

Equivalence of Extended Symbolic Finite Transducers ^{*}

Loris D’Antoni¹ and Margus Veanes²

¹ University of Pennsylvania
lorisdan@cis.upenn.edu

² Microsoft Research
margus@microsoft.com

Abstract. Symbolic Finite Transducers augment classic transducers with symbolic alphabets represented as parametric theories. Such extension enables succinctness and the use of potentially infinite alphabets while preserving closure and decidability properties. Extended Symbolic Finite Transducers further extend these objects by allowing transitions to read consecutive input elements in a single step. While when the alphabet is finite this extension does not add expressiveness, it does so when the alphabet is symbolic. We show how such increase in expressiveness causes decision problems such as equivalence to become undecidable and closure properties such as composition to stop holding. We also investigate how the automata counterpart, Extended Symbolic Finite Automata, differs from Symbolic Finite Automata. We then introduce the subclass of Cartesian Extended Symbolic Finite Transducers in which guards are limited to conjunctions of unary predicates. Our main result is an equivalence algorithm for such subclass in the single-valued case. Finally, we model real world problems with Cartesian Extended Symbolic Finite Transducers and use the equivalence algorithm to prove their correctness.

1 Introduction

Finite automata have proven to be an effective tool in a wide range of applications, from regular expressions to network packet inspection [15]. Finite transducers extend finite automata with outputs and can model functions from strings to strings such as natural language transformations [12]. Due to their closure and decidability properties, these models are widely used in practice but they have three major disadvantages: 1) their number of transitions usually “blows up” when dealing with large alphabets; 2) they cannot model infinite alphabets; and 3) transitions cannot express relations between adjacent input symbols.

Symbolic Finite Automata/Transducers [16] or SFAs/SFTs respectively, are an extension of traditional automata and transducers that attempts to solve

^{*} Loris D’Antoni’s research was partially supported by NSF Expeditions in Computing award CCF 1138996.

problems 1 and 2 above by allowing transitions to be labelled with arbitrary predicates in a specified theory. When such theory is decidable SFAs and SFTs enjoy the same properties of finite automata and transducers, such as closure under composition and decidability of equivalence (for single-valued SFTs). In [16], Symbolic Transducers or STs (SFTs with registers) are proposed in order to cope with the third problem above. STs are however undecidable with respect to most analysis problems, even emptiness.

In our previous work on the topic of analysis of string coders [4], we introduce *Extended Symbolic Finite Automata/Transducers* or ESFAs/ESFTs, that add finite lookahead to SFAs/SFTs. This extension allows to read multiple input symbols in a single transition and combine their values in the output. ESFTs can be viewed as a subclass of STs with a restricted use of registers that mimic “look-behind”. This view is used in [4] to map ESFTs directly to STs in order to overcome the problem that ESFTs are not closed under composition. In other words, it addresses the composition problem by first converting ESFTs to STs, then composing the STs, and finally converting the result back into an ESFT using a semi-decision procedure. The formal properties of ESFTs have not been fully understood yet. From the point of view of analysis, the key operations that are desired are *composition* and *equivalence* (for single-valued ESFTs). Then, for example, the functional correctness of a string (*encoder, decoder*) pair (E, D) (for example UTF8 to UTF16 encoding) can be decided by checking the equivalence of $\lambda x.D(E(x))$ with $\lambda x.x$. Other properties, such as commutativity and idempotence, also depend on composition and equivalence.

The topic that is left open in [4] is decidability of equivalence checking of ESFTs. Our main theoretical contribution in this paper is a complete classification, in terms of guard complexity and lookahead, of the cases in which the equivalence problem is decidable for ESFTs. We first show that one-equality or equivalence in the single-valued case is in general undecidable, contrasting the finite alphabet setting where lookahead does not matter [18, Theorem 2.17]. We then introduce the notion of Cartesian ESFT, in which transition guards are constrained to be conjunctions of unary predicates, and show that one-equality and equivalence are decidable for single-valued Cartesian ESFTs. This is a proper extension of the decidability result of one-equality of SFTs [16]. The key tool that we need to prove the result is Lemma 2.

We also analyze basic properties of ESFAs and show how they differ from SFAs. We prove ESFAs to be not closed under intersection and show that equivalence and universality of ESFAs are both undecidable problems.

Applications. We present four applications of our models in different areas. We first extend the result of [4] by proving the correctness of four real world string encoders. The new equivalence algorithm is a full decision procedure for the Cartesian case, unlike the semi-decision procedure in [4] that may fail to terminate in some incorrect instances. Our second and third applications are in the context of networking and present new classes of programs that can be modelled as ESFAs/ESFTs. We show how 1) ESFAs can be used for the task of deep-packet inspection, and 2) ESFTs can succinctly represent transformations

between headers of different network protocols. Our fourth case study shows the use of additional theories for the analysis of list manipulating programs.

Contributions. In summary, we offer the following contributions:

- we study the closure and decidability properties of ESFAs (Section 3.1);
- we study the equivalence problem for ESFTs (Section 3.2):
 - we prove the equivalence of single-valued ESFTs to be undecidable;
 - we present a novel algorithm for the equivalence of single-valued Cartesian ESFTs;
- we extend the negative result on ESFTs composition presented in [4] (Section 3.3); and
- we analyze the performance of the equivalence algorithm for “Cartesian” ESFTs on real examples and propose new applications for ESFAs and ESFTs (Section 4).

We finally summarize previous work and conclude (Section 5 and 6).

2 Extended Symbolic Finite Transducers

We assume a recursively enumerable (r.e.) *background universe* \mathcal{U} with built-in function and relation symbols. Definitions below are given with \mathcal{U} as an implicit parameter. We use λ -expressions for representing anonymous functions that we call λ -terms. A Boolean λ -term $\lambda x.\varphi(x)$, where x is a variable of type σ is called a σ -predicate. Our notational conventions are consistent with the definition of symbolic transducers [16]. The universe is multi-typed with \mathcal{U}^τ denoting the sub-universe of elements of type τ . We write Σ for \mathcal{U}^σ and Γ for \mathcal{U}^γ .

A *label theory* is given by a recursively enumerable set Ψ of formulas that is closed under Boolean operations, substitution, equality and if-then-else terms. A label theory Ψ is *decidable* when satisfiability for $\varphi \in \Psi$, $IsSat(\varphi)$, is decidable.

For σ -predicates φ , we assume an effective *witness* function \mathcal{W} such that, if $IsSat(\varphi)$ then $\mathcal{W}(\varphi) \in \llbracket \varphi \rrbracket$, where $\llbracket \varphi \rrbracket \subseteq \mathcal{U}^\sigma$ is the set of all values that satisfy φ ; φ is *valid*, $IsValid(\varphi)$, when $\llbracket \varphi \rrbracket = \mathcal{U}^\sigma$.

We are studying in this paper an extension of SFTs with *lookahead*, called *extended* SFTs or *ESFTs*. Originally, ESFTs were introduced in [4] for the purposes of analyzing string encoders and decoders, where a semi-decision procedure was provided for converting STs (SFTs with registers) into ESFTs.

