

# Equivalence of finite-valued streaming string transducers is decidable

Anca Muscholl

LaBRI, University of Bordeaux

Gabriele Puppis

CNRS, LaBRI

## Abstract

In this paper we provide a positive answer to a question left open by Alur and and Deshmukh in 2011 by showing that equivalence of finite-valued copyless streaming string transducers is decidable.

**2012 ACM Subject Classification** Theory of computation → Transducers

**Keywords and phrases** String transducers, equivalence, Ehrenfeucht conjecture

**Digital Object Identifier** 10.4230/LIPIcs.ICALP.2019.121

**Funding** DeLTA project (ANR-16-CE40-0007)

## 1 Introduction

Finite transducers are simple devices that allow to reason about data transformations in an effective, and even efficient way. In their most basic form they transform strings using finite control. Unlike automata, their power heavily depends on various parameters, like non-determinism, the capability of scanning the input several times, or the kind of storage they may use. The oldest transducer model, known as generalized sequential machine, extends finite automata by outputs. Inspired by an approach that applies to arbitrary relational structures [10], logic-based transformations (also called transductions) were considered by Engelfriet and Hoogetboom [13]. They showed that two-way transducers and monadic-second order (MSO) definable transductions are equivalent in the deterministic case (and even if the transduction is single-valued, which is more general than determinism). This equivalence supports thus the notion of “regular” functions, in the spirit of classical results on regular word languages from automata theory and logics due to Büchi, Elgot, Trakhtenbrot, Rabin, and others. A one-way transducer model that uses write-only registers as additional storage was proposed a few years ago by Alur and Cerný [2], and called streaming string transducer (SST). SST were shown equivalent to two-way transducers and MSO definable transductions in the deterministic setting, and again, even in the single-valued case.

In the relational case the picture is less satisfactory, as expressive equivalence is only preserved for SST and non-deterministic MSO transductions [5], which extend the original MSO transductions by existentially quantified monadic parameters. On the other hand, two-way transducers and SST are incomparable in the relational case. Between functions and relations there is however one class of transductions that exhibits a better behavior, and this is the class of finite-valued transductions. Being finite-valued means that there exists some constant  $k$  such that every input belonging to the domain has at most  $k$  outputs.

Finite-valued transductions were intensively studied in the setting of one-way and two-way transducers. For one-way transducers,  $k$ -valuedness can be checked in PTIME [18]. In addition, every  $k$ -valued one-way transducer can be effectively decomposed into a union of  $k$  unambiguous one-way transducers of exponential size [27, 25]. For both two-way transducers and SST, checking  $k$ -valuedness is in PSPACE.



© A. Muscholl and G. Puppis;

licensed under Creative Commons License CC-BY

46th International Colloquium on Automata, Languages, and Programming (ICALP 2019).

Editors: Christel Baier, Ioannis Chatzigiannakis, Paola Flocchini, and Stefano Leonardi;

Article No. 121; pp. 121:1–121:23



Leibniz International Proceedings in Informatics

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



Besides expressiveness, another fundamental question concerning transducers is the equivalence problem, that is, the problem of deciding whether two transducers define the same relation (or the same partial function if we consider the single-valued case). The equivalence problem turns out to be PSPACE-complete for deterministic two-way transducers [17], single-valued two-way transducers, as well as for single-valued SST [5]. For deterministic SST, equivalence is in PSPACE [3], but it is open whether this complexity upper bound is optimal. For arbitrary SST, and in fact even for non-deterministic one-way transducers over a unary output alphabet, equivalence is undecidable [14, 19]. The equivalence problem for  $k$ -valued one-way transducers was shown to be decidable by Culik and Karhumäki using an elegant argument based on Ehrenfeucht’s conjecture [11], and the authors noted that the same proof goes through for two-way transducers as well. The decidability status for the equivalence problem for  $k$ -valued SST was first stated as an open problem in [5]. Another open problem is whether SST and two-way transducers are equivalent in the finite-valued case, like in the single-valued case. It is worth noting, however, that in the full relational case SST and two-way transducers are incomparable. Concerning this last open question, a partial positive answer was given in [15], by decomposing any finite-valued SST with only one register into a finite union of unambiguous SST. This decomposition result also entails the decidability of the equivalence problem for the considered class.

The main result of this paper is a positive answer to the first question left open in [5]:

► **Theorem 1.** *The equivalence problem for finite-valued SST is decidable.*

We show the above result with a proof idea due to Culik and Karhumäki [11], based on the Ehrenfeucht conjecture. Our proof is much more involved, because SST produce their outputs piece-wise, in contrast to one-way and two-way transducers, that produce output linearly while reading the input. We manage to overcome this obstacle using some (mild) word combinatorics and word equations, by introducing a suitable normalization procedure for SST. We believe that our technique will also allow to solve the second problem left open in [5], which is the expressive equivalence between finite-valued SST and two-way transducers.

### Related work.

The equivalence problem for transducers has recently raised interest for more complex types of transducers in the single-valued case: Filiot and Reynier showed that equivalence of copyful, deterministic SST is decidable by showing them equivalent to HDTOL systems and applying [11], which contains the above-mentioned result as a special case. Subsequently, Benedikt et al. showed that equivalence of copyful, deterministic SST has Ackerman complexity, with a proof based on polynomial automata and ultimately on Hilbert’s basis theorem [6]. Interestingly, the use of Hilbert’s basis theorem goes back to the proof of Ehrenfeucht’s conjecture [1, 16]. A similar approach was used by Boiret et al. in [7] to show that bottom-up register automata over unordered forests have a decidable equivalence problem, see also the nice survey [8].

### Overview.

Section 2 introduces the transducer model, then Section 3 sets up the technical machinery that allows to normalize finite-valued SST. Section 4 shows the major normalization result, which holds for left quotients of SST. Finally Section 5 recalls the Ehrenfeucht-based proof for equivalence and the application to finite-valued SST. A full version of the paper is available at <https://arxiv.org/abs/1902.06973>.

## 2 Streaming string transducers

A *streaming string transducer* (SST) is a tuple  $T = (\Sigma, \Gamma, X, Q, U, I, E, F, x_{\text{out}})$ , where  $\Sigma$  and  $\Gamma$  are finite input and output alphabets,  $X$  is a finite set of registers (usually denoted  $x, x', x_1, x_2$ , etc.),  $Q$  is a finite set of states,  $U$  is a finite set of register updates, that is, functions from  $X$  to  $(X \uplus \Gamma)^*$ ,  $I, F \subseteq Q$  are subsets of states, defining the initial and final states,  $E \subseteq Q \times \Sigma \times U \times Q$  is a transition relation, describing, for each state and input symbol, the possible register updates and target states, and finally  $x_{\text{out}} \in X$  is a register for the output. Note that, compared to the original definition from [2], here we forbid for simplicity the use of final production rules, that perform an additional register update after the end of the input. This simplification is immaterial with respect to the decidability of the equivalence problem. For example, it can be enforced, without loss of generality, by assuming that all well-formed inputs are terminated by a special marker, say  $\neg$ , on which the transducer can apply a specific transition. We assume here that *all inputs of a transducer are non-empty and of the form  $u \neg$ , with  $\neg$  not occurring in  $u$ .*

### Copyless restriction and capacity.

An SST as above is *copyless* if for all register updates  $f \in U$ , every register  $x \in X$  appears at most once in the word  $f(x_1) \dots f(x_m)$ , where  $X = \{x_1, \dots, x_m\}$ . For a copyless SST, every output has length at most linear in the length of the input. More precisely, every output associated with an input  $u$  has length at most  $c|u|$ , where  $c = \max_{f \in U} \sum_{x \in X} |f(x)|_{\Gamma}$  is the maximum number of letters that the SST can add to its registers along a single transition (this number  $c$  is called *capacity* of the SST).

*Hereafter, we assume that all SST are copyless.*

### Register updates and flows.

Every register update, and in general every function  $f : X \rightarrow (X \uplus \Gamma)^*$  is naturally extended to a morphism on  $(X \uplus \Gamma)^*$ , by defining it as identity over  $\Gamma$ . When reasoning with register updates, it is sometimes possible to abstract away the specific words over  $\Gamma$ , and only consider how the contents of the registers flows into other registers. Formally, the *flow* of an update  $f : X \rightarrow (X \uplus \Gamma)^*$  is the bipartite graph that consists of two ordered sequences of nodes, one on the left and one on the right, with each node in a sequence corresponding to a specific register, and arrows that go from the node corresponding to register  $x$  to a right node corresponding to register  $x$  whenever  $x$  occurs in  $f(x)$ . For example, the flow of the update  $f$  defined by  $f(x_1) = a x_1 a a x_3$ ,  $f(x_2) = b a$ , and  $f(x_3) = x_2 b$  is the second bipartite graph in the figure on page 5.

Note that there are finitely many flows on a fixed number of registers. Moreover, flows can be equipped with a natural composition operation: given two flows  $F_1$  and  $F_2$ ,  $F_1 \cdot F_2$  is the bipartite graph obtained by glueing the right nodes of  $F_1$  with the left nodes of  $F_2$ , and by shortcutting pairs of consecutive arrows. We call *flow monoid* of an SST  $T$  the monoid of flows generated by the updates of  $T$ , with the composition operation as associative product.

### Transitions, runs, and loops.

A transition  $(q, a, f, q')$  of an SST  $T$  is conveniently denoted by the arrow  $q \xrightarrow[a]{a/f} q'$ , and the subscript  $T$  is often omitted when clear from the context. A *run* on  $w = a_1 \dots a_n$  is a sequence of transitions of the form  $q_0 \xrightarrow{a_1/f_1} q_1 \xrightarrow{a_2/f_2} \dots \xrightarrow{a_n/f_n} q_n$ . Sometimes, a run as above is equally denoted by  $q_0 \xrightarrow{w/f} q_n$ , so as to highlight the underlying input  $w$  and the

induced register update  $f = f_1 \circ \dots \circ f_n$ . A run is *initial* (resp. *final*) if it begins with an initial (resp. final) state; it is *successful* if it is both initial and final.

Given two registers  $x, x'$  and a run  $\rho: q \xrightarrow{w/f} q'$ , we say that  $x$  *flows into*  $x'$  *along*  $\rho$  if  $x$  occurs in  $f(x')$ . Note that this property depends only on the flow of the induced update  $f$ .

An SST is said to be *trimmed* if every state occurs in at least one successful run, so every state is reachable from the initial states and co-reachable from the final states. This property can be easily enforced with a polynomial-time preprocessing.

When reasoning with automata, it is common practice to use pumping arguments. Pumping will also be used here, but the notion of loop needs to be refined as to take into account the effect of register updates. Formally, a *loop* of a run  $\rho$  of an SST is any non-empty factor of  $\rho$  of the form  $\gamma: q \xrightarrow{w/f} q$ , that starts and ends in the same state  $q$ , and induces a *flow-idempotent* update, namely, an update  $f$  such that  $f$  and  $f \circ f$  have the same flow.

### Outputs and finite-valuedness.

The *output* of a successful run  $\rho: q_0 \xrightarrow{w/f} q_n$  is defined as  $\text{out}(\rho) = (f_0 \circ f)(x_{\text{out}})$ , where  $f_0(x) = \varepsilon$  for all  $x \in X$ . Sometimes, we write  $\text{out}(f)$  in place of  $\text{out}(\rho)$ . The *relation realized by an SST* is the set of pairs  $(u, v) \in \Sigma^* \times \Gamma^*$ , where  $u$  is a well-formed input (namely, terminating with  $\neg$ ) and  $v$  is the output associated with some successful run on  $u$ . An SST is *k-valued* if for every input  $u$ , there are at most  $k$  different outputs associated with  $u$ . It is *single-valued* (resp. *finite-valued*) if it is  $k$ -valued for  $k = 1$  (resp. for some  $k \in \mathbb{N}$ ). The domain of an SST  $T$ , denoted  $\text{Dom}(T)$ , is the set of input words that have some successful run in  $T$ . Two SST  $T_1, T_2$  are *equivalent*, denoted as  $T_1 \equiv T_2$ , if they realize the same relation over  $\Sigma^* \times \Gamma^*$ .

### Register valuations.

A *register valuation* is a function from  $X$  to  $\Gamma^*$ . Given a successful run  $\rho: q_0 \xrightarrow{a_1/f_1} q_1 \xrightarrow{a_2/f_2} \dots \xrightarrow{a_n/f_n} q_n$  and a position  $i \in \{0, \dots, n\}$  in it, the *register valuation at position  $i$  in  $\rho$*  is the function  $\text{val}_{\rho,i}$  that is defined inductively on  $i$  as follows:  $\text{val}_{\rho,0}(x) = \varepsilon$ , for all  $x \in X$ , and  $\text{val}_{\rho,i+1} = \text{val}_{\rho,i} \circ f_i$ . Note that  $\text{val}_{\rho,n}(x_{\text{out}})$  coincides with the final output  $\text{out}(f_1 \circ \dots \circ f_n)$  produced by  $\rho$ .

