

Transducers with origin information

Mikołaj Bojańczyk*

University of Warsaw

Abstract. Call a string-to-string function *regular* if it can be realised by one of the following equivalent models: MSO transductions, two-way deterministic automata with output, and streaming transducers with registers. This paper proposes to treat origin information as part of the semantics of a regular string-to-string function. With such semantics, the model admits a machine-independent characterisation, Angluin-style learning in polynomial time, as well as effective characterisations of natural subclasses such as one-way transducers or first-order definable transducers.

This paper is about string-to-string functions which can be described by deterministic two-way automata with output [AU70]. As shown in [EH01], this model is equivalent to MSO definable string transductions. Another equivalent model, used in [AC10], is a deterministic one-way automaton with registers that store parts of the output¹. Examples of such functions include: duplication $w \mapsto ww$; reversing $w \mapsto w^R$; a function $w \mapsto ww^R$ which maps an input to a palindrome whose first half is w ; and a function which duplicates inputs of even length and reverses inputs of odd length. As witnessed by the multiple equivalent definitions, this class of string-to-string function is robust, and therefore, following [AC10], we call it the class of *regular string-to-string functions*. Regular string-to-string functions have good closure properties. For instance, if f and g are regular, then the composition $w \mapsto f(g(w))$ is also regular, which is straightforward if the MSO definition is used, but nontrivial if the two-way automata definition is used [CJ77]. Also the concatenation $w \mapsto f(w) \cdot g(w)$ is regular, which is apparent in any of the three definitions. Equivalence of regular string-to-string functions is decidable, as was shown in [Gur82] using the two-way automata definition.

Origins. The motivation of this paper is the simple observation that the models discussed above, namely deterministic two-way automata with output, MSO definable string transductions, and automata with registers, provide more than just a function from strings to strings. In each case, one can also reconstruct

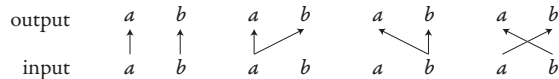
* Supported by ERC Starting Grant “Sosna”

¹ Registers are similar to attributes in attribute grammars. The equivalence of MSO definable transductions with a form of attribute grammars, in the tree-to-tree case, was shown in [BE00]. In the special case of string-to-string functions, the attribute grammars from [BE00] correspond to left-to-right deterministic automata with registers and regular lookahead.

origin information, which says how positions of the output string originate from positions in the input string. How do we reconstruct the origin of a position x in an output string? In the case of a deterministic two-way automaton, this is the position of the head when x was output. In the case of an MSO definable transducer, this is the position in which x is interpreted. In the case of an automaton with registers, this is the position in the input when the letter x was first loaded into a register. In other words, for a transducer we can consider two semantics: the *standard semantics*, where the output is a string, and the *origin semantics*, where the output is a string with origin information. The second semantics is finer in the sense that some transducers might be equivalent under the standard semantics, but not under the origin semantics.

Tracking origin information for transducers has been studied before, for instance in the programming language community, see e.g. [vDKT93]. Origin information has also been used as a technical tool in the study of tree-to-tree transducers. Examples include [EM03], where origin information is used to characterise those macro tree transducers which are MSO definable, and [LMN10], where origin information is used to get a Myhill-Nerode characterisation of deterministic top-down tree transducers. The novelty of this paper is that origin information is built into the semantics of a transducer.

Origin semantics. To illustrate the difference between the two semantics (standard and origin) of a string-to-string transducer, consider a transducer which is the identity on the string ab , and which maps other strings to the empty string. If we care about origins, then this description is incomplete, and can be instantiated in four different ways depicted below.



For example, the second diagram above describes a two-way automaton that first reads its input to determine if it is ab , and then moves its head to the first position, where it outputs both a and b .

Another example is the identity function on strings over a one letter alphabet, which can be realised by copying the input left-to-right or right-to-left. Actually the function can be realised in infinitely many different ways once origins are taken into account: consider an automaton that outputs n letters in input positions divisible by n , and then outputs the remainder under division by n in the last input position.

This paper is a study of the more refined semantics. Almost any “natural” construction for transducers will respect origin information. For instance, the translation from [EH01] which converts an MSO interpretation into a deterministic two-way automaton remains correct when the origin information is taken into account. The same holds for the other translations between the three models. In other words, one can also talk about *regular string-to-string functions with origin information*. Various closure properties, such as composition and concatenation,

are retained when origins are taken into account. Some results become easier to prove, e.g. decidability of equivalence of string-to-string transducers.

A machine independent characterisation. The main contribution of this paper is a machine independent characterisation of regular string-to-string functions with origin information, which is given in Theorem 1. The characterisation is similar to the Myhill-Nerode theorem, which says that a language L is regular if and only if it has finitely many left derivatives of the form

$$w^{-1}L \stackrel{\text{def}}{=} \{v : wv \in L\}.$$

From the usual Myhill-Nerode theorem for regular languages one obtains a canonical device, which is the minimal deterministic automaton. The situation is similar here. We define a notion of left and right derivatives for string-to-string functions with origin information, and show that a function is regular if and only if it has finitely many left and right derivatives (finitely many left derivatives is not enough, same for right derivatives). The proof of the theorem yields a canonical device, which is obtained from the function itself and not its representation as a two-way automaton, MSO transduction, or machine with registers. One use for the canonical device is testing equivalence: two devices are equivalent if and only if they yield the same canonical machine.

Another use of the canonical device is that it is easy to see when the underlying function actually belongs to a restricted class, e.g. if it can be defined by a deterministic one-way automaton with output (see Theorem 4), or by functional nondeterministic one-way automaton with output (see Theorem 3). A more advanced application is given in Theorem 5, which characterises the first-order fragment of MSO definable transducers with origin information.

Learning. One of the advantages of origin information is that it allows functions to be learned, using an Angluin style algorithm. We show that a regular string-to-string function with origin information can be learned with a number of queries that is polynomial in the size of the canonical device. The queries are of two types: the learner can ask for the output on a given input string; or the learner can propose a transducer with origin information, and in case this is not the correct one, then the teacher gives a counterexample string where the proposed transducer produces a wrong output.

