

Decidability in the Logic of Subsequences and Supersequences*

Prateek Karandikar¹ and Philippe Schnoebelen²

¹ LIAFA, University Paris Diderot, France

² Laboratoire Specification & Verification (LSV), Cachan, France

Abstract

We consider first-order logics of sequences ordered by the subsequence ordering, aka sequence embedding. We show that the Σ_2 theory is undecidable, answering a question left open by Kuske. Regarding fragments with a bounded number of variables, we show that the FO^2 theory is decidable while the FO^3 theory is undecidable.

1998 ACM Subject Classification F.4.1 Mathematical Logic, F.4.3 Formal Languages

Keywords and phrases subsequence, subword, logic, first-order logic, decidability, piecewise-testability, Simon's congruence

Digital Object Identifier 10.4230/LIPIcs.FSTTCS.2015.84

1 Introduction

A subsequence of a (finite) sequence $u = (x_1, \dots, x_\ell)$ is a sequence obtained from u by removing any number of elements. For example, if $u = (a, b, a, b, a)$ then $u' = (b, b, a)$ is a subsequence of u , a fact we denote with $u' \sqsubseteq u$. Other examples that work for any u are $u \sqsubseteq u$ (remove nothing) and $() \sqsubseteq u$.

In this paper we consider decidability and complexity questions for the first-order logic of finite sequences with the subsequence ordering as the only predicate. The notion of subsequence is certainly a fundamental one in logic, and it occurs prominently in several areas of computer science: in pattern matching (of texts, of DNA strings, etc.), in coding theory, in algorithmics, and in many other areas. We also note that sequences and their subsequences are a special case of a more general notion where a family of finite labelled structures (e.g., trees, or graphs, or ..) are compared via a notion of embedding. Closer to our own motivations, the automatic verification of unreliable channel systems and related problems generate many formulae where the subsequence ordering appears prominently [2, 4, 8, 11].

While decision methods for logics of sequences have been considered in several contexts, the corresponding logics usually do not include the subsequence predicate: they rather consider the prefix ordering, and/or membership in a regular language, and/or functions for taking contiguous subsequences or computing the length of sequences, see, e.g., [10, 7, 1].

As far as we know, Kuske's article [12] is the only one that specifically considers the decidability of the first-order logic of the subsequence ordering *per se*. The article also considers more complex orderings since these decidability questions first occurred in automated

* The first author was partially supported by Tata Consultancy Services and the CEFIPRA Raman-Charpak fellowship. This work was done when he was affiliated to Chennai Mathematical Institute, India and LSV, ENS Cachan, France.



deduction under the name of *ordered constraints solving* and they involve rather specific orderings on terms and strings [5].

Kuske considers the first-order logic of subsequences over a set of atoms A , denoted $\text{FO}(A^*, \sqsubseteq)$, and notes that the undecidability of its Σ_4 theory can be seen by reinterpreting an earlier undecidability result from [6] for the first-order logic of the lexicographic path ordering. He then shows that already the Σ_3 theory is undecidable even when A contains only two elements, and also shows that the Σ_1 theory is decidable so that the status of the Σ_2 theory remains open.

Our contribution. In this paper we show that the Σ_2 theory of the subsequence ordering is undecidable. On the positive side, we show that the FO^2 theory is decidable (but FO^3 is not). We also prove some complexity bounds for the decidable fragments: the Σ_1 theory is NP-complete and the FO^2 theory is PSPACE-hard.

Outline of the paper. The relevant definitions and basic results are given in section 2. Section 3 develops the reduction that proves undecidability for the Σ_2 and FO^3 theories. Section 4 presents a further reduction that proves undecidability for the Σ_2 theory even when constants are not allowed in the formulae. Then section 5 shows decidability for the two-variable fragment FO^2 .

Since our constructions heavily rely on concepts and results from formal language theory, we shall from now on speak of “words”, and “letters” (from an “alphabet”) rather than sequences and atoms. Note however that the logic $\text{FO}(A^*, \sqsubseteq)$ is defined for any kind of set A .

2 Basic notions

Let $A = \{a_1, a_2, \dots\}$ be a set called *alphabet*, whose elements are called *letters*. In this paper we only consider finite alphabets for ease of exposition but without any real loss of generality. A *word* is a finite sequence of letters like aac and we use u, v, \dots , to denote words, and A^* to denote the set of all words over A . Concatenation of word is written multiplicatively, and ϵ denotes the empty word. We also use regular expressions like $(ab + c)^*$ to denote regular languages (i.e., subsets of A^*). The length of a word u is denoted $|u|$ and, for $a \in A$, we let $|u|_a$ denote the number of occurrences of a in u .

We say that a word u is a *subword* (i.e., a subsequence) of v , written $u \sqsubseteq v$, when u is some $a_1 \cdots a_n$ and v can be written under the form $v_0 a_1 v_1 \cdots a_n v_n$ for some $v_0, v_1, \dots, v_n \in A^*$. We say a word u is a *factor* of a word v if there exist words v_1 and v_2 such that $v = v_1 u v_2$. For $B \subseteq A$, and $w \in A^*$, we define the projection of w onto B , denoted as $\pi_B(w)$, as the subword of w obtained by removing all letters in $A \setminus B$. For example, $\pi_{\{a,b\}}(abcacbbc) = ababb$.

We assume familiarity with basic notions of first-order logic as exposed in, e.g., [9]: bound and free occurrences of variables, etc.

In particular, for $n \in \mathbb{N}$, the fragment FO^n consists of all formulae that only use at most n distinct variables (these can have multiple occurrences inside the formula).

The fragments Σ_n and Π_n of $\text{FO}(A^*, \sqsubseteq)$ are defined inductively as follows:

- an atomic formula is in Σ_n and Π_n for all $n \in \mathbb{N}$;
- a negated formula $\neg\phi$ is in Σ_n iff ϕ is in Π_n , it is in Π_n iff ϕ is in Σ_n ;
- a conjunction $\phi \wedge \phi'$ is in Σ_n (resp., in Π_n) iff both ϕ and ϕ' are;
- For $n > 0$, an existentially quantified $\exists x\phi$ is in Σ_n iff ϕ is, it is in Π_n iff ϕ is in Σ_{n-1} ;
- For $n > 0$, a universally quantified $\forall x\phi$ is in Π_n iff ϕ is, it is in Σ_n iff ϕ is in Π_{n-1} .

Note that we do not require formulae to be in prenex normal form when defining the Σ_n and Π_n fragments: for example the formula $\forall x \exists y(x \sqsubseteq y \wedge \exists x \neg(x \sqsubseteq y))$ is simultaneously in Π_2 and FO^2 .

In this article we consider three versions of $\text{FO}(A^*, \sqsubseteq)$, the first-order logic of subsequences over A :

The pure logic: the signature consists of only one predicate symbol, “ \sqsubseteq ”, denoting the subword relation. One also uses a countable set $X = \{x, x', y, z, \dots\}$ of variables ranging over words in A^* and the usual logical symbols.

