

# Future-looking logics on data words and trees

Diego Figueira and Luc Segoufin

INRIA, ENS Cachan, LSV

**Abstract.** In a data word or a data tree each position carries a label from a finite alphabet and a data value from an infinite domain.

Over data words we consider the logic  $\text{LTL}_1^\downarrow(\mathbb{F})$ , that extends  $\text{LTL}(\mathbb{F})$  with one register for storing data values for later comparisons. We show that satisfiability over data words of  $\text{LTL}_1^\downarrow(\mathbb{F})$  is already not primitive recursive. We also show that the extension of  $\text{LTL}_1^\downarrow(\mathbb{F})$  with either the reverse modality  $\mathbb{F}^{-1}$  or with one extra register is undecidable. All those lower bounds were already known for  $\text{LTL}_1^\downarrow(\mathbb{X}, \mathbb{F})$  and our results essentially show that the  $\mathbb{X}$  modality was not necessary.

Moreover we show that over data trees similar lower bounds hold for certain fragments of  $\text{XPath}$ .

## 1 Introduction

A data word (data tree) is a word (tree) where each position carries a label from a finite alphabet and a *datum* from some infinite domain. These models have been considered in the realm of semistructured data [4], timed automata [6] and extended temporal logics [10,9,13]. In this paper we consider an infinite domain with no structure where we can only test for equality or inequality between two elements.

There have been various other logics considered to specify properties over data words and data trees. For example from the standpoint of Temporal Logics (both on data words [9] and trees [13]), of First Order Logics (see [5] for data words case and [4] for trees), or of logic based on tree patterns [8,3]. The logic  $\text{LTL}_1^\downarrow(\mathbb{X}, \mathbb{F})$  is the extension of  $\text{LTL}(\mathbb{X}, \mathbb{F})$  with the ability to use one *register* for storing a data value for later comparisons. It has been studied in [10,9] where satisfiability and expressivity issues have been addressed. In [9] it has been established that the satisfiability problem for  $\text{LTL}_1^\downarrow(\mathbb{X}, \mathbb{F})$  is decidable and not primitive recursive on data words. It is also shown in [9] that the *two way* extension  $\text{LTL}_1^\downarrow(\mathbb{X}, \mathbb{F}, \mathbb{F}^{-1})$  is undecidable over data words and, similarly, that the extension to 2 registers  $\text{LTL}_2^\downarrow(\mathbb{X}, \mathbb{F})$  is undecidable.

Here we show that even without the  $\mathbb{X}$  modality, all the lower bounds aforementioned remain valid: the satisfiability problem for  $\text{LTL}_1^\downarrow(\mathbb{F})$  over data words is not primitive recursive, while for  $\text{LTL}_1^\downarrow(\mathbb{F}, \mathbb{F}^{-1})$  and  $\text{LTL}_2^\downarrow(\mathbb{F})$  is undecidable.

Note that in the absence of data values the satisfiability problem for  $\text{LTL}(\mathbb{F})$  is NP-complete whereas it is PSPACE-complete for  $\text{LTL}(\mathbb{X}, \mathbb{F})$  [16]. Our work shows that in the presence of data values, the  $\mathbb{X}$  modality does not imply further computational complexity.

Data trees can be seen as a coding of an XML document [4,13]. Therefore XPath, the node selecting language of W3C [7], can be considered as a logic over data trees. By XPath, we refer to the 1.0 specification without all the domain specific features (arithmetic, string manipulation, etc.) As XPath is at the core of many XML standard languages (like XQuery and XSLT), deciding satisfiability of some of its fragments can be of great help during optimization stages.

A data word can be seen as a special case of an unranked ordered data tree, for instance by adding a root that is the parent of all the positions of the data word. With this consideration the fragment XPath( $\rightarrow, \rightarrow^+, =$ ) of XPath that contains only the axis `next-sibling` ( $\rightarrow$ ) and `following-sibling` ( $\rightarrow^+$ ) can be seen as a fragment of  $LTL_1^\downarrow(X, F)$ . There are nonetheless two important differences between the expressive power of XPath( $\rightarrow, \rightarrow^+, =$ ) and  $LTL_1^\downarrow(X, F)$ . The first one is that the axis  $\rightarrow^+$  corresponds to the strict future modality in LTL, denoted  $F_s$  in the sequel. Of course  $F$  can be defined using  $F_s$  but the opposite is not true in the absence of  $X$ . The second difference lies in the fact that XPath can compare data values in a way strictly more limited than  $LTL_1^\downarrow$ . We illustrate this with an example. At a position of a data word,  $LTL_1^\downarrow(F)$  can store the current data value in the register, check for a later symbol  $a$  with a data value different from the one in the register, and then further check for a symbol  $b$  with a data value matching the one of the register. In this spirit, XPath could only compare the data values at the beginning and the end of the path and could not say anything about the data value of the intermediate  $a$  symbol. Hence XPath( $\rightarrow^+, =$ ) should be seen as a fragment of  $LTL_1^\downarrow(F_s)$ , incomparable with  $LTL_1^\downarrow(F)$  in terms of expressive power.

Based on these ideas, we exhibit a syntactic fragment  $sLTL_1^\downarrow(F_s)$  of  $LTL_1^\downarrow(F_s)$  that limits the use of register comparisons and has the same expressive power than XPath( $\rightarrow^+, =$ ) on data words. We then show that satisfiability over data words of  $sLTL_1^\downarrow(F_s)$  is not primitive recursive. Similarly the same restriction on  $LTL_1^\downarrow(F_s, F_s^{-1})$  correspond to XPath( $\rightarrow^+, \leftarrow^+, =$ ) and is shown to be undecidable.

Our non primitive recursive results are proved by a reduction from the emptiness problem of *faulty counter automata*. These are counter automata that may have *incrementing errors* in their counters during the run. Non-emptiness for this class of automata was proven to be decidable and not primitive recursive [18]. Our reduction will be centered in a strategy of using data values for coding – with some limitations– a *next step* move of this automaton. We show that the strict future modality together with the limited data comparison capabilities of  $sLTL_1^\downarrow(F_s)$  are sufficient for our coding. With the extra power of  $LTL_1^\downarrow(F)$  for comparing data values, we also show that strictness of the future modality can be avoided. Similar ideas are used for our undecidable results: The extra available expressive power being sufficient for forbidding the incrementing errors and thus to code emptiness of a Minsky Counter Automaton.

*Related work.* There are known complexity results concerning the satisfiability problem of several data-aware fragments of XPath on data trees. When all axes are present, XPath is undecidable [12]. When all vertical axes are present but in the absence of any horizontal axis, the status of the decidability of

$\text{XPath}(\downarrow, \downarrow^+, \uparrow, \uparrow^+, =)$  is not yet known [2]. In this paper we show that if decidable, the complexity of this fragment cannot be primitive recursive, even in the absence of the axis  $\downarrow$  and  $\uparrow$ .

In [4] decidability over data trees of the two-variable fragment of first order logic,  $\text{FO}_2(+1, \sim)$  is established. As a direct consequence a fragment of XPath, with no occurrence of `descendant/ancestor` nor `following/previous-sibling` in the paths occurring in expressions with a data value test, is shown to be decidable. On the opposite, in the present work, we focus on satisfiability under the absence of the *successor* axis.