## 3 Normalizations

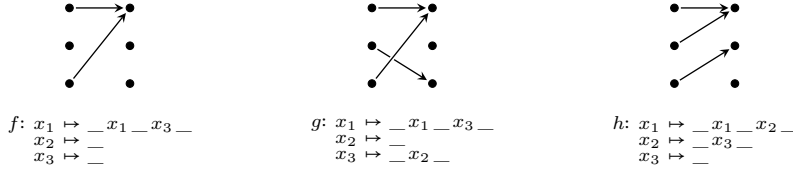
A major stumbling block in deciding equivalence of SST, as well as other crucial problems, lies in the fact that the same output can be produced by very different runs. This phenomenon already appears with much simpler transducers, e.g. with one-way transducers, where runs may produce the same output, but at different speeds. However, the phenomenon is more subtle for SST, as the output is produced piece-wise, and not sequentially: runs with same output may appear to be different in many ways, e.g. in terms of the flows of the register updates, or in terms of shifts of portions of the output. The goal of this section is to provide suitable normalization steps that remove, one at a time, the above mentioned degrees of freedom in producing the same output.

Another issue that we will be concerned with is the compatibility of the normalization steps with constructions on transducers that shortcut arbitrary long runs into a single transition. Essentially, we aim at having an effective notion of equivalence w.r.t. final outputs that works not only for transitions but also for runs.

### Normalization of flows.

In this section,  $m$  will always denote the number of registers of an SST and  $X = \{x_1, \dots, x_m\}$  the set of registers. It is convenient to equip  $X$  with a total order, say  $x_1 < \dots < x_m$ . Accordingly, we let  $\chi = x_1 \dots x_m$  be the juxtaposition of all register names, and  $f(\chi) = f(x_1) \dots f(x_m)$  for every register update  $f$ .

We say that a register update  $f$  is *non-erasing* if for every register  $x$ ,  $f(\chi)$  contains at least an occurrence of  $x$  (in fact, exactly one, since  $T$  is copyless). This can be rephrased as a property of the flow of  $f$ , where every node on the left must have an outgoing arrow. In a similar way, we say that  $f$  is *non-permuting* if registers appear in  $f(\chi)$  with their natural order and without jumps, that is,  $f(\chi) \in \Gamma^* x_1 \Gamma^* \dots \Gamma^* x_k \Gamma^*$ , for some  $k \leq m$ . As before, this can be rephrased by saying that the arrows in the flow of  $f$  must not be crossing, and the target nodes to the right must form a prefix of  $\chi$ . Below are some examples of updates with their flows: the first update  $f$  is erasing, the second update  $g$  is non-erasing but permuting, and the third update  $h$  is non-erasing and non-permuting.



We say that  $T$  is *flow-normalized* if all its register updates are non-erasing and non-permuting. Note that a flow-normalized SST with  $m$  registers can have at most  $2^m$  different flows.

► **Proposition 2.** *One can transform any SST into an equivalent flow-normalized one.*

Recall that a register valuation is a function from  $X$  to  $\Gamma^*$ . With a flow-normalized SST, one can also define a dual notion of valuation, representing ‘gaps’ between registers that shrink along the run. For this we introduce  $m + 1$  fresh variables  $y_0, y_1, \dots, y_m$ , called *gaps*. Hereafter,  $Y = \{y_0, y_1, \dots, y_m\}$  will always denote the set of gaps. We use the term *valuation* to generically denote a register/gap valuation, that is, a function from  $X \uplus Y$  to  $\Gamma^*$ .

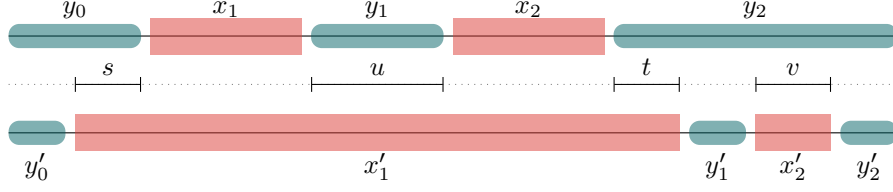
The idea is that a gap  $y_j$  represents a word that is inserted between register  $x_j$  (if  $j > 0$ ) and register  $x_{j+1}$  (if  $j < n$ ) so as to form the final output. Formally, given a word  $w \in \Gamma^* x_1 \Gamma^* \dots \Gamma^* x_k \Gamma^*$ , with  $k \leq m$ , and given two registers  $x_i, x_j$ , with  $i < j$ , we denote by  $w\langle x_i, x_j \rangle$  the maximal factor of  $w$  strictly between the unique occurrence of  $x_i$  and the unique occurrence of  $x_j$ , using the following conventions for the degenerate cases: if  $i = 0$ , then  $w\langle x_i, x_j \rangle$  is a maximal prefix of  $w$ ; if  $i > 0$  but there is no occurrence of  $x_i$ , then  $w\langle x_i, x_j \rangle = \varepsilon$ ; finally, if there is an occurrence of  $x_i$  but no occurrence of  $x_j$  in  $w$ , then  $w\langle x_i, x_j \rangle$  is a maximal suffix. Given a run  $\rho: q_0 \xrightarrow{a_1/f_1} q_1 \xrightarrow{a_2/f_2} \dots \xrightarrow{a_n/f_n} q_n$  and a position  $i$  in it, the valuation at position  $i$  of  $\rho$  is the function  $\text{val}_{\rho,i}: X \uplus Y \rightarrow \Gamma^*$  such that

- $\text{val}_{\rho,i}$  restricted to  $X$  is the register valuation at position  $i$  of  $\rho$ ,
- $\text{val}_{\rho,i}$  maps every gap  $y_j$  to the word  $(f_{i+1} \circ \dots \circ f_n)(\chi)\langle x_j, x_{j+1} \rangle$ .

By definition, the image of the word  $\zeta = y_0 x_1 y_1 \dots x_m y_m$  via the valuation  $\text{val}_{\rho,i}$  is always equal to the final output  $\text{out}(\rho)$ , for all positions  $i$ . In this sense, the sequence of valuations  $\text{val}_{\rho,0}, \text{val}_{\rho,1}, \dots, \text{val}_{\rho,n}$  can be identified with a sequence of factorizations of  $\text{out}(\rho)$ . For example, below are the factorizations of the output before and after a transition with

## 121:6 Equivalence of finite-valued streaming string transducers is decidable

register update  $f$  such that  $f(x_1) = s x_1 u x_2 t$  and  $f(x_2) = v$ , for  $s, u, t, v \in \Gamma^*$ :



This also suggests the principle that gaps, like registers, are updated along transitions via suitable morphisms, but in a symmetric way, that is, from right to left. For instance, in the above picture, the gaps  $y_0, y_1, y_2$  are updated by the function  $f^*$  such that  $f^*(y_0) = y_0 s$ ,  $f^*(y_1) = u$ , and  $f^*(y_2) = t y_1 v y_2$ . In general, the function  $f^*$ , called *gap update*, is uniquely determined by the register update  $f$ , and vice versa,  $f$  is uniquely determined by the gap update  $f^*$ . Another perhaps interesting phenomenon is that the gap update  $f^*$  is also non-erasing and non-permuting (the notion of non-permuting gap assignment is defined w.r.t. the reverse order  $y_m < \dots < y_0$ ).

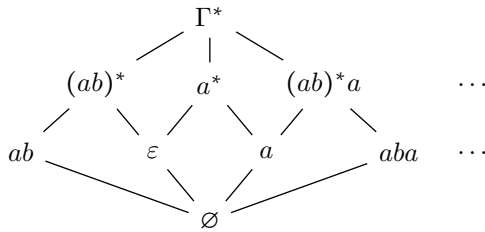
### Normalization of states.

The next normalization step splits the states of an SST in such a way that it becomes possible to associate with each state an over-approximation of the possible register/gap valuations witnessed when the state is visited along a successful run. These over-approximations are very simple languages over the output alphabet  $\Gamma$ , e.g. singleton languages like  $\{aba\}$  and periodic languages like  $\{ab\}^* \{a\}$  (often denoted  $(ab)^* a$  to improve readability). Basically our over-approximations refer to length and period constraints. The *period* of a word  $w$  is the least number  $0 < p \leq |w|$  such that  $w$  is a prefix of  $(w[1, p])^\omega$ . For example, the period of  $w = abcab$  is 3.

For a given parameter  $\alpha \in \mathbb{N}$  we define the family  $\mathcal{L}_\alpha$  that contains:

- the empty language  $\emptyset$ ,
- the singleton languages  $\{u\}$ , with  $u \in \Gamma^*$  and  $|u| \leq \alpha$ ,
- the periodic languages  $u^* v$ , with  $u \in \Gamma^+$  *primitive* (i.e.  $u = w^k$  only if  $k = 1$ ),  $|u| \leq \alpha$ , and  $v \in \Gamma^*$  strict prefix of  $u$ ,
- the universal language  $\Gamma^*$ .

The languages in  $\mathcal{L}_\alpha$ , partially ordered by containment, form a finite meet semi-lattice, where the meet is the intersection  $\cap$ . We depict here part of the lattice  $\mathcal{L}_\alpha$  for a parameter  $\alpha \geq 3$ :



The semi-lattice structure allows to derive a best over-approximation in  $\mathcal{L}_\alpha$  of any language  $L \subseteq \Gamma^*$ , that is:  $L^{\uparrow\alpha} = \cap \{L' \in \mathcal{L}_\alpha : L' \supseteq L\}$ . We will mostly use the approximation operator  $^{\uparrow\alpha}$  on singleton languages. For example, for  $\alpha = 3$ , we have  $\{aba\}^{\uparrow\alpha} = \{aba\}$ ,  $\{ababa\}^{\uparrow\alpha} = (ab)^* a$ , and  $\{abbb\}^{\uparrow\alpha} = \Gamma^*$ . Note also that if  $|w| \leq \alpha$  then  $w^{\uparrow\alpha} = \{w\}$ . A useful property is the compatibility of  $^{\uparrow\alpha}$  with concatenation, which immediately extends to compatibility with word morphisms:

► **Lemma 3.**  $(L_1 \cdot L_2)^{\uparrow\alpha} = (L_1^{\uparrow\alpha} \cdot L_2^{\uparrow\alpha})^{\uparrow\alpha}$  for every  $\alpha \in \mathbb{N}$  and  $L_1, L_2 \subseteq \Gamma^*$ .

Recall that  $X, Y$  denote, respectively, the sets of registers and gaps of a flow-normalized SST. Given a valuation  $\nu : X \uplus Y \rightarrow \Gamma^*$ , its  $\alpha$ -approximant is the function  $\nu^{\uparrow\alpha} : X \uplus Y \rightarrow \mathcal{L}_\alpha$  that maps any  $z \in X \uplus Y$  to the language  $\{\nu(z)\}^{\uparrow\alpha}$ . The set of  $\alpha$ -approximants is denoted  $\mathcal{L}_\alpha^{X \uplus Y}$ , and consists of all maps from  $X \uplus Y$  to  $\mathcal{L}_\alpha$ . Further let

$$\text{Val}_q = \{\text{val}_{\rho,i} : \rho \text{ successful run visiting } q \text{ at any position } i\}$$

be the set of possible valuations induced by an arbitrary successful run when visiting state  $q$ .

A first desirable property is that all valuations in  $\text{Val}_q$  have the *same*  $\alpha$ -approximant, which is thus determined by the state  $q$ . Formally, given a flow-normalized SST  $T$  with trimmed state space  $Q$ , we say that  $T$  *admits  $\alpha$ -approximants* if every state  $q \in Q$  can be effectively annotated with an  $\alpha$ -approximant  $A_q \in \mathcal{L}_\alpha^{X \uplus Y}$  in such a way that

$$\forall \nu \in \text{Val}_q : \quad \nu^{\uparrow\alpha} = A_q. \quad (1)$$

This condition is best understood as an invariant on lengths and periods that can be enforced on valuations of registers and gaps when visiting a particular state. For instance, when the approximant  $A_q$  guarantees a certain period, then this period will be the same for all valuations occurring at  $q$ , independently of the specific initial run that may lead to  $q$  (for registers), and of the run that may lead from  $q$  to an accepting state (for gaps).