Definition 1. An *Extended Symbolic Finite Transducer (ESFT)* with *input type* σ and *output type* γ is a tuple $A = (Q, q^0, R)$,

- Q is a finite set of *states*;
- $q^0 \in Q$ is the *initial state*;
- R is a finite set of *rules*, $R = \Delta \cup F$, where
- Δ is a set of *transitions* $r = (p, \ell, \varphi, f, q)$, denoted $p \xrightarrow[\ell]{\varphi/f} q$, where $p \in Q$ is the *start state* of r ;

- $\ell \geq 1$ is the *lookahead* of r ;
 - φ , the *guard* of r , is a σ^ℓ -predicate;
 - f , the *output* of r , is a $(\sigma^\ell \rightarrow \gamma)$ -sequence;
 - $q \in Q$ is the *continuation* state of r .
- F is a set of *finalizers* $r = (p, \ell, \varphi, f)$, denoted $p \xrightarrow[\ell]{\varphi/f} \bullet$, with components as above and where ℓ may be 0.

The *lookahead* of A is the maximum of all lookaheads of rules in R . An ESFT all of whose rules have output \square is an *Extended Symbolic Finite Automaton (ESFA)*.

A finalizer is a rule without a continuation state. A finalizer with lookahead ℓ is used when the end of the input sequence has been reached with *exactly* ℓ input elements remaining. A finalizer is a generalization of a final state. In a classical setting, finalizers can be avoided by adding a new symbol to the alphabet that is only used to mark the end of the input. In the presence of arbitrary input types, this is not always possible without affecting the theory, e.g., when the input type is \mathbb{Z} then that symbol would have to be outside \mathbb{Z} .

In the sequel let $A = (Q, q^0, R)$, $R = \Delta \cup F$, be a fixed ESFT with input type σ and output type γ . The semantics of rules in R is as follows:

$$\llbracket p \xrightarrow[\ell]{\varphi/f} q \rrbracket \stackrel{\text{def}}{=} \{ p \xrightarrow{[a_0, \dots, a_{\ell-1}]/\llbracket f \rrbracket(a_0, \dots, a_{\ell-1})} q \mid (a_0, \dots, a_{\ell-1}) \in \llbracket \varphi \rrbracket \}$$

We write $s_1 \cdot s_2$ for the concatenation of two sequences s_1 and s_2 .

Definition 2. For $u \in \Sigma^*$, $v \in \Gamma^*$, $q \in Q$, $q' \in Q \cup \{\bullet\}$, define $q \xrightarrow{u/v}_A q'$ as follows: there exists $n \geq 0$ and $\{p_i \xrightarrow{u_i/v_i} p_{i+1} \mid i \leq n\} \subseteq \llbracket R \rrbracket$ such that

$$u = u_0 \cdot u_1 \cdots u_n, \quad v = v_0 \cdot v_1 \cdots v_n, \quad q = p_0, \quad q' = p_{n+1}.$$

Let also $q \xrightarrow{\square/\square}_A q$ for all $q \in Q_A$.

Definition 3. The *transduction* of A , $\mathcal{T}_A(u) \stackrel{\text{def}}{=} \{v \mid q^0 \xrightarrow{u/v}_A \bullet\}$.

The following example represents typical (realistic) ESFTs over a label theory of linear modular arithmetic. We use the following abbreviated notation for rules, by omitting explicit λ 's. We write

$$p \xrightarrow[\ell]{\varphi(\bar{x})/[f_1(\bar{x}), \dots, f_k(\bar{x})]} q \quad \text{for} \quad p \xrightarrow[\ell]{\lambda \bar{x} \cdot \varphi(\bar{x}) / \lambda \bar{x} \cdot [f_1(\bar{x}), \dots, f_k(\bar{x})]} q,$$

where φ and f_i are terms whose free variables are among $\bar{x} = (x_0, \dots, x_{\ell-1})$.

Example 1. The example illustrates the standard encoding BASE64, that is used to transfer binary data in textual format, e.g., in emails via the protocol MIME. The digits of the encoding are chosen in the safe ASCII range of characters that remain unmodified during transport over textual media. Assume that the input

type and the output type are both BYTE, that is the set of integers between 0 and 255. *Base64encode* is an ESFT with one state and four rules:

$$\begin{aligned}
& p \xrightarrow{\frac{\text{true}/[\ulcorner b_2^7(x_0) \urcorner, \ulcorner (b_0^1(x_0) \ll 4) | b_4^7(x_1) \urcorner, \ulcorner (b_0^3(x_1) \ll 2) | b_6^7(x_2) \urcorner, \ulcorner b_0^5(x_2) \urcorner]}{3}} p \\
& p \xrightarrow{\frac{\text{true}/[]}{0}} \bullet \quad p \xrightarrow{\frac{\text{true}/[\ulcorner b_2^7(x_0) \urcorner, \ulcorner b_0^1(x_0) \ll 4 \urcorner, \ulcorner = \urcorner, \ulcorner = \urcorner]}{1}} \bullet \\
& p \xrightarrow{\frac{\text{true}/[\ulcorner b_2^7(x_0) \urcorner, \ulcorner (b_0^1(x_0) \ll 4) | b_4^7(x_1) \urcorner, \ulcorner b_0^3(x_1) \ll 2 \urcorner, \ulcorner = \urcorner]}{2}} \bullet
\end{aligned}$$

where $b_n^m(x)$ extracts bits m through n from x , e.g., $b_2^3(13) = 3$, $x|y$ is bitwise OR of x and y , $x \ll k$ is x shifted left by k bits, and $\ulcorner x \urcorner$ is the mapping

$$\ulcorner x \urcorner \stackrel{\text{def}}{=} (x \leq 25 ? x + 65 : (x \leq 51 ? x + 71 : (x \leq 61 ? x - 4 : (x = 62 ? '+' : '/'))))$$

of values between 0 and 63 into a standardized sequence of safe ASCII character codes. The last two finalizers correspond to the cases when the length of the input sequence is not a multiple of three. Observe that the length of the output sequence is always a multiple of four. The character '=' (61 in ASCII) is used as a padding character and it is not a BASE64 digit. i.e., '=' is not in the range of $\ulcorner x \urcorner$.