In [13] it is shown that the emptiness problem for alternating automata with one register over data trees is decidable. As a direct consequence it is shown that  $\text{XPath}(\downarrow, \downarrow^+, \rightarrow, \rightarrow^+, =)$ , with some restriction on the expressions with a data value test, is decidable and not primitive recursive. As a consequence of the present work, the hardness result already holds for  $\text{XPath}(\downarrow^+, \rightarrow, =)$  and for  $\text{XPath}(\rightarrow^+, =)$ .

We finally remark that, nevertheless, the satisfiability of  $\text{XPath}(\downarrow, \downarrow^+, =)$  is “only” EXPTIME-complete [11].

## 2 Preliminaries

We fix a finite set  $\Sigma$  of *labels* and an infinite set  $D$  of *data values*. The models we consider are either data words or data trees. A data word  $\sigma$  over a finite alphabet  $\Sigma$  is a non-empty word of  $(\Sigma \times D)^*$ . We assume no structure on  $D$  and  $D$  will be used only to perform tests for (in)equalities. In our examples we will always use data values from  $\mathbb{N}$  seen as a *set* of numbers. The data trees we will be using are unranked and ordered. The *domain* of an unranked ordered tree is represented by a prefix-closed set  $T$  of elements from  $\mathbb{N}^*$  such that whenever  $n(i+1) \in T$  then  $ni \in T$ . The elements of  $T$  will be called *nodes* of the tree. A data tree  $t$  over  $\Sigma, D$  is a tree domain  $T$  together with a labeling function  $\lambda$  assigning an element of  $\Sigma \times D$  to any node of  $t$ . We use the standard terminology for trees such as descendant, ancestor, sibling, etc.

*LTL with registers* Our most expressive logic for data words is  $\text{LTL}_n^\downarrow(\text{F}, \text{X}, \text{F}^{-1}, \text{X}^{-1})$ , the Linear Temporal Logic with the freeze quantifier ( $\downarrow_i$ ), test predicate ( $\uparrow_i$ ) and next (X) and future (F) temporal operators together with their inverse modalities ( $\text{X}^{-1}, \text{F}^{-1}$ ). Sentences are defined:

$$\varphi ::= a \mid \downarrow_i \varphi \mid \uparrow_i \mid \neg \varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \mathcal{O} \varphi$$

where  $a$  is a symbol from a finite alphabet  $\Sigma$ ,  $i \in \{1, \dots, n\}$ ,  $\mathcal{O}$  ranges over  $\{\text{F}, \text{X}, \text{F}^{-1}, \text{X}^{-1}\}$ . We also use the *strict future* temporal operator  $\text{F}_s$  as a shortcut for  $\text{XF}$ . We will consider fragments of this logics as simple restrictions to certain temporal operators.

Given a data word  $\sigma$ , we write  $\sigma_i$  for the  $i$ th element (duple) of the word, and  $\pi_1, \pi_2$  for the projections on  $\Sigma$  and  $D$ . A *valuation* is defined as a partial map  $v : \{1, \dots, n\} \rightarrow D$ . The satisfaction relation  $\models$  is inductively defined (we omitted  $\text{X}^{-1}$  and  $\text{F}^{-1}$  for simplicity):



### 3.1 Lower bound for $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$ .

**Theorem 2.** *Satisfiability of  $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$  on data words is not primitive recursive.*

*Proof.* We exhibit a PTIME reduction from the non-emptiness of ICA to formulas of  $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$ . Let  $C = \langle \Sigma, Q, q_0, n, \delta, F \rangle$  be an ICA.

Let  $L = \{(\text{inc}_i)_{1 \leq i \leq n}, (\text{dec}_i)_{1 \leq i \leq n}, (\text{ifzero}_i)_{1 \leq i \leq n}\}$ , and  $\hat{\Sigma} = Q \times (\Sigma \cup \{\varepsilon\}) \times L \times Q$ . We construct a formula  $\varphi_C \in \text{sLTL}_1^\downarrow(\mathbb{F}_s)$  that is satisfied by a data word iff  $C$  accepts the word. We view a run of  $C$  of the form:

$$\langle q_0, v_0 \rangle \xrightarrow{a, \text{inc}_i} \langle q_1, v_1 \rangle \xrightarrow{b, \text{dec}_j} \langle q_2, v_2 \rangle \xrightarrow{b, \text{ifzero}_i} \langle q_3, v_3 \rangle \cdots$$

as a string in  $\hat{\Sigma}$ :

$$\langle q_0, a, \text{inc}_i, q_1 \rangle \langle q_1, b, \text{dec}_j, q_2 \rangle \langle q_2, b, \text{ifzero}_i, q_3 \rangle \cdots$$

The formula  $\varphi_C$  will force that any string that satisfies it codes a run of  $C$ . In order to do this,  $\varphi_C$  must ensure that

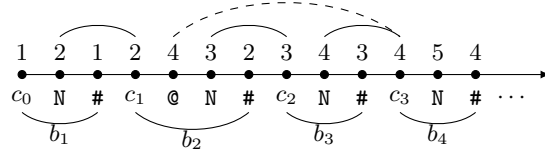
- (begin) the string starts with  $q_0$ ,
- (end) the string ends with a state of  $F$ ,
- (tran) every symbol of  $\hat{\Sigma}$  in the string corresponds to a transition of  $C$ ,
- (chain) the last component of a symbol of  $\hat{\Sigma}$  is equal to the first component of the next symbol,
- (pair) for each  $i$ , every symbol that contains  $\text{inc}_i$  occurring in the string to the left of a symbol containing  $\text{ifzero}_i$ , can be paired with a symbol containing  $\text{dec}_i$  and occurring in between the  $\text{inc}_i$  and the  $\text{ifzero}_i$ .

Before continuing let us comment on the (pair) condition. If we were coding runs of a perfect Minsky CA (ie, with no incremental errors), to the left of any position containing a  $\text{ifzero}_i$ , we would require a perfect matching between  $\text{inc}_i$  and  $\text{dec}_i$  operations in order to make sure that the value of the counter  $i$  is indeed zero at the position of the test. But as we are dealing with ICA, we only need to check that each  $\text{inc}_i$  has an associated  $\text{dec}_i$  to its right and before the test for zero, but we do not enforce the converse, that all  $\text{dec}_i$  match an  $\text{inc}_i$ . This is fortunate because this would require past navigational operators.

The first difficulty comes from the fact that (pair) is not a regular relation. The pairing will be obtained using data values. The second difficulty is to enforce (chain) without having access to the string successor relation. In order to simulate the successor relation we add extra symbols to the alphabet, with suitable associated data values.

Let  $\Sigma' = \hat{\Sigma} \cup \{\mathbb{N}, \#, \textcircled{\#}\}$ . The coding of a run consists in a succession of *blocks*. Each block is a sequence of 3 or 4 symbols, “ $c \mathbb{N} \#$ ” or “ $c \textcircled{\#} \mathbb{N} \#$ ”, with  $c \in \hat{\Sigma}$ . The data value associated to the  $c$  and  $\#$  symbols of a block is the same and uniquely determines the block: no two blocks may have the same data value. The data value associated with the symbol  $\mathbb{N}$  of a block is the data value of the next block. If a block contains a symbol  $c$  that codes a  $\text{inc}_i$  operation that is later in the string followed by a  $\text{ifzero}_i$ , then this block contains a symbol  $\textcircled{\#}$  whose data value is the same as the data value of the block containing  $\text{dec}_i$  that is paired with  $c$ .

