sigma^* \setminus a\Sigma^*$ .

Clearly,  $xa^*y$  is a chain in  $(\Sigma^*, \sqsubseteq)$ .

Let  $w \in \Sigma^*$  and  $i, \ell \in \mathbb{N}$  with  $xa^i y \sqsubseteq w \sqsubseteq xa^\ell y$ . We have to show that  $w$  belongs to  $xa^*y$ . Note that  $xy \sqsubseteq w$  since  $xy \sqsubseteq xa^i y \sqsubseteq w$ . Let  $w_1$  be the prefix of length  $|x|$  of  $w$ ,  $w_3$  be the suffix of length  $|y|$  of  $w$ , and  $w_2$  be the unique word with  $w = w_1w_2w_3$ . Since  $w$  is a subword of  $xa^\ell y$ , dropping the first  $|x|$  letters in both  $w$  and  $xa^\ell y$  preserves the subword relation. The same holds when dropping the last  $|y|$  letters, hence  $w_2 \sqsubseteq a^\ell$ , i.e.  $w_2 \in a^*$ . By similar reasoning for  $xy \sqsubseteq w_1w_2w_3$ , we can conclude that  $x \sqsubseteq w_1w_2$ . Since  $w_2 \in a^*$ , but  $x$  does not end on  $a$ ,  $x$  has to be a subword of  $w_1$  and thus  $x = w_1$ , as both words are of the same length. Symmetrically, we can show  $y = w_3$ . Consequently, we have  $w \in xa^*y$ .  $\square$

The third item of the following lemma implies, together with Theorem 3.2, that the maximal  $a$ -prefixes of two words from  $\min([x]_n)$  differ in length by at most one.

**Lemma 3.11** *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  $aa u \in \text{Perm}(B_1, \dots, B_k)$ , 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  $|m - n| \leq 1$ .*

**Proof** Since  $B_1 \neq \emptyset$ , the first claim follows from

$$\begin{aligned} \text{Perm}(B_1, \dots, B_k) &= \text{Perm}(B_1) \cdot \text{Perm}(B_2, \dots, B_k) \\ &= \bigcup_{b \in B_1} b \cdot \text{Perm}(B_1 \setminus \{b\}) \cdot \text{Perm}(B_2, \dots, B_k) \\ &= \bigcup_{b \in B_1} b \cdot \text{Perm}(B_1 \setminus \{b\}, B_2, \dots, B_k). \end{aligned}$$

Now assume  $aa u \in \text{Perm}(B_1, \dots, B_k)$ . Then, by the above,  $a \in B_1$  and the word  $au$  belongs to the set  $\text{Perm}(B_1 \setminus \{a\}, B_2, \dots, B_k)$ . If, towards a contradiction,  $B_1 \neq \{a\}$ , then  $B_1 \setminus \{a\} \neq \emptyset$ . Hence, by the first claim again,  $a \in B_1 \setminus \{a\}$ , a contradiction. Thus, indeed,  $B_1 = \{a\}$ .

Towards a contradiction, assume  $u, v \notin a\Sigma^*$ ,  $|m - n| > 1$ , and  $a^m u, a^n v \in \text{Perm}(B_1, \dots, B_k)$ . Without loss of generality, we may assume  $m \geq n + 2$ . By the second claim, we get  $B_1 = \{a\}$  from  $m \geq n + 2 \geq 2$ . Hence,

$$a^{m-1}u, a^{n-1}v \in \text{Perm}(B_1 \setminus \{a\}, B_2, \dots, B_k) = \text{Perm}(B_2, \dots, B_k).$$

By induction, we obtain  $a^{m-n}u, v \in \text{Perm}(B_{n+1}, \dots, B_k)$ . Since  $m - n \geq 2$ , the second claim implies  $B_{n+1} = \{a\}$  and therefore  $v \in a\Sigma^*$ . But this contradicts our choice of  $v$ .  $\square$

**Lemma 3.12** *Let  $\overline{B}$  be a tuple of finite nonempty sets of letters,  $x_1, x_2, y_1, y_2 \in \Sigma^*$  be words with  $x_1x_2, y_1y_2 \in \text{Perm}(\overline{B})$ , and  $a, b \in \Sigma$  letters with  $x_1ax_2 \neq y_1by_2$ .*

*Then,  $x_1a^*x_2$  and  $y_1b^*y_2$  are convex chains that intersect, at most, in  $x_1x_2$ .*

**Proof** Without loss of generality, we assume  $|x_1| \leq |y_1|$ . Since  $|x_1x_2| = |y_1y_2|$ , we get  $|x_2| \geq |y_2|$ . By Lemma 3.10,  $C_1 = x_1a^*x_2$  and  $C_2 = y_1b^*y_2$  form convex chains.

It remains to be shown that their intersection is contained in  $\{x_1x_2\} = \{y_1y_2\}$ . So let  $v \in C_1 \cap C_2$ . Then, there exist non-negative integers  $\ell$  and  $m$  with  $v = x_1a^\ell x_2 = y_1b^m y_2$ .

Since the words  $x_1x_2$  and  $y_1y_2$  are of equal length, we have  $\ell = m$ . If  $\ell = 0$ , then  $v = x_1a^0x_2 = y_1b^0y_2$  is in  $\{x_1x_2\}$  and we are done. Thus, we may assume  $\ell > 0$ .

Since  $x_1x_2$  and  $y_1y_2$  both belong to  $\text{Perm}(\overline{B})$ , we get  $|x_1x_2|_a = |y_1y_2|_a$ , implying  $0 < \ell = |a^\ell|_a = |x_1a^\ell x_2|_a - |x_1x_2|_a = |y_1b^\ell y_2|_a - |y_1y_2|_a = |b^\ell|_a$  and therefore  $a = b$ .