Base64decode is an ESFT that decodes a BASE64 encoded sequence back into the original byte sequence. *Base64decode* has also one state and four rules:

$$\begin{aligned}
& q \xrightarrow{\frac{\bigwedge_{i=0}^3 \beta_{64}(x_i) / [(\lceil x_0 \rceil \ll 2) | b_4^5(\lfloor x_1 \rfloor), (b_0^3(\lfloor x_1 \rfloor) \ll 4) | b_2^5(\lfloor x_2 \rfloor), (b_0^1(\lfloor x_2 \rfloor) \ll 6) | \lfloor x_3 \rfloor]}{4}} q \\
& q \xrightarrow{\frac{\text{true}/[]}{0}} \bullet \quad q \xrightarrow{\frac{\beta_{64}(x_0) \wedge \beta'_{64}(x_1) \wedge x_2 = '=' \wedge x_3 = '=' / [(\lceil x_0 \rceil \ll 2) | b_4^5(\lfloor x_1 \rfloor)]}{4}} \bullet \\
& q \xrightarrow{\frac{\beta_{64}(x_0) \wedge \beta_{64}(x_1) \wedge \beta''_{64}(x_2) \wedge x_3 = '=' / [(\lceil x_0 \rceil \ll 2) | b_4^5(\lfloor x_1 \rfloor), (b_0^3(\lfloor x_1 \rfloor) \ll 4) | b_2^5(\lfloor x_2 \rfloor)]}{4}} \bullet
\end{aligned}$$

The function $\lfloor y \rfloor$ is the inverse of $\ulcorner x \urcorner$, i.e., $\lfloor \ulcorner x \urcorner \rfloor = x$, for $0 \leq x \leq 63$. The predicate $\beta_{64}(y)$ is true iff y is a valid BASE64 digit, i.e., $y = \ulcorner x \urcorner$ for some x , $0 \leq x \leq 63$. The predicates $\beta'_{64}(y)$ and $\beta''_{64}(y)$ are restricted versions of $\beta_{64}(y)$. Unlike *Base64encode*, *Base64decode* does not accept all input sequences of bytes, and sequences that do not correspond to any encoding are rejected.³ \boxtimes

The following subclass of ESFTs captures transductions that behave as partial functions from Σ^* to Γ^* .

Definition 4. A function $\mathbf{f} : X \rightarrow 2^Y$ is *single-valued* if $|\mathbf{f}(x)| \leq 1$ for all $x \in X$. An ESFT A is *single-valued* if \mathcal{T}_A is single-valued.

A sufficient condition for single-valuedness is determinism. We define $\varphi \wedge \psi$, where φ is a σ^m -predicate and ψ a σ^n -predicate, as the $\sigma^{\max(m,n)}$ -predicate $\lambda(x_1, \dots, x_{\max(m,n)}). \varphi(x_1, \dots, x_m) \wedge \psi(x_1, \dots, x_n)$. We define *equivalence of f and g modulo φ* , $f \equiv_{\varphi} g$, as: $IsValid(\lambda \bar{x}. (\varphi(\bar{x}) \Rightarrow f(\bar{x}) = g(\bar{x})))$.

Definition 5. A is *deterministic* if for all $p \xrightarrow{\varphi/f} q, p \xrightarrow{\varphi'/f'} q' \in R$:

³ For more information see <http://www.rise4fun.com/Bek/tutorial/base64>.

- (a) Assume $q, q' \in Q$. If $IsSat(\varphi \wedge \varphi')$ then $q = q'$, $\ell = \ell'$ and $f \equiv_{\varphi \wedge \varphi'} f'$.
- (b) Assume $q = q' = \bullet$. If $IsSat(\varphi \wedge \varphi')$ and $\ell = \ell'$ then $f \equiv_{\varphi \wedge \varphi'} f'$.
- (c) Assume $q \in Q$ and $q' = \bullet$. If $IsSat(\varphi \wedge \varphi')$ then $\ell > \ell'$.

Intuitively, determinism means that no two rules may overlap. It follows from the definitions that if A is deterministic then A is single-valued. Both ESFTs in Example 1 are deterministic.

The *domain* of a function $\mathbf{f} : X \rightarrow 2^Y$ is $\mathcal{D}(\mathbf{f}) \stackrel{\text{def}}{=} \{x \in X \mid \mathbf{f}(x) \neq \emptyset\}$ and for an ESFT A , $\mathcal{D}(A) \stackrel{\text{def}}{=} \mathcal{D}(\mathcal{T}_A)$. When A is single-valued, and $u \in \mathcal{D}(A)$, we treat A as a partial function from Σ^* to Γ^* and write $A(u)$ for the value v such that $\mathcal{T}_A(u) = \{v\}$. For example, $Base64encode("Foo") = "Rm9v"$ and $Base64decode("QmFy") = "Bar"$.

Cartesian ESFTs. We introduce a subclass of ESFTs that plays an important role in this paper. A binary relation R over X is *Cartesian over X* if R is the Cartesian product $R_1 \times R_2$ of some $R_1, R_2 \subseteq X$. The definition is lifted to n -ary relations and σ^n -predicates for $n \geq 2$ in the obvious way. In order to decide if a satisfiable σ^n -predicate φ is Cartesian over σ , let $(a_0, \dots, a_{n-1}) = \mathcal{W}(\varphi)$ and perform the following validity check:

$$IsCartesian(\varphi) \stackrel{\text{def}}{=} \forall \bar{x} (\varphi(\bar{x}) \Leftrightarrow \bigwedge_{i < n} \varphi(a_0, \dots, a_{i-1}, x_i, a_{i+1}, \dots, a_{n-1}))$$

In other words, a σ^n -predicate φ is Cartesian over σ if φ can be rewritten equivalently as a conjunction of n independent σ -predicates.

Definition 6. An ESFT (ESFA) is *Cartesian* if all its guards are Cartesian.

Both ESFTs in Example 1 are Cartesian. $Base64encode$ trivially so, while the guards of all rules of $Base64decode$ are conjunctions of independent unary predicates. In contrast, a predicate such as $\lambda(x_0, x_1).x_0 = x_1$ is not Cartesian.

Note that $IsCartesian(\varphi)$ is decidable by using the decision procedure of the label theory. Namely, decide unsatisfiability of $\neg IsCartesian(\varphi)$.

3 ESFAs and ESFTs Properties

We prove some basic properties of ESFAs and ESFTs and show how they drastically differ from SFAs and SFTs. First, we investigate basic ESFAs properties. Secondly, we prove the undecidability of ESFT equivalence and propose a new one-equality algorithm for the subclass of Cartesian ESFTs. Finally, we present some preliminary results on ESFT composition.

3.1 ESFAs Properties

In this section, we analyse closure and decidability properties of Extended Symbolic Finite Automata. We show how ESFAs have properties similar to those of context free grammars rather than regular languages. First, we show how checking the emptiness of the intersection of two ESFA definable languages is an undecidable problem.

Theorem 1 (Domain Intersection). *Given two ESFAs A and B with lookahead 2 over quantifier free successor arithmetic and tuples, checking whether there exists an input accepted by both A and B is undecidable.*

While checking the emptiness of an ESFA is a decidable problem, it is not possible to decide whether an ESFA accepts every possible input. It follows that equivalence is also undecidable.

Theorem 2 (Emptiness, Universality and Equivalence). *Given an ESFA A and B over σ checking whether A does not accept any input is decidable while, checking whether A accepts all the sequences in σ^* or whether A and B accept the same language are both undecidable problems.*

Combining Theorems 1 and 2 with a simple construction we obtain the following closure properties.

Theorem 3 (Closure Properties). *ESFAs are closed under union, but they are not closed under complement and intersection.*

Finally, Cartesian ESFAs capture exactly the class of SFA definable languages.

