

Fragments of Fixpoint Logic on Data Words*

Thomas Colcombet¹ and Amaldev Manuel²

1 LIAFA, CNRS, Université Paris 7-Paris Diderot
thomas.colcombet@liafa.univ-paris-diderot.fr

2 MIMUW, University of Warsaw
amal@mimuw.edu.pl

Abstract

We study fragments of a μ -calculus over data words whose primary modalities are ‘go to next position’ (X^g), ‘go to previous position’ (Y^g), ‘go to next position with the same data value’ (X^c), ‘go to previous position with the same data value’ (Y^c). Our focus is on two fragments that are called the bounded mode alternation fragment (BMA) and the bounded reversal fragment (BR). BMA is the fragment of those formulas that whose unfoldings contain only a bounded number of alternations between global modalities (X^g, Y^g) and class modalities (X^c, Y^c). Similarly BR is the fragment of formulas whose unfoldings contain only a bounded number of alternations between left modalities (Y^g, Y^c) and right modalities (X^g, X^c). We show that these fragments are decidable (by inclusion in Data Automata), enjoy effective Boolean closure, and contain previously defined logics such as the two variable fragment of first-order logic and DataLTL. More precisely the definable language in each formalism obey the following inclusions that are effective.

$$\text{FO}^2 \subsetneq \text{DataLTL} \subsetneq \text{BMA} \subsetneq \text{BR} \subsetneq \nu \subseteq \text{Data Automata} .$$

Our main contribution is a method to prove inexpressibility results on the fragment BMA by reducing them to inexpressibility results for combinatorial expressions. More precisely we prove the following hierarchy of definable languages,

$$\emptyset = \text{BMA}^0 \subsetneq \text{BMA}^1 \subsetneq \dots \subsetneq \text{BMA} \subsetneq \text{BR} ,$$

where BMA^k is the set of all formulas whose unfoldings contain at most $k-1$ alternations between global modalities (X^g, Y^g) and class modalities (X^c, Y^c). Since the class BMA is a generalization of FO^2 and DataLTL the inexpressibility results carry over to them as well.

1998 ACM Subject Classification F.4.1

Keywords and phrases Data words, Data automata, Fixpoint logic

1 Introduction

Data words are words of the form $(a_1, d_1) \dots (a_n, d_n) \in (\Sigma \times \mathcal{D})^*$ where Σ is a finite set of letters and \mathcal{D} is an infinite domain of data values. Typically the alphabet Σ abstracts a finite set of actions or events and the set of data values \mathcal{D} models some sort of identity information. Thus, data words can model a number of scenarios where the information is linearly ordered and it is composed of finite as well as unbounded elements. For example the authors of [1] imagine Σ as the actions of a finite program and \mathcal{D} as process ids. Then, an execution trace of a system with unbounded instances of the program can be modeled as

* The research leading to these results has received funding from the European Union’s Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 259454.



a data word in which each action is associated with the identifier of the process which has generated it.

The paradigmatic question in the study of data words is to develop suitable models (in particular automata and logics) to specify properties of data words. Sure enough there exists a rich variety of models for specifying properties of data words that includes Data Automata [4], Register Automata [14, 9], Pebble Automata [18], Class Memory Automata [1], Class Automata [2], Walking Automata [17], Variable Automata [12], First-Order logic with two variables [4], guarded MSO logic [5], DataLTL [15], Freeze-Logics [9, 13], Logic of Repeating Values [8], XPath [10, 11], Regular expressions [16], *Data Monoids* [3] etc.

In this work we further study a modal fixpoint logic on data words that we introduced in [6]. This logic is composed of four modalities that allow to evaluate formulas on the successor, class successor (the nearest future position with the same data value), predecessor and class predecessor (nearest past position with the same data value) positions, Y^g, Y^c . In addition there is a couple of zeroary modalities that describes whether these positions coincide or not. To build the formulas, besides the usual Boolean operations, it is allowed to form the least and greatest fixpoints of formulas. In [6] it is shown that the satisfiability problem for the set of formulas that use only least fixpoints is undecidable, whereas the fragment that consists of only greatest fixpoints is subsumed by Data Automata and hence it has a decidable satisfiability problem. The main result of the work was the decidability of an alternation-free fragment of the logic that further bounds the number of change of directions in evaluating the formulas by using a generalisation of Data Automata.

Contributions

In the present paper, we aim at restricting the power of the above μ -calculus logic for data words for obtaining classes that are closed under all Boolean connectives, mirroring, and enjoy decidability of emptiness and universality. We consider two restricted fragments that achieve this goal. The first one, called BMA (for Bounded Mode Alternation) syntactically bounds the number of changes between class and global modes. The second, called BR (for Bounded Reversal), syntactically bounds the number of changes between left modalities and right modalities.

It is easy to show that BMA is contained in Data Automata. It is not very difficult to show that BMA is subsumed by BR, that is to say for every formula in BMA there is an equivalent one in BR. Our main result is the strictness of this last inclusion, i.e. that there is a formula in BR for which there is no equivalent formula in BMA. This proof uses a deep result from combinatorics called the Hales-Jewett theorem. As a proof device we use a sort of circuits called combinatorial expressions that were introduced in [7]. These expressions define functions over an infinite domain (for instance the integers). They are built by composing *gates* that are functions of two kinds, either the function has a bounded arity, or the function has a bounded domain. In [7] it is shown that certain properties (a property is a function that has a binary codomain) for instance *the given sequence of positive integers has gcd 1* or *the given sequence of integers sum to 0* cannot be computed by expressions of fixed depth. We use a variant of this theorem in this paper to show that there is a formula in BR for which there is no equivalent one in BMA. More precisely it is shown that there is a specific formula in BR such that if it has an equivalent formula in BMA, then it is possible to construct expressions of fixed depth for a particular property and since that particular property cannot be computed by fixed depth expressions, we derive a contradiction. One thing to note is that since the techniques developed in [7] are general enough to derive impossibility results for a large family of properties, correspondingly the

proof method developed here can be used to show inexpressibility results for a variety of formulas.

Now we examine the implications of our result in a larger context. As mentioned earlier, the results mentioned here have very close connection with Data Automata (DA for short). The well known feature of DA is that it subsumes the logic $\text{FO}^2(\Sigma, <, +1, \sim, +^c1)$ on data words where Σ denotes the unary predicates indicating the letters, $<$ is the linear order on positions, $+1$ is the successor relation on positions, \sim is the equivalence relation on positions with respect to the data values (i.e. $i \sim j$ if $d_i = d_j$), and $+^c1$ is the class successor relation. It is known that data languages recognisable by DA are closed under union, intersection and letter-to-letter projection, but not under complementation [4]. Since FO^2 formulas are closed under Boolean operations, it is evident that Data Automata strictly subsumes the logic FO^2 . This observation prompts the question that if there are other classes subsumed by DA that are closed under Boolean operations. The fragments introduced in the paper answer this question positively. Note only that, there are automata theoretic characterisations that are natural variants of DA for both these fragments (we only present the one for BMA).