Since  $x_1$  and  $y_1$  both are prefixes of  $x_1a^\ell x_2$  with  $|x_1| \leq |y_1|$ , the word  $x_1$  is a prefix of  $y_1$ , i.e., there is a word  $x'_1$  with  $y_1 = x_1x'_1$ . Symmetrically, we get a word  $y'_2$  with  $x_2 = y'_2y_2$ . From  $x_1x'_1a^\ell y_2 = y_1a^\ell y_2 = x_1a^\ell x_2 = x_1a^\ell y'_2y_2$ , we conclude  $x'_1a^\ell = a^\ell y'_2$  (and therefore in particular  $|x'_1| = |y'_2|$ ). Aiming at a contradiction, assume  $|x'_1| = |y'_2| \leq \ell$ . Then,  $x'_1$  is a prefix of  $a^\ell$  and similarly  $y'_2$  a suffix of  $a^\ell$ , hence  $x'_1 = y'_2 = a^k$  for some nonnegative integer  $k \in \mathbb{N}$ . But then  $y_1by_2 = y_1ay_2 = x_1a^kay_2 = x_1aa^k y_2 = x_1ax_2$ , as opposed to our assumption. Consequently  $|x'_1| = |y'_2| > \ell$ , implying that there exists  $k \in \mathbb{N}$  and a word  $w \in \Sigma^* \setminus a\Sigma^*$  such that  $x'_1 = a^\ell a^k w$  and  $y'_2 = a^k w a^\ell$ . If  $w = \varepsilon$ , then  $x'_1 = y'_2 = a^{k+\ell}$  and therefore (as above)  $y_1by_2 = y_1ay_2 = x_1a^{k+\ell}ay_2 = x_1aa^{k+\ell}y_2 = x_1ax_2$ , as opposed to our assumption. Hence,  $w = cw'$  for some letter  $c \neq a$  and some word  $w' \in \Sigma^*$ .

Note that

$$\begin{aligned} x_1a^k cw' a^\ell y_2 &= x_1x_2 \in \text{Perm}(\overline{B}) \text{ and} \\ x_1a^{\ell+k} cw' y_2 &= y_1y_2 \in \text{Perm}(\overline{B}). \end{aligned}$$

Applying Lemma 3.11(1), we obtain a tuple  $\overline{C}$  of non-empty subsets of  $\Sigma$  with  $a^k cw' a^\ell y_2, a^{\ell+k} cw' y_2 \in \text{Perm}(\overline{C})$ .

Since  $c \neq a$ , Lemma 3.11(3) implies  $|\ell + k - k| \leq 1$ , i.e.,  $\ell \leq 1$ . But  $\ell = 1$  is impossible since  $x_1ax_2 \neq y_1ay_2$ . Hence,  $\ell = 0$  and therefore  $v = x_1x_2$ .

Recall that we considered an arbitrary word  $v \in C_1 \cap C_2$  and derived  $v \in \{x_1x_2\}$ . Hence, indeed,  $C_1 \cap C_2 \subseteq \{x_1x_2\}$ . □

Note that, in the lemma above, the two words  $x_1x_2$  and  $y_1y_2$  have the same Parikh image. However, replacing the requirement  $x_1x_2, y_1y_2 \in \text{Perm}(\overline{B})$  by this weaker property does not suffice for the claim of the lemma: consider  $x_1 = aac, y_2 = caa, x_2 = y_1 = \varepsilon$ , and  $a = b$ . Then,  $x_1x_2 = aac$  and  $y_1y_2 = caa$  satisfy the modified prerequisites, but  $x_1a^*x_2 = aaca^*$  and  $y_1b^*y_2 = b^*caa = a^*caa$  are two convex chains that intersect in  $aacaa$ .

**Lemma 3.13** *Let  $u \in \Sigma^*$  and  $n \in \mathbb{N}$  such that  $[u]_n$  is infinite but not a single chain. Then,  $[u]_n$  contains two 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** Since  $[u]_n$  is infinite but not a single chain, [10, Lemma 6.2 and 6.3] implies that there are words  $x_1, x_2, y_1, y_2 \in \Sigma^*$  and letters  $a, b \in \Sigma$  such that  $x_1x_2, y_1y_2 \in \min([u]_n)$ ,  $x_1ax_2, y_1by_2 \in [u]_n$ , and  $x_1ax_2 \neq y_1by_2$ .

By Theorem 3.2, there exists a tuple  $\overline{B}$  of nonempty subsets of  $\Sigma$  such that  $x_1x_2, y_1y_2 \in \min([u]_n) \subseteq \text{Perm}(\overline{B})$ . By Lemma 3.12,  $x_1a^*x_2$  and  $y_1b^*y_2$  are convex chains whose intersection is contained in  $\{x_1x_2\}$ . By Lemma 3.1(3) they are both subsets of  $[u]_n$ , obviously containing elements from  $\min([u]_n)$ . □

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 3.14** *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** If  $t = 0$ , then  $L \parallel^{\geq t} = \Sigma^* \in \text{PT}(0)$ . Hence, it remains to consider the case  $t > 0$ .

By Lemma 3.13, there exist two 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  $i \in \{1, 2\}$ . We prove that

$$L \parallel^{< t} \subseteq \Sigma^{< g_{|\Sigma|}(n) + t}.$$

Let  $v \in \Sigma^*$  with  $|v| \geq g_{|\Sigma|}(n) + t > g_{|\Sigma|}(n)$ . Then, by Theorem 3.2 and the definition of the function  $g_{|\Sigma|}$ ,  $v \notin \min(L)$  implying  $v \notin C_1 \cap C_2$ , without loss of generality, we assume  $v \notin C_1$ . Since  $C_1 \cap \min(L) \neq \emptyset$ , the chain  $C_1$  contains some word of length  $\leq g_{|\Sigma|}(n)$ . Consequently, its word  $x_{t-1}$  number  $t - 1$  satisfies  $|x_{t-1}| < g_{|\Sigma|}(n) + t \leq |v|$ , i.e.,  $v$  cannot be a subword of  $x_{t-1}$ . Now Lemma 3.8 implies  $v \notin C_1 \parallel^{< t}$ . From  $C_1 \subseteq L$ , we now obtain  $v \in C_1 \parallel^{\geq t} \subseteq L \parallel^{\geq t}$ . Consequently,  $v \notin L \parallel^{< t}$  which proves the above claim.

Since all words in  $L \parallel^{< t}$  are “short”, we obtain  $L \parallel^{< t} \in \text{PT}(g_{|\Sigma|}(n) + t)$  and the same holds for the complement  $L \parallel^{\geq t}$  of this set.  $\square$

We can now put the above three propositions together to verify the last claim of Theorem 2.1.