Theorem 4 (Cartesian ESFA iff SFA). *SFAs and Cartesian ESFAs are equivalent in expressiveness.*

From Theorem 4 we have that Cartesian ESFAs enjoy all the properties of SFAs (regular languages) such as boolean closures and decidability of equivalence.

3.2 Equivalence of ESFTs

While the general equivalence problem of $\mathcal{T}_A = \mathcal{T}_B$ is already undecidable for very restricted classes of finite state transducers [6], the problem is decidable for SFTs in the single-valued case. More generally, one-equality of transductions (defined next) is decidable for SFTs (over decidable label theories).

Definition 7. Functions $\mathbf{f}, \mathbf{g} : X \rightarrow 2^Y$ are *one-equal*, $\mathbf{f} \stackrel{\perp}{=} \mathbf{g}$, if for all $x \in X$, if $x \in \mathcal{D}(\mathbf{f}) \cap \mathcal{D}(\mathbf{g})$ then $|\mathbf{f}(x) \cup \mathbf{g}(x)| = 1$. Let

$$\mathbf{f} \uplus \mathbf{g}(x) \stackrel{\text{def}}{=} \begin{cases} \mathbf{f}(x) \cup \mathbf{g}(x), & \text{if } x \in \mathcal{D}(\mathbf{f}) \cap \mathcal{D}(\mathbf{g}); \\ \emptyset, & \text{otherwise.} \end{cases}$$

Proposition 1. $\mathbf{f} \stackrel{\perp}{=} \mathbf{g}$ iff $\mathbf{f} \uplus \mathbf{g}$ is single-valued.

Note that $\mathbf{f} \stackrel{\perp}{=} \mathbf{f}$ iff \mathbf{f} is single-valued. Thus, one-equality is a more refined notion than single-valuedness, because an effective construction of $A \uplus B$ such that $\mathcal{T}_{A \uplus B} = \mathcal{T}_A \uplus \mathcal{T}_B$ may not always be feasible or even possible for some classes of transducers.

Definition 8. Functions $\mathbf{f}, \mathbf{g} : X \rightarrow 2^Y$ are *domain-equivalent* if $\mathcal{D}(\mathbf{f}) = \mathcal{D}(\mathbf{g})$.

Definitions 7 and 8 are lifted to (E)SFTs. For domain-equivalent single-valued transducers A and B , $A \stackrel{\perp}{=} B$ implies equivalence of A and B ($\mathcal{T}_A = \mathcal{T}_B$).

A natural question that arises is whether decidability of one-equality of SFTs generalizes to ESFTs. The answer is positive for the subclass of *Cartesian* ESFTs (that includes ESFTs in Example 1), but negative in general. We first show that one-equality of ESFTs over decidable label theories is undecidable in general.

Theorem 5 (One-Equality). *One-equality of ESFTs with lookahead 2, over quantifier free successor arithmetic and tuples is undecidable.*

Proof. We give a reduction from the Domain Intersection problem of Theorem 1. Let A_1 and A_2 be ESFTs with lookahead 2 over quantifier free successor arithmetic and tuples. We construct ESFTs A'_i , for $i \in \{1, 2\}$, as follows:

$$A'_i = (Q_{A_i}, q_{A_i}^0, \Delta_{A_i} \cup \{p \xrightarrow{\varphi/[i]} \bullet \mid p \xrightarrow{\varphi} \bullet \in F_{A_i}\})$$

So $\mathcal{T}_{A'_i}(t) = \{[i]\}$ if $t \in \mathcal{D}(A_i)$ and $\mathcal{T}_{A'_i}(t) = \emptyset$ otherwise. Let $\mathbf{f} = \mathcal{T}_{A_1} \uplus \mathcal{T}_{A_2}$. So

- $|\mathbf{f}(t)| = 0$ iff $t \notin \mathcal{D}(A_1) \cup \mathcal{D}(A_2)$;
- $|\mathbf{f}(t)| = 1$ iff $t \in \mathcal{D}(A_1) \cup \mathcal{D}(A_2)$ and $t \notin \mathcal{D}(A_1) \cap \mathcal{D}(A_2)$;
- $|\mathbf{f}(t)| = 2$ iff $t \in \mathcal{D}(A_1) \cap \mathcal{D}(A_2)$.

It follows that $A'_1 \stackrel{\perp}{=} A'_2$ iff (by Proposition 1) \mathbf{f} is single-valued iff $\mathcal{D}(A_1) \cap \mathcal{D}(A_2) = \emptyset$. Now use Theorem 1. \square

The main decidability result of the paper is Theorem 6 that extends the corresponding result for SFTs [16, Theorem 1]. We use the following definitions. A transition, $p \xrightarrow{\varphi/f} q$ where $\ell > 1$, φ is Cartesian and $\mathcal{W}(\varphi) = (a_1, \dots, a_\ell)$, is represented, given $\varphi_i = \lambda x. \varphi(a_1, \dots, a_{i-1}, x, a_{i+1}, \dots, a_\ell)$, by the following path of *split* transitions,

$$p \xrightarrow{\varphi_1/f} p_1 \xrightarrow{\varphi_2/\perp} p_2 \cdots p_{\ell-1} \xrightarrow{\varphi_\ell/\perp} q$$

where p_i for $1 \leq i < \ell$ are new *temporary* states, and the output f is postponed until all input elements have been read. Let Δ_A^s denote such *split* view of Δ_A . Here we assume that all finalizers have lookahead zero, since we do not assume ESFTs here to be deterministic.

Example 2. It is trivial to transform any ESFT into an equivalent (possibly nondeterministic) form where all finalizers have zero lookahead. Consider the ESFT *Base64encode* in Example 1. In the last two finalizers, replace \bullet with a new state p_1 and add the new finalizer $p_1 \xrightarrow{\text{true}/\perp} \bullet$. \square

Definition 9. Let A and B be Cartesian ESFTs with same input and output types and zero-lookahead finalizers. The *product* of A and B is the following

product ESFT $A \times B$. The initial state $q_{A \times B}^0$ of $A \times B$ is (q_A^0, q_B^0) . The states and transitions of $A \times B$ are obtained as the least fixed point of

$$\left. \begin{array}{l} (p, q) \in Q_{A \times B} \\ p \xrightarrow[1]{\varphi/f} p' \in \Delta_A^s \\ q \xrightarrow[1]{\psi/g} q' \in \Delta_B^s \end{array} \right\} \text{IsSat}(\varphi \wedge \psi) \xrightarrow{\Rightarrow} (p', q') \in Q_{A \times B}, \quad (p, q) \xrightarrow[1]{\varphi \wedge \psi / (f, g)} (p', q') \in \Delta_{A \times B}$$

Let $F_{A \times B}$ be the set of all rules $(p, q) \xrightarrow[0]{\text{true}/(v, w)} \bullet$ such that $p \xrightarrow[0]{\text{true}/v} \bullet \in F_A$, $q \xrightarrow[0]{\text{true}/w} \bullet \in F_B$, and $(p, q) \in Q_{A \times B}$. Finally, remove from $Q_{A \times B}$ (and $\Delta_{A \times B}$) all dead ends (non-initial states from which \bullet is not reachable).