In the algorithm, the learner uses the origin information. However, it seems that the learner's advantage from the origin information does not come at any significant cost to the teacher. Suppose that we want to learn a transducer inside a text editor, e.g. the user wants to teach the text editor that she is thinking of the transducer which replaces every $=$ by $:=$. If a user is trying to show an example of this transducer on some input, then she will probably place the cursor on occurrences of $=$ in the input, delete them, and retype $:=$, thus giving origin information to the algorithm. A user who backspaces the whole input and retypes a new version will possibly be thinking of some different transformation. It would be wasteful to ignore this additional information supplied by the user.

Thank you. I would like to thank Sebastian Maneth and the anonymous referees for their valuable feedback; Anca Muscholl, Szymon Toruńczyk and Igor Walukiewicz for discussions about the model; and Rajeev Alur for asking the question about a machine-independent characterisation of transducers.

1 Regular string to string transducers

A *string-to-string function* is any function from strings over some fixed input alphabet to strings over some fixed output alphabet, such that the empty string is mapped to the empty string. A *string-to-string function with origin information* is defined in the same way, but for every input string w it provides not only an output string $f(w)$, but also origin information, which is a function from positions in $f(w)$ to positions in w . We consider total functions, although the results can easily be adapted to partial functions. In this section we recall three equivalent models recognising string-to-string functions.

Streaming transducer. Following [AC10], a *streaming transducer* is defined as follows. It has finite *input* and *output alphabets*. There is a finite set of *control states* with a distinguished *initial state*, and a finite set of *registers*, with a distinguished *output register*. The transition function inputs a control state and an input letter, and outputs a new control state and a *register update*, which is a sequence of register operations of two possible types:

- *Concatenate.* Replace the contents of register r with rs , and replace the contents of register s by the empty string;
- *Create.* Replace the contents of register r with output letter b .

Finally, there is an *end of input function*, which maps each state to a sequence of register operations of the first type².

When given an input string, the transducer works as follows. It begins in the initial state with all registers containing the empty string. Then it processes each input letter from left to right, updating the control state and the registers according to the transition function. Once the whole input has been processed, the end of input function is applied to the last state, yielding another sequence of register operations, and finally the value of the transducer is extracted from the output register. For the origin semantics, we observe that every letter in a register is created once using an operation of type *create*, and then moved around using operations of type *concatenate*. The origin of an output letter is defined to be the input position which triggered the transition whose register update contained the appropriate *create* operation.

Observe that the register operations do not allow copying registers. This is an important restriction which guarantees, among other things, that the size of the output is linear in the size of the input.

² The end of input function is prohibited to produce new output letters so that the origin information can be assigned. Alternatively, one could assume that the positions produced by the end of input function have a special origin, “created out of nothing”.

Example 1. By composing the atomic register operations and using additional registers, we can recover additional register operations such as “add letter b to the end of register r ”, “add letter b to the beginning of register r ”, “move register r to register s , leaving r empty”. The examples use the additional operations.

Consider the function $w \mapsto ww^R$, where w^R is the reverse of w . The transducer has one control state and two registers, used to store w and w^R . When it reads an input letter a , the transducer adds a to the end of the register storing w and adds a to the beginning of the register storing w^R . The end of input update concatenates both registers, and puts the result in the first register, which is the output register.

A transducer for the duplication function is obtained in a similar way. Observe that since the register operations do not allow copying, it is still necessary to have two registers, both storing w .

Deterministic two-way automaton with output. A deterministic two-way automaton with output is like a deterministic two-way finite automaton, except that every transition is additionally labelled by a string (possibly empty) over the output alphabet. A run over an input w can be seen as a sequence of pairs $(\delta_1, x_1), \dots, (\delta_n, x_n)$ where δ_i is a transition and x_i a position in the string $\vdash w \dashv$. The transition δ_i reads the label of position x_i and the state generated by the previous transition, and chooses the new position x_{i+1} , a new state, and what will be appended to the output. The output of the automaton is the concatenation of the strings labelling the transitions $\delta_1, \dots, \delta_n$. The origin of a position in the output string that is generated by the transition δ_i is defined to be the position x_i . To make the origin well-defined, we require that every output letter is produced for transitions that have their source in input letters, and not over the markers \vdash and \dashv .

MSO transduction. Following [Tho97], a string over an alphabet A can be treated as a relational structure, whose universe is the positions of the string, and which has a binary position order predicate $x < y$ and label predicates $a(x)$ for the letters of the alphabet. To transform strings into strings, we can use MSO interpretations. An MSO *interpretation* is a function from structures over some fixed input vocabulary (set of relation names with their arities) to structures over some fixed output vocabulary, which is specified by a system of MSO formulas, as follows. There is a universe formula with one free variable over the input vocabulary, which selects the elements of the universe of the input structure that will appear in the universe of the output structure. Furthermore, for every predicate of the output vocabulary there is a formula over the input vocabulary of the same arity, which says how the predicates are defined in the output structure.

Another function from structures to structures is called *k-copying*; which maps a structure to k disjoint copies of itself, together with binary relations $1(x, y), \dots, k(x, y)$ such that $i(x, y)$ holds if y is the i -th copy of x . A *copying MSO transduction* consists of first a copying function, followed by an MSO interpretation. A string-to-string function f is called *MSO-definable* if there is some copying MSO transduction such that for every input string w , the transduction

transforms the relational structure corresponding to w into a relational structure corresponding to $f(w)$. The origin information in such a transducer is defined in the natural way: a position in the output string is interpreted in some copy of a position in the input string, the latter is defined to be the origin.

Equivalence of the models. Deterministic two-way automata with output define the same translations as copying MSO transductions in [EH01]. The same proof works if the semantics with origin information is used. Streaming transducers are shown to be equivalent to the previous two models in [AC10]; the same proof also works with the origin semantics. A string-to-string function with origin information is called a *regular string-to-string function with origin information* if it can be defined by any one of the three models mentioned above.

2 A machine independent characterisation

In this section we present a Myhill-Nerode style characterisation of regular transducers with origin information.