For instance in the example of Fig. 1 one can see four blocks  $b_1, b_2, b_3, b_4$ . Each of them starts with a symbol from  $\hat{\Sigma}$  coding a transition of the ICA and ends with a # with the same data value marking the end of the block. Inside the block, the data value of N is the same as the data value of the next block. The data value of @ corresponds to that of a future block. In this example  $c_1$  must correspond to a  $\text{inc}_i$  while  $c_3$  to  $\text{dec}_i$  and there must be a  $\text{ifzero}_i$  somewhere to the remaining part of the word (say,  $b_5$ ). Moreover  $c_2$  can't be a  $\text{ifzero}_i$  as otherwise the data value of the symbol @ would refer to a block to the left of  $c_2$ .



**Fig. 1.** Coding of an ICA run.

We now show that the coding depicted above can be enforced in  $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$ . By  $\pi_1, \pi_2, \pi_3, \pi_4$  we denote the projection of each symbol of  $\hat{\Sigma}$  into its corresponding component. To simplify the presentation we use the following abbreviations:

$$\begin{aligned} \hat{\sigma} &\equiv \bigvee_{c \in \hat{\Sigma}} c & \text{INC}(i) &\equiv \bigvee_{c \in \hat{\Sigma}, \pi_3(c) = \text{inc}_i} c & \text{INC} &\equiv \bigvee_i \text{INC}(i) \\ \text{DEC}(i) &\equiv \bigvee_{c \in \hat{\Sigma}, \pi_3(c) = \text{dec}_i} c & \text{LAST} &\equiv \hat{\sigma} \wedge \neg(\rightarrow^+[\hat{\sigma}]) & \text{IZ}(i) &\equiv \bigvee_{c \in \hat{\Sigma}, \pi_3(c) = \text{ifzero}_i} c \end{aligned}$$

The formula  $\varphi_C$  that we construct is the conjunction of all the formula mentioned below.

#### *Forcing the structure*

$\text{G}(\text{LAST} \Rightarrow \neg \mathbb{F}_s \top)$  : The string ends with the last transition,

$\bigwedge_{c \in \{N, @, \#\}} \text{G}(c \Rightarrow \neg(\downarrow \mathbb{F}_s(c \wedge \uparrow)))$  : the data value associated to each N, # and @ uniquely determines the occurrence of that symbol,

$\text{G}((\hat{\sigma} \wedge \neg \text{LAST}) \Rightarrow (\downarrow \mathbb{F}_s(N \wedge \mathbb{F}_s(\# \wedge \uparrow)))$  : each occurrence of  $\hat{\Sigma}$  (except the last one) is in a block that contains a N and then a #,

$\bigwedge_i \text{G}((\text{INC}(i) \wedge \mathbb{F}_s(\text{IZ}(i))) \Rightarrow \downarrow \mathbb{F}_s(@ \wedge \mathbb{F}_s(N \wedge \mathbb{F}_s(\# \wedge \uparrow))))$  : every  $\text{inc}_i$  block to the left of a  $\text{ifzero}_i$  must have a @ before the N,

$\text{G}((\hat{\sigma} \wedge \neg \text{INC}) \Rightarrow \neg \downarrow \mathbb{F}(@ \wedge \mathbb{F}_s(\# \wedge \uparrow)))$  : only blocks  $\text{inc}$  are allowed to have a @,

$\bigwedge_{s \in \{N, @\}} \text{G}(\hat{\sigma} \Rightarrow \neg(\downarrow \mathbb{F}_s(s \wedge \mathbb{F}_s(s \wedge \mathbb{F}_s(\# \wedge \uparrow))))$  : there is at most one occurrence of N and @ in each block,

$\bigwedge_{s \in \hat{\Sigma} \cup \{\#\}} \text{G}(\hat{\sigma} \Rightarrow \neg \downarrow \mathbb{F}_s(s \wedge \mathbb{F}_s(\# \wedge \uparrow)))$  : there is exactly one symbol # and one symbol of  $\hat{\Sigma}$  per block,

$\text{G}(N \Rightarrow (\downarrow \mathbb{F}_s(\# \wedge \mathbb{F}_s(\hat{\sigma} \wedge \uparrow))))$  : the data value of each symbol N points to a block to its right,

$\text{G}(N \Rightarrow \neg \downarrow \mathbb{F}_s(\# \wedge \mathbb{F}_s(\# \wedge \mathbb{F}_s(\hat{\sigma} \wedge \uparrow))))$  : in fact N has to point to the next block.

Once the structure has the expected shape, we can enforce the run as follows. All the formulas below are based on the following trick. In a test of the form

$\varepsilon \rightarrow^+ [N] \rightarrow^+ [\#]$  which is typically performed at a position of symbol  $\hat{\Sigma}$ , the last symbol  $\#$  must have the same data value as the initial position. Hence, because of the structure above, both must be in the same block. Hence the middle symbol  $N$  must also be inside that block. From the structure we know that the data value of this  $N$  points to the next block. Therefore by adding in the middle bracket  $[N]$  a test of the form  $(\varepsilon \rightarrow^+ [s])$  we can transfer some finite information from the current block to the next one. This gives the desired successor relation.

*Forcing a run*

(begin) 
$$\bigvee_{c \in \hat{\Sigma}, \pi_1(c) = q_0} c$$

(end) 
$$F_s(\text{LAST} \wedge \bigvee_{c \in \hat{\Sigma}, \pi_4(c) \in F} c)$$

(tran) All elements used from  $\hat{\Sigma}$  correspond to valid transitions. Let  $\hat{\Sigma}^C$  be that set of transitions of  $C$ ,

$$G\left(\bigwedge_{c \in \hat{\Sigma} \setminus \hat{\Sigma}^C} \neg c\right)$$

(chain) For every  $c \in \hat{\Sigma}$ ,

$$G\left(c \Rightarrow (\downarrow F_s(N \wedge F_s(\# \wedge \uparrow)) \wedge \bigvee_{\substack{d \in \hat{\Sigma}, \\ \pi_4(c) = \pi_1(d)}} (\downarrow F_s(d \wedge \uparrow)))\right)$$

(pair) We must first make sure that the block of the  $\textcircled{0}$  of an  $\text{inc}_k$  is matched with a block of a  $\text{dec}_k$ :

$$\bigwedge_k G(\text{INC}(k) \Rightarrow (\neg(\downarrow F_s(\textcircled{0} \wedge F_s(\# \wedge \uparrow))) \vee (\downarrow F_s(\textcircled{0} \wedge (\downarrow F_s(\text{DEC}(k) \wedge \uparrow)) \wedge F_s(\# \wedge \uparrow))))))$$

Now, each time we are in a block with  $\text{inc}_k$  that afterward has a  $\text{ifzero}_k$  then:

- 1. The block must contain an  $\textcircled{0}$  element:

$$\bigwedge_k G(\text{INC}(k) \Rightarrow ((\downarrow F_s(\textcircled{0} \wedge F_s(\# \wedge \uparrow))) \vee \neg F_s(\text{IZ}(k))))$$

- 2. The data value of that  $\textcircled{0}$  element must point to a future block *before* any occurrence of  $\text{ifzero}_k$ :

$$\bigwedge_k G(\text{INC}(k) \Rightarrow \neg(\downarrow F_s(\textcircled{0} \wedge (\downarrow F_s(\text{IZ}(k) \wedge F_s \uparrow)) \wedge F_s(\# \wedge \uparrow))))$$

This concludes the construction of  $\varphi_C$ . The correctness of the construction is quite standard and can found in the appendix.  $\square$

### 3.2 Undecidability of $\text{sLTL}_1^\downarrow(F_s, F_s^{-1})$

We now consider  $\text{sLTL}_1^\downarrow(F_s, F_s^{-1})$  over data words. The extra modality can be used to code the run of a (non faulty) Minsky CA.