Another and perhaps more important question is how to show that a given data language is not expressible in FO^2 . Note that in some cases, using the techniques on words over finite alphabets it is possible to show that a given data language is not definable in FO^2 (for instance to show that data words of even length are not definable in FO^2). We are interested in those cases where such reductions are not possible, in particular where the property given is dependent on the data values. We don't have a complete solution to this problem yet, but our method to prove inexpressibility results on BMA offers a partial answer. This is because the logic $\text{FO}^2(\Sigma, <, +1, \sim, +^c1)$, as it is shown in this paper, is equivalent to the unary fragment of a temporal logic, namely DataLTL [15], which is a strict subfragment of BMA. DataLTL is the temporal logic where usual temporal operators such as *until*, *future*, *past* etc. exist both on the linear order on positions (called the global order) as well as on the suborders formed by subsets of positions that share the same data value (called the *class orders*). For instance the temporal operator $F^g\varphi$ is true a position if there is a position in the future that satisfies the formula φ , whereas the formula $F^c\varphi$ is true at a position if there is a future position that has the same data value as the current position and that satisfies the formula φ . The unary fragment of DataLTL is the subclass of formulas that uses only the unary temporal operators (such as F^g, P^g, F^c etc). Since every such operator is expressible in FO^2 it is immediate that unary DataLTL is subsumed by the logic FO^2 . But the converse direction, which is shown in the paper, is not obvious, since it is not immediate how to translate formulas like $\exists y (a(x) \wedge b(y) \wedge x < y \wedge x \not\sim y)$. Thus inexpressibility results on the fragment BMA renders directly corresponding results on all sublogics including FO^2 and DataLTL.

Finally let us also note that the translations outlined in this paper, namely

$$\text{FO}^2 \subsetneq \text{DataLTL} \subsetneq \text{BMA} \subsetneq \text{BR} \subsetneq \nu \subseteq \text{DA},$$

constitutes an alternate proof the main result of [4] that FO^2 is subsumed by DA. The proof in [4] is a direct translation of FO^2 formulas by a intricate case analysis. Our proof, however, is modular and makes use of analogous constructions from automata theory on finite words.

Related work

As mentioned already this work is strongly related to DataLTL, FO^2 and DA. One other very popular ecosystem on data words is that of Register Automata and the associated lo-

gics such as Freeze LTL, Freeze μ -calculus, Xpath [9, 13, 8, 10, 11] etc. Our inexpressibility result implies that BMA is incomparable to Register Automata (in particular nondeterministic 1-Register Automata). Since all our modalities are expressible in terms of successor, predecessor, freeze operator and fixpoint operators, our fixpoint logic is subsumed by Freeze μ -calculus of [13]. However it should be noted that the latter logic is highly undecidable [9]. The decidable fragment of Freeze μ -calculus (and also Freeze LTL) is unidirectional (only future modalities) but our logic is naturally two-way. Finally the decidable two-way fragment of Freeze LTL, namely Simple Freeze LTL is equivalent to FO^2 and hence it is subsumed by BMA. Therefore our method of proving inexpressibility extends to this logic as well.

Structure of the document

In Section 2 we present the definition of our fixpoint logic and give some examples. In Section 3 we recall the *composition* operator (*comp*) on sets of formulas and define the fragments BMA and BR using it. Thereafter, a characterisation of the class BMA in terms of cascades of automata, that is used in the proof of the separation theorem, is given. In Section 4, first we recollect the paradigm of combinatorial expressions and state the necessary results for our purpose. Afterwards it is shown how to translate a cascade on data words with a specific structure to expressions and the separation theorem is proved. In Section 5 we conclude. Due to space constraints some proofs have been moved to the appendix.

2 μ -Calculus on Data Words

In this section, we recall the basics of the μ -calculus on data words [6].

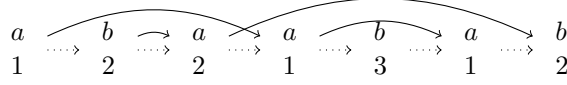
Fix an infinite set \mathcal{D} of *data values*. A *data word* is a word $u = (a_1, d_1) \cdots (a_n, d_n) \in (\Sigma \times \mathcal{D})^*$ where Σ is a finite set of *letters*. We often name *positions* the letter indices in a word. A maximal set of positions in u with the same data value is called a *class*. The set of classes in u define an equivalence relation \sim , called the *class relation*, on the set of positions of u . Given a permutation σ of \mathcal{D} , it can be applied on a data word as expected, yielding $\sigma(u) = (a_1, \sigma(d_1)) \cdots (a_n, \sigma(d_n))$. The data words u and $\sigma(u)$ have the same class relation. A *data language* is a set of data words that is invariant under such applications of the permutations of \mathcal{D} .

For our purposes, it is convenient to see data words as graphs in the following manner. To each data word $w = (a_1, d_1) \cdots (a_n, d_n) \in (\Sigma \times \mathcal{D})^*$ associate the graph $G_w = ([n], \ell, +1, +^c1)$ where $[n]$ is the set of positions $\{1, \dots, n\}$, $\ell : \Sigma \rightarrow 2^{[n]}$ is the labelling function $\ell(a) = \{i \mid a_i = a\}$, the binary relation $+1$ denotes the *successor* relation on positions, *i.e.*, $+1(i, j)$ if $j = i + 1$, and the binary relation $+^c1$ denotes the *class successor* relation on positions, *i.e.*, $+^c1(i, j)$ if $i < j$, $d_i = d_j$, and $d_m \neq d_i$ for all $i < m < j$. We call *predecessor* relation (*resp.*, *class predecessor* relation) the reverse of the successor relation (*resp.*, *class successor* relation). We implicitly identify a data word with its graph. Figure 1 shows a data word and its corresponding graph.

Seen as such graphs, data words are naturally prone to the use of temporal logics. Let $Prop = \{p, q, \dots\}$ and $Var = \{x, y, \dots\}$ be countable sets of *propositional variables* and *fixpoint variables* respectively. The μ -calculus on data words is the set of all *formulas* φ respecting the following syntax:

$$\varphi := x \mid A \mid \neg A \mid M\varphi \mid \varphi \vee \varphi \mid \varphi \wedge \varphi \mid \mu x.\varphi \mid \nu x.\varphi$$

where $M := X^g \mid X^c \mid Y^g \mid Y^c$ and $A := p \in Prop \mid \mathcal{S} \mid \mathcal{P} \mid \text{fst}^c \mid \text{fst}^g \mid \text{lst}^c \mid \text{lst}^g$



■ **Figure 1** A data word and its graph. Dotted and thick arrows denote the successor and class successor relations respectively.

The elements of \mathbf{M} are called *modalities*, and the ones of A , *atoms*. The set of zeroary modalities $\{\text{fst}^c, \text{fst}^g, \text{lst}^c, \text{lst}^g, \mathcal{P}, \mathcal{S}\}$ will be denoted by the symbol Z for the rest of the paper.