**Proposition 3.15** *Let  $L \in \text{PT}(n)$  be a language over  $\Sigma$  and  $t \in \mathbb{N}$ . Then,  $L \parallel^{\geq t} \in \text{PT}(g_{|\Sigma|}(n) + t)$ .*

**Proof** Since  $L$  is of height  $\leq n$ , there is a finite set of words  $\{x_1, \dots, x_m\}$  with  $x_i \approx_n x_j$  for all  $1 \leq i < j \leq m$  such that  $L$  is the union of the equivalence classes  $[x_i]_n$ . Since equivalence classes are disjoint, we obtain

$$L \parallel^{\geq t} = \bigcup_{1 \leq i \leq m} \bigcap [x_i]_n \parallel^{\geq g(i)}$$

where the union is taken over all functions  $g: \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, t\}$  with  $\sum_{1 \leq i \leq m} g(i) = t$ . The previous propositions show that any of the languages  $[x_i]_n \parallel^{\geq s}$  is piecewise testable of height  $\leq g_{|\Sigma|}(n) + t$ . Since the class  $\text{PT}(g_{|\Sigma|}(n) + t)$  is closed under Boolean operations, the claim follows.  $\square$

### 4 Expressive power and quantifier elimination

Having completed the language-theoretic part of this paper, we now come to its consequences in logic, i.e., we consider the threshold counting logic  $C_{\text{ptNFA}}^2$  that has two variables  $x$  and  $y$ , unary predicates for each piecewise testable language (represented by some pt-nfa), the subword order, a constant for every word, and threshold quantifiers of the form  $\exists^{\geq t}$  for  $t \in \mathbb{N}$ . The central result, Theorem 4.3, states that every language definable in this logic 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 4.1** *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  $\|\varphi\| \leq n$ .*
- (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** (1) Since  $L \in \text{PT}(n)$ , it is a finite union of equivalence classes  $[v]_n$  for  $v \in \Sigma^*$ . Such an equivalence class  $[v]_n$  can be defined by the formula

$$\varphi(x) = \bigwedge_{\substack{u \sqsubseteq v \\ |u| \leq n}} u \sqsubseteq x \wedge \bigwedge_{\substack{u \not\sqsubseteq v \\ |u| \leq n}} \neg(u \sqsubseteq x).$$

Since  $\varphi$  uses constants of length  $\leq n$ , only, we have  $\|\varphi\| \leq n$ .

(2) Now let  $\varphi(x) \in \text{FO}_{\text{ptNFA}}^2$  be a quantifier-free formula with  $\|\varphi(x)\| \leq n$ . First, suppose  $x \in L(\mathcal{A})$  is a subformula of  $\varphi(x)$ . Then, the depth of the pt-nfa  $\mathcal{A}$  is  $\leq n$ . Hence, by [17, Thm. 8],  $L(\mathcal{A}) \in \text{PT}(n)$ . By the first statement, any subformula  $x \in L(\mathcal{A})$  can be replaced by a quantifier- and automata-free formula  $\lambda(x) \in \text{FO}_{\emptyset}^2$  with  $\|\lambda(x)\| \leq n$ . Consequently, we can assume that  $\varphi(x)$  is automata-free, i.e., belongs to  $\text{FO}_{\emptyset}^2$ . Now replace subformulas of the form  $x \sqsubseteq v$  (with  $v$  a word) by

$$\bigvee_{u \sqsubseteq v} x = u,$$

such that the formula  $\varphi(x)$  becomes a Boolean combination of formulas  $u \sqsubseteq x$  and  $u = x$  with constants  $u$  of length  $\leq n$ . Note that  $\{u\}^\uparrow$  is of height  $\leq |u|$  and  $\{u\}$  is of height  $\leq |u| + 1$ . Hence,  $\varphi(x)$  defines a Boolean combination of languages from  $\text{PT}(n + 1)$ , i.e., a language from  $\text{PT}(n + 1)$ . □

**Remark** The above lemma shows that quantifier- and automata-free formulas of norm  $\leq n$  suffice to describe all piecewise testable languages of height  $\leq n$ , but any such formula is only guaranteed to define a piecewise testable language of height  $\leq n + 1$ . The bounds are tight as the following two examples demonstrate (with  $\Sigma = \{a\}$ ):

- (1) The language  $\{aaa\}a^*$  belongs to  $\text{PT}(3)$ , but cannot be defined by a formula of norm  $\leq 2$ .
- (2) The formula  $x = aaa$  of norm 3 defines the language  $\{aaa\}$  from  $\text{PT}(4) \setminus \text{PT}(3)$ .

Note that all heights appearing in Theorem 2.1 are bounded by  $(|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + m)$ . We now bound this function by a polynomial.

**Lemma 4.2** *Let  $c = 2 \cdot |\Sigma|$  and let  $m, n \in \mathbb{N}$ . Then,  $(|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + m) < (m + n + 2)^c$ .*

**Proof** If  $|\Sigma| = 1$ , we get

$$\begin{aligned} (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + m) &= 2(m + n) && \text{since } g_1(n) = n \\ &< 2(m + n + 2) \\ &\leq (m + n + 2)^2 \\ &= (m + n + 2)^c. \end{aligned}$$

If  $|\Sigma| \geq 2$ , we obtain

$$\begin{aligned} (|\Sigma| + 1) \cdot (g_{|\Sigma|}(n) + m) &\leq (|\Sigma| + 1) \cdot ((n + 2)^{|\Sigma|} + m) && \text{by [16 Thm.3.7 Eq. (3.12)]} \\ &< 2^{|\Sigma|} \cdot ((n + 2)^{|\Sigma|} + m) && \text{since } |\Sigma| \geq 2 \\ &\leq 2^{|\Sigma|} \cdot (m + n + 2)^{|\Sigma|} \\ &\leq (m + n + 2)^{2|\Sigma|} = (m + n + 2)^c. \end{aligned}$$

□

**Theorem 4.3** *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** We prove the claim by induction on the construction of the formula  $\varphi$ .

First suppose  $\varphi(x)$  is quantifier-free. Then, by Lemma 4.1(2), the language  $L_\varphi$  is piecewise testable of height  $\leq \|\varphi\| + 1 < 2^{c^{2^{|\varphi|}}}$  since  $c \geq 2$ . If the formula  $\varphi$  is a Boolean combination of formulas, the claim follows by induction since the doubly exponential function is monotone.