**Theorem 3.** *Satisfiability of  $\text{sLTL}_1^\downarrow(F_s, F_s^{-1})$  is undecidable over data words.*

*Proof.* Consider a Minsky CA  $C$ . We revisit the proof of Theorem 2. It is easy to see that we can enforce the absence of faulty increments during the run by asking that every  $\text{dec}_i$  element that appears is referenced by some previous  $\text{inc}_i$  block:  $\bigwedge_i G(\text{DEC}(i) \Rightarrow (\varepsilon = \text{+} \leftarrow [\text{Q}]))$ . In this way we make sure that there are no  $\text{dec}$  that are unrelated to a corresponding  $\text{inc}$  and then the coding is that of a *perfect* (non faulty) run.  $\square$

## 4 The case of $\text{LTL}_1^\downarrow$

In this section we lift the lower bounds of the previous section by considering the temporal operator  $F$  instead of  $F_s$ . We can only do so by removing at the same time the restriction to simple formulas. Hence the results of this section cannot be applied to  $\text{XPath}$ . Notice that  $\text{LTL}_1^\downarrow(F)$  and  $\text{sLTL}_1^\downarrow(F_s)$  are incomparable in term of expressive power. We do not know whether  $\text{sLTL}_1^\downarrow(F)$ , which is weaker than the two above mentioned logics, is already not primitive recursive. The results of this section improve the results of [9] which shows that satisfiability is not primitive recursive for  $\text{LTL}_1^\downarrow(X, F)$  and undecidable for  $\text{LTL}_1^\downarrow(X, F, F^{-1})$ .

**Theorem 4.** *Over data words,*

1. *Satisfiability of  $\text{LTL}_1^\downarrow(F)$  is decidable and not primitive recursive.*
2. *Satisfiability of  $\text{LTL}_1^\downarrow(F, F^{-1})$  is undecidable.*

*Proof.* We only prove Item 1. The proof of Item 2. is similar and can be found in the appendix.

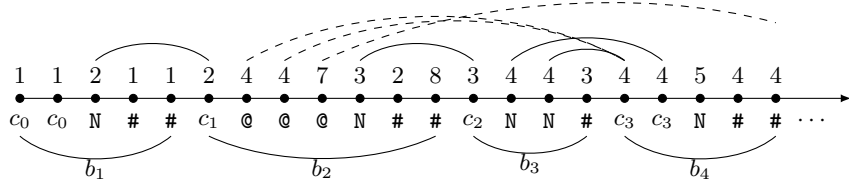
Consider an ICA  $C$  and recall the coding of runs of  $C$  used in the proof of Theorem 2. Note that all data formulas used are of the form  $\varepsilon = \alpha$  which are straightforward to code in  $\text{LTL}_1^\downarrow$ . In the construction of the formula  $\varphi_C$ , whenever we use  $\rightarrow^+$  with two different symbols on the left-hand side and the right-hand side, this axis can be equivalently replaced by the  $F$  temporal operator. This is the case in all formulas except in three places: (i) the formula saying that  $N$  should point to the next block contains  $[\#] \rightarrow^+ [\#]$ . But from the structure that is enforced, this can equivalently be replaced by  $[\#] \rightarrow^+ [N] \rightarrow^+ [\#]$ . (ii) To enforce that each block contains at most one occurrence of symbols in  $\hat{\Sigma} \cup \{N, \#, \text{Q}\}$ . (iii) To enforce that no two symbols in  $\hat{\Sigma} \cup \{N, \#, \text{Q}\}$  have the same data values.

In order to cope with (ii) and (iii), we use a slightly different coding for runs of  $C$ . This coding is the same as the one for the proof of Theorem 2 except that we allow succession of equal symbols, denoted as *group* in the sequel. Note that in a group two different occurrences of the same symbol in general may have different data values, as we can no longer enforce their distinctness. However, as we will see, we can enforce that the  $\hat{\Sigma}$  group of elements of a block have all the same data value.

Hence a block now is either a group of  $c \in \hat{\Sigma}$  followed by a group of  $N$  followed by a group of  $\#$  or the same with a group of  $\text{Q}$  in between. A coding of a run is depicted in Figure 2.

This structure is enforced by modifying the formulas of the proof of Theorem 2 as follows.





**Fig. 2.**  $LTL_1^\downarrow(F)$  cannot avoid having repeated consecutive symbols.

- (i) The formulas that limit the number of occurrences of symbols in a block are replaced by formulas *limiting the number of groups* in a block.
- (ii) The formulas requiring that no two occurrences of a same symbol may have the same data values are replaced by formulas requiring that no two occurrences of a same symbol *in different groups* may have the same data values.
- (iii) In all other formulas,  $\rightarrow^+$  is replaced by  $F$ .
- (iv) Finally, we ensure that although we may have repeated symbols inside a block, all symbols from  $\hat{\Sigma}$  have the same data value.

$$G(\hat{\sigma} \Rightarrow \neg \downarrow_1 F((\hat{\sigma} \vee \#) \wedge \neg \uparrow_1 \wedge F(\# \wedge \uparrow_1)))$$

Note that this implies that each  $N$  of a group must have the same data values as they all point to the next block. However there could still be  $@$  symbols with different data values as depicted in Fig. 2.

The new sentences now imply, for instance, that:

1. Every position from a group of  $c \in \hat{\Sigma}$  have the same data value which is later matched by a group of  $\#$ .
2. Every position from a group of  $N$  has the same data value as a position  $c \in \hat{\Sigma}$  of the *next* block (and then, it has only one data value).
3. Every position from a group  $@$  has the same data value as a position  $c \in \hat{\Sigma}$  of a block to its right. Note that the data values of two  $@$  of the same group may correspond to the data values of symbols in different blocks. This is basically the main conceptual difference with the previous proof.

The proof of correctness of the construction is left to the reader. □

Note that in the previous proof we used the fact that  $LTL_1^\downarrow(F)$ , although it has only one register, can make (in)equality tests several times throughout a path (as used in the formula of item (iv) in the proof) something that simple formulas and XPath can't do.

*Two registers.* When 2 registers are available the previous result can be adapted to code a (non faulty) Minsky CA. The proof can be found in the appendix.

**Theorem 5.** *Satisfiability of  $LTL_2^\downarrow(F)$  is undecidable over data words.*

## 5 Data trees and XPath

We now turn to data trees. An XML document can be seen as an unranked tree with *attributes* and data values in its nodes. While a data tree has only *one* data value per node, an XML document element may have *several* attributes each of which with an associated data value. All XPath fragments treated in this paper can force all elements of an XML document to have only one attribute. Therefore all hardness results in the present work hold also for the class of XML documents.