The semantic of a formula φ , over a data word w is the set of positions of w where “ φ is true” (denoted as $\llbracket \varphi \rrbracket_w$). The formal definition is given in Figure 2. The different constructs have their expected meaning, keeping in mind that the class modalities $\mathbf{X}^c, \mathbf{Y}^c, \text{fst}^c, \text{lst}^c$ have to be interpreted on the word restricted to the current data value. The modality \mathcal{S} (*resp.*, \mathcal{P}) holds at a position i if the successor and class successor i coincide (*resp.* the predecessor and class predecessor coincide).

$$\begin{array}{ll}
\llbracket \text{fst}^g \rrbracket_w = \{1\} & \llbracket \mathbf{X}^g \varphi \rrbracket_w = \llbracket \varphi \rrbracket_w - 1 \\
\llbracket \text{lst}^g \rrbracket_w = \{n\} & \llbracket \mathbf{Y}^g \varphi \rrbracket_w = \llbracket \varphi \rrbracket_w + 1 \\
\llbracket \text{fst}^c \rrbracket_w = \{i \mid \nexists j = i - c\} & \llbracket \mathbf{X}^c \varphi \rrbracket_w = \llbracket \varphi \rrbracket_w - c \\
\llbracket \text{lst}^c \rrbracket_w = \{i \mid \nexists j = i + c\} & \llbracket \mathbf{Y}^c \varphi \rrbracket_w = \llbracket \varphi \rrbracket_w + c \\
\llbracket \varphi_1 \wedge \varphi_2 \rrbracket_w = \llbracket \varphi_1 \rrbracket_w \cap \llbracket \varphi_2 \rrbracket_w & \llbracket \mathcal{S} \rrbracket_w = \{i \mid i + 1 = i + c\} \\
\llbracket \varphi_1 \vee \varphi_2 \rrbracket_w = \llbracket \varphi_1 \rrbracket_w \cup \llbracket \varphi_2 \rrbracket_w & \llbracket \mathcal{P} \rrbracket_w = \{i \mid i - 1 = i - c\} \\
\llbracket \mu x. \varphi \rrbracket_w = \bigcap \{S \subseteq [n] \mid \llbracket \varphi \rrbracket_{w[\ell(x) := S]} \subseteq S\} & \llbracket p \rrbracket_w = \ell(p) \\
\llbracket \nu x. \varphi \rrbracket_w = \bigcup \{S \subseteq [n] \mid S \subseteq \llbracket \varphi \rrbracket_{w[\ell(x) := S]}\} & \llbracket \neg p \rrbracket_w = [n] \setminus \ell(p) \\
\llbracket x \rrbracket_w = \ell(x) &
\end{array}$$

■ **Figure 2** Semantics of μ -calculus on data words $w = ([n], +1, +c, \ell)$.

Note that in this definition of the logic, negations in a formula are located at the leaves. It is nevertheless possible, as usual, to negate such formulas by pushing the negation toward the leaves, but this requires a bit of care when negating modalities. For instance, $\neg \mathbf{X}^c \varphi$ is not equivalent to $\mathbf{X}^c \neg \varphi$, but to $\text{lst}^c \vee \mathbf{X}^c \neg \varphi$. Similar arguments have to be used for all modalities. Following these ideas, we define the dual modalities $\tilde{\mathbf{X}}^g \varphi \equiv \text{lst}^g \vee \mathbf{X}^g \varphi$, $\tilde{\mathbf{Y}}^g \varphi \equiv \text{fst}^g \vee \mathbf{Y}^g \varphi$, $\tilde{\mathbf{X}}^c \varphi \equiv \text{lst}^c \vee \mathbf{X}^c \varphi$ and $\tilde{\mathbf{Y}}^c \varphi \equiv \text{fst}^c \vee \mathbf{Y}^c \varphi$. These modalities are considered dual since $\tilde{\mathbf{X}}^g \varphi \equiv \neg \mathbf{X}^g \neg \varphi, \dots$

Next we lay out some terminology and abbreviations which we will use in the subsequent sections. Let λ denote either μ or ν . Every occurrence of a fixpoint variable x in a subformula $\lambda x. \psi$ of a formula is called *bound*. All other occurrences of x are called *free*. A formula is called a *sentence* if all the fixpoint variables in φ are bound. If $\varphi(x_1, \dots, x_n)$ is a formula with free variables x_1, \dots, x_n , then by $\varphi(\psi_1, \dots, \psi_n)$ we mean the formula obtained by substituting ψ_i for each x_i in φ . As usual the bound variables of $\varphi(x_1, \dots, x_n)$ may require a renaming to avoid the capture of the free variables of ψ_i 's. For a sentence φ and a position i in the word w , we denote by $w, i \models \varphi$ if $i \in \llbracket \varphi \rrbracket_w$. The notation $w \models \varphi$ abbreviates the case when $i = 1$. The *data language* of a sentence φ , denoted as $L(\varphi)$, is the set of data

words w such that $w \models \varphi$.

By μ -fragment we mean the subset of μ -calculus which uses only μ -fixpoints. Similarly ν -fragment stands for the subset which uses only ν -fixpoints.

► **Example 1** (temporal modalities). An example of a formula would be $\varphi \mathbf{U}^g \psi$ which holds if ψ holds in the future, and φ holds in between. This can be implemented as $\mu x. \psi \vee (\varphi \wedge \mathbf{X}^g x)$. The formula $\varphi \mathbf{U}^c \psi = \mu x. \psi \vee (\varphi \wedge \mathbf{X}^c x)$ is similar, but for the fact that it refers only to the class of the current position. The formula $\mathbf{F}^g \varphi$ abbreviates $\top \mathbf{U}^g \varphi$, and its *dual* is $\mathbf{G}^g \varphi = \neg \mathbf{F}^g \neg \varphi$. The constructs $\mathbf{S}^g, \mathbf{S}^c, \mathbf{P}^g, \mathbf{P}^c, \mathbf{H}^g$ and \mathbf{H}^c , are defined analogously, using past modalities, and correspond respectively to $\mathbf{U}^g, \mathbf{U}^c, \mathbf{F}^g, \mathbf{F}^c, \mathbf{G}^g$ and \mathbf{G}^c . For instance, $\mathbf{F}^c \mathbf{P}^c \varphi$ expresses that there is a position in the class that satisfies φ and $\mathbf{F}^c \mathbf{P}^c (\varphi \wedge \tilde{\mathbf{X}}^c \mathbf{G}^c \neg \varphi \wedge \tilde{\mathbf{Y}}^c \mathbf{H}^c \neg \varphi)$ expresses that there exists exactly one position which satisfies φ in the class.

3 The bounded reversal and bounded mode alternation fragments

In this section we introduce the bounded mode alternation and bounded reversal fragments (BMA and BR) and compare these two fragments between themselves and with other logics (Theorem 5).

Before delving into the technical details let us outline the intuition behind each of the fragments. The four modalities $\mathbf{X}^g, \mathbf{Y}^g, \mathbf{X}^c$ and \mathbf{Y}^c can be divided along two axis. Based on the directions: there are the *left modalities* $\mathbf{Y}^g, \mathbf{Y}^c$, and *right modalities* $\mathbf{X}^g, \mathbf{X}^c$. Based on the modes: there are *global modalities* $\mathbf{X}^g, \mathbf{Y}^g$, and *class modalities* $\mathbf{X}^c, \mathbf{Y}^c$. The BR and BMA fragments are defined by limiting the number of alternation between this types of modalities. This is formally achieved using the operation *comp* that we define now.

Let Ψ be a set of μ -calculus formulas. Define the sets $\text{comp}^i(\Psi)$ for $i \in \mathbb{N}$ inductively

- $\text{comp}^0(\Psi) = \emptyset$,
- $\text{comp}^{i+1}(\Psi) = \{\psi(\varphi_1, \dots, \varphi_n) \mid \psi(x_1, \dots, x_n) \in \Psi, \varphi_1, \dots, \varphi_n \in \text{comp}^i(\Psi)\}$ in which the substitution $\psi(\varphi_1, \dots, \varphi_n)$ is allowed only if none of the free variables of $\varphi_1, \dots, \varphi_n$ get bound in $\psi(\varphi_1, \dots, \varphi_n)$.