Note that there is no way in the pure logic to refer to specific elements of A in the logic. However, whether a formula ϕ is true, denoted $\models_{A^*} \phi$, may depend on A (in fact, its cardinality). For example, the closed formula

$$\forall x, y(x \sqsubseteq y \vee y \sqsubseteq x),$$

stating that \sqsubseteq is a total ordering, is true if, and only if, A contains at most one letter.

The basic logic: extends the pure logic by adding all words $u \in A^*$ as constant symbols (denoting themselves). For example, assuming A contains a, b and c , one can write the following sentence:

$$\exists x(ab \sqsubseteq x \wedge bc \sqsubseteq x \wedge abc \not\sqsubseteq x)$$

which is true, as witnessed by the valuation $x \mapsto bcab$.

The extended logic: further allows all regular expressions as unary predicates (with the expected semantics). For these predicates we adopt a more natural notation, writing e.g. $x \in \text{expr}$ rather than $P_{\text{expr}}(x)$. For example, the extended logic allows writing

$$\forall x([\exists y(y \in (ab)^* \wedge x \sqsubseteq y)] \Leftrightarrow x \in (a + b)^*)$$

which states that the regular language $(a + b)^*$ is the downward closure of $(ab)^*$, i.e., the set of all subwords of its words.

When writing formulae we freely use abbreviations like $x \sqsubset y$ for $x \sqsubseteq y \wedge \neg(y \sqsubseteq x)$ and $x \supseteq y$ for $y \sqsubseteq x$. Note that equality can be defined as an abbreviation since $x \sqsubseteq y \wedge y \sqsubseteq x$ is equivalent to $x = y$. Finally, we use negated symbols as in $x \not\sqsubseteq y$ or $x \notin (ab)^*$ with obvious meaning.

When we write $\text{FO}(A^*, \sqsubseteq)$ without any qualification we refer by default to the basic logic. The pure logic is apparently a very restricted logic, where one may hardly express more than generic properties of the subword ordering like saying that (A^*, \sqsubseteq) is a total ordering, or is a lattice. However, Theorem 3.1 below shows that the pure logic is quite expressive.

We conclude this expository section with

► **Theorem 2.1.** *The truth problem for the Σ_1 fragment of $\text{FO}(A^*, \sqsubseteq)$ is NP-complete even when restricting to a fixed alphabet.*

Proof sketch. The upper bound follows from the decidability proof in [12] since it is proved there that a satisfiable quantifier-free formula $\phi(x_1, \dots, x_n)$ can be satisfied with words of size in $O(n)$ assigned to the x_i 's. Guessing linear-sized witnesses u_1, \dots, u_n and checking that $\models_{A^*} \phi(u_1, \dots, u_n)$ can be done in NP.

For the lower bound, we reduce from boolean satisfiability. Consider a boolean formula $\phi(x_1, \dots, x_n)$ over n boolean variables. We reduce it to an $\text{FO}(A^*, \sqsubseteq)$ formula in the Σ_1 fragment

$$\psi \equiv \exists z, x_1, \dots, x_n(\phi')$$

where ϕ' is obtained from ϕ by replacing each occurrence of x_i with $x_i \sqsubseteq z$ (hence replacing $\neg x_i$ with $x_i \not\sqsubseteq z$). Then, for any alphabet A with at least one letter, ϕ is satisfiable if and only if $\models_{A^*} \psi$. ◀

3 Undecidability for Σ_2

We are interested in solving the *truth problem*. This asks, given an alphabet A and a sentence $\phi \in \text{FO}(A^*, \sqsubseteq)$, whether ϕ is true in the structure (A^*, \sqsubseteq) , written $\models_{A^*} \phi$. Restricted versions of the truth problems are obtained for example by fixing A (we then speak of the truth problem *over* A) and/or by restricting to a fragment of the logic.

This section is devoted to proving the following main result.

► **Theorem 3.1** (Undecidability). *The truth problem for $\text{FO}(A^*, \sqsubseteq)$ is undecidable even when restricted to formulae in the $\Sigma_2 \cap \text{FO}^3$ fragment of the basic logic.*

This is done by encoding Post's Correspondence Problem in $\text{FO}(A^*, \sqsubseteq)$. The reduction is described in several stages.

3.1 Expressing simple properties

We start with a list of increasingly complex properties and show how to express them in the basic $\text{FO}(A^*, \sqsubseteq)$ logic. We keep track of what fragment is used, with regards to both the number of distinct variables, and the quantifier alternation depth.

Note that when we claim that a property with m free variables can be expressed in FO^n (necessarily $n \geq m$), we mean that the formula only uses at most n variables *including the m free variables*.

We let $A = \{a_1, \dots, a_\ell\}$ denote an arbitrary alphabet, use B to denote subsets of A , and a, b, \dots to denote arbitrary letters from A .

P1. “ $x \in B^*$ ” can be expressed in $\Sigma_0 \cap \text{FO}^1$: using

$$\bigwedge_{a \in A \setminus B} a \not\sqsubseteq x.$$

P2. “ $\pi_B(y) \sqsubseteq x$ ” can be expressed in $\Pi_1 \cap \text{FO}^3$: building on P1, we use

$$\forall z ((z \sqsubseteq y \wedge z \in B^*) \implies z \sqsubseteq x),$$

noting that $\pi_B(y) \sqsubseteq x$ is equivalent to $\pi_B(y) \sqsubseteq \pi_B(x)$.

P3. “ $x = \pi_B(y)$ ” can be expressed in $\Pi_1 \cap \text{FO}^3$: building on P1, P2, and using

$$\pi_B(y) \sqsubseteq x \wedge x \sqsubseteq y \wedge x \in B^*.$$

P4. “ $\pi_B(x) = \pi_B(y)$ ” can be expressed in $\Pi_1 \cap \text{FO}^3$: building on P2, and using

$$\pi_B(y) \sqsubseteq x \wedge \pi_B(x) \sqsubseteq y.$$

P5. “ $x \in aA^*$ ”, i.e., “ x starts with a ”, can be expressed in $\Sigma_2 \cap \text{FO}^3$: building on P1, and using

$$\exists z (a \sqsubseteq z \wedge [\bigwedge_{b \in A \setminus \{a\}} ba \not\sqsubseteq z] \wedge z \sqsubseteq x \wedge \pi_{A \setminus \{a\}}(x) \sqsubseteq z).$$

Here the first two conjuncts require that z contains an occurrence of a and cannot start with another letter. The last two conjuncts require that z is a subword of x which has at least all the occurrences in x of all letters other than a .

Clearly, the mirror property “ $x \in A^*a$ ” can be expressed in $\Sigma_2 \cap \text{FO}^3$ too.

P6. “ $x \notin A^*aaA^*$ ” can be expressed in $\Sigma_2 \cap \text{FO}^3$: building on P3, and using

$$\exists y \left(y = \pi_{A \setminus \{a\}}(x) \wedge \forall z \left[(aa \sqsubseteq z \wedge y \sqsubseteq z \wedge z \sqsubseteq x) \implies \bigvee_{b \in A \setminus \{a\}} aba \sqsubseteq z \right] \right).$$

Note that this is equivalent to “ x does not have aa as a factor”. Here $z \sqsubseteq x$ implies that any two occurrences of a in z must come from x . Furthermore, if these are not contiguous in x they cannot be contiguous in z in view of $y = \pi_{A \setminus \{a\}}(x) \sqsubseteq z$.