Now let  $\varphi(x) = \exists^{\geq t} y: \varphi'(x, y)$ . Our *first goal* is to express the formula  $\varphi(x)$  as a Boolean combination of formulas  $\alpha(x)$  with  $\|\alpha\| \leq \|\varphi'\|$  and  $\exists^{\geq s} y: (x\theta y \wedge \gamma(y))$  with  $s \leq t$ ,  $\theta \in \{\sqsubseteq, \supseteq, =, \parallel\}$ , and  $\|\gamma\| \leq \|\varphi'\|$ .

There exists a finite set  $A$  of formulas of the following form such that  $\varphi'(x, y)$  is a Boolean combination of formulas from  $A$ :

- formulas where at most  $x$  or  $y$ , but not both, are free
- atomic formulas  $x \sqsubseteq y$ ,  $x = y$ , and  $y \sqsubseteq x$

Note that all formulas  $\alpha$  from  $A$  satisfy  $\|\alpha\| \leq \|\varphi'\|$  since they are subformulas of  $\varphi'(x, y)$ .

For  $B \subseteq A$  set

$$\delta_B(x, y) = \bigwedge_{\beta \in B} \beta \wedge \bigwedge_{\alpha \in A \setminus B} \neg \alpha.$$

Then, there is a set  $\mathcal{B}$  of subsets of  $A$  such that  $\varphi'(x, y)$  is equivalent to

$$\bigvee_{B \in \mathcal{B}} \delta_B(x, y).$$

Since any pair of words can satisfy at most one formula  $\delta_B(x, y)$ , the formula  $\varphi(x) = \exists^{\geq t} y: \varphi'(x, y)$  is equivalent to

$$\varphi_1(x) = \bigvee \bigwedge_{B \in \mathcal{B}} \exists^{\geq t_B} y: \delta_B(x, y)$$

where the disjunction extends over all tuples  $(t_B)_{B \in \mathcal{B}}$  of natural numbers from  $\{0, 1, \dots, t\}$  that sum up to  $t$ .

So far, we expressed the formula  $\varphi(x)$  as a Boolean combination of formulas  $\exists^{\geq s} y: \delta(x, y)$  with  $s \leq t$  and  $\delta(x, y)$  a conjunction of possibly negated formulas from  $A$ . Note that any such formula is equivalent to the disjunction over all formulas

$$\begin{aligned} \exists^{\geq s_1} y: x \sqsubseteq y \wedge \delta(x, y) \\ \wedge \exists^{\geq s_2} y: x \supseteq y \wedge \delta(x, y) \\ \wedge \exists^{\geq s_3} y: x \parallel y \wedge \delta(x, y) \\ \wedge \exists^{\geq s_4} y: x = y \wedge \delta(x, y) \end{aligned}$$

where the disjunction extends over all tuples  $(s_1, s_2, s_3, s_4)$  of natural numbers from  $\{0, 1, \dots, s\}$  that sum up to  $s$ .

So far, we expressed the formula  $\varphi(x)$  as a Boolean combination of formulas  $\exists^{\geq s} y: (x\theta y \wedge \delta(x, y))$  with  $s \leq t$ ,  $\delta(x, y)$  a conjunction of possibly negated formulas from  $A$ , and  $\theta \in \{\sqsubseteq, \supseteq, =, \parallel\}$ .

We now consider one such formula. Since  $\delta(x, y)$  is a conjunction of possibly negated formulas from  $A$ , we can write it 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  $x \sqsubseteq y$ ,  $x \sqsupseteq y$ , 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

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

Thus, we reached our first goal: we expressed the formula  $\varphi(x)$  as a Boolean combination

- (1) of formulas  $\alpha(x)$  with  $\|\alpha\| \leq \|\varphi'\|$  and
- (2) of formulas  $\exists^{\geq s} y: (x\theta y \wedge \gamma(y))$  with  $s \leq t$ ,  $\theta \in \{\sqsubseteq, \sqsupseteq, =, \parallel\}$ , and  $\|\gamma\| \leq \|\varphi'\|$ .

Since the class  $PT(n)$  is closed under Boolean operations, it suffices to show that any such formula defines a piecewise testable language of height  $< 2^{c^{2\|\varphi'\|}}$ . By the induction hypothesis, this is clear for formulas from (1) since  $\|\varphi'\| \leq \|\varphi\|$ .

Our *second and final goal* is to show that it also holds for formulas from (2). So let  $s \leq t$ ,  $\theta \in \{\sqsubseteq, \sqsupseteq, =, \parallel\}$ , and  $\gamma(y)$  be a formula with  $\|\gamma\| \leq \|\varphi'\|$  and consider the formula  $\exists^{\geq s} y: (x\theta y \wedge \gamma(y))$ .

We consider the language

$$L = \{w \in \Sigma^* \mid S \models \gamma(w)\}$$

that, by the induction hypothesis, is piecewise testable of height  $< 2^{c^{2\|\varphi'\|}}$ . Now we have to consider the four possible values of  $\theta$  separately.

- 1. Let  $\theta = \sqsubseteq$ . Then, the formula  $\exists^{\geq s} y: (x \sqsubseteq y \wedge \gamma(y))$  is equivalent to

$$\begin{aligned} &\gamma(x) \wedge \exists^{\geq s+1} y: (x \sqsubseteq y \wedge \gamma(y)) \\ &\vee \neg\gamma(x) \wedge \exists^{\geq s} y: (x \sqsubseteq y \wedge \gamma(y)). \end{aligned}$$

Consequently, the set of words satisfying  $\exists^{\geq s} y: (x \sqsubseteq y \wedge \gamma(y))$  equals

$$(L \cap L\downarrow^{\geq s+1}) \cup (L\downarrow^{\geq s} \setminus L).$$

From Theorem 2.1(2), we obtain

$$L\downarrow^{\geq s+1}, L\downarrow^{\geq s} \in \text{PT}\left(\left(|\Sigma| + 1\right) \cdot \left(g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + 1\right)\right).$$