Factorised output. Suppose that f is a string-to-string function with origin information and output alphabet B . A *factorised input* is a tuple of strings w_1, \dots, w_n over the input alphabet, which is meant to describe an input string factorised into n blocks. Given such a factorised input, define an *output block of type i* to be a maximal connected subset of positions in the output $f(w_1 \cdots w_n)$ that originates in w_i . Define the *factorised output* corresponding to a factorised input w_1, \dots, w_n , denoted by

$$f(w_1 | \dots | w_n) \in (\{1, \dots, n\} \times B^+)^*.$$

to be the sequence of output blocks read from left to right, with each block described by its type and corresponding part of the output. In particular, if we concatenate all of the strings coming from B^+ , we obtain the output string $f(w_1 \cdots w_n)$. When $n = 3$, instead of numbers 1, 2, 3 we use “left”, “middle” and “right” to indicate types of blocks. We use fraction-style notation for output blocks, with the lower part indicating the type, and the upper part describing the output. For instance, if f is the duplicating function, then

$$f(ab|cd|e) = \overset{ab}{\text{left}} \overset{cd}{\text{middle}} \overset{e}{\text{right}} \overset{ab}{\text{left}} \overset{cd}{\text{middle}} \overset{e}{\text{right}}.$$

Some input blocks might be empty, as in the following example:

$$f(ab||e) = \overset{ab}{\text{left}} \overset{e}{\text{right}} \overset{ab}{\text{left}} \overset{e}{\text{right}}.$$

If some of the input blocks are underlined, then in the output we just keep the information that there is a nonempty output block, but we do not store the actual output strings which originate in the underlined blocks. For example,

$$f(\underline{ab}|cd|e) = \overset{cd}{\text{left}} \overset{cd}{\text{middle}} \overset{cd}{\text{right}} \overset{cd}{\text{left}} \overset{cd}{\text{middle}} \overset{cd}{\text{right}}.$$

Note that we will never have two consecutive blocks of the same type, e.g. left left, in the factorised output, since blocks are maximal. In particular, for underlined input blocks we lose track of how long their corresponding output blocks are.

Derivatives. Define a *two-sided derivative* of string-to-string function with origin information f to be any function of the form

$$f_{u_w} \stackrel{\text{def}}{=} v \mapsto f(\underline{u}|v|\underline{w}),$$

for some choice of strings u and w over the input alphabet. *Left derivatives* and *right derivatives* are the special cases of the two-sided derivative when either u or w is empty, i.e. they are functions of the forms:

$$f_{u_} \stackrel{\text{def}}{=} v \mapsto f(\underline{u}|v) \qquad f_{_w} \stackrel{\text{def}}{=} v \mapsto f(v|\underline{w}).$$

Example 2. Let f be the function $w \mapsto w^R w$. Then

$$f_{u_w}(v) = \text{right} \overset{v^R}{\text{middle}} \text{left} \overset{v}{\text{middle}} \text{right}$$

for every nonempty strings u or w . When the string u is empty, then the left block disappears, likewise when w is empty then the right blocks disappear. In particular, this function has four possible values for the two-sided derivative. There are two possible values for the left derivative $f_{u_}$, namely the functions

$$v \mapsto \overset{v^R v}{\text{right}} \qquad v \mapsto \overset{v^R}{\text{right}} \overset{v}{\text{left}} \text{right}.$$

Example 3. Let f be the function which is the identity on strings of even length, and which maps strings of odd length to the empty string. This function has three possible left derivatives $f_{v_}$, depending on whether v is empty, nonempty and even length, or odd length. Below is the derivative for the last case.

$$w \mapsto \begin{cases} \overset{w}{\text{left}} \text{right} & \text{if } w \text{ has odd length} \\ \epsilon & \text{otherwise} \end{cases}$$

Example 4. Here is a function with finitely many right derivatives, but infinitely many left derivatives. Consider first the function which scans its input from left to right, and outputs only those letters whose position is a prime number

$$f(a_1 \cdots a_n) = w_1 \cdots w_n \quad \text{where } w_i = \begin{cases} a_i & \text{if } i \text{ is a prime number} \\ \epsilon & \text{otherwise.} \end{cases}$$

This particular function has infinitely many right derivatives, since

$$f_{\cdot w}(v) = \begin{cases} f^{(v)}_{\text{left right}} & \text{if there is a prime number in } \{|v| + 1, \dots, |vw|\} \\ f^{(v)}_{\text{left}} & \text{otherwise.} \end{cases}$$

However finitely many right derivatives can be obtained by making the last position to be output unconditionally, i.e. in the string-to-string function

$$g(a_1 \cdots a_n) = f(a_1 \cdots a_{n-1})a_n.$$

In this case, g has only two right derivatives, namely

$$v \mapsto \begin{matrix} g^{(v)} \\ \text{left} \end{matrix} \quad v \mapsto \begin{matrix} f^{(v)} \\ \text{left right} \end{matrix}.$$

The function has infinitely many left derivatives $g_{v\cdot}$ because the criterion “ i is a prime number” needs to be replaced by “ $i + |v|$ is a prime number”.

To present our machine independent characterisation, we need a notion of regularity for functions from tuples of strings to a finite set. Define the *language encoding* of a function $f : (A^*)^n \rightarrow C$, with C finite, to be

$$\{w_1 \# w_2 \# \cdots \# w_n \# f(w_1, \dots, w_n) : w_1, \dots, w_n \in A^*\} \subseteq (A \cup C \cup \{\#\})^*.$$

assuming $\#$ is a symbol outside $A \cup C$. The function f is called a *regular colouring* if its language encoding is regular. Among several models of automata reading tuples of strings, regular colourings correspond to the weakest model, called recognisable. For instance, the equality function, seen as a colouring of string pairs by “equal” or “not equal”, is not a regular colouring.

Theorem 1 (Machine independent characterisation). *For a string-to-string function f with origin information, the following conditions are equivalent*

1. f is regular;
2. f has finitely many left derivatives and finitely many right derivatives;
3. for every letter a in the input alphabet, the following is a regular colouring

$$(v, w) \mapsto f(\underline{v}|a|\underline{w}).$$

The function $(v, a, w) \mapsto f(\underline{v}|a|\underline{w})$, where v, w are words and a is a letter over the input alphabet, is called the *characteristic function* of f .