► **Remark 3.2.** Note that the “ $y = \pi_{A \setminus \{a\}}(x)$ ” subformula in P6 uses one variable apart from y and x . We use the same variable name z that is used later in the formula, so that the formula is in FO^3 . We similarly reuse variable names whenever possible in later formulae.

P7. “ $x \notin A^*BBA^*$ ” can be expressed in $\Sigma_2 \cap \text{FO}^3$: as in P6 with

$$\exists y \left(y = \pi_{A \setminus B}(x) \wedge \forall z \bigwedge_{a, a' \in B} \left[(aa' \sqsubseteq z \wedge y \sqsubseteq z \wedge z \sqsubseteq x) \implies \bigvee_{b \in A \setminus B} aba' \sqsubseteq z \right] \right).$$

Note that this is equivalent to “ x has no factor in BB ”.

P8. “ $|\pi_B(x)| = 2$ ” can be expressed in $\Sigma_0 \cap \text{FO}^1$: using

$$\left(\bigvee_{a, a' \in B} aa' \sqsubseteq x \right) \wedge \bigwedge_{a, a', a'' \in B} aa'a'' \not\sqsubseteq x.$$

3.2 Expressing regular properties

Building on the previous formulae, our next step is to show how any regular property can be expressed in the basic logic by using an enlarged alphabet.

► **Lemma 3.3.** For any regular $L \subseteq A^*$ there is an extended alphabet $A' \supseteq A$ and a formula $\phi_L(x)$ in $\Sigma_2 \cap \text{FO}^3$ over A' such that for all $u \in A'^*$, $u \in L$ if and only if $\models_{A'^*} \phi_L(u)$.

Proof. Let $\mathcal{A} = (Q, A, \delta, I, F)$ be a NFA recognising L so that $u \in L$ iff \mathcal{A} has an accepting run on input u . We define $\phi_L(x)$ so that it states the existence of such a run, i.e., we put $\phi_L(x) \equiv \exists y \psi_{\mathcal{A}}(x, y)$ where $\psi_{\mathcal{A}}(x, y)$ expresses that “ y is an accepting run of \mathcal{A} over x ”.

Let $A' \stackrel{\text{def}}{=} A \cup Q$, assuming w.l.o.g. that A and Q are disjoint. A run $q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} q_n$ of \mathcal{A} can be seen as a word $q_0 a_1 q_1 a_2 \dots a_n q_n$ in A'^* . We now define $\psi_{\mathcal{A}}(x, y)$ as the conjunction $\psi_1(x, y) \wedge \psi_2(x, y)$, with

$$\begin{aligned} \psi_1 &\equiv (y \text{ has no factor from } AA) \wedge (y \text{ has no factor from } QQ) \\ &\wedge \left(\bigvee_{q \in I} y \text{ begins with } q \right) \wedge \left(\bigvee_{q \in F} y \text{ ends with } q \right) \wedge (\pi_A(y) = x), \\ \psi_2 &\equiv \forall z \left(\begin{aligned} &(x \sqsubseteq z \wedge z \sqsubseteq y \wedge z \text{ has exactly two occurrences of letters from } Q) \\ &\implies \left(\bigvee_{q, q' \in Q} \bigvee_{a, a' \in A} qaa'q' \sqsubseteq z \vee \bigvee_{(q, a, q') \in \delta} qaq' \sqsubseteq z \right) \end{aligned} \right). \end{aligned}$$

Here ψ_1 reuses simple properties from the previous subsection and states that y is a word alternating between Q (states of \mathcal{A}) and A (proper letters), starting with an initial state of \mathcal{A}

and ending with an accepting state, hence has the required form $q_0 a_1 \dots a_n q_n$. Furthermore, $\pi_A(y) = x$ ensures that y has the form of an accepting run over x . Note that it also ensures $x \in A^*$.

With ψ_2 , one further ensures that the above y respects the transition table of \mathcal{A} , i.e., that $(q_{i-1}, a_i, q_i) \in \delta$ for $i = 1, \dots, n$. Indeed, assume $z \in A'^*$ satisfies $x \sqsubseteq z \sqsubseteq y$ and contains two occurrences from Q . Thus z is $a_1 \dots a_i q_i a_{i+1} a_{i+2} \dots a_j q_j a_{j+1} a_{j+2} \dots a_n$ for some $1 \leq i < j \leq n$. If now $j > i + 1$ then z contains $q_i a_{i+1} a_{i+2} q_j$ as a subword and the disjunction after the implication is fulfilled. However, if $j = i + 1$, the only way to fulfil the disjunction is to have $(q_{j-1}, a_j, q_j) \in \delta$.

Finally, $\psi_A(x, y)$ exactly states that y is an accepting run for x and $\models_{A'^*} \phi_L(u)$ holds iff $u \in L$. One easily checks that ψ_1 is in $\Sigma_2 \cap \text{FO}^3$, ψ_2 is in $\Pi_1 \cap \text{FO}^3$, so that ψ_A and ϕ_L are in $\Sigma_2 \cap \text{FO}^3$. We reuse variables wherever possible to ensure that only three variables are used (see remark 3.2). For example, the implementation of “ y has no factor from QQ ” from P7 needs two other variables, and here we use x and z for it. ◀

3.3 Encoding Post’s Correspondence Problem

It is now easy to reduce Post’s Correspondence Problem to the truth problem for the basic $\text{FO}(A^*, \sqsubseteq)$ logic.

Suppose we have a PCP instance \mathcal{P} consisting of pairs $(u_1, v_1), \dots, (u_n, v_n)$ over the alphabet Γ . We let $N = \{1, \dots, n\}$, consider the alphabet $A \stackrel{\text{def}}{=} \Gamma \cup N$, and define

$$\phi_{\mathcal{P}} \equiv \exists x, x' \left(\begin{array}{l} x \in (1u_1 + \dots + nu_n)^+ \wedge x' \in (1v_1 + \dots + nv_n)^+ \\ \wedge \pi_N(x) = \pi_N(x') \wedge \pi_{\Gamma}(x) = \pi_{\Gamma}(x') \end{array} \right). \quad (1)$$

Clearly, $\phi_{\mathcal{P}}$ is true iff the PCP instance has a solution.

It remains to check that $\phi_{\mathcal{P}}$ is indeed a formula in the Σ_2 fragment: this relies on Lemma 3.3 for expressing membership in two regular languages, and the P4 properties for ensuring that x and x' contain the same indexes from N and the same letters from Γ . Finally, we note that $\phi_{\mathcal{P}}$ is also a FO^3 formula.

4 Undecidability for the pure logic

In this section we give a stronger version of the undecidability for the Σ_2 fragment.