The set of formulas $\text{comp}(\Psi)$ is defined as $\text{comp}(\Psi) = \bigcup_{i \in \mathbb{N}} \text{comp}^i(\Psi)$. For a formula $\psi \in \text{comp}(\Psi)$, the *comp-height* of ψ in $\text{comp}(\Psi)$ is the least i such that ψ is in $\text{comp}^i(\Psi)$.

We are now ready to define the BR and BMA fragments of the μ -calculus. For a set of modalities M , define *formulas*(M) to be the set of formulas that uses only the modalities M apart from the zeroary modalities.

► **Definition 2.** The BMA and the BR fragments of μ -calculus are respectively:

$$\begin{aligned} \text{BMA} &= \text{comp}(\text{formulas}(\{\mathbf{X}^g, \mathbf{Y}^g\}) \cup \text{formulas}(\{\mathbf{X}^c, \mathbf{Y}^c\})), \\ \text{and BR} &= \text{comp}(\text{formulas}(\{\mathbf{X}^g, \mathbf{X}^c\}) \cup \text{formulas}(\{\mathbf{Y}^g, \mathbf{Y}^c\})). \end{aligned}$$

Further, BMA^k denotes the subset of BMA with comp-height k . Similarly BR^k stands for the subset of BR with comp-height k .

► **Example 3.** Let

$$\begin{aligned} \varphi_1 &= \nu x. (\mathbf{X}^c x \vee \mathbf{X}^g \mu y. (q \wedge \mathbf{Y}^c y)), & \varphi_2 &= \nu x. (\mathbf{X}^c \text{!st}^g \vee \mathbf{X}^c \mathbf{Y}^g x), \\ \varphi_3 &= \mu x. ((\nu y. q \vee \mathbf{X}^c y) \vee \mathbf{X}^g x \vee \mathbf{Y}^g x), & \text{and } \varphi_4 &= \mu x. (\mathbf{X}^c \mathbf{X}^g x \vee p). \end{aligned}$$

The formula φ_1 is in BR^2 and in BMA^3 . The formula φ_2 is neither in BR nor in BMA. The formula φ_3 is in BMA^2 but not in BR. The formula φ_4 is in BR^1 but not in BMA.

► **Example 4.** Consider the language L_k that contains the data words such that, by applying k -times the sequence of the global successor followed by the class successor, one reaches a position labeled with letter a . The language L is the union of all L_k for k ranging over all non-negative integers. The language L_k is defined by φ_k and L by φ defined as follows:

$$\varphi_k = \overbrace{X^g X^c \dots X^g X^c}^{k\text{-times}} a, \quad \text{and} \quad \varphi = \mu x. (X^g X^c x \vee a).$$

The formula φ_k is in BR^1 and in BMA^{2k} . The formula φ is in BR^1 , but not in BMA . Later in Section 4 we will prove that a variant of L is not definable by any formula in BMA .

Let us now state the main theorem of this section, namely the inclusions between the fragments of the μ -calculus in terms of the data languages defined. Below DataLTL denotes the temporal logic on data words consisting of the modalities $\{\mathcal{S}, \mathcal{P}, X^g, Y^g, X^c, Y^c, U^g, S^g, U^c, S^c\}$, $u\text{DataLTL}$ is the unary sublogic consisting of the modalities $\{\mathcal{S}, \mathcal{P}, X^g, X^c, Y^g, Y^c, F^c, F^g, P^g, P^c\}$ and ν denotes the fragment of the μ -calculus containing only the greatest fixpoints (ν 's).

► **Theorem 5.** *The following inclusions hold for definable languages,*

$$\text{FO}^2(\Sigma, <, +1, \sim, +^c1) = u\text{DataLTL} \subsetneq \text{DataLTL} \subsetneq \text{BMA} \subseteq \text{BR} \subsetneq \nu \subseteq \text{DA}.$$

3.1 Characterizing BMA as cascades of automata

Next we give a characterization of BMA in terms of cascades of finite state automata. It is classical that composition (*comp*) corresponds to the natural operation of composing sequential transducers that compute subset of subformulas that are true at each position. Given a μ -calculus formula φ over words, we can see it as a transducer that reads the input, and labels every position with one extra bit of information denoting the truth value of the formula φ at that position. Under this view, the composition of formulas corresponds to applying the transducers in sequence: the first transducer reads the input, and adds some extra labelling on it. Then a second transducer reads the resulting word, and processes it in a similar way, etc... If we push this view further, we can establish exact correspondences between the class BMA , and suitable cascades of transducers. Furthermore, the *comp*-height of the formula matches the number of transducers involved in the cascade.

First we introduce some notation. Given a data word $w = (a_1, d_1) \dots (a_n, d_n)$ the *string projection* of w , denoted by $\text{str}(w)$, is the word $a_1 \dots a_n$. For a class $S = \{i_1, \dots, i_m\}$ the *class projection* corresponding to S , denoted as $\text{str}(w|_S)$, is the finite word $a_{i_1} \dots a_{i_m}$. For a word $u = b_1 \dots b_n$, the relabelling of w by u is the data word $(b_1, d_1) \dots (b_n, d_n)$. Similarly the relabelling of the class S in w by $b_1 \dots b_m$ is the data word $(a'_1, d_1) \dots (a'_n, d_n)$ where $a'_i = b_j$ if $i = i_j$ and a_i otherwise.

The *marking* of a position i in the data word w , in notation $m(i)$, is the subset of zeroary modalities Z satisfied by i . The *marked string projection* of w , denoted as $\text{mstr}(w)$, is the word $(a_1, m(i)) \dots (a_n, m(n))$ over the alphabet $\Sigma \times 2^Z$. Given a class $S = \{i_1, \dots, i_n\}$ the *marked class projection* of S , denoted as $\text{mstr}(w|_S)$, is the finite word $(a_{i_1}, m(i_1)) \dots (a_{i_n}, m(i_n))$.

A *functional letter-to-letter transducer* $\mathcal{A} : \Sigma^* \rightarrow \Gamma^*$ over words is a nondeterministic finite state letter-to-letter transducer such that every input word $w \in \Sigma^*$ has at most one output word $\mathcal{A}(w) \in \Gamma^*$.

We next disclose two forms of transductions possible by a word transducer on data words. Let $\mathcal{A} : (\Sigma \times 2^Z)^* \rightarrow \Gamma^*$ be a functional letter-to-letter transducer. The automaton \mathcal{A} acts as a *global transducer* when it runs on the marked string projection $\text{mstr}(w)$ of the input data word $w \in (\Sigma \times \mathcal{D})^*$. If the run succeeds then the unique output data

word $w' \in (\Gamma \times \mathcal{D})^* = \mathcal{A}(w)$ (by abuse of notation) is the relabelling of w with the word $\mathcal{A}(\text{mstr}(w))$. Automaton \mathcal{A} is a *class transducer* when for each class S in the input data word w , a copy of the automaton \mathcal{A} runs on the marked class projection $\text{mstr}(w|_S)$. If all the runs succeed then the unique output data word $\mathcal{A}(w)$ (by abuse of notation) is the relabelling of each class of S of w by $\text{mstr}(w|_S)$.