Proof (rough sketch). The implication from 1 to 2 is shown by using deterministic two-way automata with output. For the implication from 2 to 3, one observes that the functions $v \mapsto f_{v\cdot}$ and $w \mapsto f_{\cdot w}$ are regular colorings, and that $f(\underline{v}|a|\underline{w})$ is uniquely determined by $f_{v\cdot}$, a and $f_{\cdot w}$. For the implication from 3 to 1, one shows that an arbitrary string-to-string function with origin information can be uniquely reconstructed based on its characteristic function, and when the characteristic function happens to be a regular coloring then this reconstruction can be done by a finite state device.

Since a string-to-string function is uniquely determined by its characteristic function, instead of studying string-to-string functions, one can study their characteristic functions. This is the case in the learning algorithm from Section 3, and the studies of subclasses of transducers in Sections 4 and 5. The characteristic function can be computed based on a representation as a transducer model, e.g. from a copying MSO transduction. In particular, Theorem 1 gives a conceptually simple equivalence check for origin semantics: compute the characteristic functions and test if they are equal. The complexity of this algorithm, especially in the case when the function is given by streaming transducers, is left open.

As shown in Example 4, it is not enough to require finitely many derivatives of one kind, say right derivatives, since a function might have finitely many derivatives of one kind, but infinitely many derivatives of the other kind³.

3 Learning

This section shows that transducers with origin information can be learned. We first recall the Angluin algorithm for regular languages, which will be used as a black box in our learning algorithm for learning transducers. The setup for the Angluin algorithm is as follows. A teacher knows a regular language. A learner wants to learn this language, by asking two kinds of queries. In a *membership query*, the learner gives a string and the teacher responds whether this string is in the language. In an *equivalence query*, the learner proposes a candidate for the teacher's language, and the teacher either says that this candidate is correct, in which case the protocol is finished by learner's success, or otherwise the teacher returns a *counterexample*, which is a string in the symmetric difference between the candidate and teacher's languages.

Angluin proposed an algorithm [Ang87], in which the learner learns the language by asking a number of queries which is polynomial in the minimal deterministic automaton for the teacher's language, and the size of the counterexamples given during the interaction. Theorem 2 shows that a variant of this algorithm works for regular string-to-string transducers with origin information. In the case of transducers, the membership query becomes a *value query*, where the learner gives a string and the teacher responds with the output of the transducer on that string. In the equivalence query, the counterexample becomes a string where the transducer proposed by the learner gives a different value than the transducer of the teacher. In both the value query and in the counterexample, the teacher also provides the origin information.

Theorem 2. *A regular string-to-string function with origin information can be learned using value and equivalence queries in polynomial time (both number of queries and computation time) in terms of the number of left and right derivatives, and the size of the counterexamples given by the teacher.*

³ It does follow from the theorem that a function with finitely many left and right derivatives has finitely many two-sided derivatives. This is because every regular string-to-string function has finitely many two-sided derivatives.

4 Order-preserving transducers

In this section, we present two characterisations of subclasses of transducers. For semantics without origins, [FGRS13] shows how to decide if a deterministic two-way transducer is equivalent to a nondeterministic one-way transducer, while [WK94] shows how to decide (in polynomial time) if a nondeterministic one-way transducer is equivalent to a deterministic one-way transducer. This section shows analogous results for the origin semantics. Unlike [FGRS13, WK94], the characterisations for the origin semantics are self-evident, which shows how changing the semantics (and therefore changing the problem) makes some technical problems go away. A more difficult characterisation, about first-order definable transducers, is presented in the next section.

In the following theorem, a string-to-string function with origin information is called *order preserving* if for every input positions $x < y$, every output position corresponding to x is before every output position corresponding to y .

Theorem 3. *For a regular string-to-string function with origin information f , the following conditions are equivalent.*

1. f is order-preserving.
2. $f(\underline{v}|\underline{w})$ is one of ϵ , left, right or left right for all input strings v, w .
3. f is recognised by a streaming transducer with lookahead which has only one register, and which only appends output letters to that register.
4. f is recognised by a nondeterministic one-way automaton with output, which has exactly one run over every input string.

Proof. The implication from item 1 to item 2 follows straight from the definition. For the implication from item 2 to 3, we observe that if condition 1 is satisfied, then the transducer constructed in the proof of Theorem 1 will only have one register, and it will only append letters to that register during the run. For the implication from item 3 to item 4, we observe that a nondeterministic one-way automaton with output can guess, for each position of the input, what the lookahead will say. Since the lookahead is computed by a deterministic right-to-left automaton, this leads to a unique run on every input string. The implication from item 4 to item 1 also follows straight from the definition.

Observe that the condition in item 2 can be decided, even in polynomial time, when the characteristic function of the transducer is known.

We can further restrict the model by requiring that the transducer in item 3 does not use any lookahead, or equivalently, by requiring that the automaton in item 4 be deterministic. This restricted model is characterised in the following theorem.

Theorem 4. *Let f be a regular string-to-string function which satisfies any of the equivalent conditions in Theorem 3. Then f is defined by a left-to-right deterministic automaton with output if and only if all input strings u, v, w satisfy*

$$f(u|\underline{v}) = f(u|\underline{w}).$$

Proof. The left-to-right implication is immediate. For the right-to-left implication, we observe that the assumption implies that

$$f(\underline{u}|a|\underline{v})$$

does not depend on v , but only on f_{u_-} and the letter a . Furthermore, since f satisfies the assumptions from Theorem 3, the above value is of the form

$$\text{left } \overset{x}{\text{right}},$$

where each block is possibly missing. After reading input u , the automaton stores in its control state the derivative f_{u_-} . When it reads a letter a , it updates its control state, and outputs the string w , which depends only on the control state and input letter a .

5 First-order definable transducers

In this section we consider first-order definable transducers. Recall that when coding a string as a relational structure, we have a predicate for the order. We underline this because, unlike for MSO, for first-order logic order is more powerful than successor. The notion of first-order definability makes sense for:

- *languages*: there is a first-order formula that is true in the strings from the language and false in strings from outside the language.
- *regular colourings*: the language encoding is first-order definable.
- *string-to-string functions with origin information*: the same definition as for MSO-definable ones, except that set quantification is disallowed.