► **Theorem 4.1** (Undecidability for the pure logic). *The truth problem for $\text{FO}(A^*, \sqsubseteq)$ is undecidable even when restricted to formulae in the Σ_2 fragment of the pure logic.*

The proof is by constructing a Σ_2 formula $\psi(x_1, \dots)$ in the pure logic that defines all the letters and constant words we need to reuse the reduction from the previous section.

Kuske solves the problem in the special case of a formula using only $\{\epsilon, a, b, ab, ba, aa, bb, aba, bab\}$ as constants [12]. We provide a more generic construction whereby all words (up to a fixed length) can be defined in a single Σ_2 formula. One inherent difficulty is that it is impossible to properly define constant words in the pure logic. Of course, with the pure logic one can only define properties up to a bijective renaming of the letters, so $\psi(x_1, \dots)$ will only define letters and words up to renaming. But a more serious problem is that we can only define properties *invariant by mirroring* as we now explain.

For a word $u = a_1 a_2 \dots a_\ell$, we let \tilde{u} denote its mirror image $a_\ell \dots a_2 a_1$.

► **Lemma 4.2** (Invariance by mirroring). *If $\psi(x_1, \dots, x_n)$ is a formula in the pure logic and u_1, \dots, u_n are words in A^* , then $\models_{A^*} \phi(u_1, \dots, u_n)$ if, and only if, $\models_{A^*} \phi(\tilde{u}_1, \dots, \tilde{u}_n)$.*

Proof Sketch. By structural induction on ϕ , noting that the only atomic formulae in the pure logic have the form $x \sqsubseteq y$, and that $u \sqsubseteq v$ iff $\tilde{u} \sqsubseteq \tilde{v}$ for any $u, v \in A^*$. \blacktriangleleft

4.1 Defining letters and short constant words

We now define $\psi(x_1, \dots)$. In our construction ψ has the form $\psi_1 \wedge \psi_2 \wedge \dots \wedge \psi_{13}$ and features a large number of free variables. We describe the construction in several stages, explaining what valuation of its free variables can make ψ true. We start with

$$\forall y(z \sqsubseteq y) \tag{\psi_1}$$

$$\wedge \bigwedge_{1 \leq i \neq j \leq n} x_i \not\sqsubseteq x_j \tag{\psi_2}$$

$$\wedge \bigwedge_{i=1}^n \forall y[y \sqsubseteq x_i \implies (x_i \sqsubseteq y \vee y \sqsubseteq z)] \tag{\psi_3}$$

Here ψ_1 implies $z = \epsilon$, then ψ_2 implies $x_i \neq \epsilon$ so that ψ_3 requires that each x_i is a single letter and furthermore x_1, \dots, x_n must be different letters as required by ψ_2 .

We continue with:

$$\wedge \bigwedge_{i=1}^n (x_i \sqsubseteq x_i^2 \wedge x_i^2 \not\sqsubseteq x_i \wedge \forall y[y \sqsubseteq x_i^2 \implies (y \sqsubseteq x_i \vee x_i^2 \sqsubseteq y)]) \tag{\psi_4}$$

Note that n new free variables, x_1^2, \dots, x_n^2 are involved. First ψ_4 requires that any x_i^2 has at least two letters (it must contain x_i strictly). But it also requires that any subword of x_i^2 is ϵ or x_i or x_i^2 , thus x_i^2 has length 2 and can only be $x_i x_i$.

In the same style we introduce new free variables x_1^3, \dots, x_n^3 and x_1^4, \dots, x_n^4 and require that x_i^3 equals $x_i x_i x_i$, and that x_i^4 equals $x_i x_i x_i x_i$ with:

$$\wedge \bigwedge_{i=1}^n (x_i^2 \sqsubseteq x_i^3 \wedge x_i^3 \not\sqsubseteq x_i^2 \wedge \forall y[y \sqsubseteq x_i^3 \implies (y \sqsubseteq x_i^2 \vee x_i^3 \sqsubseteq y)]) \tag{\psi_5}$$

$$\wedge \bigwedge_{i=1}^n (x_i^3 \sqsubseteq x_i^4 \wedge x_i^4 \not\sqsubseteq x_i^3 \wedge \forall y[y \sqsubseteq x_i^4 \implies (y \sqsubseteq x_i^3 \vee x_i^4 \sqsubseteq y)]) \tag{\psi_6}$$

We introduce new free variables $\{y_{i,j}\}_{1 \leq i \neq j \leq n}$ and conjuncts:

$$\wedge \bigwedge_{1 \leq i \neq j \leq n} \forall y (y \sqsubseteq y_{i,j} \implies y \sqsubseteq z \vee x_i \sqsubseteq y \vee x_j \sqsubseteq y) \tag{\psi_8}$$

$$\wedge \bigwedge_{1 \leq i \neq j \leq n} (x_i \sqsubseteq y_{i,j} \wedge x_j \sqsubseteq y_{i,j} \wedge x_i^2 \not\sqsubseteq y_{i,j} \wedge x_j^2 \not\sqsubseteq y_{i,j}) \tag{\psi_9}$$

$$\wedge \bigwedge_{1 \leq i \neq j \leq n} (y_{i,j} \not\sqsubseteq y_{j,i}) \tag{\psi_{10}}$$

Here ψ_8 requires that any $y_{i,j}$ only contains letters among x_i and x_j , and ψ_9 requires that it contains exactly one occurrence of x_i and one of x_j . So that $y_{i,j}$ is either $x_i x_j$ or $x_j x_i$. With ψ_{10} we require that $y_{j,i}$ is, among $x_i x_j$ and $x_j x_i$, the word not assigned to $y_{i,j}$.

Now, in view of Lemma 4.2, it is impossible to fix e.g. $y_{i,j} = x_i x_j$. However we can force all $y_{i,j}$ to have “the same orientation”. Let i, j, k be three different indexes in $\{1, \dots, n\}$ and consider the following formula

$$\xi_{i,j,k} \equiv \exists t \left[\begin{array}{l} \forall y (y \sqsubseteq t \implies y \sqsubseteq z \vee x_i \sqsubseteq y \vee x_j \sqsubseteq y \vee x_k \sqsubseteq y) \tag{\xi_1} \\ \wedge x_i^2 \sqsubseteq t \wedge x_i^3 \not\sqsubseteq t \wedge x_j \sqsubseteq t \wedge x_j^2 \not\sqsubseteq t \wedge x_k \sqsubseteq t \wedge x_k^2 \not\sqsubseteq t \tag{\xi_2} \\ \wedge y_{i,j} \sqsubseteq t \wedge y_{j,i} \sqsubseteq t \wedge y_{i,k} \sqsubseteq t \wedge y_{k,i} \not\sqsubseteq t \wedge y_{j,k} \sqsubseteq t \wedge y_{k,j} \not\sqsubseteq t \tag{\xi_3} \end{array} \right]$$

We claim that, in conjunction with the earlier ψ -conjuncts, $\xi_1 \wedge \xi_2 \wedge \xi_3$ requires $t = x_i x_j x_i x_k$ or $t = x_k x_i x_j x_i$: indeed by ξ_1 , t only contains letters among $\{x_i, x_j, x_k\}$, then by ξ_2 , t contains exactly 2 occurrences of x_i and exactly one occurrence each of x_j and x_k , then by ξ_3 , t has $x_i x_j$ and $x_j x_i$ as subwords, so the single occurrence of x_j is between the two occurrences of x_i and, by ξ_3 again, the occurrence of x_k is outside the two x_i occurrences. Finally, satisfying $\xi_{i,j,k}$ requires $y_{i,k}$ and $y_{j,k}$ to have the same orientation.