Note that

$$\begin{aligned} \left(|\Sigma| + 1\right) \cdot \left(g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + 1\right) &< \left(2^{c^{2\|\varphi'\|}} + 3\right)^c && \text{by Lemma 4.2} \\ &< \left(2^{c^{2\|\varphi'\|}} + 2^{c^{2\|\varphi'\|}}\right)^c && \text{since } c \geq 2, \|\varphi'\| \geq 1 \\ &= \left(2 \cdot 2^{c^{2\|\varphi'\|}}\right)^c \\ &< \left(2^{c^{\text{bin}(s)}} \cdot 2^{c^{2\|\varphi'\|}}\right)^c \\ &= \left(2^{c^{\text{bin}(s)} + c^{2\|\varphi'\|}}\right)^c \\ &\leq \left(2^{c^{\text{bin}(s)} + 2\|\varphi'\|}\right)^c && \text{since } c \geq 2, \text{bin}(s), \|\varphi'\| \geq 1 \\ &= 2^{c^{1 + \text{bin}(s)} + 2\|\varphi'\|} \\ &\leq 2^{c^{2\|\varphi\|}} \end{aligned}$$

where the last inequality holds since  $\|\varphi\| = \text{bin}(t) + \|\varphi'\|$ ,  $\text{bin}(s) \geq 1$ , and  $s \leq t$ .



It follows that  $L \downarrow^{\geq s+1}$  and  $L \downarrow^{\geq s}$  both are piecewise-testable of height  $< 2^{c^{2\|\varphi\|}}$ . Since this also holds for the language  $L$  and since  $\text{PT}(2^{c^{2\|\varphi\|}} - 1)$  is closed under Boolean combinations, this settles the case  $\theta = \sqsubseteq$ .

- Now let  $\theta = \sqsupseteq$ . Similarly to above, the set of words satisfying  $\exists^{\geq s} y : (x \sqsupseteq y \wedge \gamma(y))$  is a Boolean combination of the languages  $L$ ,  $L \uparrow^{\geq s+1}$ , and  $L \uparrow^{\geq s}$ . By Theorem 2.1(1), the latter two languages both belong to

$$\text{PT}(g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s).$$

Note that

$$g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) \geq 2^{c^{2\|\varphi'\|}} \geq 16 > 4$$

since  $c \geq 2$  and  $\|\varphi'\| \geq 1$ . It follows that

$$\begin{aligned} g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s &< g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s + g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) - 4 && \text{since } g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) > 4 \\ &\leq 2 \cdot (g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s - 2) \\ &\leq (|\Sigma| + 1) \cdot (g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s - 2) && \text{since } |\Sigma| \geq 1 \\ &< (2^{c^{2\|\varphi'\|}} + s)^c && \text{by Lemma 4.2} \\ &\leq (2^{c^{2\|\varphi'\|}} + 2^{c^{\text{bin}(s)}})^c && \text{since } c \geq 2 \\ &\leq (2^{c^{\text{bin}(s)+2\|\varphi'\|}})^c \\ &\leq 2^{c^{2\|\varphi\|}} \end{aligned}$$

where the last equality follows from  $\|\varphi\| = \text{bin}(t) + \|\varphi'\|$  and  $s \leq t$ . Thus, we showed that  $L \uparrow^{\geq s+1}$  and  $L \uparrow^{\geq s}$  both are of height  $< 2^{c^{2\|\varphi\|}}$ . Since this also holds for the language  $L$  and since  $\text{PT}(2^{c^{2\|\varphi\|}} - 1)$  is closed under Boolean operations, this settles the case  $\theta = \sqsupseteq$ .

- Next consider the case  $\theta = \parallel$ . By Theorem 2.1(3), the set of words satisfying  $\exists^{\geq s} y : (x \parallel y \wedge \gamma(y)) = L \parallel^{\geq s}$  belongs to  $\text{PT}(g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s)$ . The claim follows from  $g_{|\Sigma|}(2^{c^{2\|\varphi'\|}}) + s < 2^{c^{2\|\varphi\|}}$  as we verified in the previous case.

- It remains to consider the (trivial) case that  $\theta$  is equality. Then, the set of words satisfying  $\exists^{\geq s} y : (x = y \wedge \gamma(y))$  equals

- $\Sigma^* \in \text{PT}(0)$  if  $s = 0$ ,
- $L \in \text{PT}(2^{c^{2\|\varphi'\|}} - 1)$  if  $s = 1$ , and
- $\emptyset \in \text{PT}(0)$  if  $s > 1$ .

Consequently, it is always of height  $< 2^{c^{2\|\varphi\|}}$ .

Thus, we reached our second and final goal.

In summary, we proved that the set of words satisfying  $\varphi(x)$  is a Boolean combination of piecewise testable languages of height  $< 2^{c^{2\|\varphi\|}}$  and therefore belongs to this class as well.

This finishes the inductive proof of the theorem. □

Since piecewise testable languages of bounded height can be defined by quantifier-free formulas from  $\text{FO}_{\neq}^2$ , we obtain the following quantifier-elimination result (that, differently from the theorem above, applies also to formulas with two free variables).

**Corollary 4.4** *Let  $c = 2 \cdot |\Sigma|$ . Every  $C_{\text{ptNFA}}^2$ -formula  $\varphi$  is equivalent to some quantifier- and automata-free formula  $\psi \in \text{FO}_{\emptyset}^2$  with  $\|\psi\| < 2^{c^{2\|\varphi\|}}$ .*

For first-order formulas  $\varphi$ , this result can be found in [10, Cor. 7.4 & Thm. 7.5].

**Proof** Let  $\varphi(x, y)$  be some formula from  $C_{\text{ptNFA}}^2$ . As in the previous proof, it is a Boolean combination of formulas with a single free variable of norm  $\leq \|\varphi\|$  and of the formulas  $x \sqsubseteq y$ ,  $x = y$ , and  $x \sqsupseteq y$ . By Theorem 4.3, any formula  $\alpha(x)$  with a single free variable defines a piecewise testable language  $L$  of height  $< 2^{c^{2\|\varphi\|}}$ . By Lemma 4.1(1), this language can be defined by a quantifier- and automata-free formula  $\alpha'(x)$  from  $\text{FO}_{\emptyset}^2$  with  $\|\alpha'\| < 2^{c^{2\|\varphi\|}}$ . Replacing, in the Boolean combination  $\varphi(x, y)$ , all occurrences of  $\alpha(x)$  with  $\alpha'(x)$ , we obtain a quantifier- and automata-free formula  $\psi(x, y)$  that is equivalent to  $\varphi(x, y)$  and satisfies  $\|\psi\| < 2^{c^{2\|\varphi\|}}$ .  $\square$