We lift the definition of transductions to product ESFTs. A pair-state $(p, q) \in Q_{A \times B}$ is *aligned* if all transitions from (p, q) have outputs (f, g) such that $f \neq \perp$ and $g \neq \perp$. The relation $\xrightarrow[\text{aligned}]{/}_{A \times B}$ is defined analogously to ESFTs.

Lemma 1 (Product). *For all aligned $(p, q) \in Q_{A \times B}$, $u \in \Sigma^*$, $v, w \in \Gamma^*$:*

$$(p, q) \xrightarrow[\text{aligned}]{u/(v, w)}_{A \times B} \bullet \Leftrightarrow p \xrightarrow[u/v]{/}_A \bullet \wedge q \xrightarrow[u/w]{/}_B \bullet.$$

We define also, for all $u \in \Sigma^*$, $\mathcal{T}_{A \times B}(u) \stackrel{\text{def}}{=} \{(v, w) \mid q_{A \times B}^0 \xrightarrow[u/(v, w)]{\text{aligned}} \bullet\}$ and $\mathcal{D}(A \times B) \stackrel{\text{def}}{=} \mathcal{D}(\mathcal{T}_{A \times B})$. Lemma 1 implies that $\mathcal{D}(A \times B) = \mathcal{D}(A) \cap \mathcal{D}(B)$ and $A \not\equiv B$ iff there exists u and $v \neq w$ such that $(v, w) \in \mathcal{T}_{A \times B}(u)$.

Next we prove an *alignment* lemma that allows us to either effectively eliminate all non-aligned pair-states from $A \times B$ without affecting $\mathcal{T}_{A \times B}$ or else to establish that $A \not\equiv B$. A product ESFT is *aligned* if all pair-states in it are aligned.

Lemma 2 (Alignment). *If $A \stackrel{\perp}{=} B$ then there exists an aligned product ESFT that is equivalent to $A \times B$. Moreover, there is an effective procedure that either constructs it or else proves that $A \not\equiv B$, if the label theory is decidable.*

Proof. The product $A \times B$ is incrementally transformed by eliminating non aligned pair-states from it. Each iteration preserves equivalence. Using DFS, initialize the search *frontier* to be $\{q_{A \times B}^0\}$. Pick (and remove) a state (p, q) from the frontier and consider all transitions starting from it. The main two cases are the following:

1. If there are transitions from (p, q) where both the A -output f and the B -output g are $(\sigma^\ell \rightarrow \gamma)$ -sequences with equal lookahead (say $\ell = 2$):

$$(p, q) \xrightarrow[1]{\varphi/(f, g)} (p_1, q_1) \xrightarrow[1]{\psi/(\perp, \perp)} (p_2, q_2)$$

replace the path with the following combined transition with lookahead 2

$$(p, q) \xrightarrow[2]{\lambda(x_0, x_1) \cdot \varphi(x_0) \wedge \psi(x_1) / (f, g)} (p_2, q_2).$$

and add (p_2, q_2) to the frontier unless (p_2, q_2) has already been visited. Note that $(p_2, q_2) \in Q_A \times Q_B$ and thus (p_2, q_2) is aligned.

2. Assume there are transitions where the A -output f is a $(\sigma^k \rightarrow \gamma)$ -sequence and the B -output g is a $(\sigma^\ell \rightarrow \gamma)$ -sequence ($k \neq \ell$, say $k = 2$ and $\ell = 1$):

$$(p, q) \xrightarrow{\varphi/(f,g)} (p_1, q_1) \xrightarrow{\psi/(\perp, g_1)} (p_2, q_2)$$

So p_1 is temporary while q_1 is not.

Decide if f can be split into two independent $(\sigma \rightarrow \gamma)$ -sequences f_1 and f_2 such that for all $a_1 \in \llbracket \varphi \rrbracket$ and $a_2 \in \llbracket \psi \rrbracket$, $\llbracket f \rrbracket(a_1, a_2) = \llbracket f_1 \rrbracket(a_1) \cdot \llbracket f_2 \rrbracket(a_2)$. To do so, choose h_1 and h_2 such that $f = \lambda(x, y).h_1(x, y) \cdot h_2(x, y)$ (note that the total number of such choices is $|f| + 1$ where $|f|$ is the length of the output sequence), let $f_1 = \lambda x.h_1(x, \mathcal{W}(\psi))$, $f_2 = \lambda x.h_2(\mathcal{W}(\varphi), x)$ and check validity of the *split predicate*

$$\forall x y ((\varphi(x) \wedge \psi(y)) \Rightarrow f(x, y) = f_1(x) \cdot f_2(y))$$

If there exists a valid split predicate then pick such f_1 and f_2 , and replace the above path with

$$(p, q) \xrightarrow{\varphi/(f_1, g)} (p'_1, q'_1) \xrightarrow{\psi/(f_2, g_1)} (p_2, q_2)$$

where (p'_1, q'_1) is a new *aligned* pair-state added to the frontier.

Suppose that splitting fails. We show that $A \not\stackrel{\perp}{=} B$, by way of contradiction. Assume $A \stackrel{\perp}{=} B$.

Since splitting fails, the following *dependency* predicates are satisfiable:

$$\begin{aligned} D1 &= \lambda(x, x', y). \varphi(x) \wedge \varphi(x') \wedge \psi(y) \wedge f(x, y) \neq f(x', y) \\ D2 &= \lambda(x, y, y'). \varphi(x) \wedge \psi(y) \wedge \psi(y') \wedge f(x, y) \neq f(x, y') \end{aligned}$$

Let $(a_1, a'_1, a_2) = \mathcal{W}(D1)$ and $(e_1, e_2, e'_2) = \mathcal{W}(D2)$. Assume that $A \stackrel{\perp}{=} B$. We proceed by case analysis over $|f|$. We know that $|f| \geq 1$, or else splitting is trivial.

- (a) Assume first that $|f| = 1$. Let

$$[b] = \llbracket f \rrbracket(a_1, a_2), [b'] = \llbracket f \rrbracket(a'_1, a_2), [d] = \llbracket f \rrbracket(e_1, e_2), [d'] = \llbracket f \rrbracket(e_1, e'_2).$$