We continue the construction of ψ with:

$$\wedge \bigwedge_{1 \leq i \neq j \leq n} \bigwedge_{k \notin \{i,j\}} \xi_{i,j,k} \quad (\psi_{11})$$

As just explained, this will force all $y_{i,j}$'s to have the same orientation, i.e., any satisfying assignment will have $y_{i,j} = x_i x_j$ for all i, j , or $y_{i,j} = x_j x_i$ for all i, j .

4.2 Defining long constant words

Once we have defined all words of length 2 (up to mirroring) over the alphabet $\{x_1, \dots, x_n\}$ (up to renaming), it is easier to systematically define all words of length 3, 4, etc. Actually, we only use constant words of length at most 4 for the formula $\phi_{\mathcal{P}}$ from section 3.

The general strategy relies on a technical lemma we now explain. For $n \in \mathbb{N}$ we say that two words u and v are n -equivalent, written $u \sim_n v$, if u and v have the same set of subwords of length up to n . Thus \sim_n is the piecewise-testability congruence introduced by Simon, see [16, 15].

► **Lemma 4.3.** *Let $n \geq 2$, and let u and v be words of length $n + 1$ with $u \neq v$. Then $u \not\sim_n v$.*

Proof. See appendix. ◀

We can thus introduce new variables $y_{i,j,k}$ and $y_{i,j,k,m}$ for all $i, j, k, m \in \{1, \dots, n\}$ (allowing repetitions of indexes) and require $y_{i,j,k} = x_i x_j x_k$ and $y_{i,j,k,m} = x_i x_j x_k x_m$, up to mirroring but with the same orientation for all the y_{i_1, \dots, i_ℓ} 's. Then we complete the construction of ψ with the following conjuncts:

$$\wedge \bigwedge_{1 \leq i, j, k \leq n} \text{“formula defining } y_{i,j,k}\text{”} \quad (\psi_{12})$$

$$\wedge \bigwedge_{1 \leq i, j, k, m \leq n} \text{“formula defining } y_{i,j,k,m}\text{”}. \quad (\psi_{13})$$

In order to require that, for example, $y_{1,5,2} = x_1 x_5 x_2$, it is enough to:

- enumerate all words of length upto 2, and for each say whether it is or is not a subword of $y_{1,5,2}$ ($y_{1,5,2} \sqsupseteq y_{1,5,2} \wedge x_1^2 \not\sqsupseteq y_{1,5,2} \wedge \dots$),
- and require that $y_{1,5,2}$ has length 3, by saying that every subword of $y_{1,5,2}$ is itself or is one of the words of length upto 2, and that $y_{1,5,2}$ is distinct from all these words.

The correctness of the construction is guaranteed by Lemma 4.3.

Once all 3-letter words have been defined, we can use them to define 4-letter words (and if needed, 5-letter words, and so on) similarly, with correctness following from Lemma 4.3.

Finally, we let $\phi'_{\mathcal{P}}$ be obtained from the formula $\phi_{\mathcal{P}}$ —see Eq. (1) page 89— by replacing every constant letter $a_i \in A$ by the variable x_i , and every constant word $a_{i_1} \dots a_{i_\ell} \in A^*$ by the variable y_{i_1, \dots, i_ℓ} (we use z for the constant word ϵ , and x_i^2 for the constant word $x_i x_i$).

Now we define $\psi_{\mathcal{P}}$ with

$$\psi_{\mathcal{P}} \equiv \exists Z (\psi_1 \wedge \dots \wedge \psi_{13} \wedge \phi'_{\mathcal{P}})$$

where $Z = \{z, x_1, \dots, x_n, x'_1, \dots, x'_n, x_1^2, x_1^3, x_1^4, \dots, y_{1,1}, \dots, y_{i_1, \dots, i_\ell}, \dots\}$ collects all the free variables we used in $\psi_1 \wedge \dots \wedge \psi_{13}$.

Noting that each ψ_i as well as $\phi'_{\mathcal{P}}$ is a Σ_2 formula, we get that the resulting $\psi_{\mathcal{P}}$ is a Σ_2 formula in the pure logic that is true in (A^*, \sqsupseteq) iff the PCP instance \mathcal{P} is positive. This concludes the proof of Theorem 4.1.

4.3 Undecidability for a fixed alphabet

The above Theorem 4.1 applies to the truth problem *for unbounded alphabet*, i.e., where we ask whether $\models_{A^*} \phi$ for given A and ϕ . In this proof, the alphabet A depends on the PCP instance \mathcal{P} since it includes symbols for the states of the regular automata that define the languages $(1u_1 + \dots + nu_n)^+$ and $(1v_1 + \dots + nv_n)^+$ in Eq. (1), and further includes symbols in $N = \{1, \dots, n\}$.

It is possible to further show undecidability of the Σ_2 fragment even *for a fixed alphabet* A as we now explain. For this we consider a variant of Post's Correspondence Problem:

► **Definition 4.4.** The *variant PCP problem* asks, given an alphabet Γ , pairs $(u_1, v_1), \dots, (u_n, v_n)$ over Γ , and an extra word $w \in \Gamma^*$, whether there exists a sequence i_1, \dots, i_ℓ over $\{1, \dots, n\}$ such that $w u_{i_1} \dots u_{i_\ell} = v_{i_1} \dots v_{i_\ell}$.

► **Lemma 4.5.** *There is a fixed Γ and a fixed sequence of pairs over Γ for which the variant PCP problem (with only w as input) is undecidable.*

Proof Sketch. One adapts the standard undecidability proof for PCP. Instead of reducing from the question whether a given TM halts, one reduces from the question whether a fixed TM accepts a given input. Note that in the case of a universal TM, the problem is undecidable. Fixing the TM will lead to a fixed sequence of pairs $(u_1, v_1), \dots, (u_n, v_n)$, and the input of the TM will provide the w parameter of the problem. ◀

► **Theorem 4.6 (Undecidability for fixed alphabet).** *There exists a fixed alphabet A such that the truth problem for the pure logic $\text{FO}(A^*, \sqsubseteq)$ is undecidable even when restricted to formulae in Σ_2 .*

Proof Sketch. We adapt the proof of Theorems 3.1 and 4.1 by reducing from the variant PCP problem with fixed Γ and sequence of pairs. The encoding formula can be

$$\equiv \exists x, x' \left(\begin{array}{l} x \in \Gamma'^* \cdot (1u_1 + \dots + nu_n)^+ \wedge x' \in \rho(1v_1 + \dots + nv_n)^+ \\ \wedge \pi_{\Gamma'}(x) = \hat{w} \wedge \pi_N(x) = \pi_N(x') \wedge \pi_{\Gamma \cup \Gamma'}(x) = \pi_{\Gamma \cup \Gamma'}(x') \end{array} \right) \quad (2)$$

to be compared with Eq. (1). Here we use $\Gamma' = \{\hat{a}, \hat{b}, \dots\}$, a renamed copy of $\Gamma = \{a, b, \dots\}$, to be able to extract the w prefix in x . The word \hat{w} is simply w from the variant PCP instance with all letters from Γ replaced by corresponding letters from Γ' . We then need to extend the language $(1v_1 + \dots + nv_n)$ for x' so that letters from Γ' can be used in place of the corresponding letters from Γ . This is done by applying a simple transduction $\rho \stackrel{\text{def}}{=} \left(\bigcup_{a \in \Gamma} \begin{bmatrix} a \\ a \end{bmatrix} \cup \begin{bmatrix} \hat{a} \\ a \end{bmatrix} \right)^*$.