*The logic XPath.* XPath is a two-sorted language, with *path* expressions  $(\alpha, \beta, \dots)$  and *node* expressions  $(\varphi, \psi, \dots)$ . These are defined by mutual recursion:

$$\begin{aligned} \alpha &::= \downarrow | \downarrow^+ | \uparrow | \uparrow^+ | \rightarrow | \rightarrow^+ | \leftarrow | \leftarrow^+ | \varepsilon | \alpha\beta | \alpha \cup \beta | \alpha[\varphi] \\ \varphi &::= \sigma | \langle \alpha \rangle | \neg\varphi | \varphi \wedge \psi | \alpha = \beta | \alpha \neq \beta \quad (\sigma \in \Sigma) \end{aligned}$$

A path expression essentially describes a traversal of the tree by using the axis: **child** ( $\downarrow$ ), **descendant** ( $\downarrow^+$ ), the **next-sibling** ( $\rightarrow$ ), **following-sibling** ( $\rightarrow^+$ ) and their inverses, with the ability to test at any stage for node conditions. Let  $t$  be a data tree with domain  $T$  and labeling function  $\lambda$ . The semantics of XPath is defined by induction in the usual intuitive way, we only give here some important steps:

$$\begin{aligned} \llbracket \downarrow \rrbracket^t &= \{(x, xi) \mid xi \in T\} & \llbracket \alpha\beta \rrbracket^t &= \{(x, z) \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^t \wedge (y, z) \in \llbracket \beta \rrbracket^t\} \\ \llbracket \langle \alpha \rangle \rrbracket^t &= \{x \in T \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^t\} & \llbracket \alpha[\varphi] \rrbracket^t &= \{(x, y) \in \llbracket \alpha \rrbracket^t \mid y \in \llbracket \varphi \rrbracket^t\} \\ \llbracket \alpha^+ \rrbracket^t &= \text{the transitive closure of } \llbracket \alpha \rrbracket^t \\ \llbracket \alpha = \beta \rrbracket^t &= \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^t, (x, z) \in \llbracket \beta \rrbracket^t, \pi_2(\lambda(y)) = \pi_2(\lambda(z))\} \\ \llbracket \alpha \neq \beta \rrbracket^t &= \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^t, (x, z) \in \llbracket \beta \rrbracket^t, \pi_2(\lambda(y)) \neq \pi_2(\lambda(z))\} \end{aligned}$$

A key property for using results of Section 3 is that simple formulas can be translated to XPath and back. The proof of this result is straightforward by induction on the formula.

**Proposition 1.** *Over data words,  $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$  and  $\text{XPath}(\rightarrow^+, =)$  have the same expressive power. The same hold for  $\text{sLTL}_1^\downarrow(\mathbb{F}_s, \mathbb{F}_s^{-1})$  and  $\text{XPath}(\leftarrow^+, \rightarrow^+, =)$ . Moreover, in both cases, the transformation from  $\text{sLTL}_1^\downarrow$  to XPath takes polynomial time while it takes exponential time in the other direction.*

We remark that there is a big difference between  $\text{XPath}(\rightarrow^+, =)$  over data words and  $\text{XPath}(\downarrow^+, =)$  over data trees. Indeed  $\text{XPath}(\downarrow^+, =)$  is closed under bisimulation and hence it cannot assume that the tree is a vertical path. As the string structure was essential in the proof of Theorem 2, the non primitive recursiveness of  $\text{XPath}(\rightarrow^+, =)$  over data words does not lift to  $\text{XPath}(\downarrow^+, =)$  over data trees. Actually one can show that satisfiability of  $\text{XPath}(\downarrow^+, =)$  is ExpTime-complete [11]. However if one consider the logic  $\text{XPath}(\downarrow^+, \rightarrow, =)$  then the axis  $\rightarrow$  can be used to enforce a vertical path and therefore it follows from Theorem 2 and Proposition 1 that:

**Corollary 1.** *Satisfiability of  $\text{XPath}(\downarrow^+, \rightarrow, =)$  on data trees is not primitive recursive.*

Similarly in  $\text{XPath}(\downarrow^+, \uparrow^+, =)$  one can simulate a string by going down to a leaf using  $\downarrow^+$  and then use the path from that leaf to the root as a string using  $\uparrow^+$ . Hence we also have:

**Corollary 2.** *Satisfiability of  $\text{XPath}(\downarrow^+, \uparrow^+, =)$  on data trees is not primitive recursive.*

Note that the decidability of  $\text{XPath}(\downarrow^+, \rightarrow, =)$  and of  $\text{XPath}(\downarrow^+, \uparrow^+, =)$  is still an open problem.

**Open problem:** Is  $\text{XPath}(\downarrow^+, \uparrow^+, =)$  decidable over data trees?

It would be interesting to know whether the strictness of the axis  $\downarrow^+$  is necessary in the above two results. This boils down to know whether  $\text{sLTL}_1^\downarrow(\mathbb{F})$  is already not primitive recursive over data words. Note that the proof of Theorem 4 use in a essential way the possibility to make (in)equality tests several times throughout a path. This is exactly what cannot be expressed in  $\text{sLTL}_1^\downarrow(\mathbb{F})$ .

**Open problem:** Is satisfiability of  $\text{sLTL}_1^\downarrow(\mathbb{F})$  primitive recursive over data words?

We conclude with the following result which is a simple consequence of Theorem 3 and Proposition 1:

**Corollary 3.** *Satisfiability of  $\text{XPath}(\leftarrow^+, \rightarrow^+, =)$  and of  $\text{XPath}(\downarrow^+, \uparrow^+, \rightarrow, =)$  over data trees is undecidable.*

## 6 Discussion

By [9] it is known that satisfiability of  $\text{LTL}_1^\downarrow(\mathbb{X}, \mathbb{F})$  with infinite data words is undecidable, as well as  $\text{LTL}_2^\downarrow(\mathbb{X}, \mathbb{F})$ . The proof of Theorem 4 can be extended to code runs of ICA over infinite data words, which is known to be undecidable, to show that these results hold in the absence of  $\mathbb{X}$ .

**Theorem 6.** *On infinite data words, the satisfiability problem of  $\text{LTL}_1^\downarrow(\mathbb{F})$  is undecidable.*

### Summary of results

In the table below we summarize the main results and some of the consequences we have mentioned. In this table,  $\overline{\text{PR}}$  stands for not primitive recursive, decidability unknown while  $\text{R} \setminus \text{PR}$  stands for decidable and not primitive recursive.

Logic	Complexity	Details
$\text{LTL}_1^\downarrow(\mathbb{F})$	$\text{R} \setminus \text{PR}$	Thm 4 & [9]
$\text{LTL}_1^\downarrow(\mathbb{F}, \mathbb{F}^{-1})$	undecidable	Thm 4
$\text{LTL}_2^\downarrow(\mathbb{F})$	undecidable	Thm 5
$\text{sLTL}_1^\downarrow(\mathbb{F}_s)$	$\text{R} \setminus \text{PR}$	Thm 2 & [9]
$\text{sLTL}_1^\downarrow(\mathbb{F}_s, \mathbb{F}_s^{-1})$	undecidable	Thm 3
$\text{XPath}(\downarrow^+, \rightarrow, =)$	$\overline{\text{PR}}$	Cor 1
$\text{XPath}(\downarrow^+, \uparrow^+, =)$	$\overline{\text{PR}}$	Cor 2
$\text{XPath}(\downarrow^+, \uparrow^+, \rightarrow, =)$	undecidable	Cor 3