The proposition below shows that it is always possible to refine any SST  $T$  so as to admit  $\alpha$ -approximants, for any parameter  $\alpha$ . The proof for  $\alpha$ -approximants that concern only registers could be understood as unfolding the SST  $T$ , and merging nodes corresponding to any two inputs  $u$  and  $v$ , with  $u$  prefix of  $v$ , whenever the induced  $\alpha$ -approximants at  $u$  and  $v$  are the same for every register  $x$ . In general, the resulting SST can be seen a covering of the original SST  $T$ , in the sense formalized by Sakarovitch and de Souza in [24]:  $T'$  is a *covering* of  $T$  if the states of  $T'$  can be mapped homomorphically to states of  $T$ , while preserving transitions and the distinction into initial and final states, and, moreover, the outgoing transitions of every state of  $T'$  map one-to-one to outgoing transitions of a corresponding state of  $T$ . This implies that the successful runs of  $T$  and those of  $T'$  are in one-to-one correspondence.

► **Proposition 4.** *Let  $T$  be a flow-normalized SST, and let  $\alpha \in \mathbb{N}$ . One can construct an equivalent flow-normalized SST  $T'$  that admits  $\alpha$ -approximants and that is a covering of  $T$ .*

*Notation.* Whenever an SST admits  $\alpha$ -approximants  $A_q$  as above, it is convenient to denote its states by triples of the form  $(q, A_X, A_Y)$ , where  $A_X$  (resp.  $A_Y$ ) is the restriction of the  $\alpha$ -approximant  $A_q$  of state  $q$  to registers (resp. gaps).

Note that the smaller the parameter  $\alpha$ , the weaker is the property required for  $\alpha$ -approximants (in particular, for  $\alpha = 0$  the lattice  $\mathcal{L}_\alpha$  collapses to  $\emptyset$  and  $\Gamma^*$ ). Choosing  $\alpha$  to be at least the capacity of the SST is already a reasonable choice, as it gives a nice characterization of equivalence of transitions w.r.t. the produced outputs (cf. Lemma 5 below). However, we will see that it is desirable to have even finer approximants, in such a way that our results will be compatible with left quotients of SST, that shortcut arbitrary long runs into single transitions. We postpone the technical details to Section 4, and only provide a rough intuition underlying the choice of the appropriate parameter  $\alpha$ . We will choose  $\alpha$  much larger than the capacity of the SST, so that, by pumping arguments, one can show that, for every state  $q$  and every parameter  $\beta \geq \alpha$ , the  $\beta$ -approximant cannot be strictly smaller than the  $\alpha$ -approximant on *all* valuations from  $\text{Val}_q$ .



### Normalization of transitions.

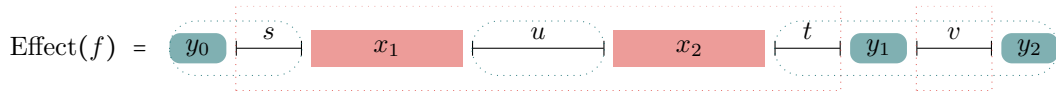
We finally turn to studying a notion of equivalence on transitions that is similar to the two-sided Myhill-Nerode equivalence on words. We will only compare transitions that consume the same input letter and link the same pair of states. Instead of using words as two-sided contexts, we will use initial and final runs that can be attached to the considered transitions in order to form successful runs, and instead of comparing membership in a language, we will compare the effect on the produced outputs.

Consider two transitions  $\tau_1 : q \xrightarrow{a/f_1} q'$  and  $\tau_2 : q \xrightarrow{a/f_2} q'$ . We say that  $\tau_1$  and  $\tau_2$  are *equivalent* if for every initial run  $\rho$  leading to  $q$  and every final run  $\sigma$  starting in  $q'$ , the outputs  $\text{out}(\rho\tau_1\sigma)$  and  $\text{out}(\rho\tau_2\sigma)$  are equal. We often refer to  $(\rho, \sigma)$  as a *context* for  $\tau_1, \tau_2$ .

In general, two transitions of an SST having the same source, target and  $\Sigma$ -label might turn out to be non-equivalent, and still produce the same output within specific contexts. However, Lemma 5 below shows that this is not the case with  $\alpha$ -approximants at hand, provided that  $\alpha$  is at least the capacity of the SST. More precisely, we will show that the equivalence of two transitions  $\tau_1 : (q, A_X, A_Y) \xrightarrow{a/f_1} (q', A'_X, A'_Y)$  and  $\tau_2 : (q, A_X, A_Y) \xrightarrow{a/f_2} (q', A'_X, A'_Y)$ , where  $f_1, f_2$  have the same flow, only depends on the  $\alpha$ -approximants  $A_X, A'_Y$  that annotate the source and target states. This will imply that  $\tau_1, \tau_2$  either always produce the same output or always produce different outputs, independently of the surrounding contexts. To prove the statement, we have to consider register valuations induced by initial runs, and symmetrically gap valuations induced by final runs. It helps to introduce the following:

*Notation.* Given an initial run  $\rho$ ,  $\text{val}_{\bullet\rho}$  is the register valuation induced at the end of  $\rho$ ; symmetrically,  $\text{val}_{\bullet\sigma}$  is the gap valuation induced at the beginning of a final run  $\sigma$ .

We also recall a consequence of the flow normalization: the effect on the final output of an update  $f$  that occurs in a successful run can be described by a word  $\text{Effect}(f)$  over the alphabet  $X \uplus Y \uplus \Gamma$ , defined as  $\text{Effect}(f) = y_0 f(x_1) y_1 \dots f(x_m) y_{m+1}$ . Note that in  $\text{Effect}(f)$  each register (resp. gap) occurs exactly once, according to the order  $x_1 < \dots < x_m$  (resp.  $y_0 < \dots < y_m$ ) — the occurrences of registers and gaps, however, may not be strictly interleaved. In  $\text{Effect}(f)$ , an occurrence of  $x_i \in X$  represents an abstract valuation for register  $x_i$  before applying the update  $f$ , while an occurrence of  $y_j \in Y$  represents an abstract valuation for the gap  $y_j$  after applying  $f$ . In particular, note that the  $x$ 's and the  $y$ 's refer to valuations induced at different positions of a run. The maximal factors of  $\text{Effect}(f)$  that are entirely over  $\Gamma$  represent the words that need to be added in order to get the register valuation after  $f$ , or equally the gap valuation before  $f$ . For instance, by reusing the example update  $f$  from page 6, where  $f(x_1) = s x_1 u x_2 t$  and  $f(x_2) = v$ , the effect of  $f$  is described by the word  $\text{Effect}(f) = y_0 s x_1 u x_2 t y_1 v y_2$ , suggestively depicted as



(here the lengths of the blocks labeled with variables are immaterial). Note that the  $y_j$  above are in fact the  $y'_j$  from the picture at page 6. In the above figure we have also highlighted with dotted rectangles the factors that represent gap valuations before the update (e.g.  $y_0 s$ ), and register valuations after the update (e.g.  $s x_1 u x_2 t$ ).

Given a valuation  $\nu$  and two approximants  $A \in \mathcal{L}_\alpha^X$  and  $A' \in \mathcal{L}_\alpha^Y$ , one for register valuations and the other for gap valuations, we write  $\nu \in A \uplus A'$  to mean that  $\nu(x) \in A(x)$  and  $\nu(y) \in A'(y)$  for all  $x \in X$  and  $y \in Y$ .



► **Lemma 5.** *Let  $T$  be a trimmed flow-normalized SST. Given two transitions  $\tau_i : q \xrightarrow{a/f_i} q'$ , with  $i \in \{1, 2\}$ , a context  $(\rho, \sigma)$  for them, and the  $\alpha$ -approximants  $A = ((\text{val}_{\rho, \bullet})^{\uparrow \alpha})|_X$  and  $A' = ((\text{val}_{\bullet, \sigma})^{\uparrow \alpha})|_Y$ , with  $\alpha \in \mathbb{N}$ , the following holds:*

1. *If  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds on all valuations  $\nu \in A \uplus A'$ , then  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$ .*
2. *If  $\tau_1, \tau_2$  have the same flow,  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$ , and  $\alpha \geq c$ , where  $c$  is the capacity of  $T$ , then  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds on all valuations  $\nu \in A \uplus A'$ .*

Recall that in an SST that admits  $\alpha$ -approximants, states are of the form  $(q, A_X, A_Y)$ , and we have  $((\text{val}_{\rho, \bullet})^{\uparrow \alpha})|_X = A_X$  (resp.  $((\text{val}_{\bullet, \sigma})^{\uparrow \alpha})|_Y = A_Y$ ) for every initial run  $\rho$  that ends in  $(q, A_X, A_Y)$  (resp. for every final run  $\sigma$  that starts in  $(q, A_X, A_Y)$ ). By pairing this with Lemma 5, we immediately obtain the following corollary:

► **Corollary 6.** *Let  $T$  be a trimmed flow-normalized SST. One can decide in polynomial time whether two given transitions  $\tau_1, \tau_2$  of  $T$  with the same flow are equivalent. Moreover, if  $T$  has capacity  $c$  and admits  $\alpha$ -approximants for some  $\alpha \geq c$ , then only two cases can happen:*

1. *either  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$  for every context  $(\rho, \sigma)$  (so  $\tau_1, \tau_2$  are equivalent),*
2. *or  $\text{out}(\rho \tau_1 \sigma) \neq \text{out}(\rho \tau_2 \sigma)$  for every context  $(\rho, \sigma)$  (so  $\tau_1, \tau_2$  are not equivalent).*

Another important consequence is the following theorem, that normalizes finite-valued SST in order to bound the maximum number of transitions linking the same pair of states and consuming the same input letter. This number is called *edge ambiguity* for short.

► **Theorem 7.** *Let  $T$  be a  $k$ -valued, flow-normalized SST that has  $m$  registers, capacity  $c$ , and that admits  $\alpha$ -approximants, for some  $\alpha \geq c$ . One can construct an equivalent SST  $T'$ , with the same states and the same registers as  $T$ , that has edge ambiguity at most  $k \cdot 2^m$ .*

**Proof.** By Corollary 6,  $T$  has at most  $k$  pairwise non-equivalent transitions with the same input letter, the same source and target states, and the same flow. Moreover equivalence of such transitions can be decided. We can then normalize  $T$  by removing in each equivalence class all but one transitions with the same flow. Since  $T$  has at most  $2^m$  flows, the normalization results in an equivalent SST  $T'$  with edge ambiguity at most  $k \cdot 2^m$ . ◀

The next section is devoted to prove a very similar result as above, but for all SST that can be obtained by shortcutting runs into single transitions, and that thus have arbitrary large capacity. This will be the main technical ingredient for establishing the decidability of the equivalence problem for  $k$ -valued SST.

## 4 Shortcut construction

Here we focus on a transformation of relations that absorbs the first input letter when this is equal to a specific element, say  $a \in \Sigma \setminus \{-\}$ . Such a transformation maps any relation  $R$  to the relation  $R_a = \{(u, v) : (au, v) \in R\}$ . Observe that  $R = R_{\varepsilon} \cup \bigcup_{a \in \Sigma} R_a$ , where  $R_{\varepsilon} = R \cap (\{-\} \times \Gamma^*)$ .

It is easy to see that the class of relations realized by SST is effectively closed under the transformation  $R \mapsto R_a$ . To prove this closure property, it is convenient to restrict, without loss of generality, to SST with *transient initial states*, namely, SST where no transition reaches an initial state. Under this assumption, the closure property also preserves the state space (though some states may become useless), the set of registers, the property of being  $k$ -valued, as well as the  $\alpha$ -approximants, if they are admitted by the original SST. However, the transformation does not preserve the capacity, which may increase.

## 121:10 Equivalence of finite-valued streaming string transducers is decidable

► **Lemma 8.** *Given a flow-normalized SST  $T$  with transient initial states, and given a letter  $a \in \Sigma \setminus \{-\}$ , one can construct an SST  $T_a$  with transient initial states such that*

- $\text{Dom}(T_a) = \{u \in \Sigma^* : au \in \text{Dom}(T)\},$
- $T_a$  on input  $u$  produces the same outputs as  $T$  on input  $au$ .

Moreover,  $T_a$  has the same states and the same registers as  $T$ ; if  $T$  has capacity  $c$ , then  $T_a$  has capacity  $2c$ ; if  $T$  admits  $\alpha$ -approximants, then so does  $T_a$  (via the same annotation).

The above construction can be applied inductively to compute an SST for any *left quotient*  $R_u = \{(v, w) : (uv, w) \in R\}$  of  $R$ . For  $u = a_1 \dots a_n \in (\Sigma \setminus \{-\})^*$ , we let  $T_u = (\dots (T_{a_1})_{a_2} \dots)_{a_n}$ .