► **Definition 6.** A *cascade of class and global transducers over data words* (hereafter simply cascade) \mathcal{C} is a sequence $\langle \Sigma = \Sigma_0, \mathcal{A}_1, \Sigma_1, \dots, \Sigma_{n-1}, \mathcal{A}_n, \Sigma_n \rangle$ such that $\mathcal{A}_1, \dots, \mathcal{A}_n$ is a sequence of class and global transducers over data words and for each i , the transducer \mathcal{A}_i has input alphabet $\Sigma_{i-1} \times 2^Z$ and output alphabet Σ_i . Sets Σ_0, Σ_n are respectively the *input* and *output* alphabets of the cascade \mathcal{C} and n is the *height* of the cascade.

The cascade \mathcal{C} has a *successful run* on a given data word w if there is a sequence of data words $w_0 = w, w_1, \dots, w_{n-1}, w_n$ such that each transducer \mathcal{A}_i has a successful run on w_{i-1} outputting the data word w_i . The data word w_n is the output of the cascade \mathcal{C} , in notation $\mathcal{C}(w) = w_n$. The language accepted by the cascade \mathcal{C} , denoted as $L(\mathcal{C})$, is the set of all data words w on which \mathcal{C} has a successful run.

Two cascades \mathcal{C}_1 and \mathcal{C}_2 can be composed to form the cascade $\mathcal{C}_1 \circ \mathcal{C}_2$ if the output alphabet of \mathcal{C}_1 and the input alphabet of \mathcal{C}_2 are the same. Composition of cascades is the natural analogue of composition of formulas, this is expressed by the following proposition.

► **Proposition 7.** *Let L be a set of data words. Then the following statements are equivalent.*

1. L is definable by a formula in BMA of comp-height k .
2. L is recognisable by a cascade of height k .

4 Separation of the fragments BMA and BR

In this section we prove the main theorem of the paper, namely the separation of the fragments of BMA and BR. More precisely it is shown that there is a formula in BR that has no equivalent formula in BMA. We start by presenting our technical tool, namely combinatorial expressions [7].

4.1 Combinatorial expressions

Put simply, combinatorial expressions are *circuits* over a *data domain* \mathcal{E} . For our purposes it is sufficient to assume that \mathcal{E} is a set that contains all the usual data types such as Booleans, integers, finite words etc. We form expressions by composing variables (denoted by X, Y etc.) and functions (denoted by f, g etc.) whose domains and ranges are explicitly specified. A variable X has *range* $E \subseteq \mathcal{E}$, denoted as $X : E$, if it takes values from the set E . We say a function $f : E_1 \times \dots \times E_k \rightarrow F$, where $E_1, \dots, E_k, F \subseteq \mathcal{E}$, has *arity* k , *domain* $E_1 \times \dots \times E_k$ and *range* F . The expressions are built using two specific classes of functions (called gates), namely:

- *binary functions* — when $k \leq 2$, and,
- *finitary functions* — when each of E_1, \dots, E_k is finite.

For example the addition on integers $+: \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$ is a binary function, whereas the Boolean disjunction of k inputs $\vee : \{0, 1\}^k \rightarrow \{0, 1\}$ is a finitary function.

► **Definition 8.** *Combinatorial expressions* are defined inductively;

- a variable $X : E$ is a combinatorial expression with range E , and *depth* 0.

- if $f : E_1 \times \dots \times E_k \rightarrow F$ is a binary or a finitary function, and t_1, \dots, t_k are combinatorial expressions with ranges E_1, \dots, E_k and depths d_1, \dots, d_k respectively, then $f(t_1, \dots, t_k)$ is a combinatorial expression with range F and depth $\max(d_1, \dots, d_k) + 1$.

Let $t(\bar{X})$ be a combinatorial expression that contains (possibly vacuously) the variables $\bar{X} = X_1 : E_1, \dots, X_n : E_n$. For the valuation $\bar{a} = a_1, \dots, a_n$, where $a_i \in E_i$ for each i , of the variables \bar{X} , the value of the expression t , denoted as $t(\bar{a})$, is defined inductively; if t is a variable X_i then $t(\bar{a}) = a_i$, and if $t = f(t_1, \dots, t_k)$ then $t(\bar{a}) = f(t_1(\bar{a}), \dots, t_k(\bar{a}))$. Assume $F \subseteq \mathcal{E}$ is the range of the expression t . Naturally t defines a map $\llbracket t \rrbracket : \bar{a} \rightarrow t(\bar{a})$ from the set $E_1 \times \dots \times E_n$ to the set F . Given a map $m : E_1 \times \dots \times E_n \rightarrow F$, where $E_1, \dots, E_n, F \subseteq \mathcal{E}$, we say the map is *recognised* by an expression t if $\llbracket t \rrbracket = m$. A particular case is when the range of the map m is restricted to a set of size two (without loss of generality $\{0, 1\}$); in which case we say that t recognises the *property* $\{a_1, \dots, a_n : m(a_1, \dots, a_n) = 1\}$.

► **Example 9.** Each map $f : E^n \rightarrow F$, for some $E, F \subseteq \mathcal{E}, n \in \mathbb{N}$, has an expression of depth $\lceil \log n \rceil + 1$ recognising it. Let $cat : E^* \times E^* \rightarrow E^*$ be the concatenation operation on words over the alphabet E and let $t(X_1 : E, \dots, X_n : E)$ be an expression of depth $\lceil \log n \rceil$ that consists of only the function cat and that computes the concatenation of the inputs. Let $u : E^* \rightarrow F$ be a binary function on words over E such that $u(e_1 \dots e_n) = f(e_1, \dots, e_n)$. The map f is recognised by expression $u(t(X_1 : E, \dots, X_n : E))$.

► **Example 10.** Consider the set P_n of n -tuples (u_1, \dots, u_n) of words in $\{0, 1\}^*$ that all have equal length. The property P_n is recognised by the expression

$$t = \bigwedge (el(X_1, X_2), \dots, el(X_1, X_n), el(X_2, X_3), \dots, el(X_2, X_n), \dots, el(X_{n-1}, X_n))$$

where \bigwedge is the Boolean conjunction on $n \cdot (n-1)/2$ inputs and $el : A^* \times A^* \rightarrow \{0, 1\}$ is the function on words defined as $el(u, v) = 1$ iff the words u and v are of the same length. The function \bigwedge is finitary and the function el is binary. The expression t has depth 2.

In the previous example, regardless of the value of n the expression t has a constant depth. But there exists properties for which it is not the case.

► **Definition 11.** Let V_n be the set of n -tuples (u_1, \dots, u_n) of words over the alphabet $\{0, 1\}$ such that:

1. the words u_1, \dots, u_n are of the same length, and;
2. there exists a position $1 \leq i \leq |u_1|$ such that the i th letter of each of u_1 to u_n is 1.

It is shown in [7] that,

► **Theorem 12.** *There is no expression of depth at most k that recognises the property $V_{2^{k+1}}$.*

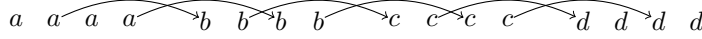
4.2 Application

We now apply the above theorem to derive our inexpressibility results. The idea is to define a data language B_n that corresponds to the property V_n and to show that if there is a BMA-formula of comp-height k recognising B_n then there is a combinatorial expression of depth $\mathcal{O}(k)$ (precise bound disclosed later) recognising the property V_n . This claim along with the Theorem 12 implies a lower bound on the comp-height of formulas defining the language B_n .