## References

1. Rajeev Alur and Thomas A. Henzinger. A really temporal logic. *J. ACM*, 41(1):181–204, 1994.
2. Michael Benedikt, Wenfei Fan, and Floris Geerts. XPath satisfiability in the presence of DTDs. *J. ACM*, 55(2), 2008.
3. Henrik Björklund, Wim Martens, and Thomas Schwentick. Optimizing conjunctive queries over trees using schema information. In *MFCS*, pages 132–143, 2008.
4. Mikolaj Bojańczyk, Claire David, Anca Muscholl, Thomas Schwentick, and Luc Segoufin. Two-variable logic on data trees and XML reasoning. In *PODS*, pages 10–19, 2006.
5. Mikolaj Bojańczyk, Anca Muscholl, Thomas Schwentick, Luc Segoufin, and Claire David. Two-variable logic on words with data. In *LICS*, pages 7–16, 2006.
6. Patricia Bouyer, Antoine Petit, and Denis Thérien. An algebraic approach to data languages and timed languages. *Inf. Comput.*, 182(2):137–162, 2003.
7. J. Clark and S. DeRose. XML path language (XPath). Website, 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
8. Claire David. Complexity of data tree patterns over XML documents. In *MFCS*, pages 278–289, 2008.
9. Stéphane Demri and Ranko Lazić. LTL with the freeze quantifier and register automata. In *LICS*, pages 17–26. IEEE Computer Society Press, 2006.
10. Stéphane Demri, Ranko Lazić, and David Nowak. On the freeze quantifier in constraint LTL: Decidability and complexity. In *TIME*, pages 113–121, 2005.
11. Diego Figueira. Satisfiability of downward XPath with data equality tests. In *PODS*, Providence, Rhode Island, USA, 2009. ACM Press.
12. Floris Geerts and Wenfei Fan. Satisfiability of XPath queries with sibling axes. In *DBPL*, volume 3774, pages 122–137. Springer, 2005.
13. Marcin Jurdziński and Ranko Lazić. Alternating automata on data trees and XPath satisfiability. *CoRR*, abs/0805.0330, 2008.
14. Alexei Lisitsa and Igor Potapov. Temporal logic with predicate lambda-abstraction. In *TIME*, pages 147–155, 2005.
15. Marvin L. Minsky. *Computation: finite and infinite machines*. Prentice-Hall, Inc., 1967.
16. A. Nakamura and H. Ono. On the size of refutation Kripke models for some linear modal and tense logics. *SL*, 39(4):325–333, 1980.
17. J. Ouaknine and J. Worrell. On Metric temporal logic and faulty Turing machines. In *FoSSACS*, volume 3921, pages 217–230, 2006.
18. Philippe Schnoebelen. Verifying lossy channel systems has nonprimitive recursive complexity. *Information Processing Letters*, 83(5):251–261, September 2002.

## A Missing proofs

### A.1 Correctness of the construction of the proof of Theorem 2

Given an ICA  $C$ , the proof exhibited a formula  $\varphi_C \in \text{XPath}(\rightarrow^*, =)$  that is satisfied by a data word iff  $C$  accepts the word. We show here that the construction is correct.

**From data word to an accepting ICA run** Let  $w$  be a data word satisfying  $\varphi_C$ . Let  $B_1 B_2 \cdots B_m$  be the decomposition of  $w$  into blocks as described above. Let  $c_i = \langle q_i, a_i, \ell_i, q_{i+1} \rangle$  be the element of  $B_i$  in  $\hat{\Sigma}$ . We show that  $w' = a_1 a_2 \cdots a_m$  is accepted by  $C$ . We construct a run of  $C$  on  $w'$  such that  $C$  is in state  $q_i$  before reading letter  $a_i$  of  $w'$  and execute the transitions  $c_i$ . By (**tran**) each one of these are valid transitions of the ICA, and (**chain**) ensures that the sequence is consistent. (**begin**) takes care of the initialization condition and (**end**) of the accepting condition. It remains to construct the valuation  $v_1, \dots, v_{m+1}$  of the counter at each step such that we have  $\langle q_i, v_i \rangle \xrightarrow{a_i, \ell_i} \langle q_{i+1}, v_{i+1} \rangle$  as desired.

As expected  $v_1$  is the zero valuation. We construct the rest of the valuations by induction, simulating a perfect counter machine and introducing incrementing errors in a lazy way, whenever necessary. At step  $i$ , given a transition  $\langle q_i, a_i, \ell_i, q_{i+1} \rangle$  and the current valuation  $v_i$  we set  $v_{i+1}$  such that  $\langle q_i, v_i \rangle \xrightarrow{a_i, \ell_i} \langle q_{i+1}, v_{i+1} \rangle$ .

- If  $\ell_i$  is **inc<sub>k</sub>**, then  $v_{i+1} := v_i[k \mapsto v_i(k) + 1]$  and clearly  $\langle q_i, v_i \rangle \xrightarrow{a_i, \ell_i} \langle q_{i+1}, v_{i+1} \rangle$ .
- If  $\ell_i$  is **dec<sub>k</sub>** then:
  - If  $v_i(k) > 0$  then  $v_{i+1} := v_i[k \mapsto v_i(k) - 1]$ .
  - If  $v_i(k) = 0$ , then  $v_{i+1} := v_i$ .

In both cases one can verify that  $\langle q_i, v_i \rangle \xrightarrow{a_i, \ell_i} \langle q_{i+1}, v_{i+1} \rangle$ , the transition being with an incrementing error of 1 in the second case.

- If  $\ell_i$  is **ifzero<sub>k</sub>**, we set  $v_{i+1} := v_i$ . To ensure that  $\langle q_i, v_i \rangle \xrightarrow{a_i, \ell_i} \langle q_{i+1}, v_{i+1} \rangle$  we must show that necessarily  $v_i(k) = 0$ . This is a consequence of condition **pair**. This condition enforces that each **inc<sub>k</sub>** must be paired with a unique **dec<sub>k</sub>** to its right and before the **ifzero<sub>k</sub>** test. Hence the lazy strategy above ensures that each increment is decremented later on and before the zero test.

**From an accepting ICA run to data word** Given an accepting run of the ICA  $\langle q_1, v_1 \rangle \xrightarrow{a_1, \ell_1} \langle q_2, v_2 \rangle \xrightarrow{a_2, \ell_2} \dots \xrightarrow{a_m, \ell_m} \langle q_{m+1}, v_{m+1} \rangle$  we construct a data word  $w$  verifying  $\varphi_C$  as follows. Let  $c_i$  be  $\langle q_i, a_i, \ell_i, q_{i+1} \rangle$ .  $w$  consists of a concatenation of  $m$  blocks  $B_1 B_2 \cdots B_m$  each one defined:

- If  $i = m$ ,  $B_i$  is the word  $c_i$  with data value  $m$ .
- If  $i < m$ , and  $\ell_i \in \{\text{dec}_k, \text{ifzero}_k\}$ , then set  $B_i$  to  $c_i \#$  with respective data values  $i, i + 1, i$ .

$$\langle c_i, i \rangle \langle \text{N}, i + 1 \rangle \langle \#, i \rangle$$

- If  $i < m$  and  $\ell_i$  is  $\text{inc}_k$ . Consider the minimal  $d > i$  such that  $v_{d+1}(k) = v_i(k)$  and set  $B_i$  to
  - $c_i \text{ @ N \#}$  with respective data values  $i, d, i + 1, i$  if  $d$  exists
  - $c_i \text{ N \#}$  with respective data values  $i, i + 1, i$  if not.