Theorem 5. *The following conditions are equivalent for a regular string-to-string function f with origin information.*

1. *it is definable by a first-order string-to-string transduction.*
2. *the colourings $w \mapsto f_{w_-}$ and $w \mapsto f_{-w}$ are first-order definable.*
3. *for every letters a, b , the following is a first-order definable colouring*

$$(u, v, w) \mapsto f(\underline{u}|a|\underline{v}|b|\underline{w})$$

Before proving the theorem, we observe that condition in item 2 is effective. Using a straightforward extension of the the Schützenberger-McNaughton-Papert Theorem, one can decide if a regular colouring is first-order definable. By applying the decision procedure to the functions $w \mapsto f_{w_-}$ and $w \mapsto f_{-w}$, we can decide if a regular string-to-string function with origin semantics is first-order definable. It is unclear if this sheds any light for the analogous question for semantics without origins.

Without origin information, a variant of first-order definable transducers was considered in [MSTV06], namely the transducers which are first-order definable in the sense of Theorem 5 and simultaneously order preserving in the sense

of Theorem 3. For instance, the doubling transducer $w \mapsto ww$ is first-order definable in the sense of Theorem 5, but not in the sense of [MSTV06], because it is not order preserving. By testing for both Condition 2 from Theorem 5 and Condition 2 of Theorem 3, we get an effective characterisation of the origin version of the transducers from [MSTV06].

6 Further work

Preliminary work indicates that the ideas in this paper extend to MSO-definable tree-to-tree transducers; this should be followed up. Another direction for further study is the computational complexity of equivalence with respect to origin semantics; in particular finding models for which equivalence is polynomial time.

References

- [AC10] Rajeev Alur and Pavol Cerný. Expressiveness of streaming string transducers. In *FSTTCS*, pages 1–12, 2010.
- [Ang87] Dana Angluin. Learning regular sets from queries and counterexamples. *Inf. Comput.*, 75(2):87–106, 1987.
- [AU70] Alfred V. Aho and Jeffrey D. Ullman. A characterization of two-way deterministic classes of languages. *J. Comput. Syst. Sci.*, 4(6):523–538, 1970.
- [BE00] Roderick Bloem and Joost Engelfriet. A comparison of tree transductions defined by monadic second order logic and by attribute grammars. *J. Comput. Syst. Sci.*, 61(1):1–50, 2000.
- [CJ77] Michal Chytil and Vojtech Jákl. Serial composition of 2-way finite-state transducers and simple programs on strings. In *ICALP*, pages 135–147, 1977.
- [EH01] Joost Engelfriet and Hendrik Jan Hoogeboom. MSO definable string transductions and two-way finite-state transducers. *ACM Trans. Comput. Log.*, 2(2):216–254, 2001.
- [EM03] Joost Engelfriet and Sebastian Maneth. Macro tree translations of linear size increase are MSO definable. *SIAM J. Comput.*, 32(4):950–1006, 2003.
- [FGRS13] Emmanuel Filiot, Olivier Gauwin, Pierre-Alain Reynier, and Frédéric Servais. From two-way to one-way finite state transducers. In *LICS*, pages 468–477. IEEE Computer Society, 2013.
- [Gur82] Eitan M. Gurari. The equivalence problem for deterministic two-way sequential transducers is decidable. *SIAM J. Comput.*, 11(3):448–452, 1982.
- [LMN10] Aurélien Lemay, Sebastian Maneth, and Joachim Niehren. A learning algorithm for top-down XML transformations. In Jan Paredaens and Dirk Van Gucht, editors, *PODS*, pages 285–296. ACM, 2010.
- [MSTV06] Pierre McKenzie, Thomas Schwentick, Denis Thérien, and Heribert Vollmer. The many faces of a translation. *J. Comput. Syst. Sci.*, 72(1):163–179, 2006.
- [Tho97] W. Thomas. Languages, automata, and logic. In *Handbook of Formal Language Theory*, volume III, pages 389–455. Springer, 1997.
- [vDKT93] Arie van Deursen, Paul Klint, and Frank Tip. Origin tracking. *J. Symb. Comput.*, 15(5/6):523–545, 1993.
- [WK94] Andreas Weber and Reinhard Klemm. Economy of description for single-valued transducers. In *STACS*, pages 607–618, 1994.

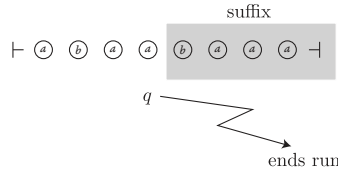
A Proof of Theorem 1

In this part of the appendix, we prove Theorem 1, which says that for a string-to-string function with origin information, the following conditions are equivalent

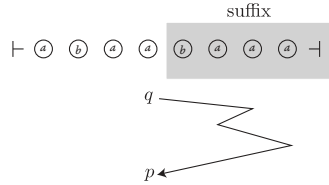
1. f is regular string-to-string transducer with origin information;
2. f has finitely many left derivatives and finitely many right derivatives;
3. for every a in the input alphabet, the following is a regular colouring

$$(v, w) \mapsto f(\underline{v}|a|\underline{w}).$$

Implication from 1 to 2. To show that there are finitely many derivatives in a regular transducer, suppose that f is a function that is recognised by a two-way deterministic automaton with output. When the head of this automaton enters a suffix of the input from the left in some state q , then several things can happen: it might not return from that suffix as in the following picture



or it might return from that suffix in another state p as in the following picture



In either case, whether it returns or not, the automaton can output an empty or nonempty string. Define the *suffix type* of a string to be the partial function

$$Q \rightarrow (Q \cup \{\text{"no return"}\}) \times \{\text{"empty"}, \text{"nonempty"}\},$$

which says what happens for each state q . It is not difficult to see that $f(x|y)$ depends only on the suffix type of y . Since there are finitely many suffix types, there are also finitely many right derivatives. For the left derivatives, an almost symmetric argument works. The symmetry is broken because a two way automaton begins in the leftmost position, so the notion of prefix type also needs to account for the first time that a prefix is exited by the head.

Implication from 2 to 3.