Note that the above corollary implies in particular that the logics  $C_{\text{ptNFA}}^2$  and  $\text{FO}_{\emptyset}^2$  are equally expressive (a description of this expressive power in terms of subword-piecewise testable relations can be found in [10, Thm. 7.2(ii)]). It bounds the *norm* of the resulting formula  $\psi$  in terms of the norm of  $\varphi$  (which, in turn, is bounded by the size of  $\varphi$ ). Since  $\psi$  is automata- and quantifier-free, its norm equals the maximal length of a constant appearing in  $\psi$ , i.e., all words in  $\psi$  are of length at most doubly exponential in  $|\varphi|$ . Hence, the number of distinct atomic formulas in  $\psi$  is at most triply exponential. It follows that the size of  $\psi$  is at most fivefold exponential in the norm (and therefore the size) of  $\varphi$ .

This explosion of size is not surprising since, in  $\psi$ , we are not allowed to use quantification (let alone threshold counting quantification) nor piecewise testable predicates. The following result shows that disallowing threshold counting quantification alone already results in an exponential increase in formula size.

**Theorem 4.5** *For  $\Sigma = \{a\}$ , the logic  $C_{\emptyset}^2$  is exponentially more succinct than  $\text{FO}_{\emptyset}^2$ . More precisely, there is a sequence  $(\varphi_n(x))_{n \in \mathbb{N}}$  of formulas in  $C_{\emptyset}^2$  of size  $O(n)$  such that, for every sequence of equivalent formulas  $(\psi_n(x))_{n \in \mathbb{N}}$  from  $\text{FO}_{\emptyset}^2$ ,  $\psi_n$  is of size  $\Omega(2^n)$ .*

We do not know whether the same result holds for non-singleton alphabets.

**Proof** Let  $\Sigma = \{a\}$ . For  $n \in \mathbb{N}$  consider the formula

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

Note that  $|\varphi_n| = O(n)$  since the thresholds  $2^n$  and  $2^n + 1$  are encoded in binary. Furthermore,  $a^{2^n - 1}$  is the only word satisfying this formula.

Now let  $\psi_n(x)$  be a formula from  $\text{FO}_{\emptyset}^2$  that is equivalent to  $\varphi_n(x)$ . First note that  $a^m \sqsubseteq z$  (where  $z$  is any variable) is equivalent to the formula

$$\alpha_m(z) = \begin{cases} z = z & \text{if } m = 0 \\ \exists z' : (z' \sqsubseteq z \wedge z' \neq z \wedge \alpha_{m-1}(z')) & \text{otherwise} \end{cases}$$

(where  $z'$  is the other variable). Thus, replacing all subformulas  $a^m \sqsubseteq z$  and  $z \sqsubseteq a^m$  by  $\alpha_m(z)$  and  $\neg \alpha_{m+1}(z)$ , respectively, we eliminate all constants from  $\psi_n(x)$ . This replacement results in a linear increase in formula size, only (note that the size of the word  $a^m$  is  $m$ ). So, from now on, we can assume that  $\psi_n(x)$  is constant-free.

Now, let  $\psi(x)$  be any constant-free formula from  $\text{FO}_{\emptyset}^2$  of quantifier-rank  $d$  and let  $k, \ell \in \mathbb{N}$  with  $k, \ell \geq d$ . Then, by induction on  $d$ , one can show that  $S \models \psi(a^k) \iff S \models \psi(a^\ell)$ .

Since  $a^{2^n - 1}$  is the only word satisfying  $\psi_n(x)$ , the quantifier-rank of  $\psi_n(x)$  is  $\geq 2^n - 1$ . Hence, the size of  $\psi_n(x)$  is exponential in  $n$ .  $\square$

### 5 Complexity of the $C_{\text{ptNFA}}^2$ -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 4.4.<sup>4</sup> 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 [3].

**Lemma 5.1** *Let  $\varphi(x) = \exists^{\geq t} y: \psi(x, y)$  be a formula from  $C_{\text{ptNFA}}^2$ . Let  $c = 2 \cdot |\Sigma|$ ,  $N \in \mathbb{N}$  with  $2^{c^{2\|\varphi\|}} \leq N$ , and  $u \in \Sigma^*$  with  $|u| < N$ . Then,  $S \models \varphi(u)$  iff there are  $t$  words  $v$  of length  $< N^{2c}$  such that  $S \models \psi(u, v)$ .*

**Proof** We have to show that, whenever  $\varphi(u)$  holds, then there are  $t$  short words  $v$  such that  $\psi(u, v)$  holds (the other implication is trivial).

So assume there are at least  $t$  words in the language  $L := \{v \in \Sigma^* \mid S \models \psi(u, v)\}$ .

By Corollary 4.4, there exists a quantifier- and automata-free formula  $\psi'(x, y) \in \text{FO}_{\emptyset}^2$  equivalent to  $\psi(x, y)$  such that  $\|\psi'\| < 2^{c^{2\|\psi\|}} < 2^{c^{2\|\varphi\|}} \leq N$ . Since  $|u| < N$ , also the norm of the quantifier- and automata-free formula  $\psi'(u, y)$  is  $< N$ . Note that  $L$  is defined by this formula. Hence, by Lemma 4.1(2),  $L$  is piecewise testable of height  $\leq N$ . Since  $L$  contains at least  $t$  words, the definition of the function  $g_{|\Sigma|}$  together with the convexity of all equivalence classes implies that  $L$  contains mutually distinct words  $v_1, \dots, v_t$  of length  $< g_{|\Sigma|}(N) + t \leq (N + t + 2)^c$  (by Lemma 4.2). We have  $|\text{bin}(t)| \leq \|\varphi\|$  which implies  $t \leq N$ . Hence,  $(N + t + 2)^c \leq (2N + 2)^c$  which is smaller than  $N^{2c}$  since  $N \geq 16$ . Thus, we have  $|v_i| < N^{2c}$  for all  $1 \leq i \leq t$ . Consequently, we found  $t$  “short” witnesses for  $\psi(u, y)$ .  $\square$

**Proposition 5.2** *There is an alternating algorithm that, on input of a formula  $\varphi(x, y) \in C_{\text{ptNFA}}^2$  and words  $u$  and  $v$ , decides whether  $S \models \varphi(u, v)$ . This alternating algorithm runs in time doubly exponential in  $\|\varphi(u, v)\|$  and uses  $O(\|\varphi\|)$  alternations.*

**Proof** Before we come to the actual proof, we explain the idea underlying our approach. First, from  $\varphi, u, v$ , and  $N$ , we could compute a propositional formula (whose atomic propositional formulas are atomic formulas from  $C_{\text{ptNFA}}^2$ ) that is equivalent to  $\varphi(u, v)$ . This is possible since, by the previous lemma, we can restrict quantification in  $\varphi$  to words of bounded length. To serve as the basis of an alternating algorithm, we need in addition that the propositional formula is in negation normal form (i.e., at most atomic formulas are negated).