In the end, we only use two fixed regular languages, and thus a fixed alphabet A . Note however that encoding the input w will require using constant words of unbounded lengths. Here we rely on the fact that our reduction from basic to pure logic can define constant words of arbitrary length in the Σ_2 fragment. ◀

5 Decidability for the FO^2 fragment

In this section we show that for finite alphabets, the truth problem for the 2-variable fragment $\text{FO}^2(A^*, \sqsubseteq)$ is decidable. The proof was first sketched by Kuske [13].

5.1 Rational relations

We recall the basics of rational relations. See [3, Chap. 3] or [14, Chap. 4] for more details.

For finite alphabets A and B , the *rational relations* between A^* and B^* are defined as the subsets of $A^* \times B^*$ recognised by asynchronous transducers. The set of rational relations between A^* and B^* is exactly the closure of the finite subsets of $A^* \times B^*$ under union, concatenation, and Kleene star.

For example, it is easy to see that the subword relation, seen as a subset of $A^* \times A^*$ is a rational relation [3, Example III.5.9], and that the strict subword relation is rational too:¹

$$\sqsubseteq = \left(\bigcup_{a \in A} \begin{bmatrix} a \\ \epsilon \end{bmatrix} \cup \begin{bmatrix} a \\ a \end{bmatrix} \right)^* , \quad \sqsubset = \sqsubseteq \cdot \left(\bigcup_{a \in A} \begin{bmatrix} a \\ \epsilon \end{bmatrix} \right) \cdot \sqsubseteq .$$

Define now the *incomparability relation* over A^* , denoted \perp , by $u \perp v$ iff $u \not\sqsubseteq v \wedge v \not\sqsubseteq u$.

► **Lemma 5.1.** *The incomparability relation over A^* is a rational relation.*

Proof. We cannot simply use the fact that \sqsubseteq and \sqsupseteq are rational relations since rational relations are not closed under intersection. The way out is to express incomparability as a union $\perp = T_1 \cup T_2$ of rational relations, using the following equivalence

$$u \perp v \text{ iff } \overbrace{(u \not\sqsubseteq v \wedge |u| \leq |v|)}^{(u,v) \in T_1} \vee \overbrace{(v \not\sqsubseteq u \wedge |v| \leq |u|)}^{(u,v) \in T_2} . \quad (3)$$

The equivalence holds since $|u| > |v|$ implies $u \not\sqsubseteq v$.

We show (see Coro. 5.3) that T_1 is rational. A symmetric reasoning shows that T_2 is rational. This concludes since the union of two rational relations is rational. ◀

In the following proof, we write $w(0 : -i]$ to denote the prefix of length $|w| - i$ of an arbitrary word w (assuming $0 \leq i \leq |w|$).

► **Lemma 5.2.** *$(u, v) \in T_1$ iff there exists an integer ℓ , a factorisation $u = a_1 a_2 \dots a_\ell a u'$ of u , and a factorisation $v = v_1 a_1 v_2 a_2 \dots v_\ell a_\ell v v'$ of v such that*

- $a_1, \dots, a_\ell \in A$ and $v_1, \dots, v_\ell \in A^*$ are such that a_i does not occur in v_i for all $i = 1, \dots, \ell$,
- $a, b \in A$ are two letters with $a \neq b$, and
- $u', v' \in A^*$ are two suffixes with $|u'| = |v'|$.

Proof. The (\Leftarrow) direction is clear: the listed conditions guarantee $|u| \leq |v|$ and $u \not\sqsubseteq v$.

To see the (\Rightarrow) direction, we assume $(u, v) \in T_1$ and write $u = a_1 \dots a_n$, with $n = |u|$, knowing that $n > 0$ since $u \not\sqsubseteq v$. We say that $i \in \{0, \dots, n\}$ is *good* if $u(0 : -i] \sqsubseteq v(0 : -i]$, and *bad* otherwise. Clearly, n is good and 0 is bad. Let $m > 0$ be the smallest good index: it is easy to check that taking $\ell = n - m$, $a = a_{\ell+1}$ and $u' = a_{\ell+2} \dots a_n$ proves the claim. ◀

► **Corollary 5.3.** *T_1 is a rational relation.*

Proof. Lemma 5.2 directly translates as

$$T_1 = \left(\bigcup_{a \in A} \left[\bigcup_{b \neq a} \begin{bmatrix} b \\ \epsilon \end{bmatrix} \right]^* \cdot \begin{bmatrix} a \\ a \end{bmatrix} \right)^* \cdot \left(\bigcup_a \bigcup_{b \neq a} \begin{bmatrix} b \\ a \end{bmatrix} \right) \cdot \left(\bigcup_{a, a'} \begin{bmatrix} a' \\ a \end{bmatrix} \right)^* .$$

¹ When writing such regular expressions we use the vector notation $\begin{bmatrix} y \\ x \end{bmatrix}$ to denote (x, y) . Note that the domain and the range of the relation correspond to the bottom and, resp., the top, lines of the vectors. We use \cdot to mean concatenation.

5.2 Decidability for FO²

Let $\mathcal{R} \stackrel{\text{def}}{=} \{=, \sqsubset, \sqsupset, \perp\}$ consists of the following four relations on A^* : equality, strict subword relation, its inverse, and incomparability. These four relations form a partition of $A^* \times A^*$, i.e., for all $u, v \in A^*$, exactly one of $u = v$, $u \sqsubset v$, $u \sqsupset v$, and $u \perp v$ holds.

For any $R \in \mathcal{R}$ and language $L \subseteq A^*$, we define the *preimage* of L by R , denoted $R^{-1}(L)$, as being the language $\{x \in A^* : \exists y \in L : (x, y) \in R\}$. We saw in section 5.1 that each relation $R \in \mathcal{R}$ is rational: we deduce that $R^{-1}(L)$ is regular whenever L is. Furthermore, using standard automata-theoretic techniques, a description of the preimage $R^{-1}(L)$ can be computed effectively from a description of L .