Thus $b \neq b'$ and $d \neq d'$. Since (p, q) is aligned, and (p_1, q_1) is reachable and alive (by construction of $A \times B$, \bullet is reachable from (p_1, q_1)), there exists $\alpha, \beta \in \Sigma^*$, $u_1, u_2, v_1, v_2, v_3, v_4 \in \Gamma^*$, such that, by *IsSat*($D1$),

$$\left. \begin{array}{l} p_0 \xrightarrow{\alpha/u_1} p \xrightarrow{[a_1, a_2]/[b]} p_2 \xrightarrow{\beta/u_2} \bullet \\ q_0 \xrightarrow{\alpha/v_1} q \xrightarrow{[a_1]/\llbracket g \rrbracket(a_1)} q_1 \xrightarrow{[a_2] \cdot \beta/v_2} \bullet \end{array} \right\} \begin{array}{l} (A \stackrel{\perp}{=} B) \quad u_1 \cdot [b] \cdot u_2 = \\ \quad \quad \quad v_1 \cdot \llbracket g \rrbracket(a_1) \cdot v_2 \end{array}$$

$$\left. \begin{array}{l} p_0 \xrightarrow{\alpha/u_1} p \xrightarrow{[a'_1, a_2]/[b']} p_2 \xrightarrow{\beta/u_2} \bullet \\ q_0 \xrightarrow{\alpha/v_1} q \xrightarrow{[a'_1]/\llbracket g \rrbracket(a'_1)} q_1 \xrightarrow{[a_2] \cdot \beta/v_2} \bullet \end{array} \right\} \begin{array}{l} (A \stackrel{\perp}{=} B) \quad u_1 \cdot [b'] \cdot u_2 = \\ \quad \quad \quad v_1 \cdot \llbracket g \rrbracket(a'_1) \cdot v_2 \end{array}$$

By $b \neq b'$, $|v_1| \leq |u_1| < |v_1 \cdot \llbracket g \rrbracket(a_1)| = |v_1| + |g|$. Also, by *IsSat*(D2),

$$\left. \begin{array}{l} p_0 \xrightarrow{\alpha/u_1} p \xrightarrow{[e_1, e_2]/[d]} p_2 \xrightarrow{\beta/u_2} \bullet_A \\ q_0 \xrightarrow{\alpha/v_1} q \xrightarrow{[e_1]/\llbracket g \rrbracket(e_1)} q_1 \xrightarrow{[e_2] \cdot \beta/v_3} \bullet_B \end{array} \right\} \begin{array}{l} (A \stackrel{1}{=} B) \quad u_1 \cdot [d] \cdot u_2 = \\ \quad \quad \quad v_1 \cdot \llbracket g \rrbracket(e_1) \cdot v_3 \end{array}$$

$$\left. \begin{array}{l} p_0 \xrightarrow{\alpha/u_1} p \xrightarrow{[e_1, e'_2]/[d']} p_2 \xrightarrow{\beta/u_2} \bullet_A \\ q_0 \xrightarrow{\alpha/v_1} q \xrightarrow{[e_1]/\llbracket g \rrbracket(e_1)} q_1 \xrightarrow{[e'_2] \cdot \beta/v_4} \bullet_B \end{array} \right\} \begin{array}{l} (A \stackrel{1}{=} B) \quad u_1 \cdot [d'] \cdot u_2 = \\ \quad \quad \quad v_1 \cdot \llbracket g \rrbracket(e_1) \cdot v_4 \end{array}$$

By $d \neq d'$, $|v_1 \cdot \llbracket g \rrbracket(e_1)| = |v_1| + |g| \leq |u_1|$. But $|u_1| < |v_1| + |g|$. \dagger
(b) The case where $|f| > 2$ is similar to (a).

The remaining cases are similar and effectively eliminate all non-aligned pair-states from $A \times B$ or else establish that $A \stackrel{1}{=} B$. \boxtimes

Assume $A \times B$ is aligned and let $\lceil A \times B \rceil$ be the following *product SFT* (product ESFT all of whose transitions have lookahead 1) over the input type σ^* . For each $p \xrightarrow{\lambda \bar{x} \cdot \varphi(x_0, x_1, \dots, x_{\ell-1})/(f, g)}_{\ell} q$ in $\Delta_{A \times B}$ let y be a variable of sort σ^* and let φ_1 be the σ^* -predicate

$$\lambda y. \varphi(y[0], y[1], \dots, y[\ell-1]) \wedge \text{tail}^\ell(y) = \square \bigwedge_{i < \ell} \text{tail}^i(y) \neq \square$$

where $y[i]$ is the term that accesses the i 'th head of y and $\text{tail}^i(y)$ is the term that accesses the i 'th tail of y . Lift f to the $(\sigma^* \rightarrow \gamma)$ -sequence $f_1 = \lambda y. f(y[0], y[1], \dots, y[\ell-1])$ and lift g similarly to g_1 . Add the rule $p \xrightarrow{\varphi_1/(f_1, g_1)}_1 q$ as a rule of $\lceil A \times B \rceil$. Thus, the domain type of $\mathcal{T}_{\lceil A \times B \rceil}$ is $(\Sigma^*)^*$ while the range type is $2^{\Gamma^*} \times \Gamma^*$. For $u = [u_0, u_1, \dots, u_n] \in (\Sigma^*)^*$, let $\llbracket u \rrbracket \stackrel{\text{def}}{=} u_0 \cdot u_1 \cdots u_n$ in Σ^* .

Lemma 3 (Grouping). *Assume $A \times B$ is aligned. For all $u \in \Sigma^*$ and $v, w \in \Gamma^*$: $(v, w) \in \mathcal{T}_{A \times B}(u)$ iff $\exists z (u = \llbracket z \rrbracket \wedge (v, w) \in \mathcal{T}_{\lceil A \times B \rceil}(z))$.*

Proof. The type lifting does not affect the semantics of the label-theory specific transformations. \boxtimes

Note that, $[[a_1, a_2], [a_3]]$ and $[[a_1], [a_2, a_3]]$ may be distinct inputs of the lifted product, while both correspond to the same flattened input $[a_1, a_2, a_3]$ of the original product. Intuitively, the internal subsequences correspond to input alignment boundaries of the two ESFTs A and B .

So, in particular, grouping preserves the property: there exists an input u and outputs $v \neq w$ such that $(v, w) \in \mathcal{T}_{A \times B}(u)$. We use the following lemma that is extracted from the main result in [16, Proof of Theorem 1].

Lemma 4 (SFT One-Equality [16]). *Let C be a product SFT over a decidable label theory. The problem of deciding if there exist u and $v \neq w$ such that $(v, w) \in \mathcal{T}_C(u)$ is decidable.*

We can now prove the main decidability result of this paper.

Theorem 6 (Cartesian ESFT One-Equality). *One-equality of Cartesian ESFTs over decidable label theories is decidable.*

Proof. Let A and B be Cartesian ESFTs. Construct $A \times B$. By the Product lemma 1, $\mathcal{D}(A \times B) = \mathcal{D}(A) \cap \mathcal{D}(B)$ and $A \not\equiv B$ iff there exist u and $v \neq w$ such that $(v, w) \in \mathcal{T}_{A \times B}(u)$. By using the Alignment lemma 2, construct aligned product SFT C such that $\mathcal{T}_C = \mathcal{T}_{A \times B}$ or else determine that $A \equiv B$. Now lift C to $\lceil C \rceil$, and by using the Grouping lemma 3, $A \not\equiv B$ iff there exist u and $v \neq w$ such that $(v, w) \in \mathcal{T}_{\lceil C \rceil}(u)$. Finally, observe that adding the sequence operations for accessing the head and the tail of sequences in the lifting construction do, by themselves, not affect decidability of the label theory, apply Lemma 4. \square

3.3 Composition of ESFTs

In this section we show some preliminary results (mainly negative) on ESFT composition. In particular, ESFTs and Cartesian ESFTs are not closed under composition.