The last and most technical step consists in proving that edge ambiguity can be uniformly bounded in every SST  $T_u$ , provided that the initial SST  $T$  is finite-valued, flow-normalized, and admits  $\alpha$ -approximants for a large enough  $\alpha$ . More precisely, we aim at establishing that, for  $\alpha$  much larger than the capacity of  $T$ , the notion of  $\alpha$ -approximant, besides satisfying Equation (1), also satisfies the following property:

$$\begin{aligned} \forall \beta \geq \alpha \quad \exists \rho : \quad & (\text{val}_{\rho, \bullet})^{\uparrow \beta} = (\text{val}_{\rho, \bullet})^{\uparrow \alpha} \\ \forall \beta \geq \alpha \quad \exists \sigma : \quad & (\text{val}_{\bullet, \sigma})^{\uparrow \beta} = (\text{val}_{\bullet, \sigma})^{\uparrow \alpha}. \end{aligned} \tag{2}$$

In this case we say that the  $\alpha$ -approximants are *tight*.

Intuitively, the above property can be explained as follows. When considering an SST with states annotated with  $\alpha$ -approximants, it may happen that for some larger parameter  $\beta$  some initial runs induce register valuations at a state  $(q, A_X, A_Y)$  whose  $\beta$ -approximants are strictly included in  $A_X$  (e.g. possibly entailing new periodicities). These runs should be thought of as exceptional cases, and there is a way of pumping them so as to restore the equality between  $A_X$  and the induced  $\beta$ -approximant.

Let us first see how tight approximants are used. The theorem below assumes that there is an SST  $T'$  with tight approximants (later we will show how to compute such an SST), and bounds the edge ambiguity of the SST  $T'_u$  that realizes a left quotient of  $T'$ .

► **Theorem 9.** *Let  $T'$  be a  $k$ -valued, flow-normalized SST realizing  $R$ , with transient initial states and tight  $\alpha$ -approximants. For every  $u \in \Sigma^*$ , one can construct an SST  $T'_u$  realizing  $R_u$ , with the same states and the same registers as  $T'$ , and with edge ambiguity at most  $k \cdot 2^m$ .*

**Proof.** The crux is to show that the SST  $T'_u$  obtained from Lemma 8 has at most  $k$  pairwise non-equivalent transitions with the same flow (for any given source/target state and label). Once this is proven, one can proceed as in the proof of Theorem 7, by removing all but one transition with the same flow in each equivalence class. By way of contradiction, assume that  $T'_u$  has  $k + 1$  pairwise non-equivalent transitions  $\tau_1, \dots, \tau_{k+1}$  with the same flow. Since  $T'$  admits tight  $\alpha$ -approximants, by Lemma 8 we know that the source and target state, respectively, of the previous transitions are annotated with tight  $\alpha$ -approximants, say  $(A_X, A_Y)$  and  $(A'_X, A'_Y)$ , respectively.

We begin by applying the first claim of Lemma 5, implying that the equation  $\text{Effect}(f_i) = \text{Effect}(f_j)$  is violated for some valuation  $\nu \in A_X \uplus A'_Y$ . Then, we let  $\beta = \max(\alpha, |u|c)$  and use Equations (1) and (2) to get a context  $(\rho, \sigma)$  such that  $(\text{val}_{\rho, \bullet})^{\uparrow \beta} = A_X$  and  $(\text{val}_{\bullet, \sigma})^{\uparrow \beta} = A'_Y$ . Finally, knowing that  $\beta$  is at least the capacity of  $T'_u$ , we apply the second claim of Lemma 5 to get  $\text{out}(\rho \tau_i \sigma) \neq \text{out}(\rho \tau_j \sigma)$ , thus witnessing non-equivalence of all pairs of transitions  $\tau_i, \tau_j$  at the same time. This contradicts the assumption that  $T'_u$  (and hence  $T'$ ) is  $k$ -valued. ◀

Now, let  $T$  be a flow-normalized SST with  $m$  registers, capacity  $c$ , and trimmed state space  $Q$ . Below, we show how to compute, with the help of Proposition 4, an SST  $T'$

equivalent to  $T$  that admits tight approximants. For simplicity, we will mostly focus on register valuations induced by initial runs, even though similar results can be also stated for gap valuations induced by final runs. We begin by giving a few technical results based on pumping arguments. We say that register  $x$  is *productive* along  $\rho$  if the update induced by  $\rho$  maps  $x$  to a word that contains at least one letter from  $\Gamma$ . We also recall that a loop of a run needs to induce a flow-idempotent update.

► **Lemma 10.** *If  $\rho = \rho_1 \gamma \rho_2$  is an initial run of  $T$ , with  $\gamma$  loop, then for every  $n > 0$  the pumped run  $\rho^{(n)} = \rho_1 \gamma^n \rho_2$  induces valuations  $\text{val}_{\rho^{(n)}} \bullet$  mapping any register  $x$  to a word of the form  $u_0 v_1^{n-1} u_1 \dots v_{2m}^{n-1} u_{2m}$ , where  $u_0, \dots, u_{2m}, v_1, \dots, v_{2m} \in \Gamma^*$  depend on  $\rho$  and  $x$ , but not on  $n$ . Moreover, we have  $v_i \neq \varepsilon$  for some  $i$  if there is a register  $x'$  that is productive along  $\gamma$  and that flows into  $x$  along  $\rho_2$ .*

Given a tuple of pairwise disjoint loops  $\bar{\gamma} = \gamma_1, \dots, \gamma_\ell$  in a run  $\rho$ , we write  $\rho' \triangleright_{\bar{\gamma}} \rho$  when  $\rho'$  is obtained from  $\rho$  by simultaneously pumping  $n$  times every loop  $\gamma_i$ , for some  $n > 0$ . When using this notation, we often omit the subscript  $\bar{\gamma}$ ; in this case we tacitly assume that  $\bar{\gamma}$  is uniquely determined from  $\rho$ . In this way, when writing, for instance,  $\rho', \rho'' \triangleright \rho$ , we will know that  $\rho', \rho''$  are obtained by pumping the same loops of  $\rho$ . We also say that a property on runs holds for all but finitely many  $\rho' \triangleright \rho$  if it holds on runs  $\rho'$  that are obtained from  $\rho$  by pumping  $n$  times the loops in a fixed tuple  $\bar{\gamma}$ , for all  $n > n_0$  and for a sufficiently large  $n_0$ .

► **Lemma 11.** *Let  $\rho$  be an initial run and  $x$  a register. If  $\text{val}_{\rho \bullet}(x)$  has length (resp. period) larger than  $\alpha = mc|Q|2^{3 \cdot 2^m}$ , then for every  $\beta \geq \alpha$  and for all but finitely many  $\rho' \triangleright \rho$ ,  $\text{val}_{\rho' \bullet}(x)$  has length (resp. period) larger than  $\beta$ .*

Using the previous lemmas and the fact that the type of quantification “for all but finitely many runs” commutes with conjunctions (e.g. those used to enforce properties on each register  $x \in X$ ), we obtain that  $\alpha$ -approximants are tight for sufficiently large  $\alpha$ :

► **Proposition 12.** *Let  $T'$  be the SST admitting  $\alpha$ -approximants that is obtained from  $T$  using Proposition 4, for any  $\alpha \geq mc|Q|2^{3 \cdot 2^m}$ , where  $Q$  is the set of states of  $T$ . The  $\alpha$ -approximants of  $T'$  are tight.*

## 5 Equivalence algorithm

The equivalence algorithm for  $k$ -valued SST follows a classical approach of Culik and Karhumäki [11] that is based on so-called *test sets*. A test set for two SST  $T_1, T_2$  over input alphabet  $\Sigma$  is a set  $F \subseteq \Sigma^*$  such that  $T_1, T_2$  are equivalent if and only if they are equivalent over  $F$ . The main contribution of [11] is to show that *finite* test sets exist and be computed effectively for  $k$ -valued one-way transducers. The key ingredient of their proof is to show the existence of a test set that works for *all* transducers with fixed number of states. An essential observation is that for  $k$ -valued one-way, or even two-way, transducers one can assume that the edge ambiguity is at most  $k$ . The reason for this is simply that the output is generated sequentially. For SST the situation is far more complex because the output is generated piecewise. The purpose of the normalizations performed in Section 3 was precisely to restore the property of bounded edge ambiguity.

In a nutshell, the existence of a test set for transducers is a consequence of Ehrenfeucht’s conjecture, whereas the effectiveness is based on the resolution of word equations due to Makanin (see e.g. the survey [12]).

Ehrenfeucht’s conjecture was originally stated as a conjecture about formal languages: for every language  $L \subseteq \Sigma^*$ , there is a finite subset  $F \subseteq L$  such that for all morphisms

## 121:12 Equivalence of finite-valued streaming string transducers is decidable

$f, g : \Sigma^* \rightarrow \Delta^*$ ,  $f(w) = g(w)$  for every  $w \in L$  if and only if  $f(w) = g(w)$  for every  $w \in F$ . Such a set  $F$  is called a *test set* for  $L$ .

There is an equivalent formulation of Ehrenfeucht's conjecture in terms of a compactness property of word equations [21]. Let  $\Sigma$  and  $\Omega$  be two alphabets, where the elements in  $\Omega$  are called unknowns. A word equation is a pair  $(u, v) \in \Omega^* \times \Omega^*$ , and a solution is a morphism  $\sigma : \Omega^* \rightarrow \Sigma^*$  such that  $\sigma(u) = \sigma(v)$ . Ehrenfeucht's conjecture is equivalent to saying that any system of equations over a finite set  $\Omega$  of unknowns has a finite, equivalent subsystem, where equivalence means that the solution sets are the same. The latter compactness property was proved in [1, 16] by encoding words by polynomials and using Hilbert's basis theorem.

In view of Propositions 2, 4, and 12, we can restrict without loss of generality to SST that are flow-normalized and that admit tight approximants. Hereafter, we shall tacitly assume that all transducers are of this form. Given some integers  $k, n, m$ , and  $e$ , let  $\mathcal{C}_k(n, m, e)$  be the class of  $k$ -valued SST with at most  $n$  states,  $m$  registers, and edge-ambiguity at most  $e$ . Note that if  $T$  is  $k$ -valued, then by Theorem 7 it belongs to  $\mathcal{C}_k(n, m, e)$ , where  $n, m$  are the number of states and registers of  $T$  and  $e = k \cdot 2^m$ . Similarly, by Lemma 8 and Theorem 9, every left quotient  $T_u$  also belongs to  $\mathcal{C}_k(n, m, e)$ .

Now, let us fix  $k, n, m, e$  and consider an arbitrary SST  $T$  from  $\mathcal{C}_k(n, m, e)$ . Following [11] we first build an abstraction of  $T$  by replacing each maximal factor from  $\Gamma^*$  occurring in some update function of  $T$ , by a distinct unknown from  $\Omega$ . The SST  $\Delta(T)$  obtained in this way is called a *schema*; its outputs are words over  $\Omega$ . Note that the assumption of bounded edge ambiguity is essential here to get a uniform bound on the number of unknowns required for a schema. Clearly, there are only finitely many schemas of SST in  $\mathcal{C}_k(n, m, e)$ . We denote by  $\phi_T : \Omega \rightarrow \Gamma^*$  the partial mapping (concretization) that associates with each unknown the corresponding word from  $\Gamma^*$  as specified by the updates of  $T$ .

We can rephrase the equivalence  $T_1 \equiv T_2$  of two arbitrary SST from  $\mathcal{C}_k(n, m, e)$  as an infinite “system” of word equations<sup>1</sup>  $\mathcal{S} = \bigwedge_{u \in \Sigma^*} \bigvee_{\pi} S_{\pi}$  over set of unknowns  $\Omega \uplus \Omega'$ . The unknowns from  $\Omega$  are used for the schema  $\Delta(T_1)$ , whereas those from  $\Omega'$  are used for  $\Delta(T_2)$ ; in particular,  $\phi_{T_1} : \Omega \rightarrow \Gamma^*$  and  $\phi_{T_2} : \Omega' \rightarrow \Gamma^*$ . The disjunctions in  $\mathcal{S}$  are finite, with  $\pi$  ranging over the possible schemas  $\Delta_1, \Delta_2$  (for  $T_1$  and  $T_2$ , respectively) and the possible partitions of the set of runs of  $\Delta_1$  and  $\Delta_2$  over the input  $u$ , into at most  $k$  groups (one for each possible output). Finally,  $S_{\pi}$  is a (finite) system of word equations, stating the equality of the words from  $\Omega^* \cup \Omega'^*$  that belong to the same group according to  $\pi$ .