This construction assigns the unique data value  $i$  to each block  $B_i$ , and set the data value of each symbol  $\text{N}$  to the data value of the next block  $B_{i+1}$ . The constraints (**begin**), (**end**), (**tran**), (**chain**) are obviously true. It remains to verify the constraints on  $\text{@}$ : they must have a unique data value, must point to a corresponding **dec** instruction that occur before any corresponding zero test, and no two  $\text{@}$  may point to the same **dec** instruction. Consider a position  $i$  such that  $\ell_i$  is  $\text{inc}_k$  with a subsequent zero test  $\text{ifzero}_k$  at position  $j > i$ . Because the zero test is correct, the increment of counter  $k$  made at step  $i$  must be decremented at some position between  $i$  and  $j$ . Hence there is a  $i < d < j$  such that  $v_{d+1}(k) = v_i(k)$ . By construction the data value of the  $\text{@}$  symbol of  $B_i$  is the minimal such  $d$ . By minimality,  $c_d$  must be a  $\text{dec}_k$  instruction. Assume now that at position  $i' < i$  there is also a  $\text{inc}_k$  instruction. The data value of the  $\text{@}$  symbol of  $B_{i'}$  cannot be  $d$ . Because if this would be the case then  $v_{i'}(k) = v_{d+1}(k) = v_i(k)$  and  $d$  would not be minimal for the position  $i'$ . Hence (**pair**) is also satisfied and  $w \models \varphi_C$ .

## A.2 Proof of Theorem 4.2

Recall the statement.

**Theorem 4.2** Over data words, Satisfiability of  $\text{LTL}_1^\downarrow(\text{F}, \text{F}^{-1})$  is undecidable.

*Proof (sketch).* Consider a Minsky Counter Automaton  $C$ . We revisit the previous proof. We modify the coding of a run of  $C$  by also inserting a group of symbols  $\text{@}$  in each **dec** block. Using  $\text{F}^{-1}$ , we add formulas ensuring that *every* block containing a **dec** instruction also contains a group of  $\text{@}$  symbols such that their data values refer to previous blocks containing a matching **inc** instruction. We also make sure that two  $\text{@}$  from different blocks have different data values. All this can be done by taking the  $\text{F}^{-1}$  counterpart of the formulas we constructed with  $\text{F}$ .

We show that the new formula  $\varphi_C$  is satisfied by a data word iff  $C$  accept the word. Constructing a data word that satisfies  $\varphi_C$  from a word accepted by  $C$  is done as in the proof of Theorem 2.

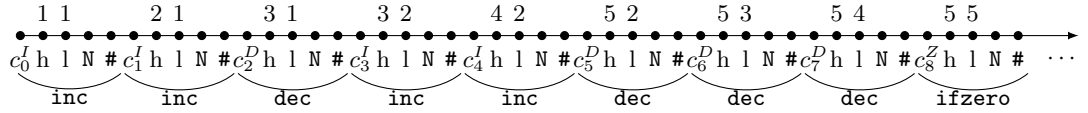
For the other direction one needs to show that from a word that satisfies  $\varphi_C$ , the corresponding word with the obvious run and obvious valuations of the counter is accepting for  $C$ . One can verify that the new extra conditions do imply that no null counter is ever decreased and each zero test is correct, based on the fact that for every  $\text{dec}_i$  there must be at least one  $\text{inc}_i$  to its left, and for every  $\text{inc}_i$  there must have a related  $\text{dec}_i$  before any  $\text{ifzero}_i$  test.

### A.3 Proof of Theorem 5

First recall the statement of Theorem 5:

**Theorem 5** Satisfiability of  $\text{LTL}_2^\downarrow(\mathcal{F})$  is undecidable over data words.

*Proof.* We adapt the proof of part 1. of Theorem 4 using ideas already present in [9,14] for coding Minsky CA with only forward temporal operators. Fix a two counter machine  $C = \langle \Sigma, Q, q_0, 2, \delta, F \rangle$ . The main idea is to modify the coding of Theorem 4 by adding in each block a symbol  $h_1$  whose data value is supposed to be the one of the last  $\text{inc}_1$  that has previously occurred, a symbol  $l_1$  whose data value is supposed to be the one of the last  $\text{inc}_1$  that has not been decreased yet and similarly with  $h_2$  and  $l_2$  for the second counter. A zero test of counter 1 can then only occur in a block where the data values of  $h_1$  and  $l_1$  match. The coding is depicted in Figure 3.



**Fig. 3.** Coding of a 1-counter machine run. Only data values of the  $h$  and  $l$  symbols are represented for the sake of clarity. Only one symbol per group is also represented. Elements depicted as  $c_i^I$ ,  $c_i^D$  and  $c_i^Z$  correspond, respectively, to increment, decrement and zero testing symbols of  $\hat{\Sigma}$ .

Based on these ideas, the formula  $\varphi_C$  that we construct makes sure that:

- (i) *The structure.* We enforce a structure that is a sequence of blocks of the form “ $\langle q, w, l, q' \rangle h_1 l_1 h_2 l_2 N \#$ ”, with the possibility that there are repeated consecutive symbols. This can be done in the same way we did before. As for the previous codings, we demand that in each block, for every element in the first group of  $\hat{\Sigma}$  symbols there is an element in the last group of symbols  $\#$  with the same data value. This is enforced with a formula of  $\text{LTL}_1^\downarrow(\mathcal{F})$  as before. With the help of the second register we can now enforce that each group of symbols has the same data value, for all symbol  $s$  we have the formula:

$$\mathbf{G}(\hat{\sigma} \Rightarrow \downarrow_1 \neg \mathbf{F}(s \wedge \downarrow_2 \mathbf{F}(s \wedge \neg \uparrow_2 \wedge \mathbf{F}(\# \wedge \uparrow_1))))$$

- (ii) Each  $\langle q, w, l, q' \rangle$  occurring in the string is in  $\delta$ . For the first one  $q = q_0$ , and for the last one  $q' \in F$ . For any  $\langle q, w, l, q' \rangle$  and  $\langle q'', w', l', q'' \rangle$  in consecutive blocks,  $q' = q''$ . This can be checked just as in the previous codings.
- (iii) In the initial block, the data values of  $h_1$  and  $l_1$  are the same and the data values of  $h_2$  and  $l_2$  are the same (we only show the formula for  $h_1$  and  $l_1$ ):

$$\downarrow_1 \neg \mathbf{F}(h_1 \wedge \downarrow_2 \mathbf{F}(l_1 \wedge \neg \uparrow_2 \wedge \mathbf{F}(\# \wedge \uparrow_1)))$$

- (iv) In any block immediately after a **ifzero** instruction, the data values of  $h_1, l_1, h_2, l_2$  are identical to those of the **ifzero** instruction. We only give the formula for **ifzero**<sub>1</sub> instruction with the  $h_1$  symbol.

$$\neg F(\text{IZ}(1) \wedge \downarrow_1 F(h_1 \wedge \downarrow_2 F(\psi \wedge F(\# \wedge \uparrow_1))))$$

$$\text{where } \psi \equiv N \wedge \downarrow_1 F(h_1 \wedge \neg \uparrow_2 \wedge F(\# \wedge \uparrow_1))$$

- (v) In any block immediately after an **inc**<sub>1</sub> instruction,  $h_1$  is not in the same class as any preceding  $h_1$  (stated in the formula below) and  $l_1, h_2, l_2$  are in the same class as in the previous block (this can be stated using a formula similar to the one for for item (iv)).

$$\neg F(h_1 \wedge \downarrow_1 F(\text{INC}(1) \wedge \downarrow_2 \wedge F(\varphi \wedge F(\# \wedge \uparrow_2))))$$

$$\text{where } \varphi \equiv N \wedge \downarrow_2 F(\uparrow_2 \wedge F(h_1 \wedge \uparrow_1 \wedge F(\# \wedge \uparrow_2)))$$

We have of course corresponding formulas for blocks with **inc**<sub>2</sub> instructions.

- (vi) In any **dec**<sub>1</sub> block,  $h_1$  and  $l_1$  have different data values.

$$\neg F(\text{DEC}(1) \wedge \downarrow_1 F(h_1 \wedge \downarrow_2 F(l_1 \wedge \uparrow_2 \wedge F(\# \wedge \uparrow_1))))$$

We have a corresponding formula for blocks with a **dec**<sub>2</sub> instruction.