Lemma 1. *Let w_1, \dots, w_n and w, v be input strings. Then*

$$f_{w_-} = f_{v_-} \quad \text{implies} \quad f(\underline{w}|w_1|\dots|w_n) = f(\underline{v}|w_1|\dots|w_n)$$

Proof. By using the origin information, the value

$$f(\underline{w}|w_1|\cdots|w_n)$$

can be obtained from the value

$$f(\underline{w}|w_1\cdots w_n) \stackrel{\text{def}}{=} f_{w_-}(w_1\cdots w_n).$$

Lemma 2. *The functions $w \mapsto f_{w_-}$ and $w \mapsto f_{-v}$ are regular colourings.*

Proof. The images of the two functions are finite by the assumption 2, which says that there are finitely many left and right derivatives. By symmetry, we only study the left derivatives f_{w_-} . By Lemma 1, it follows that

$$f_{w_-} = f_{v_-} \quad \text{implies} \quad f(\underline{w}|a|u) = f(\underline{v}|a|u)$$

holds for every input letter a and every input string u . Since $f(\underline{w}|a|u)$ uniquely determines $f(\underline{wa}|u)$, it follows that

$$f_{w_-} = f_{v_-} \quad \text{implies} \quad f_{wa_-} = f_{va_-}$$

for every input letter a . This means that the set of left derivatives can be equipped with a transition function as in a deterministic left-to-right automaton, so that after reading a string w from the state f_{ϵ_-} , the automaton ends up in state f_{w_-} .

Thanks to Lemma 1 and its symmetric version for right derivatives,

$$(v, w) \mapsto f(\underline{v}|a|w)$$

factors through the function

$$(v, w) \mapsto (f_{v_-}, f_{-w}),$$

meaning that when a is fixed then equal results for the second function imply equal results for the first function. The second function is a regular colouring by Lemma 2. A function which factors through a regular colouring must itself be a regular colouring, which finishes the proof of item 3 in Theorem 1.

Implication from 3 to 1. In the proof, we use a two-way model called a *lookaround transducer*, which is defined as follows. The control is given by two deterministic automata: a *past automaton*, which is left-to-right deterministic, and a *future automaton*, which is right-to-left deterministic. There is a set of *registers*, with a designated *output register*. The registers are updated by a *register update function*, which inputs a state of the past automaton, an input letter, and a state of the future automaton, and outputs a sequence of register operations.

The output of the transducer on a string over the input alphabet is defined as follows. Define the *type* of a position i in the input string to be: the state of the past automaton after doing a left-to-right pass from the beginning of the

string up to and including position $i - 1$; the label of position i ; the state of the future automaton after doing a right-to-left pass from the end of the string up to and including position $i + 1$. To each type, the register update function assigns a register update. The automaton begins with all registers empty, and then it executes the register updates corresponding to the types of all positions in the string, read from left to right. After all of these register updates are executed, the value of the function is found in the output register.

Lookaround transducers define exactly the regular transducers, also under the origin semantics. For instance, it is not difficult to see that a lookaround transducer can be simulated in MSO.

To prove item 1 in Theorem 1, we construct a lookaround transducer based on the assumption that the characteristic function of f is a regular colouring. We begin by describing the registers. The transducer has a register for every left block in every partial output $f(\underline{x}|y)$. A partial output $f(v|\underline{w})$ can be interpreted as a register valuation, which is defined on the left blocks of $f(\underline{v}|\underline{w})$ and undefined on all other registers. The transducer is designed to satisfy the following invariant: when it has finished processing a prefix v of an input vw , then its register valuation is $f(v|\underline{w})$.

It remains to define the past and future automata, as well as the register update function. Our assumption, namely item 3, says that for every letter a in the input alphabet, the function

$$(v, w) \mapsto f(\underline{v}|a|\underline{w})$$

is a regular colouring. This means that there is a left-to-right deterministic automaton, which we can choose to be the past automaton, and a right-to-left deterministic automaton, which we can choose to be the future automaton, such that $f(\underline{v}|a|\underline{w})$ depends only on the state of the past automaton after reading v , the input letter a , and the state of the future automaton after reading w from right-to-left. The following lemma shows that the register update function can be defined to satisfy the invariant.

Lemma 3. *Based on $f(\underline{v}|a|\underline{w})$, one can construct a sequence of register operations which transforms the register valuation corresponding to $f(v|\underline{aw})$ into the register valuation corresponding to $f(va|\underline{w})$*

Proof. For every left block b in $f(va|\underline{w})$ there is a corresponding sequence b' of left and middle blocks in $f(v|\underline{a}|\underline{w})$. The register update required by the lemma is defined according to this sequence: for every left block b in $f(va|\underline{w})$, its value is defined to be the concatenation, according to the sequence b' , of the values of the left and middle blocks in $f(v|\underline{aw})$ and $f(\underline{v}|a|\underline{w})$, respectively.

B Learning

In this part of the appendix, we prove Theorem 2, which says that a regular string-to-string transducer with origin information can be learned using value

and equivalence queries in polynomial time (both number of queries and computation time) in terms of the number of left and right derivatives, and the size of the counterexamples given by the teacher.

By Theorem 1, learning a transducer with origin information f is the same as learning the characteristic function

$$(v, w) \mapsto f(\underline{v}|a|\underline{w}).$$

We simply use the Angluin algorithm for the language encoding of the characteristic function. To be more exact, in order to get the optimal complexity we use a variant of the language encoding, with the string w written in reverse, namely the language

$$L_f \stackrel{\text{def}}{=} \{v\#a\#w^R\#f(\underline{v}|a|\underline{w})\}.$$

It is not difficult to see that the minimal deterministic automaton of the language L_f is polynomial in the parameters from the statement of the lemma. The reason why the string w is reversed is that the automaton for the right derivatives reads its input from right to left.

Therefore, we can apply the Angluin algorithm to learn the language L_f . The only technical issue is that the queries for learning L_f need to be translated into the queries for learning f . A membership query

$$u \stackrel{?}{\in} L_f$$

corresponds to a value query for the transducer f , as follows. If the string w does not have the right format – exactly three appearances of $\#$, with exactly one letter between the second and third appearance – then the learner can immediately answer “no” without bothering the teacher. Otherwise, the learner extracts the (v, a, w) stored in u , and asks for the value of $f(vaw)$. Using the origin information, learner computes the value $f(\underline{v}|a|\underline{w})$, and can thus determine if u belongs to L_f . The correspondence between equivalence queries and counterexamples is done in a similar fashion.

C First-order definable transducers

In this part of the appendix, we prove Theorem 5, which says that the following conditions are equivalent for a regular string-to-string transducer f with origin information.