In the following we consider FO² formulae using only x and y as variables. We allow formulae to have regular predicates of the form $x \in L$ for fixed regular languages L (i.e., we consider the extended logic). Furthermore, we consider a variant of the logic where we use the binary relations \sqsubset , $=$ and \perp instead of \sqsubseteq . This will be convenient later. The two variants are equivalent, even when restricting to FO^m or Σ_m fragments: in one direction we observe that $x \sqsubseteq y$ can be defined with $x \sqsubset y \vee x = y$, in the other direction one defines $x \sqsubset y$ with $x \sqsubseteq y \wedge y \not\sqsubseteq x$ and $x \perp y$ with $x \not\sqsubseteq y \wedge y \not\sqsubseteq x$. We also use $x \sqsupset y$ as shorthand for $y \sqsubset x$.

► **Lemma 5.4.** *Let $\phi(x)$ be an FO² formula with at most one free variable. Then there exists a regular language $L_\phi \subseteq A^*$ such that $\phi(x)$ is equivalent to $x \in L_\phi$. Furthermore, a description for L_ϕ can be computed effectively from ϕ .*

Proof. By structural induction on $\phi(x)$. If $\phi(x)$ is an atomic formula of the form $x \in L$, the result is immediate. If $\phi(x)$ is an atomic formula that uses a binary predicate R from \mathcal{R} , the fact that it has only one free variable means that $\phi(x)$ is a trivial $x = x$, or $x \sqsubset x$, or \dots , so that L_ϕ is A^* or \emptyset .

For compound formulae of the form $\neg\phi'(x)$ or $\phi_1(x) \vee \phi_2(x)$, we use the induction hypothesis and the fact that regular languages are closed under boolean operations.

There remains the case where $\phi(x)$ has the form $\exists y \phi'(x, y)$. We first replace any subformulae of ϕ' having the form $\exists x \psi(x, y)$ or $\exists y \psi(x, y)$ with equivalent formulae of the form $y \in L_\psi$ or $x \in L_\psi$ respectively, for appropriate languages L_ψ , using the induction hypothesis. Thus we may assume that ϕ' is quantifier-free. We now rewrite ϕ' by pushing all negations inside with the following meaning-preserving transformations:

$$\neg\neg\psi \rightarrow \psi \qquad \neg(\psi_1 \vee \psi_2) \rightarrow \neg\psi_1 \wedge \neg\psi_2 \qquad \neg(\psi_1 \wedge \psi_2) \rightarrow \neg\psi_1 \vee \neg\psi_2$$

and then eliminating negations completely with:

$$\neg(z \in L) \rightarrow z \in (A^* \setminus L) \qquad \neg(z_1 R_1 z_2) \rightarrow z_1 R_2 z_2 \vee z_1 R_3 z_2 \vee z_1 R_4 z_2$$

where R_1, R_2, R_3, R_4 are relations such that $\mathcal{R} = \{R_1, R_2, R_3, R_4\}$. Thus, we may now assume that ϕ' is a positive boolean combination of atomic formulae. We write ϕ' in disjunctive normal form, that is, as a disjunction of conjunctions of atomic formulae. Observing that $\exists y(\phi_1 \vee \phi_2)$ is equivalent to $\exists y \phi_1 \vee \exists y \phi_2$, we assume w.l.o.g. that ϕ' is just a conjunction of atomic formulae. Any atomic formula of the form $x \in L$, for some L , can be moved outside the existential quantification, since $\exists y(x \in L \wedge \psi)$ is equivalent to $x \in L \wedge \exists y \psi$. All atomic formulae of the form $y \in L$ can be combined into a single one, since regular languages are closed under intersection.

Finally we may assume that $\phi'(x, y)$ is a conjunction of a single atomic formula of the form $y \in L$ (if no such formula appears, we can write $y \in A^*$), and some combination of atomic formulae among $x \sqsubset y$, $x \sqsupset y$, $x = y$, and $x \perp y$. If at least two of these appear, then

their conjunction is unsatisfiable, and so $\phi(x)$ is equivalent to $x \in \emptyset$. If none of them appear, $\exists y(y \in L)$ is equivalent to $x \in A^*$ (or to $x \in \emptyset$ if L is empty). If exactly one of them appears, say $x R y$, then $\exists y(y \in L \wedge xRy)$ is equivalent to $x \in L_\phi$ for $L_\phi = R^{-1}(L)$, which is regular as observed earlier. ◀

► **Theorem 5.5.** *The truth problem for $\text{FO}^2(A^*, \sqsubseteq)$ is decidable.*

Proof. Lemma 5.4 provides a recursive procedure for computing the set of words that make $\phi(x)$ true. When ϕ is a closed formula, this set is A^* or \emptyset depending on whether ϕ is true or not. ◀

5.3 Hardness for FO^2

The main question left open in this paper is the complexity of the decidable FO^2 theory. The recursive procedure described in Lemma 5.4 is potentially non-elementary since nested negations lead to nested complementations of regular languages.

Our preliminary attempts suggest that the question is difficult. At the moment we can only demonstrate the following lower bound.

► **Theorem 5.6.** *Truth checking for the basic logic, restricting to FO^2 sentences which only use letters (that is, words of length 1) as constants, is PSPACE-hard.*

Proof. We reduce from TQBF, the truth problem for quantified boolean formulae. W.l.o.g. a given instance of TQBF has the form $\phi' = \exists p_1 \forall p_2 \dots \exists p_{2n-1} \forall p_{2n} \phi$.

Consider the alphabet A with $4n$ letters, T_i and F_i for each $1 \leq i \leq 2n$. A word $w \in A^*$ is intended to encode a (partial) boolean valuation V_w of the variables p_1, \dots, p_{2n} : if T_i appears in w , $V_w(p_i) = \text{true}$, and if F_i appears in w , $V_w(p_i) = \text{false}$. We do not consider “inconsistent” words, in which both T_i and F_i appear. Observe that if x and y represent partial valuations and $x \sqsubseteq y$, then V_y extends V_x . Conversely, any valuation extending V_x can be represented by a suitable y' with $x \sqsubseteq y'$.

For each i , let $\varphi_i(w)$ be a formula that says “the domain of V_w is $\{x_1, \dots, x_i\}$ ”:

$$\bigwedge_{1 \leq j \leq i} ((T_j \sqsubseteq w \vee F_j \sqsubseteq w) \wedge \neg(T_j \sqsubseteq w \wedge F_j \sqsubseteq w)) \wedge \bigwedge_{i < j \leq 2n} (T_j \not\sqsubseteq w \wedge F_j \not\sqsubseteq w)$$

We now translate the given TQBF instance ϕ' into an FO^2 sentence ψ' in our logic:

$$\begin{aligned} \psi' = & \exists x(\varphi_1(x) \wedge \forall y((\varphi_2(y) \wedge x \sqsubseteq y) \implies \exists x(\varphi_3(x) \wedge y \sqsubseteq x \wedge \dots \\ & \wedge \exists x(\varphi_{2n-1}(x) \wedge y \sqsubseteq x \wedge \forall y((\varphi_{2n}(x) \wedge x \sqsubseteq y) \implies \psi))) \end{aligned}$$

where ψ is obtained from ϕ by replacing each p_i with $T_i \sqsubseteq y$.