- (vii) In any **ifzero**<sub>1</sub> block,  $h_1$  and  $l_1$  have the same data value. This can be expressed with a formula as above. We have a corresponding formula for **ifzero**<sub>2</sub> blocks.

- (viii) In any block immediately after a **dec**<sub>1</sub> block  $B$ ,

–  $h_1$  is in the same class as the previous  $h_1$ , and

–  $l_1$  is in the same class as the  $h_1$  occurring in a block immediately after the rightmost block (occurring before  $B$ ) with a  $h_1$  that is in the same class as the  $l_1$  of  $B$ .

–  $h_2$  and  $l_2$  are in the same class as their corresponding symbol of  $B$ .

The first and third conditions can be expressed using formulas similar to those above. The difficult part is to enforce the second condition. For this we introduce some macros:

$$\text{InBlock}_i(\varphi) \equiv \downarrow_i F(\varphi \wedge F(\# \wedge \uparrow_i))$$

$$\text{NextBlock}_i(\varphi) \equiv \downarrow_i F(N \wedge F(\# \wedge \uparrow_i) \wedge$$

$$\downarrow_i F(\hat{\sigma} \wedge \varphi \wedge F(\# \wedge \uparrow_i)))$$

$$\text{InNextBlock}_i(\varphi) \equiv \text{NextBlock}_i(\text{InBlock}_i(\varphi))$$

The second condition can then be stated:

$$\neg F(\hat{\sigma} \wedge \downarrow_2 F(h_1 \wedge \downarrow_1 F(N \wedge F(\# \wedge \uparrow_2) \wedge \downarrow_2 F(\hat{\sigma} \wedge \uparrow_2 \wedge F(h_1 \wedge \neg \uparrow_1 \wedge \downarrow_2 \varphi))))))$$

where  $\varphi \equiv F(\text{DEC}(i) \wedge \text{InBlock}_2(l_1 \wedge \uparrow_1) \wedge \text{InNextBlock}_1(l_1 \wedge \neg \uparrow_2))$

We have a corresponding formula for blocks following a **dec**<sub>2</sub> instruction.

□



#### A.4 Proof of Proposition 1

**Proposition 1** Over data words,  $\text{sLTL}_1^\downarrow(\mathbb{F}_s)$  and  $\text{XPath}(\rightarrow^+, =)$  have the same expressive power. The same hold for  $\text{sLTL}_1^\downarrow(\mathbb{F}_s, \mathbb{F}_s^{-1})$  and  $\text{XPath}(\leftarrow^+, \rightarrow^+, =)$ . Moreover, in both cases, the transformation from  $\text{sLTL}_1^\downarrow$  to  $\text{XPath}$  takes polynomial time while it takes exponential time in the other direction.

*Proof (sketch).* The translation from  $\text{sLTL}_1^\downarrow$  to  $\text{XPath}$  is straightforward by induction. One then naturally obtain equality tests of the form  $\varepsilon = \alpha$  or  $\varepsilon \neq \alpha$  depending on whether the unique possible negation is present or not.

The other translation is also by induction. The non trivial part concerns the translation of node expressions of the form  $\alpha = \beta$ . Over data words, one of the path expressions  $\alpha$  or  $\beta$  must end first, say  $\alpha$ . In this case the node expression can be decomposed into a disjunction of (exponentially many) expressions that first start with a path expression that *merge*  $\alpha$  with the part  $\beta_1$  of  $\beta$  that is common to with  $\alpha$ , followed by a test of the form  $\varepsilon = \beta_2$ , where  $\beta = \beta_1\beta_2$ . Tests of the form  $\varepsilon = \beta_2$  are the immediate to translate into  $\text{sLTL}_1^\downarrow$ .

## B Semantic of XPath

Recall that  $\text{XPath}$  is a two-sorted language, with *path* expressions  $(\alpha, \beta, \dots)$  and *node* expressions  $(\varphi, \psi, \dots)$ . These are defined by mutual recursion:

$$\begin{aligned} \alpha &::= \downarrow \mid \downarrow^+ \mid \uparrow \mid \uparrow^+ \mid \rightarrow \mid \rightarrow^+ \mid \leftarrow \mid \leftarrow^+ \mid \varepsilon \mid \alpha\beta \mid \alpha \cup \beta \mid \alpha[\varphi] \\ \varphi &::= \sigma \mid \langle \alpha \rangle \mid \neg\varphi \mid \varphi \wedge \psi \mid \alpha = \beta \mid \alpha \neq \beta \quad (\sigma \in \Sigma) \end{aligned}$$

A path expression essentially describes a traversal of the tree by using the axis: **child** ( $\downarrow$ ), **descendant** ( $\downarrow^+$ ), the **next-sibling** ( $\rightarrow$ ), **following-sibling** ( $\rightarrow^+$ ) and their inverses, with the ability to test at any stage for node conditions. Let  $t$  be a data tree with domain  $T$  and labeling function  $\lambda$ . The semantics of  $\text{XPath}$  is defined by induction as follows:

$$\begin{aligned} \llbracket \downarrow \rrbracket^t &= \{(x, xi) \mid xi \in T\} \\ \llbracket \downarrow^+ \rrbracket^t &= \{(x, xi) \mid xi \in T\} \\ \llbracket \uparrow \rrbracket^t &= \{(xi, x) \mid xi \in T\} \\ \llbracket \leftarrow \rrbracket^t &= \{(xi, x(i-1)) \mid x(i-1) \in T\} \\ \llbracket \rightarrow \rrbracket^t &= \{(xi, x(i+1)) \mid x(i+1) \in T\} \\ \llbracket \alpha^+ \rrbracket^t &= \text{the transitive closure of } \llbracket \alpha \rrbracket^t \\ \llbracket \varepsilon \rrbracket^t &= \{(x, x) \mid x \in T\} \\ \llbracket \alpha\beta \rrbracket^t &= \{(x, z) \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^t \wedge (y, z) \in \llbracket \beta \rrbracket^t\} \\ \llbracket \alpha \cup \beta \rrbracket^t &= \llbracket \alpha \rrbracket^t \cup \llbracket \beta \rrbracket^t \\ \llbracket \alpha[\varphi] \rrbracket^t &= \{(x, y) \in \llbracket \alpha \rrbracket^t \mid y \in \llbracket \varphi \rrbracket^t\} \\ \llbracket \sigma \rrbracket^t &= \{x \in T \mid \pi_1(\lambda(x)) = \sigma\} \\ \llbracket \langle \alpha \rangle \rrbracket^t &= \{x \in T \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^t\} \end{aligned}$$

$$\begin{aligned} \llbracket \neg \varphi \rrbracket^t &= T \setminus \llbracket \varphi \rrbracket^t \\ \llbracket \varphi \wedge \psi \rrbracket^t &= \llbracket \varphi \rrbracket^t \cap \llbracket \psi \rrbracket^t \\ \llbracket \alpha = \beta \rrbracket^t &= \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^t, \\ &\quad (x, z) \in \llbracket \beta \rrbracket^t, \pi_2(\lambda(y)) = \pi_2(\lambda(z))\} \\ \llbracket \alpha \neq \beta \rrbracket^t &= \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^t, \\ &\quad (x, z) \in \llbracket \beta \rrbracket^t, \pi_2(\lambda(y)) \neq \pi_2(\lambda(z))\} \end{aligned}$$