1. it is definable by a first-order string-to-string transduction.
2. the colorings $w \mapsto f_{w_}$ and $w \mapsto f_{_w}$ are first-order definable.
3. for every letters a, b , the following is a first-order definable coloring

$$(u, v, w) \mapsto f(\underline{u}|a|\underline{v}|b|\underline{w})$$

Implication from 1 to 2. By symmetry, it suffices to show that $w \mapsto f_{w_-}$ is first-order definable. Since there are finitely many possible values of f_{w_-} , there is a finite set U of test strings such that

$$f_{v_-} = f_{w_-} \quad \text{iff} \quad f_{v_-}(u) = f_{w_-}(u) \text{ for every } u \in U.$$

It is therefore sufficient to show that for every test string, one can compute in first-order logic, given w , the value of f_{w_-} on the test string. This is done in the following lemma.

Lemma 4. *If f is first-order definable, then for every string v the function $w \mapsto f(\underline{w}|v)$ is first-order definable.*

Proof. Define the *colored version* of an alphabet to be two disjoint copies: one called the black version, and one called the red version. Define f' to be a transducer which inputs a string over the colored version of the input alphabet, and works the same way as f , except that it outputs strings over the colored version of the output alphabet, and the color of an output letter is inherited from the corresponding input letter. Consider a function which inputs a string w over the original input alphabet, and outputs a black version of w , followed by a red version of v . We denote this function by $w \mapsto wv$. It is not difficult to see that both functions described above are first-order definable transducers, and since these are closed under composition, it follows that $w \mapsto f'(wv)$ is first-order definable. For every possible value x of $f(\underline{w}|v)$, one can write a first-order formula φ_x such that

$$f(\underline{w}|v) = x \quad \text{iff} \quad f'(wv) \models \varphi_x \quad \text{for every } w.$$

The property on the right side of the equivalence can be checked by a first-order formula working on w .

Implication from 3 to 1. For the moment, we skip the implication from 2 to 3, which is the most difficult⁴. We begin by showing that, under the assumption of item 3, the characteristic function of the transducer is first-order definable. For every letter b of the input alphabet, function

$$(v, w') \mapsto f(v|a|w'b)$$

must be first-order definable since it factors through the function

$$(v, w') \mapsto f(v|a|w'|b|\epsilon),$$

which is first-order definable. It follows that, as long as we know that w ends with b , then we can use first-order logic to obtain $f(v|a|w)$. If w is nonempty, then it has finitely many possibilities for the last letter, and therefore, we can use

⁴ The implication from 1 to 3 can be proved in the same way as in Lemma 4. Therefore the theorem with only items 1 and 3 would be easier to prove. Such a weaker theorem would still give an effective characterisation of first-order definable transduction.

first-order logic to obtain $f(v|a|w)$ as long as we know that w is nonempty. The same argument works when v is nonempty. The remaining case is when both v and w are empty, which can be detected in first-order logic.

We are now ready to define the first-order transduction. Define the *contribution* of a position in an input to be the subsequence (a string) of the output which originates from that position. The contribution depends only on $f(v|a|w)$, where v is the part of the input before the position, a is the label of the position, and w is the part of w after the position.

Let N be the maximal length of a contribution, ranging over all finitely many choices of the value of $f(v|a|w)$. The first-order interpretation defining f will copy each position of the input at most N times.

Given an input $a_1 \cdots a_n$, the first-order interpretation works as follows.

- **The universe formula.** For $i \in \{1, \dots, N\}$, the i -th copy of a position x in the input string belongs to the universe of the output string if and only if i is at most the length of the contribution of the x -th letter in the input string. For fixed i , this can be determined by a first-order formula with a free variable x , because the characteristic function is first-order definable.
- **The label formulas.** The label formulas are defined in the same way as the domain formula, only using the label of the i -th contribution.
- **The order formula.** For every $i, j \in \{1, \dots, N\}$, we need a first-order formula with two free variables x and y which says if, in the output, the i -th letter in the contribution of position x comes before the j -th letter in the contribution of position y . Assuming that x is before y , this information is entirely determined by partial output

$$f(\underline{a_1 \cdots a_{x-1}} | a_x | \underline{a_{x+1} \cdots a_{y-1}} | a_y | \underline{a_{y+1} \cdots a_n}),$$

which can be obtained using first-order logic thanks to the assumption. The case when x comes after y is symmetric, and in the special case when $x = y$ the formula simply returns the value of $i \leq j$.

Implication from 2 to 3. We say that a regular coloring g of n -tuples of strings is *aperiodic on coordinate i* if for every strings

$$w_1, w_2, \dots, w_{i-1}, x, y, z, w_{i+1}, \dots, w_n,$$

the following function is ultimately constant

$$i \mapsto g(w_1, \dots, w_{i-1}, xy^i z, w_{i+1}, \dots, w_n).$$

A straightforward consequence of the Schützenberger-McNaughton-Papert theorem is that a regular colouring is first-order definable if and only if it is aperiodic on every coordinate. Therefore, item 3 will follow once we show that for every letters a, b the function

$$(u, v, w) \mapsto (\underline{u} | a | \underline{v} | b | \underline{w})$$

is aperiodic on every coordinate. We begin with the first coordinate. Let then x, y, z, v, w be strings. We need to show that the function

$$i \mapsto f(xy^i z | a | \underline{v} | b | \underline{w})$$

is ultimately constant. For fixed a, v, b, w , the function above factors through

$$i \mapsto f_{xy^i z},$$

which must be ultimately constant by the assumption in item 2, and therefore it must also be ultimately constant. A symmetric argument shows that the function from the statement of the lemma is aperiodic on the third coordinate. We are left with the second coordinate. Let then u, x, y, z, w be strings. We need to show that the function

$$i \mapsto f(\underline{u} | a | \underline{xy^i z} | b | \underline{w}) \tag{1}$$

is ultimately constant.

Lemma 5. *Suppose that*

$$f(\underline{uv} | \underline{w}) = f(\underline{u} | \underline{w}) = f(\underline{u} | \underline{vw})$$

Then the function

$$i \mapsto f(\underline{u} | \underline{v^i} | \underline{w})$$

is ultimately constant.