Given $\mathbf{f}: X \rightarrow 2^Y$ and $\mathbf{x} \subseteq X$, $\mathbf{f}(\mathbf{x}) \stackrel{\text{def}}{=} \bigcup_{x \in \mathbf{x}} \mathbf{f}(x)$. Given $\mathbf{f}: X \rightarrow 2^Y$ and $\mathbf{g}: Y \rightarrow 2^Z$, $\mathbf{f} \circ \mathbf{g}(x) \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{f}(x))$. This definition follows the convention in [5], i.e., \circ applies first \mathbf{f} , then \mathbf{g} , contrary to how \circ is used for standard function composition. The intuition is that \mathbf{f} corresponds to the relation $R_{\mathbf{f}}: X \times Y$, $R_{\mathbf{f}} \stackrel{\text{def}}{=} \{(x, y) \mid y \in \mathbf{f}(x)\}$, so that $\mathbf{f} \circ \mathbf{g}$ corresponds to the binary relation composition $R_{\mathbf{f}} \circ R_{\mathbf{g}} \stackrel{\text{def}}{=} \{(x, z) \mid \exists y (R_{\mathbf{f}}(x, y) \wedge R_{\mathbf{g}}(y, z))\}$.

Definition 10. *A class of transducer C is closed under composition iff for every \mathcal{T}_1 and \mathcal{T}_2 that are C -definable $\mathcal{T}_1 \circ \mathcal{T}_2$ is also C -definable.*

Theorem 7 (Composition). *The following statements are true: ESFTs are not closed under composition; There exists two Cartesian ESFTs which composition is not ESFT definable; Cartesian ESFTs are not closed under composition.*

We now show that in general the composition of two ESFTs cannot be effectively computed.

Theorem 8 (Undecidability of Composition Computation). *Given two ESFTs with lookahead 2 over quantifier free successor arithmetic and tuples whose composition f is ESFT definable, it is undecidable to compute the ESFT corresponding to f .*

4 Experiments and Applications

In this section we show how several practical applications can be modelled and verified using ESFAs and ESFTs. We first use ESFTs to prove the correctness of some real world string encoders and decoders. We then show how ESFAs

and ESFTs can be useful in the context of deep packet inspection and network protocol transformations. Finally we propose ESFTs as a tool for the analysis of list manipulating programs. All our experiments are run using the tool BEK⁴.

Analysis of String Encoders. A string encoder E transforms input strings in a given format A into output strings in a different format B . A decoder D inverts such transformation. The formats A and B usually use different alphabets (character sets). The first half of Table 1 shows examples of common string encoders/decoders and their respective lookahead sizes. $E \circ D$ ($D \circ E$) denotes the sequential compositions of the encoder with the decoder (decoder with the encoder). We compute such compositions using the semi-decision procedure of [4].

The correctness of UTF8 encoding was already investigated in [4] using a semi-decision procedure for one-equality. We use the algorithm proposed in Section 3.2 to confirm such result and we prove the correctness of three new encoders: BASE64, BASE32 and BASE16. The second half of Table 1 shows the running times of the analyses. The column $E \circ D \stackrel{\perp}{=} I$ ($D \circ E \stackrel{\perp}{=} I$) shows the cost of checking whether $E \circ D$ ($D \circ E$) is one-equal to the identity transducer I . Composition times (typically 1-2 ms) are included in the measurements.

We want to stress that during our experiments we identified wrong implementations of the UTF8 encoder/decoder in which the algorithm of Section 3.2 correctly detected that one-equality fails, while the semi-decision procedure used in [4] did not terminate.

Deep Packet Inspection. Fast identification of network traffic patterns is of vital importance in network routing, firewall filtering and intrusion detection. This task is addressed with the name “deep packet inspection” (DPI) [15]. Due to performance constraints, DPI must be performed in a single pass over the input. The simplest approach is to use DFAs and NFAs to identify patterns. These representations are either not succinct or not streamable. Extended Finite Automata (XFA) [15] make use of registers to reduce the state space while preserving determinism and therefore deterministic ESFAs can be seen as a subclass of XFAs that are able to deal with finite lookahead. Deterministic ESFA can also represent the alphabet symbolically, which enables a new level of succinctness. We believe that deterministic ESFAs can help achieve further succinctness in particular problem instances. To support this hypothesis we observe that several examples shown in [15, Figure 2,3] can be represented as deterministic ESFAs with few transitions. For example the language $\sim/\backslash\text{ncmd}[\sim/\backslash\text{n}\{200}\$$ can be succinctly captured by a deterministic ESFA with one transition!

Network Protocol Conversions. Deep packet inspection can be naturally extended by adding data manipulation. As in the previous setting we are inter-

	Lookahead		Analysis (ms)	
	E	D	$E \circ D \stackrel{\perp}{=} I$	$D \circ E \stackrel{\perp}{=} I$
UTF8:	2	4	16	24
BASE64:	3	5	53	19
BASE32:	5	8	8	12
BASE16:	1	2	2	1

Table 1. Analysed encoders (E) and decoders (D), their lookaheads, and analysis times.

⁴ <http://www.rise4fun.com/Bek>.

ested in deterministic ESFTs which can commit their output at every transition without seeing the rest of the input. Deterministic ESFTs can be used to compute logs of network traffic or translate headers of one protocol into another. As an example, a simple translation from an IPv4 header to an IPv6 header⁵ can be easily implemented with a deterministic ESFT with less than 50 transitions. The same transformation using an SFT would require over 100000 transitions.

Verification of List Manipulating Programs. ESFTs can be used for verification of *list manipulating programs* as they naturally model sequential pattern matching. The ML guards $x1::x2::xs \rightarrow (x1+x2)::(f2\ xs)$ and $x1::x2::x3::xs \rightarrow (x1+x2+x3)::(f3\ xs)$, respectively belonging to the functions $f_2, f_3 : list\ int \rightarrow list\ int$, can be naturally expressed as ESFT transitions. Therefore f_2 and f_3 can be modelled as ESFTs. We can then use the one-equality algorithm of Section 3.2 to prove that $f_2(f_2(f_2\ l)) \stackrel{1}{=} f_3(f_3\ l)$ in less than 1 ms.

5 Related Work

Symbolic finite transducers (SFTs) and BEK were originally introduced in [7] with a focus on security analysis of sanitizers. The formal foundations and the theoretical analysis of the underlying SFT algorithms, in particular, an algorithm for one-equality of SFTs, modulo a decidable background theory is studied in [16]. Symbolic Transducers (STs) that allow the use of registers are also defined in [16]. Full equivalence of finite state transducers is undecidable [6], and already so for very restricted fragments [8]. In the single-valued case, decidability was established in [13], and extended to the finite-valued case in [3, 17].