For the proof we rely on data words with a special structure that encode a sequence of words. Let $v_1, \dots, v_n \in \Sigma^*$ be words of identical and even length, say $2\ell \in \mathbb{N}$. A data word

$w \in (\Sigma \times \mathcal{D})^*$ is a *coding* of the words $v_1, \dots, v_n \in \Sigma^*$, denoted as $w = \text{coding}(v_1, \dots, v_n)$, if $w = w_1 \cdots w_n$ with $v_1 = \text{str}(w_1), \dots, v_n = \text{str}(w_n)$ and the class relation is the set of tuples $(k \cdot 2\ell + 2i, (k+1) \cdot 2\ell + 2i - 1)$ for $0 \leq k < n-1, 1 \leq i \leq \ell$; the position $k \cdot 2\ell + 2i$ is the i th even position in the block w_{k+1} and $(k+1) \cdot 2\ell + 2i - 1$ is the i th odd position in the block w_{k+2} . Coding is only defined for words of identical even length and hereafter whenever we say *coding*(v_1, \dots, v_n) it is understood that v_1, \dots, v_n are of identical even length.

A data word w is a n -coding (or simply a coding when the value n is clear from the context) if it is the coding of some words $v_1, \dots, v_n \in \Sigma^*$. We write n -Codings for the set of all n -codings.



■ **Figure 3** The coding of the words $aaaa, bbbb, cccc, dddd \in \{a, b, c, d\}^*$.

Next we introduce some gates and expressions that we use in the proofs. If w is the coding of $u_1, \dots, u_n \in \Sigma^*$ then $\text{mstr}(w) = m_1(u_1) \cdot m_2(u_2) \cdots m_2(u_{n-1}) \cdot m_3(u_n)$ for binary gates $m_1, m_2, m_3 : \Sigma^* \rightarrow (\Sigma \times 2^Z)^*$ such that: (1) for all words $u = a_1 \cdots a_{2\ell} \in \Sigma^*, 2\ell > 2$

$$m_1(u) = (a_1, x_1) \cdots (a_{2\ell}, x_{2\ell}) \text{ where } x_i = \begin{cases} \{\text{fst}^g, \text{fst}^c, \text{lst}^c\} & \text{if } i = 1, \\ \{\text{fst}^c, \text{lst}^c\} & \text{if } i \text{ is odd and } i \neq 1, \\ \{\text{fst}^c\} & \text{if } i \text{ is even.} \end{cases}$$

$$m_2(u) = (a_1, x_1) \cdots (a_{2\ell}, x_{2\ell}) \text{ where } x_i = \begin{cases} \{\text{lst}^c\} & \text{if } i \text{ is odd,} \\ \{\text{fst}^c\} & \text{if } i \text{ is even.} \end{cases}$$

$$m_3(u) = (a_1, x_1) \cdots (a_{2\ell}, x_{2\ell}) \text{ where } x_i = \begin{cases} \{\text{lst}^c\} & \text{if } i \text{ is odd,} \\ \{\text{fst}^c, \text{lst}^c\} & \text{if } i \text{ is even and } i \neq 2\ell, \\ \{\text{fst}^c, \text{lst}^c, \text{lst}^g\} & \text{if } i = 2\ell. \end{cases}$$

(2) For each word $ab \in \Sigma^2$, $m_1(ab) = (a, \{\text{fst}^c, \text{fst}^g, \text{lst}^c\})(b, \{\text{fst}^c, \mathcal{S}\})$, $m_2(ab) = (a, \{\text{lst}^c, \mathcal{P}\})(b, \{\text{fst}^c, \mathcal{S}\})$, $m_3(ab) = (a, \{\text{lst}^c, \mathcal{P}\})(b, \{\text{fst}^c, \text{lst}^c, \text{lst}^g\})$. (3) For words of odd length the functions m_1, m_2, m_3 are fixed arbitrarily.

Let $\text{ise} : \Sigma^* \rightarrow \{0, 1\}$ be the binary gate defined as $\text{ise}(w) = 1$ precisely when $w \in \Sigma^*$ is not the empty word. Let $\text{bI} : \Sigma^* \times \{0, 1\} \rightarrow \Sigma^*$ be the binary function $\text{bI}(x, 1) = x$ and $\text{bI}(x, 0) = \epsilon$. For variables $\bar{X} = X_1 : \Sigma^*, \dots, X_n : \Sigma^*$, let $\text{NE}(\bar{X})$ be the expression $\bigwedge(\text{ise}(X_1), \dots, \text{ise}(X_n))$ of depth 2 that recognises the property that none of the input words is the empty word. Sometimes we use these gates and expressions over other alphabets, and then it is understood that the domains of the functions are appropriately defined.

Next we prove that class transductions and global transductions on n -codings can be defined by expressions of fixed height (irrespective of n). To summarise the intuition, let $w = w_1 \cdots w_n$ be the coding of the words $u_1, \dots, u_n \in \Sigma^*$ such that $\text{str}(w_i) = u_i$. Assume $\mathcal{A} : (\Sigma \times 2^Z)^* \rightarrow \Gamma^*$ is a class transducer that has a successful run on w and let $\mathcal{A}(w) = w' = w'_1 \cdots w'_n \in (\Gamma \times \mathcal{D})^*$ where w'_i is a relabelling of w_i . Observe that the only other positions in the class of a position in w_i appear either in w_{i-1} or w_{i+1} . Therefore to compute $\text{str}(w'_i)$ it suffices to know the words u_{i-1}, u_i, u_{i+1} and hence there is an expression that takes as inputs u_{i-1}, u_i, u_{i+1} and outputs the word $\text{str}(w'_i)$.

► **Lemma 13.** *For each class transducer $\mathcal{A} : (\Sigma \times 2^Z)^* \rightarrow \Gamma^*$ and each $n \in \mathbb{N}$ there exist combinatorial expressions $e_1(\bar{X}), \dots, e_n(\bar{X})$, where $\bar{X} = X_1 : \Sigma^*, \dots, X_n : \Sigma^*$, of depth 7 such that for all n -tuple $\bar{u} = (u_1, \dots, u_n)$ of words in Σ^* of identical even length $\text{coding}(e_1(\bar{u}), \dots, e_n(\bar{u})) = \mathcal{A}(\text{coding}(\bar{u}))$.*

Next we prove a similar claim for global transducers. The idea is as follows. Assume $\mathcal{A} : (\Sigma \times 2^Z)^* \rightarrow \Gamma^*$ is a global transducer and let $w = w_1 \cdots w_n$ be the coding of the words $u_1, \dots, u_n \in \Sigma^*$ such that $\text{str}(w_i) = u_i$. Assume that \mathcal{A} has a successful run on w and let $\mathcal{A}(w) = w' = w'_1 \cdots w'_n \in (\Gamma \times \mathcal{D})^*$ where w'_i is a relabelling of w_i . To compute $\text{str}(w'_i)$ it suffices to know the word u_i and the pair (p, q) of states of the automaton \mathcal{A} which are respectively the state of the automaton \mathcal{A} before and after reading the word $\text{mstr}(u_i)$ on the unique run on $\text{mstr}(w)$. Among these, the pair (p, q) can be computed a finitary function that aggregates the set of all possible partial runs of \mathcal{A} on each of the words u_1, \dots, u_n and hence an expression of fixed height can compute the word $\text{str}(w'_i)$.