The following lemma was stated in [11] for  $k$ -valued one-way transducers, but it holds as well for two-way transducers and for SST (even copyful, with a proper definition for  $\mathcal{C}_k(n, m, e)$ ):

► **Lemma 13.** *Given two SST  $T_1, T_2$  from  $\mathcal{C}_k(n, m, e)$ , the system  $\mathcal{S} = \bigwedge_{u \in \Sigma^*} \bigvee_{\pi} S_{\pi}$  has  $\phi_{T_1} \uplus \phi_{T_2}$  as solution if and only if  $T_1 \equiv T_2$ .*

As shown in [11], the Ehrenfeucht conjecture can be used to show that any infinite system  $\mathcal{S}$  as in Lemma 13 is equivalent to some *finite* sub-system  $\mathcal{S}_N = \bigwedge_{u \in \Sigma^{\leq N}} \bigvee_{\pi} S_{\pi}$ . This gives:

► **Lemma 14.** *Given  $n, m, e \in \mathbb{N}$ , there is  $N \in \mathbb{N}$  such that  $\Sigma^{\leq N}$  is a test set for every pair of SST  $T_1, T_2$  from  $\mathcal{C}_k(n, m, e)$ .*

Using Theorem 7 and Lemma 14 we can derive immediately the existence of a finite test set for any two  $k$ -valued SST. The last question is how to compute such a test set effectively. For this we will use the shortcut construction provided in Section 4.

---

<sup>1</sup> Formally,  $\mathcal{S}$  depends on  $n$  and  $k$ , but for simplicity we leave out the indices.

► **Lemma 15.** *Assume that the formulas  $\mathcal{S}_N$  and  $\mathcal{S}_{N+1}$  are equivalent, i.e., they have the same solutions. Then  $\Sigma^{\leq N}$  is a test set for any pair of SST from  $\mathcal{C}_k(n, m, e)$ .*

**Proof.** Let  $T_1 \equiv_r T_2$  denote equivalence of  $T_1$  and  $T_2$  relativized to  $\Sigma^{\leq r}$ . The goal is to prove that  $\Sigma^{\leq N}$  is a test set, namely, for all  $r > N$  and all  $T_1, T_2 \in \mathcal{C}_k(n, m, e)$ ,  $T_1 \equiv_r T_2$  holds if and only if  $T_1 \equiv_N T_2$ . Clearly, for any  $r \geq 0$ ,  $T_1 \equiv_{r+1} T_2$  is equivalent to  $T_{1,a} \equiv_r T_{2,a}$  for every  $a \in \Sigma$ , and  $T_1 \equiv_0 T_2$  (the latter being abbreviated as  $(*)$  below). Moreover, by Theorem 9, we have  $T_{1,a}, T_{2,a} \in \mathcal{C}_k(n, m, e)$ . This enables the following proof by induction on  $r$ :

$$\begin{array}{ccccc}
 T_1 \equiv_{r+1} T_2 & \Leftrightarrow & T_{1,a} \equiv_r T_{2,a} \quad (\forall a \in \Sigma) \quad \text{and } (*) & & T_1 \equiv_N T_2. \\
 & & \Updownarrow \text{ (ind. hyp.)} & & \Updownarrow \\
 & & T_{1,a} \equiv_N T_{2,a} \quad (\forall a \in \Sigma) \quad \text{and } (*) & \Leftrightarrow & T_1 \equiv_{N+1} T_2 \quad \blacktriangleleft
 \end{array}$$

Using Makanin's algorithm for solving word equations (and even for deciding the existential theory of word equations, see e.g. [12] for a modern presentation) we obtain:

► **Proposition 16.** *Given  $n, m, e \in \mathbb{N}$ , there is  $N \in \mathbb{N}$  such that  $\Sigma^{\leq N}$  is a test set for every pair of SST from  $\mathcal{C}_k(n, m, e)$ , and such an  $N$  can be effectively computed.*

**Proof.** By Lemma 14 we know that  $N$  exists, and Makanin's algorithm allows to determine whether  $\mathcal{S}_N, \mathcal{S}_{N+1}$  are equivalent, so to determine  $N$  by Lemma 15. ◀

We finally obtain the main result:

► **Theorem 1.** *The equivalence problem for finite-valued SST is decidable.*

Of course, Theorem 1 does not come with any complexity upper bound, mainly because of the Ehrenfeucht conjecture. The only known lower bound is PSPACE-hardness, which holds even for single-valued SST over unary output alphabets, and follows from a simple reduction from universality of NFA.

Quite surprisingly, the exact complexity of equivalence is not known even for *deterministic* SST, where the problem is known to be between NLOGSPACE and PSPACE [3]. We also recall that equivalence of deterministic SST with *unary output* can be checked in PTIME using invariants [4]. Finally, we recall that the currently best upper bound for solving word equations is PSPACE [23] (with even linear space requirement, as shown in [20]).

## 6 Conclusions

Our paper answers to a question left open in [5], showing that the equivalence problem for finite-valued SST is decidable. We followed a proof for one-way transducers due to Culik and Karhumäki [11], that is based on the Ehrenfeucht conjecture. The main contribution of the paper is to provide the technical development that allows to follow the proof scheme of [11]. We believe that this development will also allow to obtain stronger results. We conjecture that finite-valued SST can be effectively decomposed into finite unions of unambiguous SST. This would entail that in the finite-valued setting, two-way transducers and SST have the same expressive power, as is the case for single-valued transducers. If this holds with elementary complexity, then the equivalence of single-valued SST (or two-way transducers) could also be solved with elementary complexity. We believe that the complexity is indeed elementary, and leave this for future work.

## References

- 1 M.H. Albert and J. Lawrence. A proof of Ehrenfeucht's conjecture. *Theor. Comput. Sci.*, 41(1):121–123, 1985.
- 2 Rajeev Alur and Pavel Cerný. Expressiveness of streaming string transducer. In *IARCS Annual Conference on Foundation of Software Technology and Theoretical Computer Science (FSTTCS'10)*, volume 8 of *LIPIcs*, pages 1–12. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2010.
- 3 Rajeev Alur and Pavol Cerný. Streaming transducers for algorithmic verification of single-pass list-processing programs. In *POPL'11*. ACM, 2011.
- 4 Rajeev Alur, Loris D'Antoni, Jyotirmoy Deshmukh, Mukund Raghothaman, and Yifei Yuan. Regular functions and cost register automata. In *Proc. of Annual ACM/IEEE Symposium on Logic in Computer Science (LICS 2013)*, pages 13–22. IEEE, 2013.
- 5 Rajeev Alur and Jyotirmoy Deshmukh. Nondeterministic streaming string transducers. In *International Colloquium on Automata, Languages and Programming (ICALP'11)*, volume 6756 of *LNCS*. Springer, 2011.
- 6 Michael Benedikt, Timothy Duff, Aditya Sharad, and James Worrell. Polynomial automata: Zeroness and applications. In *Annual ACM/IEEE Symposium on Logic in Computer Science (LICS'17)*, pages 1–12. IEEE, 2017.
- 7 Adrien Boiret, Radosław Piórkowski, and Janusz Schmude. Reducing transducer equivalence to register automata problems solved by "hilbert method". In *IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS'18)*, volume 122 of *LIPIcs*, pages 48:1–48:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018.
- 8 Mikolaj Bojańczyk. The Hilbert method for transducer equivalence. *ACM SIGLOG News*, January 2019.
- 9 Thomas Colcombet. Factorisation forests for infinite words. In *International Symposium on Fundamentals of Computation Theory (FCT'07)*, number 4639 in *LNCS*, pages 226–237. Springer, 2007.
- 10 Bruno Courcelle and Joost Engelfriet. *Graph Structure and Monadic Second-Order Logic - A Language-Theoretic Approach*, volume 138 of *Encyclopedia of mathematics and its applications*. Cambridge University Press, 2012.
- 11 Karel Culik II and Juhani Karhumäki. The equivalence of finite valued transducers (on HDTOL languages) is decidable. *Theor. Comput. Sci.*, 47:71–84, 1986.
- 12 Volker Diekert. Makanin's algorithm. In M. Lothaire, editor, *Algebraic combinatorics on words*, volume 90 of *Encyclopedia of mathematics and its applications*, chapter 12, pages 387–442. Cambridge University Press, 2002.
- 13 Joost Engelfriet and Hendrik Jan Hooeboom. MSO definable string transductions and two-way finite-state transducers. *ACM Trans. Comput. Log.*, 2(2):216–254, 2001.
- 14 Patrick C. Fischer and Arnold L. Rosenberg. Multi-tape one-way nonwriting automata. *J. Comput. and System Sci.*, 2:88–101, 1968.
- 15 Paul Gallot, Anca Muscholl, Gabriele Puppis, and Sylvain Salvati. On the decomposition of finite-valued streaming string transducers. In *Annual Symposium on Theoretical Aspects of Computer Science (STACS'17)*, volume 66 of *LIPIcs*, pages 34:1–34:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2017.
- 16 Victor S. Guba. Equivalence of infinite systems of equations in free groups and semigroups to finite subsystems. *Mat. Zametki*, 40(3):688–690, 1986.
- 17 Eitan M. Gurari. The equivalence problem for deterministic two-way sequential transducers is decidable. *SIAM Journal of Computing*, 448–452, 1982.
- 18 Eitan M. Gurari and Oscar H. Ibarra. A note on finite-valued and finitely ambiguous transducers. *Math. Syst. Theory*, 16(1):61–66, 1983.
- 19 Oscar H. Ibarra. The unsolvability of the equivalence problem for e-free NGSM's with unary input (output) alphabet and applications. *SIAM J. of Comput.*, 7(4):524–532, 1978.



- 20 Artur Jez. Word equations in nondeterministic linear space. In *Proc. International Colloquium on Automata, Languages, and Programming (ICALP'17)*, volume 80 of *LIPICs*, pages 95:1–95:13. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2017.
- 21 Juhani Karhumäki. The Ehrenfeucht conjecture: a compactness claim for finitely generated free monoids. *Theor. Comput. Sci.*, 29:285–308, 1984.
- 22 Manfred Kufleitner. The height of factorization forests. In *International Symposium on Mathematical Foundations of Computer Science (MFCS'08)*, volume 5162 of *LNCS*, pages 443–454. Springer, 2008.
- 23 Wojciech Plandowski. Satisfiability of word equations with constants is in PSPACE. *JACM*, 51(3):483–496, 2004.
- 24 Jacques Sakarovitch and Rodrigo de Souza. On the decomposition of  $k$ -valued rational relations. In *Annual Symposium on Theoretical Aspects of Computer Science (STACS'08)*, volume 1 of *LIPICs*, pages 621–632. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2008.
- 25 Jacques Sakarovitch and Rodrigo de Souza. Lexicographic decomposition of  $k$ -valued transducers. *Theory Comput. Sci.*, 47:758–785, 2010.
- 26 Imre Simon. Factorization forests of finite height. *Theor. Comput. Sci.*, 72(1):65–94, 1990.
- 27 Andreas Weber. Decomposing a  $k$ -valued transducer into  $k$  unambiguous ones. *RAIRO-ITA*, 30(5):379–413, 1996.



## A Appendix

► **Proposition 2.** *One can transform any SST into an equivalent flow-normalized one.*

**Proof.** Let  $T = (\Sigma, \Gamma, X, Q, U, I, E, F, x_{\text{out}})$  be an SST. We need to construct an SST  $T'$  that simulates every run of  $T$  by guessing which registers in the current valuation contribute to form the final output, and in which precise order, by appropriately modifying the register updates so as to enforce non-erasing and non-permuting behaviours.

Formally, given any suffix  $\rho : q \xrightarrow{u/f} q'$  of a successful run of  $T$ , we define a partial bijection  $\pi_\rho : X \rightarrow X$  as follows: for every register  $x'$ , if  $x'$  is the  $i$ -th register occurring in  $f(x_{\text{out}})$ , then  $\pi_\rho(x') = x_i$ , otherwise, if  $x'$  does not occur in  $f(x_{\text{out}})$ , then  $\pi_\rho$  is undefined on  $x'$ . Any permutation of the form  $\pi_\rho$  can be thought of as a renaming of registers that contribute to the final output. By construction, the range of such a renaming is always an initial interval of the registers, i.e.  $\text{Rng}(\pi_\rho) = \{x_1, \dots, x_k\}$  for some  $k \leq m$ . For the sake of brevity, hereafter we call *renaming* any function of the above form, that is, any bijection from a subset  $\{x_{i_1}, \dots, x_{i_k}\}$  of  $X$  to  $\{x_1, \dots, x_k\}$ , for  $0 \leq k \leq m$ . We also let  $\|\pi\| = |\text{Dom}(\pi)|$  ( $= |\text{Rng}(\pi)|$ ) for any renaming  $\pi$ .