ESFTs were introduced in [4] as a succinct and more analysable representation of a subclass of symbolic transducers (STs). The main result in [4] is a register elimination technique that provides a way to construct (product) ESFTs from (product) symbolic transducers (STs). While this technique provides a semi-decision procedure for one-equality checking (by using grouping, Lemma 3) of non-Cartesian ESFTs, it does not provide a full decision procedure for one-equality of the Cartesian case. The procedure in [4] fails to decide *alignment* of ESFTs, that is the key lemma (Lemma 2) used in the main decidability result of Theorem 6, that is a proper extension of the decidability result of one-equality of SFTs [16, Theorem 1]. We also show that one-equality is undecidable in the non-Cartesian case, Theorem 5, that is in sharp contrast to the theory of classical automata, where the non-Cartesian case is irrelevant (from the point of view of decidability) due to the standard form [18, Theorem 2.17].

Extended Top-Down Tree Transducers [11] (ETTTs) are commonly used in natural language processing. ETTTs also allow finite lookahead on transformation from trees to trees, but only support finite alphabets. The special case in which the input is a string (unary tree) is equivalent to ESFTs over finite alphabets. This paper focuses on ESFTs over any decidable theory. We leave as future work extending the model to tree transformations.

⁵ More information at <http://www.cs.washington.edu/research/networking/napt/>

Symbolic finite transducers with lookback k (k -SLTs) [2] have a sliding window of size $k+1$ that allows, in addition to the current input character, references of up to k previous characters. SLTs use only final states, because it is unclear how to support nonfinal states in the context of learning. Thus, domain intersection (that is undecidable for ESFTs with lookahead 2) is trivial for SLTs.

In recent years there has been considerable interest in automata using infinite alphabets [14], starting with the work on register automata [9]. Finite words over an infinite alphabet are often called data words. This line of work focuses on fundamental questions about decidability, complexity, and expressiveness on classes of automata on one hand and fragments of logic on the other hand.

Streaming transducers [1] provide another recent symbolic extension of finite transducers where the label theories are restricted to be total orders, in order to maintain decidability of equivalence. Streaming transducers are largely orthogonal to SFTs or the extension of ESFTs, as presented in the current paper. For example, streaming transducers do not allow arithmetic, but can reverse the input, which is not possible with ESFTs.

The correctness of UTF8 encoder and decoder was proven in [4] using two semi-decision procedures for equivalence and composition. In this paper we show that the composition of UTF8 encoder and decoder can be expressed as a Cartesian ESFT and can be formally analyzed with the one-equality algorithm introduced in this paper. We do the same for three encoders of which the correctness was not proven before: BASE64, BASE32, BASE16.

Extended Finite Automata (XFA) are introduced in [15] for network packet inspection. XFAs are a succinct representation of DFAs that use registers and allow programs over the registers. ESFAs are orthogonal to XFAs in two ways: 1) XFAs only support finite alphabets; and 2) XFAs aim at representing *most* DFAs succinctly, while ESFAs only capture the languages that use finite lookahead. We have not investigated the application of ESFAs to network packet inspection in detail, but we think that they can help achieving a further level of succinctness. History-based finite automata [10] are another extension of DFAs that have been introduced for encoding regular expressions in the context of network intrusion detection systems, they use a single register (bit-vector) to keep track of history. The register is used together with the input character to determine when a transition is enabled.

6 Conclusion

We showed fundamental negative and positive results about several classical decision problems of ESFAs and ESFTs, establishing a sharp boundary between decidability (the Cartesian case with any decidable background) and undecidability (the non-Cartesian case with a background of successor arithmetic). While the main motivation came from typical *static analysis* problems using ESFTs, an equally important application of the Cartesian case is for efficient *code generation*. Namely, the conjuncts of a Cartesian predicate can be compiled and normalized into separate unary predicates that may for example use BDDs for

efficient and unique set representation when dealing with bit-vectors, as in the context of strings coders. Identifying classes of ESFTs that are closed under composition, as well as extending ESFTs to trees are left as open problems.

References

1. R. Alur and P. Cerný. Streaming transducers for algorithmic verification of single-pass list-processing programs. In *POPL'11*, pages 599–610. ACM, 2011.
2. M. Botincan and D. Babic. Sigma*: symbolic learning of input-output specifications. In *POPL'13*, pages 443–456. ACM, 2013.
3. K. Culic and J. Karhumäki. The equivalence of finite-valued transducers (on HD-TOL languages) is decidable. *Theoretical Computer Science*, 47:71–84, 1986.
4. L. D’Antoni and M. Veanes. Static analysis of string encoders and decoders. In R. Giacobazzi, J. Berdine, and I. Mastroeni, editors, *VMCAI 2013*, volume 7737 of *LNCS*, pages 209–228. Springer, 2013.
5. Z. Fülöp and H. Vogler. *Syntax-Directed Semantics: Formal Models Based on Tree Transducers*. EATCS. Springer, 1998.
6. T. Griffiths. The unsolvability of the equivalence problem for A -free nondeterministic generalized machines. *J. ACM*, 15:409–413, 1968.
7. P. Hooimeijer, B. Livshits, D. Molnar, P. Saxena, and M. Veanes. Fast and precise sanitizer analysis with Bek. In *USENIX*, August 2011.
8. O. Ibarra. The unsolvability of the equivalence problem for Efree NGSMS with unary input (output) alphabet and applications. *SIAM Journal on Computing*, 4:524–532, 1978.
9. M. Kaminski and N. Francez. Finite-memory automata. *TCS*, 134(2):329–363, 1994.
10. S. Kumar, B. Chandrasekaran, J. Turner, and G. Varghese. Curing regular expressions matching algorithms from insomnia, amnesia, and acalculia. In *ANCS 2007*, pages 155–164. ACM/IEEE, 2007.
11. A. Maletti, J. Graehl, M. Hopkins, and K. Knight. The power of extended top-down tree transducers. *SIAM J. Comput.*, 39(2):410–430, June 2009.
12. M. Mohri. Finite-state transducers in language and speech processing. *Comput. Linguist.*, 23(2):269–311, June 1997.
13. M. P. Schützenberger. Sur les relations rationnelles. In *GI Conference on Automata Theory and Formal Languages*, volume 33 of *LNCS*, pages 209–213, 1975.
14. L. Segoufin. Automata and logics for words and trees over an infinite alphabet. In *CSL*, pages 41–57, 2006.
15. R. Smith, C. Estan, S. Jha, and S. Kong. Deflating the big bang: fast and scalable deep packet inspection with extended finite automata. *SIGCOMM '08*, pages 207–218. ACM, 2008.
16. M. Veanes, P. Hooimeijer, B. Livshits, D. Molnar, and N. Bjorner. Symbolic finite state transducers: Algorithms and applications. In *POPL'12*, pages 137–150. ACM, 2012.
17. A. Weber. Decomposing finite-valued transducers and deciding their equivalence. *SIAM Journal on Computing*, 22(1):175–202, February 1993.
18. S. Yu. Regular languages. In G. Rozenberg and A. Salomaa, editors, *Handbook of Formal Languages*, volume 1, pages 41–110. Springer, 1997.