► **Lemma 14.** *For each global transducer $\mathcal{A} : (\Sigma \times 2^Z)^* \rightarrow \Gamma^*$ and each $n \in \mathbb{N}$ there exist combinatorial expressions $e_1(\bar{X}), \dots, e_n(\bar{X})$, where $\bar{X} = X_1 : \Sigma^*, \dots, X_n : \Sigma^*$, of depth 5 such that for all n -tuple $\bar{u} = (u_1, \dots, u_n)$ of words in Σ^* of identical even length $\text{coding}(e_1(\bar{u}), \dots, e_n(\bar{u})) = \mathcal{A}(\text{coding}(\bar{u}))$.*

The above two lemmas can be generalised to a similar claim on cascades by induction (on the height of the cascade).

► **Lemma 15.** *For a cascade $\mathcal{C} = \langle \mathcal{A}_1, \dots, \mathcal{A}_k \rangle$ with input alphabet Σ , and each $n \in \mathbb{N}$ there exist combinatorial expressions $e_1(\bar{X}), \dots, e_n(\bar{X})$, where $\bar{X} = X_1 : \Sigma^*, \dots, X_n : \Sigma^*$, of depth at most $7k$ such that for all n -tuple $\bar{u} = (u_1, \dots, u_n)$ of words in Σ^* of identical even length $\text{coding}(e_1(\bar{u}), \dots, e_n(\bar{u})) = \mathcal{C}(\text{coding}(\bar{u}))$.*

Next we define a data language B_n that corresponds to the property V_n .

For a word $w = a_1 a_2 \dots a_l \in \{0, 1\}^*$ we let $\text{pad}(w) = 1a_1 1a_2 \cdots 1a_l$. We will also use pad as a binary gate. A *bridge* in a data word w is a sequence of positions along a path that consists of alternating class successor and global successor edges. Formally the sequence of positions i_1, \dots, i_n forms a bridge in w if there exists a sequence of successor and class successor edges e_1, \dots, e_{n-1} in w such that for each $1 \leq j < n$, $e_j = (i_j, i_{j+1})$ and for each $1 \leq j < n-1$, e_j is a successor edge iff e_{j+1} is a class successor edge. A bridge is a -labelled, for $a \in \Sigma$, if all the positions in the bridge are labelled by the letter a .

► **Definition 16.** Let $B_n \subseteq (\{0, 1\} \times \mathcal{D})^*$ be the set of all data words w such that w has a 1-labelled bridge i_1, \dots, i_{2n-1} (connected by a path of $2n-2$ edges), and

1. all positions to the left of i_1 are first positions of classes,
2. all positions to the right of i_{2n-1} are last positions of classes, and
3. the path corresponding to the bridge starts with a class successor edge.

Define the data language $B = \bigcup_{n=1}^{\infty} B_n$.

The language B_n is defined by the BMA formula (also in unary-DataLTL) of comp-height $2n+1$,

$$\text{F}^g (\text{H}^g \text{fst}^c \wedge (1\text{X}^c 1\text{X}^g)^n \text{G}^g \text{lst}^c) \quad \text{where } 1\text{X}^g \varphi \text{ stands for the formula } (1 \wedge \text{X}^g \varphi), \quad (1) \\ \text{and } 1\text{X}^c \varphi \text{ for } (1 \wedge \text{X}^c \varphi).$$

Similarly the language B is defined by the BR formula

$$\text{fst}^c \text{U}^g (\mu x. (1\text{X}^c 1\text{X}^g x \vee 1\text{X}^c 1\text{X}^g \text{G}^g \text{lst}^c)). \quad (2)$$

► **Proposition 17.** *Let (u_1, \dots, u_n) be a tuple of words of identical length over the alphabet $\{0, 1\}$. Then the following are equivalent.*

1. $(u_1, \dots, u_n) \in V_n$.



■ **Figure 4** The data word w corresponding to the words $a_1a_2, b_1b_2, c_1c_2, d_1d_2$, and a bridge of length 7 in w .

2. The data word $w = \text{coding}(\text{pad}(u_1), \dots, \text{pad}(u_n))$ is in the language B_n .

For a data language $L \subseteq (\Sigma \times \mathcal{D})^*$ we write $L^c = \{w \in (\Sigma \times \mathcal{D})^* \mid w \notin L\}$ for the complement of L . The data language $L \subseteq (\Sigma \times \mathcal{D})^*$ separates the data languages $L_1, L_2 \subseteq (\Sigma \times \mathcal{D})^*$ if $L_i \cap L = \emptyset$ and $L_{1-i} \subseteq L$ for some $i \in \{0, 1\}$. A cascade \mathcal{C} (respectively a formula φ) separates the data languages L_1, L_2 if $L(\mathcal{C})$ (respectively $L(\varphi)$) separates L_1, L_2 .

► **Lemma 18.** *If there is a cascade \mathcal{C} of height k that separates the data languages $L_1 = B_n \cap n\text{-Codings}$, $L_2 = (B_n)^c \cap n\text{-Codings}$ then there is a combinatorial expression of depth $7k + 4$ recognising the property V_n .*

Proof. Assume that \mathcal{C} is a cascade of height k separating the languages L_1, L_2 . Since cascades (of height k) are closed under complementation, without loss of generality assume that $L(\mathcal{C}) \supseteq L_1$ and $L(\mathcal{C}) \cap L_2 = \emptyset$. Therefore the cascade \mathcal{C} produces an output on a data word $n\text{-Codings} \ni w \in (\{0, 1\} \times \mathcal{D})^*$ if and only if w is in the language B_n . Let $e_1(\bar{X}), \dots, e_n(\bar{X})$, for $\bar{X} = X_1 : \{0, 1\}^*, \dots, X_n : \{0, 1\}^*$, be the combinatorial expressions of depth at most $7k$, guaranteed by the Lemma 15 such that for all n -tuple $\bar{u} = (u_1, \dots, u_n)$ of words in $\{0, 1\}^*$ of identical even length, $\text{coding}(e_1(\bar{u}), \dots, e_n(\bar{u})) = \mathcal{C}(\text{coding}(\bar{u}))$.

Let $\text{pad}(\bar{X})$ stand for the vector of expressions $\text{pad}(X_1), \dots, \text{pad}(X_n)$. We claim that the expression

$$e = \bigwedge (NE(e_1(\text{pad}(\bar{X})), \dots, e_n(\text{pad}(\bar{X})), t(X_1, \dots, X_n)) ,$$

where t is the expression from Example 10 for the alphabet $\{0, 1\}$ that checks if all the input words are of the same length, computes the property V_n . The expression e has depth at most $7k + 4$. To show the claim it is enough to verify that for a tuple $\bar{u} = (u_1, \dots, u_n)$ of words from $\{0, 1\}^*$ of equal length, none of the words $v_1 = e_1(\text{pad}(\bar{u})), \dots, v_n = e_n(\text{pad}(\bar{u}))$ is the empty word if and only if $\bar{u} \in V_n$. By Lemma 15, the words v_1 to v_n are nonempty iff \mathcal{C} accepts the data word $w = \text{coding}(\text{pad}(\bar{u}))$. By assumption, the data word w is accepted by the cascade \mathcal{C} iff $w \in B_n$. By Lemma 17, the data word w is in the language B_n iff \bar{u} is in the property V_n . Hence the claim is proved. ◀

We are now ready for the main theorem;

► **Theorem 19 (Separation).** *Let $N = 7k + 4$.*

1. The data languages $L_1 = B_{2^{N+1}} \cap (2^N + 1)\text{-Codings}$ and $L_2 = (B_{2^{N+1}})^c \cap (2^N + 1)\text{-Codings}$ are not separable by a formula in BMA of comp-height k .
2. The data language $B_{2^{N+1}}$ is not definable by a formula in BMA of comp-height k .
3. Class of BMA definable languages form a hierarchy under composition height; more precisely for every k there exists a BMA-formula φ with comp-height k that has no equivalent formula of comp-height $k - 1$.
4. The class of BMA definable languages is strictly subsumed by the class of BR definable languages.