The normalized SST is defined as  $T' = (\Sigma, \Gamma, X, Q', U', I', E', F', x_1)$ , where:

- $Q' = Q \times R$ , where  $R$  is the set of all renamings,
- $U'$  contains all *non-erasing and non-permuting* updates of the form  $f[\pi \rightarrow \pi']$ , for  $f \in U$ ,  $\pi, \pi' \in R$ , where  $\pi$  is defined precisely on those registers that occur in  $f \circ (\pi')^{-1}(\{1, \dots, \|\pi'\|\})$ , and  $f[\pi \rightarrow \pi']$  is defined by

$$f[\pi \rightarrow \pi'](x_i) = \begin{cases} \pi \circ f \circ (\pi')^{-1}(x_i) & \text{if } i \leq \|\pi'\| \\ x_{i+\|\pi\|-\|\pi'\|} & \text{if } i > \|\pi'\| \text{ and } i + \|\pi\| - \|\pi'\| \leq m, \\ \varepsilon & \text{if } i > \|\pi'\| \text{ and } i + \|\pi\| - \|\pi'\| > m, \end{cases}$$

- $I' = I \times R$ ,
  - $E'$  contains all transition rules of the form  $(q, \pi) \xrightarrow{a/f'} (q', \pi')$ , with  $q \xrightarrow{a/f} q'$  transition rule in  $E$  and  $f' = f[\pi \rightarrow \pi']$ ,
  - $F' = F \times \{\pi_{\text{out}}\}$ , where  $\pi_{\text{out}}$  is the renaming defined only on  $x_{\text{out}}$  and mapping it to  $x_1$ .
- It is routine to show that  $T'$  is flow-normalized and equivalent to  $T$ . ◀

► **Lemma 3.**  $(L_1 \cdot L_2)^{\uparrow\alpha} = (L_1^{\uparrow\alpha} \cdot L_2^{\uparrow\alpha})^{\uparrow\alpha}$  for every  $\alpha \in \mathbb{N}$  and  $L_1, L_2 \subseteq \Gamma^*$ .

**Proof.** The left-to-right containment follows easily by monotonicity of  $\uparrow^\alpha$ . The converse containment boils down to proving that for every  $L \in \mathcal{L}_\alpha$  and  $w \in \Gamma^*$ ,  $L \supseteq \{w\} \cdot L_2$  implies  $L \supseteq \{w\} \cdot L_2^{\uparrow\alpha}$  (one can then take the conjunction of the latter implication over all  $w \in L_1$ , and prove in this way that  $L \supseteq L_1 \cdot L_2$  implies  $L \supseteq L_1 \cdot L_2^{\uparrow\alpha}$ , finally, using symmetric arguments, one derives that  $L \supseteq L_1 \cdot L_2$  implies  $L \supseteq L_1^{\uparrow\alpha} \cdot L_2^{\uparrow\alpha}$ ).

If  $L$  is empty, a singleton, or the universal language  $\Gamma^*$ , or if  $L_2$  is empty, then the considered implication holds trivially. So, we consider the case where  $L$  is a periodic language of the form  $u^*v$ , with  $u$  primitive and  $v$  prefix of  $u$ , and  $L_2$  is non-empty. Since  $L$  contains at least one word with  $w$  as prefix, we know that  $w \in u^*v'$ , for some  $v'$  prefix of  $u$ . Similarly, for every word  $w' \in L_2$ ,  $L$  must contain at least one word with  $w'$  as suffix, and hence  $L_2 \subseteq v''u^*v$  for some  $v''$  suffix of  $u$  such that  $v'v'' = u$ . This implies  $L_2^{\uparrow\alpha} \subseteq v''u^*v$ , and hence

$$L = u^*v \supseteq u^*v'v''u^*v \supseteq \{w\} \cdot L_2^{\uparrow\alpha}. \quad \blacktriangleleft$$

► **Proposition 4.** *Let  $T$  be a flow-normalized SST, and let  $\alpha \in \mathbb{N}$ . One can construct an equivalent flow-normalized SST  $T'$  that admits  $\alpha$ -approximants and that is a covering of  $T$ .*

**Proof.** Let  $T = (\Sigma, \Gamma, X, Q, U, I, E, F, x_1)$  be a flow-normalized SST with a trimmed state space  $Q$ , and let  $\alpha \in \mathbb{N}$ . The desired SST that admits  $\alpha$ -approximants is defined as  $T' = (\Sigma, \Gamma, X, Q', U, I', E', F', x_1)$ , where

- $Q' = Q \times \mathcal{L}_\alpha^X \times \mathcal{L}_\alpha^Y$  — namely, the states of  $T'$  are obtained by annotating the states of  $T$  with  $\alpha$ -approximants of register valuations and gap valuations,
- $I' = I \times \{A_\varepsilon\} \times \mathcal{L}_\alpha$ , with  $A_\varepsilon(x) = \varepsilon$  for all  $x \in X$  — namely, the initial states of  $T'$  have  $\alpha$ -approximants for register valuations initialized with the empty word  $\varepsilon$ ,
- $F' = F \times \mathcal{L}_\alpha^X \times \{A_\varepsilon\}$ , with  $A_\varepsilon(y) = \varepsilon$  for all  $y \in Y$  — namely, the final states of  $T'$  have  $\alpha$ -approximants for gap valuations initialized with  $\varepsilon$ ,
- $E'$  consists of transitions of the form  $(q, A_X, A_Y) \xrightarrow{a/f} (q', A'_X, A'_Y)$ , where  $q \xrightarrow{a/f} q'$  is a transition in  $E$ ,  $A'_X(x) = (A_X(f(x)))^{\uparrow\alpha}$  for all registers  $x \in X$ , and  $A'_Y(y) = (A_Y(f^*(y)))^{\uparrow\alpha}$  for all gaps  $y \in Y$ , with  $f^*$  gap update determined by  $f$ . Here,  $A_X(f(x))$  means substituting the languages from  $\mathcal{L}_\alpha$  associated with registers  $x' \in X$  into  $f(x)$ , and similarly for  $f^*(y)$ . Intuitively, the approximation  $A'_X(x)$  of a target register valuation is obtained by considering the effect of the update  $f$  on the register  $x$  when the source valuation ranges over  $A_X(x)$ , and symmetrically for a gap  $y$ .

It is clear from the above definitions that  $T'$  is a covering of  $T$ . This implies that the successful runs of  $T$  are precisely the successful runs of  $T'$  devoid of the  $\alpha$ -approximants, and hence  $T'$  is flow-normalized and equivalent to  $T$ .

It remains to prove that  $T'$  admits  $\alpha$ -approximants. This boils down to proving that for every successful run  $\rho$  of  $T'$  that visits a state  $(q, A_X, A_Y)$  at position  $i$ , we have  $A_X = (\text{val}_{\rho,i}^{\uparrow\alpha})|_X$  and  $A_Y = (\text{val}_{\rho,i}^{\uparrow\alpha})|_Y$ , where  $|_X$  (resp.  $|_Y$ ) denotes the restriction of a function to the set  $X$  (resp.  $Y$ ). For this, we recall that if  $\text{val}_{\rho,i}$  and  $\text{val}_{\rho,i+1}$  are two consecutive valuations w.r.t. a register update  $f$ , then

$$\text{val}_{\rho,i+1}|_X = \text{val}_{\rho,i}|_X \circ f \quad \text{and} \quad \text{val}_{\rho,i}|_Y = \text{val}_{\rho,i+1}|_Y \circ f^*.$$

By Lemma 3, we derive

$$(\text{val}_{\rho,i+1}^{\uparrow\alpha})|_X = ((\text{val}_{\rho,i}^{\uparrow\alpha})|_X \circ f)^{\uparrow\alpha} \quad \text{and} \quad (\text{val}_{\rho,i}^{\uparrow\alpha})|_Y = ((\text{val}_{\rho,i+1}^{\uparrow\alpha})|_Y \circ f^*)^{\uparrow\alpha}.$$

Thanks to this, using simple inductions on  $i$ , one can verify that  $A|_X = (\text{val}_{\rho,i}^{\uparrow\alpha})|_X$ , and symmetrically that  $A|_Y = (\text{val}_{\rho,i}^{\uparrow\alpha})|_Y$ .  $\blacktriangleleft$

► **Lemma 5.** *Let  $T$  be a trimmed flow-normalized SST. Given two transitions  $\tau_i : q \xrightarrow{a/f_i} q'$ , with  $i \in \{1, 2\}$ , a context  $(\rho, \sigma)$  for them, and the  $\alpha$ -approximants  $A = ((\text{val}_{\rho,\bullet})^{\uparrow\alpha})|_X$  and  $A' = ((\text{val}_{\sigma,\bullet})^{\uparrow\alpha})|_Y$ , with  $\alpha \in \mathbb{N}$ , the following holds:*

1. *If  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds on all valuations  $\nu \in A \uplus A'$ , then  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$ .*
2. *If  $\tau_1, \tau_2$  have the same flow,  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$ , and  $\alpha \geq c$ , where  $c$  is the capacity of  $T$ , then  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds on all valuations  $\nu \in A \uplus A'$ .*

**Proof.** The proof exploits the fact that, since  $T$  is flow-normalized, for every successful run  $\rho \tau_i \sigma$ , for both  $i = 1$  and  $i = 2$ , the substitution in  $\text{Effect}(f_i)$  of every variable  $x \in X$  (resp.  $y \in Y$ ) with the word  $\text{val}_{\rho,\bullet}(x)$  (resp.  $\text{val}_{\sigma,\bullet}(y)$ ) gives precisely the output  $\text{out}(\rho \tau_i \sigma)$ . This implies that

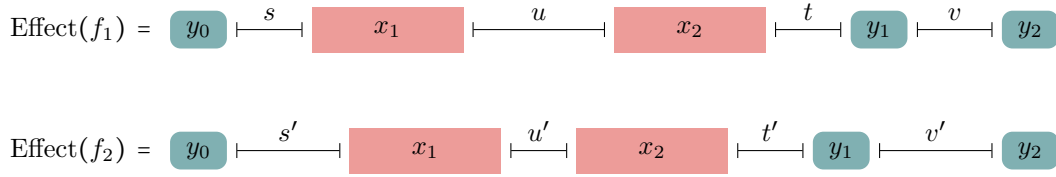
$$\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma) \quad \text{if and only if} \quad (\text{val}_{\rho,\bullet}) \uplus (\text{val}_{\sigma,\bullet}) \models \text{Effect}(f_1) = \text{Effect}(f_2) \quad (*)$$

We prove the first claim, which holds for any arbitrary parameter  $\alpha \in \mathbb{N}$ . By construction, we have  $\text{val}_{\rho,\bullet}(x) \in A(x)$ , for all  $x \in X$ , and  $\text{val}_{\sigma,\bullet}(y) \in A'(y)$ , for all  $y \in Y$ . The previous

## 121:18 Equivalence of finite-valued streaming string transducers is decidable

property  $(*)$  immediately implies that  $\tau_1$  and  $\tau_2$  produce the same output within the context  $(\rho, \sigma)$  if the equation  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds for all valuations  $\nu \in A \uplus A'$ .

To prove the second claim, we assume that  $\tau_1$  and  $\tau_2$  have the same flow, we fix a context  $(\rho, \sigma)$  for them, and we let  $A = (\text{val}_{\rho, \bullet})^{\uparrow \alpha}$  and  $A' = (\text{val}_{\bullet, \sigma})^{\uparrow \alpha}$  for some  $\alpha \geq c$ , where  $c$  is the capacity of  $T$ . We need to prove that the equation  $\text{Effect}(f_1) = \text{Effect}(f_2)$  holds for all valuations  $\nu \in A \uplus A'$ . We proceed by equating the two words  $\text{Effect}(f_1)$  and  $\text{Effect}(f_2)$ , and we study how the various blocks inside  $\text{Effect}(f_1)$  and  $\text{Effect}(f_2)$  (i.e. variables and maximal factors over  $\Gamma$ ) are aligned. The reader may refer to the figure below, which gives an example of possible alignments:



It is important to note that in any word  $\text{Effect}(f)$  the variables from  $X$  occur in the standard order  $x_1 < \dots < x_m$ , and similarly for the variables from  $Y$ . In general, it may happen that, due to updates that concatenate registers together, the  $x$ 's and the  $y$ 's are not strictly interleaved one with the other (as an example, see  $\text{Effect}(f_1)$  in the figure above). Here however, since  $\tau_1$  and  $\tau_2$  were assumed to have the same flow, we know that the interleaving of the  $x$ 's and the  $y$ 's is the same in  $\text{Effect}(f_1)$  and  $\text{Effect}(f_2)$ .