The formula ψ' uses the two variables x and y alternately, to build up suitable valuations with the appropriate alternation of \exists and \forall . It is easy to see that ϕ' is true if and only if ψ' is true.

Finally, it was not necessary to assume that ϕ' had a strict alternation of \exists and \forall , but it makes the presentation of the proof simpler. ◀

6 Concluding remarks

We considered the first-order logic of the subsequence ordering and investigated decidability and complexity questions. It was known that the Σ_3 theory is undecidable and that the Σ_1 theory is decidable. We settled the status of the Σ_2 fragment by showing that it has an

undecidable theory, even when restricting to formulae using no constants. To remain in the Σ_2 fragment, our reduction encoded language-theoretic problems rather than undecidable number-theoretic logical fragments as is more usual.

We also showed that the FO^2 theory of the subsequence ordering is decidable using automata-theoretic techniques. The FO^2 fragment is quite interesting. We note that it encompasses modal logics where the subsequence ordering correspond to one step (or its reverse) as used in the verification of unreliable channel systems.

Finally, we provided some new complexity results like Theorems 2.1 and 5.6.

We can list a few interesting directions suggested by this work. First, on the fundamental side, the main question left open is the precise complexity of the FO^2 theory.

Regarding applications, it would be interesting to see how the decidability results can be extended to slightly richer logics (perhaps with some extra functions or predicates, or some additional logical constructs) motivated by specific applications in automated reasoning or program verification.

Acknowledgements. We thank Dietrich Kuske who outlined the proof of Theorem 5.5.

References

- 1 P. A. Abdulla, M. Faouzi Atig, Yu-Fang Chen, L. Holík, A. Rezine, P. Rümmer, and J. Stenman. String constraints for verification. In *Proc. CAV 2014*, volume 8559 of *Lecture Notes in Computer Science*, pages 150–166. Springer, 2014.
- 2 P. A. Abdulla, A. Collomb-Annichini, A. Bouajjani, and B. Jonsson. Using forward reachability analysis for verification of lossy channel systems. *Formal Methods in System Design*, 25(1):39–65, 2004.
- 3 J. Berstel. *Transductions and Context-Free Languages*. B. G. Teubner, Stuttgart, 1979.
- 4 P. Bouyer, N. Markey, J. Ouaknine, Ph. Schnoebelen, and J. Worrell. On termination and invariance for faulty channel machines. *Formal Aspects of Computing*, 24(4–6):595–607, 2012.
- 5 H. Comon. Solving symbolic ordering constraints. *Int. J. Foundations of Computer Science*, 1(4):387–412, 1990.
- 6 H. Comon and R. Treinen. Ordering constraints on trees. In *Proc. CAAP '94*, volume 787 of *Lecture Notes in Computer Science*, pages 1–14. Springer, 1994.
- 7 V. Ganesh, M. Minnes, A. Solar-Lezama, and M. C. Rinard. Word equations with length constraints: What’s decidable? In *Proc. HVC 2012*, volume 7857 of *Lecture Notes in Computer Science*, pages 209–226. Springer, 2013.
- 8 Ch. Haase, S. Schmitz, and Ph. Schnoebelen. The power of priority channel systems. *Logical Methods in Comp. Science*, 10(4:4), 2014.
- 9 J. Harrison. *Handbook of Practical Logic and Automated Reasoning*. Cambridge University Press, 2009.
- 10 P. Hooimeijer and W. Weimer. StrSolve: solving string constraints lazily. *Autom. Softw. Eng.*, 19(4):531–559, 2012.
- 11 P. Karandikar and Ph. Schnoebelen. Generalized Post embedding problems. *Theory of Computing Systems*, 56(4):697–716, 2015.
- 12 D. Kuske. Theories of orders on the set of words. *RAIRO Theoretical Informatics and Applications*, 40(1):53–74, 2006.
- 13 D. Kuske. Private email exchanges, April 2014.
- 14 J. Sakarovitch. *Elements of Automata Theory*. Cambridge University Press, 2009.

- 15 J. Sakarovitch and I. Simon. Subwords. In M. Lothaire, editor, *Combinatorics on words*, volume 17 of *Encyclopedia of Mathematics and Its Applications*, chapter 6, pages 105–142. Cambridge Univ. Press, 1983.
- 16 I. Simon. Piecewise testable events. In *Proc. 2nd GI Conf. on Automata Theory and Formal Languages*, volume 33 of *Lecture Notes in Computer Science*, pages 214–222. Springer, 1975.

A Proof of Lemma 4.3

Assume $|u| = |v| = n + 1$ and $u \neq v$ as in the statement of the Lemma.

We say that a word w *distinguishes* u and v if w is a subword of exactly one of u and v . We have to prove that there exists such a distinguisher w with $|w| \leq n$.

Writing a word $w \in A^*$ under the form $w = a_1^{n_1} \dots a_k^{n_k}$ where each a_i is a letter so that $a_i \neq a_{i+1}$ for all $i = 1, \dots, k - 1$ and $n_i \geq 1$ for all $i = 1, \dots, k$ is called the *block factorisation* of w . Here k is the number of blocks in w . We now consider several cases:

- Assume that u has only one block. Then $u = a^{n+1}$ for some $a \in A$, and some one-letter word distinguishes u and v . The same reasoning applies if v has only one block.
- Assume that u and v have at least two blocks each, and there is some letter $a \in A$ such that $|u|_a \neq |v|_a$. Then a^k distinguishes u and v for some $k \leq n$.
- We are left to deal with cases where u and v have at least two blocks, and have the same Parikh image, that is, $|u|_a = |v|_a$ for every $a \in A$.

Assume now that u has exactly two blocks. Then $u \in a^+b^+$ for some $a, b \in A$ with $a \neq b$. Since v has the same number of a 's and b 's but differs from u , we must have $ba \sqsubseteq v$. But $ba \not\sqsubseteq u$, so ba is a distinguisher (here we use the assumption that $n \geq 2$).

- Finally assume that u has at least three blocks. Pick a block B of u which is neither the first nor the last, and let a be the unique letter belonging to B . Let $\ell = |u|_a$ and write u as $u = s_0as_1a \dots as_\ell$. Then

$$|s_0| + \dots + |s_\ell| = (n + 1) - \ell.$$

At least two of the numbers $|s_0|, \dots, |s_\ell|$ are strictly positive, since the two blocks immediately to the left and right of B both exist, and both do not have a . Thus for all i , $|s_i| < (n + 1) - \ell$.

Since $|v|_a = \ell$, we can write $v = t_0at_1a \dots at_\ell$. We assume $u \sim_n v$ and obtain a contradiction. For each i such that $0 \leq i \leq \ell$, consider the word $z_i = a^i s_i a^{\ell-i}$. We have $|z_i| \leq n$, and $z_i \sqsubseteq u$. Since $u \sim_n v$, we have $z_i \sqsubseteq v$. Since both z_i and v have exactly ℓ occurrences of a , we have $s_i \sqsubseteq t_i$. This holds for all i , so $u \sqsubseteq v$. But $|u| = |v|$, so $u = v$, which is a contradiction.