However, this approach has the following two problems.

First, the length bound from Lemma 5.1 is doubly exponential. Hence, the propositional formula for  $\exists^{\geq 1} y: \psi(u, y)$  is the disjunction over all formulas  $\psi(u, v)$  with  $v$  a word of doubly exponential length. But computing this formula requires triply exponential time. The solution to this first problem is that the propositional formula is not calculated explicitly. Instead, its evaluation is simulated by a procedure that takes, as arguments, a formula  $\alpha(x, y)$ , two words  $w_x$  and  $w_y$ , and a natural number  $N$  and returns the truth value of  $\alpha(w_x, w_y)$  if all quantifications are bounded by values that depend on  $N$ .

Since we do not compute the propositional formula, we cannot compute its negation normal form afterwards. Nor can we transform the  $C_{\text{ptNFA}}^2$ -formula into negation normal form since

<sup>4</sup> For first-order logic, the use of Corollary 4.4 can be replaced by the corresponding statements from [10, Cor. 7.4 & Thm. 7.5].

this logic does not allow universal quantifiers. This problem is solved by considering not just one procedure as above, but two (one for formulas  $\alpha$  occurring positively, one for negative occurrences).

Formally (and now the actual proof starts), we use the following recursive procedures  $\text{check}_P$  and  $\text{check}_N$  whose parameters are

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

- (1) If  $\alpha$  is an atomic formula, then decide whether  $\alpha(w_x, w_y)$  holds. This can be done by a nondeterministic algorithm in time linear in  $|w_x| + |w_y| + |\alpha|$ . If so, the procedure  $\text{check}_P$  returns true and false otherwise. The procedure  $\text{check}_N$  returns the negation of these values.
- (2) If  $\alpha = \beta \vee \gamma$ , then a call of  $\text{check}_P(\alpha, w_x, w_y, N)$  returns true iff at least one of  $\text{check}_P(\beta, w_x, w_y, N)$  and  $\text{check}_P(\gamma, w_x, w_y, N)$  returns true. Dually,  $\text{check}_N(\alpha, w_x, w_y, N)$  returns true iff both,  $\text{check}_N(\beta, w_x, w_y, N)$  and  $\text{check}_N(\gamma, w_x, w_y, N)$ , return true.
- (3) If  $\alpha = \neg\beta$ , then  $\text{check}_P(\alpha, w_x, w_y, N)$  returns true iff  $\text{check}_N(\beta, w_x, w_y, N)$  returns true and dually for  $\text{check}_N(\alpha, w_x, w_y, N)$ .
- (4) Let  $\alpha = \exists^{\geq t} y: \psi(x, y)$ .

Then,  $\text{check}_P(\alpha, w_x, w_y, N)$  returns true iff, for some set  $T$  of  $t$  words of length  $< N^{2c}$ , the call of  $\text{check}_P(\psi, w_x, w'_y, N^{2c})$  returns true for all words  $w'_y \in T$ . Thus, the evaluation of  $\text{check}_P$  consists of two phases: an existential phase (in which a set  $T$  is guessed, i.e. the computation branches into a sub-computation for each choice of  $T$ ), followed by a universal phase (in which for each sub-computation, i.e. for each choice of  $T$ , it is checked whether  $T$  is a set of  $t$  solutions).

Dually,  $\text{check}_N(\alpha, w_x, w_y, N)$  returns true iff, for some set  $T$  of  $t - 1$  words of length  $< N^{2c}$ , the call of  $\text{check}_N(\psi, w_x, w'_y, N^{2c})$  returns true for all words  $w'_y$  of length  $< N^{2c}$  that do *not* belong to  $T$ . As before, the call of  $\text{check}_N(\alpha, w_x, w_y, N)$  too consists of an existential phase followed by a universal phase (considering, instead of the guessed set  $T$ , its complement wrt. the set of words of length  $< N^{2c}$ ).

Now let  $\varphi(x, y)$  be a formula from  $C_{\text{ptNFA}}^2$ ,  $u, v \in \Sigma^*$ , and  $N_0 \in \mathbb{N}$  with  $|u|, |v| < N_0$  and  $2^{c^{2\|\varphi\|}} \leq N_0$ . By induction on the size of  $\varphi$  and using Lemma 5.1, one obtains that  $\mathcal{S} \models \varphi(u, v)$  iff  $\text{check}_P(\varphi, u, v, N_0)$  returns true iff  $\text{check}_N(\varphi, u, v, N_0)$  returns false.

Now let  $\psi = \varphi(u, v)$ . Then,  $\mathcal{S} \models \varphi(u, v)$  iff  $\text{check}_P(\psi, \varepsilon, \varepsilon, 2^{c^{2\|\psi\|}})$  returns true.

We now analyse the runtime of an execution of a call of  $\text{check}_P(\psi, \varepsilon, \varepsilon, 2^{c^{2\|\psi\|}})$ . First, the value of the parameter  $N$  is bounded by

$$(2^{c^{2\|\psi\|}})(2c)^d \leq 2^{c^{4\|\psi\|}} = 2^{c^{4\|\varphi(u,v)\|}}$$

where  $d \leq \|\psi\|$  is the quantifier depth of  $\psi$ . Consequently, the recursive execution considers only words of this doubly exponential length. Further, when handling a quantifier  $\exists^{\geq t}$ , it considers a set of at most  $t$  words of this doubly exponential length. Since  $t$  is at most exponential in the size of  $\psi$ , the alternating algorithm runs in at most doubly exponential time.