Of course, since  $\text{out}(\rho \tau_1 \sigma) = \text{out}(\rho \tau_2 \sigma)$  the occurrences of the first and the last variables,  $y_0$  and  $y_m$ , are aligned exactly, as they represent the same extremal gaps. For the remaining variables, which we generically denote  $z_1, \dots, z_\ell$ , we proceed by splitting the equation  $\text{Effect}(f_1) = \text{Effect}(f_2)$ , devoid of the extremal variables, into sub-equations that involve fewer variables, and reason by induction. Hereafter,  $L = R$  denotes an equation over the variables  $z_1, \dots, z_\ell$ , that occur exactly once on each side of the equation and with the same order. Moreover, the factors of  $L$  and  $R$  over  $\Gamma$  have length at most  $c$ , the capacity of the SST.

Suppose that  $L = s z_1 L'$  and  $R = s' z_1 R'$ , with  $s, s' \in \Gamma^*$  and  $L', R' \in \Gamma^* z_2 \Gamma^* \dots \Gamma^* z_\ell$ . Assume from now on that  $z_1$  is a register (gaps are treated symmetrically).

If  $s = s'$ , then the occurrences of  $z_1$  are perfectly aligned, and so  $L' = R'$ . We can then use induction. Note that in this case there is no restriction on  $z_1$ . In particular, the approximant  $A(z_1)$  for the valuation of  $z_1$  induced by the initial run  $\rho$  can well be  $\Gamma^*$ . Moreover, any solution of the equation  $L = R$ , where we replace the valuation for  $z_1$  by an arbitrary word from its approximant  $A(z_1)$ , is again a solution.

Otherwise, if  $s \neq s'$ , then either  $s$  is a prefix of  $s'$ , or the other way around. Suppose by symmetry that  $s' = s w$  for some  $w \in \Gamma^+$ , hence  $|w| \leq c$ . We get the equation  $w z_1 = z_1 w'$ , where  $w'$  is a conjugate of  $w$ , i.e.  $w = w_1 w_2$  and  $w' = w_2 w_1$  for some  $w_1, w_2 \in \Gamma^*$ . It follows that in every solution of  $L = R$ , the value of  $z_1$  must range over the periodic language  $w^* w_1$ . If we consider the register valuation  $\nu = \text{val}_{\rho, \bullet}$  induced by the initial run  $\rho$ , then we have  $\nu(z_1) \in w^* w_1$ .

Now, we assume without loss of generality that  $w$  is primitive. Since the length of  $w$  is at most  $c$  and since  $\alpha \geq c$ , we know that  $w^* w_1$  is an  $\alpha$ -approximant. Moreover, since  $\nu(z_1) \in w^* w_1$  and  $A = (\text{val}_{\rho, \bullet})^{\uparrow \alpha}$ , we get  $A(z_1) \subseteq w^* w_1$  (in particular,  $A(z_1)$  can be either a singleton or the periodic language  $w^* w_1$  itself). This implies that the equation  $s z_1 w' = s' z_1$  holds for any word from  $A(z_1)$ .

For the remaining variables, we observe that, up to any valuation that satisfies  $\text{Effect}(f_1) = \text{Effect}(f_2)$ , we get a new equation  $w' L' = R'$  in a fewer number of variables. From there by applying induction we get the desired claim.  $\blacktriangleleft$

► **Lemma 8.** *Given a flow-normalized SST  $T$  with transient initial states, and given a letter  $a \in \Sigma \setminus \{-\}$ , one can construct an SST  $T_a$  with transient initial states such that*

- $\text{Dom}(T_a) = \{u \in \Sigma^* : au \in \text{Dom}(T)\},$
- $T_a$  on input  $u$  produces the same outputs as  $T$  on input  $au$ .

Moreover,  $T_a$  has the same states and the same registers as  $T$ ; if  $T$  has capacity  $c$ , then  $T_a$  has capacity  $2c$ ; if  $T$  admits  $\alpha$ -approximants, then so does  $T_a$  (via the same annotation).

**Proof.** The construction is rather straightforward and boils down to shortcutting the first  $a$ -labeled transition in every successful run. Given  $T = (\Sigma, \Gamma, X, Q, U, I, E, F, x_1)$ , we define  $T_a = (\Sigma, \Gamma, X, Q, U', I, E', F, x_1)$ , where  $U' = U \cup U \circ U$ ,  $\circ$  denotes the functional composition, and  $E'$  contains the following transitions:

- $q_0 \xrightarrow{b/f \circ f'} q'$ , if  $E$  contains some transitions  $q_0 \xrightarrow{a/f} q \xrightarrow{b/f'} q'$ , with  $q_0$  initial state and  $b \in \Sigma$ ,
- $q \xrightarrow{b/f} q'$ , if  $E$  contains a transition  $q \xrightarrow{b/f} q'$ , with  $q$  is not initial and  $b \in \Sigma$ .

Note that, thanks to the assumption that every input of an SST ends with the special marker  $\neg$ , there is no final state in  $T$  that is a successor of an initial state along an  $a$ -labeled transition (unless of course  $a = \neg$ , which we assumed to be not the case). This essentially means that any initial  $a$ -labeled transition can be absorbed into the subsequent transitions, as precisely done in the above construction.

It is routine to check that  $T_a$  satisfies the desired claims. Here we only show that  $T_a$  admits the same  $\alpha$ -approximants as  $T$ . This follows from the fact that every successful run  $\rho$  of  $T_a$  of the form

$$\rho : (q_0, A_{0,X}, A_{0,Y}) \xrightarrow{T_a^{b/f}} (q_1, A_{1,X}, A_{1,Y}) \xrightarrow{T_a^{u/g}} (q_2, A_{2,X}, A_{2,Y})$$

can be turned to a successful run  $\rho'$  of  $T$  of the form

$$\rho' : (q_0, A_{0,X}) \xrightarrow{T^{a/f_a}} (q', A'_X, A'_Y) \xrightarrow{T^{b/f_b}} (q_1, A_{1,X}, A_{1,Y}) \xrightarrow{T^{u/g}} (q_2, A_{2,X}, A_{2,Y})$$

where  $f_a \circ f_b = f$ . In particular, we have  $\text{val}_{\rho,0} = \text{val}_{\rho',0}$  and  $\text{val}_{\rho,i} = \text{val}_{\rho',i+1}$  for all positions  $i > 0$ , and hence both  $T$  and  $T_a$  satisfy Equation (1).  $\blacktriangleleft$

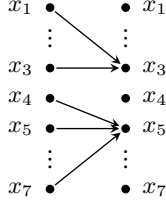
► **Lemma 10.** *If  $\rho = \rho_1 \gamma \rho_2$  is an initial run of  $T$ , with  $\gamma$  loop, then for every  $n > 0$  the pumped run  $\rho^{(n)} = \rho_1 \gamma^n \rho_2$  induces valuations  $\text{val}_{\rho^{(n)}} \bullet$  mapping any register  $x$  to a word of the form  $u_0 v_1^{n-1} u_1 \dots v_{2m}^{n-1} u_{2m}$ , where  $u_0, \dots, u_{2m}, v_1, \dots, v_{2m} \in \Gamma^*$  depend on  $\rho$  and  $x$ , but not on  $n$ . Moreover, we have  $v_i \neq \varepsilon$  for some  $i$  if there is a register  $x'$  that is productive along  $\gamma$  and that flows into  $x$  along  $\rho_2$ .*

**Proof.** We begin by observing a useful property of updates induced by loops:

▷ **Claim 17.** If  $\gamma$  is a loop, then the list of registers  $x_1, \dots, x_m$  can be partitioned into intervals  $X_1 < \dots < X_k$  such that for every  $1 \leq i \leq k$ , there is  $x' \in X_i$  so that every  $x \in X_i$  flows into  $x'$  along  $\gamma$ .

## 121:20 Equivalence of finite-valued streaming string transducers is decidable

Proof of claim. It suffices to verify that idempotent flows always have shapes similar to the flow below, where  $X_1 = \{x_1, x_2, x_3\}$  and  $X_2 = \{x_4, x_5, \dots, x_7\}$ :



◁

The next claim allows to simplify the statement of the lemma by assuming that  $\rho_2$  is empty. Its proof is straightforward, since  $T$  is non-erasing and non-permuting.

▷ **Claim 18.** If  $\rho = \rho_1 \rho_2$  is an initial run and  $x_1, \dots, x_k$  flow into  $x$  along  $\rho_2$ , then  $\text{val}_{\rho \bullet}(x)$  contains the factors  $\text{val}_{\rho_1 \bullet}(x_1), \dots, \text{val}_{\rho_1 \bullet}(x_k)$  in this precise order, possibly interleaved by other words that depend only on  $\rho_2$ . Moreover, if any of the  $x_i$ 's is productive along  $\rho_1$ , then so is  $x$  along  $\rho$ .

It now remains to prove that:

▷ **Claim 19.** If  $\rho = \rho_1 \gamma$  is an initial run, with  $\gamma$  loop, then for every  $n > 0$  the pumped run  $\rho^{(n)} = \rho_1 \gamma^n$  induces valuations  $\text{val}_{\rho^{(n)} \bullet}$  mapping any register  $x$  to a word of the form  $v_1^{n-1} u v_2^{n-1}$ , where  $u, v_1, v_2 \in \Gamma^*$  depend on  $\rho$  and  $x$ , but not on  $n$ . Moreover,  $v_1$  or  $v_2$  is non-empty if  $x$  is productive along  $\gamma$ .

We use claim 17 and for simplicity we work on the example provided there. Let us assume on the example that the updates are as follows (recall that  $\mathcal{T}$  is non-permuting):

- $f(x_3) = t_1 x_1 t_2 x_2 t_3 x_3 t_4$ ,
- $f(x_5) = t_5 x_4 t_6 x_5 t_7 x_6 t_8 x_7 t_9$ ,
- $f(x_j) = t'_j$ , for all remaining  $j \neq 3, 5$ .

Let  $\nu = \text{val}_{\rho_1 \bullet}$  be the register valuation induced by the prefix  $\rho_1$ . The valuation  $\text{val}_{\rho^{(n)} \bullet} = \nu \circ f^n$  maps e.g.

- $x_3$  to  $(t_1 t'_1 t_2 t'_2 t_3)^{n-1} (t_1 \nu(x_1) t_2 \nu(x_2) t_3) \nu(x_3) (t_4)^n$ ,
- $x_5$  to  $(t_5 t'_4 t_6)^{n-1} (t_5 \nu(x_4) t_6) \nu(x_5) (t_7 \nu(x_6) t_8 \nu(x_7) t_9) (t_7 t'_6 t_8 t'_7 t_9)^{n-1}$ .

The claim is satisfied e.g. for  $x = x_3$  by setting  $v_1 = t_1 t'_1 t_2 t'_2 t_3$ ,  $u = (t_1 \nu(x_1) t_2 \nu(x_2) t_3) \nu(x_3) t_4$ , and  $v_2 = t_4$ . ◀

► **Lemma 11.** Let  $\rho$  be an initial run and  $x$  a register. If  $\text{val}_{\rho \bullet}(x)$  has length (resp. period) larger than  $\alpha = mc|Q|2^{3 \cdot 2^m}$ , then for every  $\beta \geq \alpha$  and for all but finitely many  $\rho' \supseteq \rho$ ,  $\text{val}_{\rho' \bullet}(x)$  has length (resp. period) larger than  $\beta$ .

**Proof.** The first step consists in identifying the appropriate loops  $\gamma_1, \dots, \gamma_\ell$  inside the initial run  $\rho$ . More precisely, we need to factorize  $\rho$  as

$$\rho = \rho_0 \gamma_1 \rho_1 \dots \gamma_\ell \rho_\ell.$$

where  $\gamma_1, \dots, \gamma_\ell$  are loops, in such a way that every register with large enough induced valuation is productive along at least one loop. In fact, for technical reasons related to periodicity, we need to also guarantee that the selected loops only contribute for a bounded portion to the valuation of a register, precisely, with at most  $c|Q|2^{3 \cdot 2^m}$  letters.

For every register  $z$  and every position  $i$  of  $\rho$ , let  $X_{i,z}$  be the set of registers that flow into  $z$  along the suffix of  $\rho$  that starts at position  $i$ . Further let  $N_{i,z} = \sum_{x \in X_{i,z}} |\text{val}_{\rho,i}(x)|$ , and

let  $D_{i,j,z} = N_{j,z} - N_{i,z}$  for all  $i \leq j$ . To find a productive loop for  $z$  between positions  $i \leq j$ , it suffices to have a large enough value  $D_{i,j,z}$ :

▷ **Claim 20.** If  $D_{i,j,z} > c|Q|2^{3 \cdot 2^m}$ , then  $\rho$  contains a loop  $\gamma$  between positions  $i$  and  $j$ , and there is a register  $x$  that is productive along  $\gamma$  and that flows into  $z$  along the suffix of  $\rho$  that follows  $\gamma$ .