Before proving the lemma, we show how it completes the proof of the implication from item 2 to item 3 in Theorem 5, and therefore also completes the proof of the theorem itself. Our goal is to show that the function (1) is ultimately constant. We will show that for sufficiently large j , the function

$$i \mapsto f(\underline{u} | a | \underline{xy^j} | \underline{y^i} | \underline{y^j z} | b | \underline{w}) \tag{2}$$

is ultimately constant. This will imply that the function

$$i \mapsto f(\underline{u} | a | \underline{xy^{i+2j}} | b | \underline{w})$$

is ultimately constant, and therefore also (1) is ultimately constant. We claim that the function (2) factors through the following functions

$$i \mapsto f(\underline{u} | a | \underline{xy^j} | \underline{y^{i+j}} | z | b | \underline{w}) \tag{3}$$

$$i \mapsto f(\underline{uax} | \underline{y^j} | \underline{y^i} | \underline{y^j z} | b | \underline{w}) \tag{4}$$

$$i \mapsto f(\underline{uax} | \underline{y^{j+i}} | \underline{y^j z} | b | \underline{w}). \tag{5}$$

Indeed, the value of (2) is obtained from (4) by replacing the left part with the first three parts in (3), and replacing the right part with the last three parts in (5). The function (3) factors through

$$i \mapsto f_{\underline{y^{i+j} z} | b | \underline{w}}$$

and is therefore ultimately constant by the assumption that $w \mapsto f_{\underline{w}}$ is aperiodic. For the same reason, the function (5) is ultimately constant. We are only left with showing that (4) is ultimately constant. If j is large enough, then by aperiodicity of $w \mapsto f_{\underline{w}}$, we see that

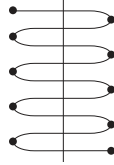
$$f_{uaxy^j} = f_{uaxy^{j+1}} \quad \text{and} \quad f_{\underline{y^j zbw}} = f_{\underline{y^{j+1} zbw}}$$

which implies that the assumptions of Lemma 5 for

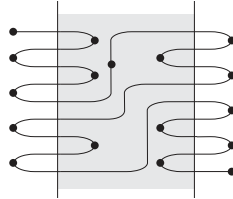
$$u = uaxy^j \quad v = y \quad w = y^j zbw.$$

The conclusion of the lemma shows that (4) is ultimately constant.

Proof (of Lemma 5). Consider the partial output $f(\underline{u}|w)$, which consists of left and right blocks in alternation. This output can be viewed as a graph (which consists of a single directed path), call it G_0 , which is illustrated in the following picture.



The vertices (black dots in the picture) in the left column correspond to left blocks, the vertices in the right column correspond to right blocks. There is a directed edge from a block to the following block; the edges in the picture are implicitly directed so that the path goes from top to bottom. Now consider the graph, call it G_1 , which corresponds to the partial output $f(\underline{u}|v|w)$, which is illustrated in the following picture.



The middle blocks are the vertices in the grey area. Thanks to the assumption

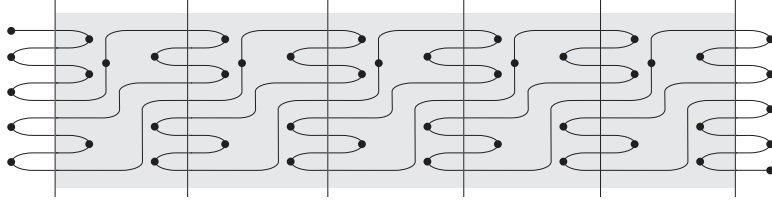
$$f(\underline{u}|vw) = f(\underline{u}|w),$$

the left blocks are visited in the same order as in G_0 . Using the other assumption, the right blocks are visited in the same order as in G_0 . However, the combined order on both left and right blocks might be different in G_0 and G_1 .

Finally, consider the graph, call it G_i , which corresponds to

$$f(\underline{u}|\overbrace{v|v|\cdots|v}^{i \text{ times}}|w),$$

which is illustrated below for $i = 5$.



Notice that the nodes in the left part of the graph, which correspond to the left blocks, are the same in every graph G_i , the same holds for the right blocks; also the order on left blocks and the order on right blocks are the same. The middle blocks of $f(\underline{u}|\underline{v}^i|\underline{w})$ correspond to maximal paths which are entirely contained in the grey area. For a left or right block x and a number i , consider a vertex $\sigma_i(x)$ and a boolean value $\tau_i(x)$, defined as follows.

- $\sigma_i(x)$ is the first left or right block after x in the graph G_i (if it exists).
- $\tau_i(x)$ says if the path in G_i from x to $\sigma_i(x)$ passes through a vertex in the grey area.

The sequence of blocks in $f(\underline{u}|\underline{v}^i|\underline{w})$ consists of the left and right blocks listed according to the function σ_i , with a middle block appended after those blocks x for which $\tau_i(x)$ says yes. To prove the lemma, it therefore suffices to show that for large enough i , the function σ_i is always the same, likewise for τ_i .

Define $loop_i$ to be the same as σ_i , but with its domain restricted to blocks x such that x and $\sigma_i(x)$ have the same type (meaning both are left blocks, or both are right blocks).

Lemma 6. *The function $loop_i$ is the same for large enough i , likewise for τ_i .*

Proof. If in the graph G_i there is a path which connects two blocks on the same type and does not pass through blocks of the other type, then same path is present in G_j . This means that if $i < j$ and $loop_i(x)$ is defined, then $loop_i(x) = loop_j(x)$. Therefore, the functions $loop_i$ stabilise eventually. A similar argument works for τ_i .

The following lemma finishes the proof.

Lemma 7. *The function σ_i can be uniquely determined from $loop_i$.*

Proof. As we have observed, the ordering on left blocks does not depend on i , likewise for right blocks. In other words, there is a unique ordering $<_L$ on left blocks and a unique ordering $<_R$ on right blocks. The function σ_i is a successor function (i.e. a function that maps every element, except one, to a successor so that a linear order is formed) which satisfies the following conditions:

- the ordering induced by σ_i on left blocks is $<_L$.
- the ordering induced by σ_i on right blocks is $<_R$.
- σ_i extends $loop_i$.
- $\sigma - loop_i$ maps left blocks to right blocks and vice versa.

It is not difficult to see that for every $loop_i$, there is at most one such function.

This completes the proof of Lemma 5.