Further note that the execution alternates between universal and existential states only linearly often. □

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

**Theorem 5.3** *The  $C_{\text{ptNFA}}^2$ -theory of  $\mathcal{S}$  belongs to  $\text{STA}(*, 2^{2^{\text{poly}(n)}}, O(n))$ , i.e., can be decided in doubly exponential alternating time with linearly many alternations.*

## 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 using doubly exponential space, more precisely, it belongs to  $\text{STA}(*, 2^{2^{\text{poly}(n)}}, O(n))$ . This extends a result from [10] 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 (if only considering the space bound). Our proof relies on two independent aspects: the consideration of the height of definable languages (which is a direct continuation from [10]) and an adaptation of Ferrante and Rackoff’s method [3].

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  $\Sigma^*$  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 [9, Thm. 5.5] can be extended to show the following result: The  $\text{FO}_{\text{NFA}}^2$ -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  $C_{\text{NFA}}^2$ -theory becomes decidable using the arguments from [16] (where the result is demonstrated in case  $\mathcal{C}$  contains the cover relation, only).  
It is not clear for which sets  $\mathcal{C}$  the  $C_{\text{ptNFA}}^2$ -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  $\text{FO}_{\emptyset}^2$ -theory.
- Addition of regular predicates: By [16], the  $C_{\text{NFA}}^2$ -theory is decidable, but the only known algorithm is non-elementary. On the other hand, the  $C_{\text{ptNFA}}^2$ -theory is decidable using elementary space. It is not clear whether there are other classes of nfAs  $\mathcal{A} \subseteq \text{NFA}$  such that the  $C_{\mathcal{A}}^2$ - or  $\text{FO}_{\mathcal{A}}^2$ -theory are decidable in elementary space.

**Acknowledgements** We thank the reviewer for the thorough reading and commenting of the submitted version; these remarks led to several simplification of proofs and also improved the readability.

**Funding** Open Access funding enabled and organized by Projekt DEAL.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Baumann, P., Ganardi, M., Thinniyam, R. S., Zetsche, G.: Existential definability over the subword ordering. In: STACS'22, volume 219 of LIPIcs, pp. 7:1–7:15. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2022)
2. Berman, L.: The complexity of logical theories. *Theoret. Comput. Sci.* **11**, 71–77 (1980)
3. Ferrante, J., Rackoff, Ch.: *The Computational Complexity of Logical Theories Lecture Notes in Mathematics*, vol. 718. Springer, Berlin (1979)
4. Finkel, A., Schnoebelen, Ph.: Well-structured transition systems everywhere! *Theoret. Comput. Sci.* **256**, 63–92 (2001)
5. Higman, G.: Ordering by divisibility in abstract algebras. *Proc. Lond. Math. Soc.* **2**, 326–336 (1952)
6. Halfon, S., Schnoebelen, Ph., Zetsche, G.: Decidability, complexity, and expressiveness of first-order logic over the subword ordering. In: LICS'17, pp. 1–12. IEEE Computer Society (2017)
7. Ježek, J., McKenzie, R.: Definability in substructure orderings. I: finite semilattices. *Algebra Univers* **61**(1), 59–75 (2009)
8. Kudinov, O.V., Selivanov, V.L.: Undecidability in the homomorphic quasiorder of finite labelled forests. *J. Log. Comput.* **17**(6), 1135–1151 (2007)
9. Karandikar, P., Schnoebelen, Ph.: Decidability in the logic of subsequences and supersequences. In: FSTTCS'15, Leibniz International Proceedings in Informatics (LIPIcs) vol. 45, pp. 84–97. Leibniz-Zentrum für Informatik (2015)
10. Karandikar, P., Schnoebelen, Ph.: The height of piecewise-testable languages and the complexity of the logic of subwords. *Log. Methods Comput. Sci.* **15**(2) (2019)
11. Kuske, D., Schwarz, Ch.: Complexity of counting first-order logic for the subword order. In: MFCS'20, Leibniz International Proceedings in Informatics (LIPIcs) vol. 170, pp. 56:1–56:13. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2020)
12. Kudinov, O.V., Selivanov, V.L., Yartseva, L.V.: Definability in the subword order. In: CiE'10, Lecture Notes in Computer Science, vol. 6158, pp. 246–255. Springer (2010)
13. Kudinov, O.V., Selivanov, V.L., Zhukov, A.V.: Definability in the h-quasiorder of labeled forests. *Ann. Pure Appl. Logic* **159**(3), 318–332 (2009)
14. Kuske, D.: Theories of orders on the set of words. *Theoret. Inf. Appl.* **40**, 53–74 (2006)
15. Kuske, D.: The subtrace order and counting first-order logic. In: CSR'20, Lecture Notes in Computer Science, vol. 12159, pp. 289–302. Springer (2020)
16. Kuske, D., Zetsche, G.: Languages ordered by the subword order. In: FoSSaCS'19, Lecture Notes in Computer Science, vol. 11425, pp. 348–364. Springer (2019)
17. Masopust, T.: Piecewise testable languages and nondeterministic automata. In: MFCS'16, Leibniz International Proceedings in Informatics (LIPIcs) vol. 58, pp. 67:1–67:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2016)
18. Masopust, T., Thomazo, M.: On Boolean combinations forming piecewise testable languages. *Theoret. Comput. Sci.* **682**, 165–179 (2017)
19. Schnoebelen, Ph.: Personal communication, February (2020)
20. Simon, I.: Hierarchies of events with dot-depth one. Ph.D. thesis, University of Waterloo (1972)
21. Simon, I.: Piecewise testable events. In: Automata Theory and Formal Languages, Lecture Notes in Comp. Science vol. 33, pp. 214–222. Springer (1975)
22. Sakarovitch, J., Simon, I.: Subwords. In: *Combinatorics on Words*, chapter 6, pp. 105–144. Addison-Wesley, Boston (1983)
23. Thinniyam, R.S.: Definability of recursive predicates in the induced subgraph order. In: 7th Indian Conference on Logic and Its Applications (ICLA'17), Lecture Notes in Computer Science, vol. 10119, pp. 211–223. Springer (2017)
24. Thinniyam, R.S.: Defining recursive predicates in graph orders. *Log. Methods Comput. Sci.* **14**(3) (2018)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.