*Proof of claim.* Let  $\sigma$  be the factor of  $\rho$  between positions  $i$  and  $j$ . Since  $T$  is copyless with capacity  $c$  and  $D_{i,j,z} > c|Q|2^{3 \cdot 2^m}$ , there are  $N > |Q|2^{3 \cdot 2^m}$  transitions between  $i$  and  $j$  along which some register in  $X_{k,z}$  is productive. Among these transitions, there are  $n > 2^{3 \cdot 2^m}$  that start with the same source state, say  $q$ . Let  $i_1 < \dots < i_n$  be the positions where the latter transitions start.

Next, consider the flows  $F_j$  of the updates induced between positions  $i_j$  and  $i_{j+1}$ , for all  $j = 1, \dots, n$ . Recall that flows are naturally equipped with an associative product, forming a monoid  $M$  of size at most  $m^2$ . By the Factorization Forest theorem [26, 9, 22], there is a factorization tree for the sequence  $F_1 \dots F_n$  that has height at most  $3|M|$  and such that every inner node with more than two successors has all children labeled by the same idempotent flow.

Since  $n > 2^{3 \cdot 2^m} \geq 2^{3 \cdot |M|}$ , there is at least one idempotent flow  $F_j$ . This proves that  $\rho$  contains a loop  $\gamma$  between positions  $i$  and  $j$ . Moreover, there is a register  $x$  that is productive along  $\gamma$  and that flows into  $z$  along the suffix of  $\rho$  that follows  $\gamma$ . ◁

We construct the desired factorization of  $\rho$  by induction as follows. We maintain a position  $i$  in  $\rho$ , representing the endpoint of the processed prefix of  $\rho$ , and a set  $Z$  of registers for which we still need to find corresponding productive loops. The position  $i$  is initialized to 0, and the set  $Z$  to the set of registers  $z$  such that  $|\text{val}_{\rho_\bullet}(z)| > mc|Q|2^{3 \cdot 2^m}$ . We then look at the first position  $j > i$  such that  $D_{i,j,z} > c|Q|2^{3 \cdot 2^m}$ , for some  $z \in Z$  (the construction terminates as soon as  $Z$  becomes empty). By Claim 20, we know that the factor of  $\rho$  between positions  $i$  and  $j$  contains a loop  $\gamma$ , and there is a register  $x$  that is productive along  $\gamma$  and flows into  $z$  along the suffix that follows  $\gamma$ . Moreover, thanks to the above eager strategy, the number of output letters that appear inside  $g(x)$ , where  $g$  is the update induced by  $\gamma$ , is at most  $c|Q|2^{3 \cdot 2^m}$ . We can thus declare  $\gamma$  to be one of the loops of our factorization, and accordingly set  $i$  to  $j$  and remove  $z$  from  $Z$ . Note that the following invariant is preserved: for all  $z \in Z$ ,  $|\text{val}_{\rho_\bullet}(z)| > mc|Q|2^{3 \cdot 2^m} - D_{0,i,z}$ . Because at each iteration the value of  $D_{0,i,z}$  increases by at most  $c|Q|2^{3 \cdot 2^m}$ , and because at most  $m$  iterations are possible, this shows that the construction can be carried over correctly.

We are now ready to prove the lemma. For the property concerning the lengths of the register valuations, suppose that  $|\text{val}_{\rho_\bullet}(x)| > \alpha = mc|Q|2^{3 \cdot 2^m}$ . By the previous constructions, there is a loop  $\gamma_i$  with a productive register  $x'$  that flows into  $x$  along the suffix  $\rho_i \gamma_{i+1} \dots \gamma_\ell \rho_\ell$ . By Lemma 10, the valuations induced at the end of the pumped runs

$$\rho^{(n)} = \rho_0 \gamma_1 \rho_1 \dots \gamma_i^n \rho_i \dots \gamma_\ell \rho_\ell$$

map  $x$  to arbitrarily long words. Moreover, the same can be said of the lengths of the valuations of  $x$  that are induced by runs obtained by pumping *simultaneously*, and by the same amount  $n$ , all loops  $\gamma_1, \dots, \gamma_\ell$ . This proves that, for every  $\beta \geq \alpha$  and for all but finitely many  $\rho' \sqsupseteq \rho$ ,  $|\text{val}_{\rho'_\bullet}(x)| > \beta$ .

We can use a similar argument to prove the property concerning the periods. Suppose that  $\text{val}_{\rho_\bullet}(x)$  has period  $p > \alpha = mc|Q|2^{3 \cdot 2^m}$ . In particular,  $|\text{val}_{\rho_\bullet}(x)| > \alpha$ . As before, there is a loop  $\gamma_i$  with a productive register  $x'$  that flows into  $x$  along the suffix  $\rho_i \gamma_{i+1} \dots \gamma_\ell \rho_\ell$ .

## 121:22 Equivalence of finite-valued streaming string transducers is decidable

Moreover, by the previous constructions we know that the effect of the loop  $\gamma_i$  on the final valuation of  $x$  is to add at most  $c|Q|2^{3 \cdot 2^m}$  letters. Let us consider runs that are obtained by pumping simultaneously all loops  $\gamma_1, \dots, \gamma_\ell$  inside  $\rho$ :

$$\rho^{(n)} = \rho_0 \gamma_1^n \rho_1 \dots \gamma_i^n \rho_i \dots \gamma_\ell^n \rho_\ell.$$

By Lemma 10 (plus Claim 18), the valuations induced at the end of the pumped runs  $\rho^{(n)}$  map  $x$  to words of the form

$$\text{val}_{\rho^{(n)} \bullet}(x) = u_0 v_1^{n-1} u_1 \dots u_{t-1} v_t^{n-1} u_t.$$

for some  $t \leq 2m\ell$ , where  $u_0, \dots, u_t, v_1, \dots, v_t \in \Gamma^*$  depend only on  $\rho$  and  $\bar{\gamma}$ ,  $|v_i| \leq \alpha$  for all  $i \leq t$ , and  $|v_i| > 0$  for some  $i \leq t$ . In particular, the above words contain arbitrarily long repetitions of non-empty words.

Now, let  $p_n$  be the period of  $\text{val}_{\rho^{(n)} \bullet}(x)$ , for all  $n > 0$ . Recall that  $p_1 = p > \alpha$ . We aim at showing that the periods  $p_n$  get arbitrarily large. Suppose, by way of contradiction, that  $p_n$  is uniformly bounded for all  $n > 0$ . Then  $p_n$  must be a constant, say  $p_n = p'$ , for infinitely many  $n$ . We also recall from the previous arguments that  $\text{val}_{\rho^{(n)} \bullet}(x)$  has arbitrarily long repetitions of words of length  $r_1 = |v_1|, \dots, r_k = |v_k|$ , with all  $r_j \leq \alpha$  and at least one  $r_j > 0$ . By Fine-Wilf's theorem, this implies that the period of  $\text{val}_{\rho^{(n)} \bullet}(x)$ , for infinitely many  $n$ , is

$$p'' = \gcd\{p', r_i\}_{r_i > 0} < \alpha < p.$$

We can transfer this property to the original word  $\text{val}_{\rho \bullet}(x)$ , by observing that  $\text{val}_{\rho \bullet}(x)$  can be obtained from any of the previous words  $\text{val}_{\rho^{(n)} \bullet}(x)$  by removing some occurrences of factors of lengths  $r_1, \dots, r_k$ . As those lengths are multiples of the period  $p''$ , the latter operation does not change the period of the entire word. Hence,  $\text{val}_{\rho \bullet}(x)$  must also have period  $p'' < p$ , which is however a contradiction.

This proves that  $p_n$  gets arbitrarily large for  $n > 0$ . In particular, for every  $\beta \geq \alpha$  and for all but finitely many runs  $\rho' \sqsupseteq \rho$ , the word  $\text{val}_{\rho' \bullet}(x)$  has period larger than  $\beta$ . ◀

► **Proposition 12.** *Let  $T'$  be the SST admitting  $\alpha$ -approximants that is obtained from  $T$  using Proposition 4, for any  $\alpha \geq mc|Q|2^{3 \cdot 2^m}$ , where  $Q$  is the set of states of  $T$ . The  $\alpha$ -approximants of  $T'$  are tight.*

**Proof.** As usual, by symmetry we can focus only on register valuations induced by initial runs. We fix, once and for all, two parameters  $\alpha, \beta$ , with  $\alpha \geq mc|Q|2^{3 \cdot 2^m}$  and  $\beta \geq \alpha$ . For the sake of readability, we also introduce the shorthands  $A_\rho = (\text{val}_{\rho \bullet})^{\uparrow \alpha}$  and  $B_\rho = (\text{val}_{\rho \bullet})^{\uparrow \beta}$ , for any initial run  $\rho$  of  $T'$ . Since  $\beta \geq \alpha$ , we have  $B_\rho(x) \subseteq A_\rho(x)$ . We need to prove that there is an initial run  $\rho'$  of  $T'$  such that, for all registers  $x$ ,  $B_{\rho'}(x) = A_{\rho'}(x)$ .

We will in fact prove a slightly stronger claim, that is: for all initial runs  $\rho$  or  $T'$ , for all but finitely many runs  $\rho' \sqsupseteq \rho$ , and for all registers  $x$ ,  $B_{\rho'}(x) = A_{\rho'}(x)$ . Towards this we analyse the possible cases when  $B_\rho(x)$  could be strictly contained in  $A_\rho(x)$ , for any initial run  $\rho'$ . By definition of  $\beta$ -approximant, this could only happen when  $B_\rho(x)$  contains only one word, or when it is a language of the form  $u^*v$ . In the former case we say for short that  $B_{\rho'}(x)$  is a *singleton*; in the latter case we say that  $B_{\rho'}$  is a *periodic language*. We then prove that, for every register  $x$  and every initial run  $\rho$  of  $T'$ :

- if  $B_\rho(x)$  is a singleton strictly included in  $A_\rho(x)$ , then, for all but finitely many runs  $\rho' \sqsupseteq \rho$ ,  $B_{\rho'}(x)$  is not a singleton;
- if  $B_\rho(x)$  is a periodic language strictly included in  $A_\rho(x)$ , then, for all but finitely many runs  $\rho' \sqsupseteq \rho$ ,  $B_{\rho'}(x)$  is not a periodic language (and thus neither a singleton).



Note that the quantification “for all but finitely many”, like universal quantification, commutes with the conjunction over the registers  $x$ . Therefore, the above two properties, paired with the previous arguments, suffice to prove the desired claim.

Now, fix a register  $x \in X$  and an initial run  $\rho$  of  $T'$ , and suppose that  $B_\rho(x)$  is a singleton or a periodic language strictly contained in  $A_\rho(x)$ .

If  $B_\rho(x)$  is a singleton, say  $B_\rho(x) = \{u\}$ , then, since  $u^{\uparrow\alpha} = A_\rho(x) \supsetneq \{u\}$ , we know that  $|u| > \alpha$ . Recall that  $T'$  is a covering of  $T$ , and in particular that the (initial) runs of  $T'$  are bijectively related to the (initial) runs of  $T$ . Let  $\tilde{\rho}$  be the initial run of  $T$  that corresponds to  $\rho$ . By Lemma 11, we get that, for all but finitely many runs  $\tilde{\rho}' \succeq \tilde{\rho}$  of  $T$ ,  $\text{val}_{\tilde{\rho}'\bullet}(x)$  has length even larger than  $\beta$ . By exploiting again the bijection between runs of  $T$  and runs of  $T'$ , we get that, for all but finitely many runs  $\rho' \succeq \rho$  of  $T'$ , the word  $\text{val}_{\rho'\bullet}(x)$  has length larger than  $\beta$ , and hence  $B_{\rho'}(x)$  cannot be a singleton.

If  $B$  is a periodic language of the form  $u^*v$ , then we get  $|u| > \alpha$ , and hence the period of  $\text{val}_{\rho'\bullet}(x)$  is larger than  $\alpha$ . Using the correspondence between runs of  $T$  and runs of  $T'$  and exploiting Lemma 11, exactly as we did before, we get that, for all but finitely many runs  $\rho' \succeq \rho$  of  $T'$ ,  $\text{val}_{\rho'\bullet}(x)$  has period larger than  $\beta$ , and hence  $B_{\rho'}$  is not a periodic language. ◀