Proof. (1.) Proof by contradiction. Assume that the data languages L_1, L_2 are separable by a BMA formula φ of comp-height k . This implies that there is cascade of height k separating

L_1, L_2 . By Lemma 18 there is an expression of depth N recognising the property V_{2^N+1} . This is in contradiction with Theorem 12.

(2.) Follows from (1).

(3.) From (2.) and the Equation (1), B_{2^N+1} is definable by a BMA formula of comp-height $2 \cdot (2^N + 1) + 1$ but not by any formula of comp-height k . Therefore (†) the set of languages defined by BMA^k is strictly contained in the set of languages defined by $\text{BMA}^{2 \cdot (2^N + 1) + 1}$.

It only remains to separate the languages definable by BMA^k and the languages definable by BMA^{k+1} , for all k . We prove this claim by contradiction. Assume that (★) there is some $m \in \mathbb{N}$ such that for every formula in BMA^{m+1} there is an equivalent formula in BMA^m . We claim that for every formula in BMA^{m+2} there is an equivalent formula in BMA^m as well. To prove the claim, let $\chi = \psi(\varphi_1, \dots, \varphi_n)$ be an arbitrary formula in BMA^{m+2} such that $\psi \in \text{BMA}^1$ and $\varphi_1, \dots, \varphi_n \in \text{BMA}^{m+1}$. By assumption (★) there exist formulas $\varphi'_1, \dots, \varphi'_n \in \text{BMA}^m$ equivalent to the formulas $\varphi_1, \dots, \varphi_n$ respectively. Therefore the formula $\chi' = \psi(\varphi'_1, \dots, \varphi'_n)$ is equivalent to the formula χ and is in BMA^{m+1} . Applying the assumption (★) again there is a formula $\chi'' \in \text{BMA}^m$ equivalent to χ' and hence also to χ , and hence the claim is proved. Extending this argument, by induction on k , it can be shown that for every formula in BMA^{m+k} there is an equivalent formula in BMA^m . This is in contradiction with the statement (†). Hence the statement is proved.

(4.) We claim that the data language B is not definable by any BMA formula. For the sake of contradiction, assume that there is a BMA formula φ of comp-height k recognising the language B and let \mathcal{C} be the cascade of height k corresponding to φ . We claim that the cascade \mathcal{C} separates the languages L_1 and L_2 . Clearly by definition of the language B , $L_1 \subseteq B$. We need to show that $L_2 \cap B = \emptyset$ and it suffices to prove that for every $w \in (2^N+1)\text{-Codings}$ if $w \in B$ then $w \notin (B_{2^N+1})^c$. Since any coding w in $(2^N+1)\text{-Codings}$ either belongs to B_{2^N+1} or does not belong to B , it follows that if $w \in B$ then $w \notin (B_{2^N+1})^c$. Therefore the cascade \mathcal{C} separates the languages L_1 and L_2 which contradicts (1.) and hence the claim follows. On the other hand, since B is definable by a formula in BR (Equation 2), the statement is proved. ◀

5 Conclusion

In this paper we studied the some fragments of μ -calculus over data words. We disclosed two fragments that are: the Bounded Reversal fragment (BR) and the Bounded Mode Alternation fragment (BMA) and proved they are separate. BR and BMA happen to form Boolean algebras making them very natural, and relatively expressive logics over data words. We also establish the relationship with earlier logics like FO^2 or DataLTL .

References

- 1 Henrik Björklund and Thomas Schwentick. On notions of regularity for data languages. *Theor. Comput. Sci.*, 411(4-5):702–715, 2010.
- 2 M. Bojańczyk and S. Lasota. An extension of data automata that captures xpath. In *Logic in Computer Science (LICS), 2010*, pages 243–252, July 2010.
- 3 Mikołaj Bojańczyk. Data monoids. In *STACS*, pages 105–116, 2011.
- 4 Mikołaj Bojańczyk, Claire David, Anca Muscholl, Thomas Schwentick, and Luc Segoufin. Two-variable logic on data words. *ACM Trans. Comput. Log.*, 12(4):27, 2011.
- 5 Thomas Colcombet, Clemens Ley, and Gabriele Puppis. On the use of guards for logics with data. In *MFCS*, volume 6907 of *LNCS*, pages 243–255. Springer, 2011.

- 6 Thomas Colcombet and Amaldev Manuel. Generalized data automata and fixpoint logic. In *FSTTCS 2014*, volume 29 of *LIPICs*, pages 267–278, 2014.
- 7 Thomas Colcombet and Amaldev Manuel. Combinatorial expressions and lower bounds. In *STACS*, 2015 (To appear).
- 8 S. Demri, D. Figueira, and M. Praveen. Reasoning about data repetitions with counter systems. In *Logic in Computer Science (LICS)*, 2013, pages 33–42, June 2013.
- 9 Stéphane Demri and Ranko Lazić. LTL with the freeze quantifier and register automata. *ACM Transactions on Computational Logic*, 10(3), April 2009.
- 10 D. Figueira. Alternating register automata on finite data words and trees. *Logical Methods in Computer Science*, 8(1), 2012.
- 11 D. Figueira. Decidability of downward XPath. *ACM Transactions on Computational Logic*, 13(4), 2012.
- 12 O. Grumberg, O. Kupferman, and S. Sheinvald. Variable automata over infinite alphabets. In *Language and Automata Theory and Applications*, pages 561–572. Springer, 2010.
- 13 M. Jurdziński and R. Lazic. Alternating automata on data trees and xpath satisfiability. *ACM Trans. Comput. Log.*, 12(3):19, 2011.
- 14 Michael Kaminski and Nissim Francez. Finite-memory automata. *Theor. Comput. Sci.*, 134(2):329–363, 1994.
- 15 Ahmet Kara, Thomas Schwentick, and Thomas Zeume. Temporal logics on words with multiple data values. In *FSTTCS*, volume 8 of *LIPICs*, pages 481–492, 2010.
- 16 L. Libkin and D. Vrgoc. Regular expressions for data words. In *LPAR*, volume 7180 of *Lecture Notes in Computer Science*, pages 274–288. SPRINGER, 2012.
- 17 A. Manuel, A. Muscholl, and G. Puppis. Walking on data words. In *Computer Science Theory and Applications*, volume 7913 of *LNCS*, pages 64–75. Springer, 2013.
- 18 F. Neven, T. Schwentick, and V. Vianu. Finite state machines for strings over infinite alphabets. *ACM Transactions on Computational Logic*, 5(3):403–435, 2004.