



A generalization of Cobham's theorem to automata over real numbers[☆]

Bernard Boigelot^{*}, Julien Brusten¹

Institut Montefiore, B28, Université de Liège, B-4000 Liège Sart-Tilman, Belgium

ARTICLE INFO

Keywords:

Automata
Real numbers
Mixed real-integer arithmetic
Cobham's theorem

ABSTRACT

This article studies the expressive power of finite-state automata recognizing sets of real numbers encoded positionally. It is known that the sets that are definable in the first-order additive theory of real and integer variables $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$ can all be recognized by weak deterministic Büchi automata, regardless of the encoding base $r > 1$. In this article, we prove the reciprocal property, i.e., a subset of \mathbb{R} that is recognizable by weak deterministic automata in every base $r > 1$ is necessarily definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. This result generalizes to real numbers the well-known Cobham's theorem on the finite-state recognizability of sets of integers. Our proof gives interesting insight into the internal structure of automata recognizing sets of real numbers, which may lead to efficient data structures for handling these sets.

© 2009 Elsevier B.V. All rights reserved.

1. Introduction

The verification of infinite-state systems, in particular the reachability analysis of systems modeled as finite-state machines extended with unbounded variables, has prompted the development of symbolic data structures for representing the sets of values that have to be handled during state-space exploration [2].

A simple representation strategy consists in using finite-state automata: The values in the considered domain are encoded as words over a given finite alphabet; a set of values is thus encoded as a language. If this language is regular, then a finite-state automaton that accepts it forms a representation of the set [23].

This approach has many advantages: Regular languages are closed under all usual set-theory operators (intersection, union, complement, Cartesian product, projection, etc.), and automata are easy to manipulate algorithmically. Deterministic automata can also be reduced to a canonical form, which simplifies comparison operations between sets.

The expressive power of automata is also well suited for verification applications. In the case of programs manipulating unbounded integer variables, it is known for a long time that the sets of integers that can be recognized by a finite-state automaton using the positional encoding of numbers in a base $r > 1$ correspond to those definable in an extension of Presburger arithmetic, i.e., the first-order additive theory of the integers $\langle \mathbb{Z}, +, < \rangle$ [10]. Furthermore, the well-known Cobham's theorem characterizes the sets that are representable by automata in all bases $r > 1$ as being exactly those that are Presburger-definable [11,8].

In order to analyze systems relying on integer and real variables, such as timed or hybrid automata, automata-based representations of numbers can be generalized to real values [3]. From a theoretical point of view, this amounts to moving

[☆] A preliminary version of this work appears in the proceedings of the 34th International Colloquium on Automata, Languages and Programming, (ICALP'07). This work is supported by the *Interuniversity Attraction Poles* program *MoVES* of the Belgian Federal Science Policy Office, and by the grant 2.4530.02 of the Belgian Fund for Scientific Research (F.R.S.-FNRS).

^{*} Corresponding author. Tel.: +32 43662970; fax: +32 43662620.

E-mail addresses: boigelot@montefiore.ulg.ac.be (B. Boigelot), brusten@montefiore.ulg.ac.be (J. Brusten).

URLs: <http://www.montefiore.ulg.ac.be/~boigelot> (B. Boigelot), <http://www.montefiore.ulg.ac.be/~brusten> (J. Brusten).

¹ Research fellow ("aspirant") of the Belgian Fund for Scientific Research (F.R.S.-FNRS).

from finite-word to infinite-word automata, which is not problematic. It has been shown that the sets of reals that can be recognized by infinite-word automata in a given encoding base are those definable in an extension of the first-order additive theory of real and integer variables $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$ [6].

In practice though, handling infinite-word automata can be difficult, especially if set complementation needs to be performed. It is however known that, for representing the sets definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$, the full expressive power of Büchi automata is not required, and that the much simpler subclass of *weak deterministic* automata is sufficient [4]. The advantage is that, from an algorithmic perspective, handling weak automata is similar to manipulating finite-word automata.

A natural question is then to characterize precisely the expressive power of weak deterministic automata representing sets of real numbers. For a given encoding base $r > 1$, it is known that the representable sets form a base-dependent extension of $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. This covers, in particular, all the sets definable in $\langle \mathbb{R}, \mathbb{Z}, +, <, P_r \rangle$, where P_r is a predicate that checks whether its argument is a power of r [7].

This article is aimed at characterizing the subsets of \mathbb{R} that can be represented as weak deterministic automata in multiple bases. Our central result is to show that, for two relatively prime bases r_1 and r_2 , the sets that are simultaneously recognizable in bases r_1 and r_2 can be defined in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. As a corollary, such sets are then representable in any base $r > 1$.

The intuition behind our proof is the following. First, we reduce the problem to characterizing the representable subsets of $[0, 1]$. We then introduce the notion of interval boundary points, as points with special topological properties, and establish that a set representable in multiple bases can only contain finitely many such points. Finally, we show that this property implies that S is definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. The argument used for this last step provides a description of the internal structure of automata representing sets definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. This result may help one to develop efficient data structures for handling such sets.

2. Representing sets of numbers with automata

In this section, we briefly present the automata-based representations of sets of integer and real values.

2.1. Number decision diagrams

Let $r > 1$ be an integer base. A natural number $x \in \mathbb{N}$ can be encoded positionally in base r by finite words $b_{p-1}b_{p-2} \dots b_1b_0$ over the alphabet $\Sigma_r = \{0, 1, \dots, r-1\}$, such that $x = \sum_{i=0}^{p-1} b_i r^i$. Negative values are encoded by their r 's-complement, i.e., the encodings of $x \in \mathbb{Z}$ with $x < 0$ are formed by the last p digits of the encodings of $r^p + x$. The length p of the encodings of a number $x \in \mathbb{Z}$ is not fixed, but must be non-zero and large enough for $-r^{p-1} \leq x < r^{p-1}$ to hold. As a consequence, the most significant digit of encodings, called the *sign digit*, is equal to $r-1$ for strictly negative numbers, and to 0 for positive numbers. This digit can always be repeated arbitrarily many times without influencing the encoded value. Each integer thus admits an infinite number of distinct encodings, differing only in the number of repetitions of their sign digit.

This encoding scheme maps a subset S of \mathbb{Z} onto a language L over Σ_r . This language contains all the encodings of the elements of S , i.e., we have $u \cdot w \in L$, with $u \in \Sigma_r$ and $w \in \Sigma_r^*$, iff $u \cdot u \cdot w \in L$. If the language L is regular, then a finite-state automaton that accepts it is called a *Number Decision Diagram (NDD)*, and is said to represent, or recognize, the set S . NDDs can be generalized to representing subsets of \mathbb{Z}^n , i.e., sets of vectors, for any $n > 0$ [10,22,2].

It has been shown [10,20,8] that the subsets of \mathbb{Z} recognizable by NDDs in a base $r > 1$ are exactly those that can be defined in the first-order theory $\langle \mathbb{Z}, +, <, V_r \rangle$ where $V_r(x)$ is the function mapping an integer $x > 0$ to the greatest power of r dividing it. Moreover, the sets that are recognizable by NDDs in every base $r > 1$ have been characterized by Cobham [11] as being exactly those that are definable in $\langle \mathbb{Z}, +, < \rangle$, i.e., Presburger arithmetic [16]. This result has been extended to subsets of \mathbb{Z}^n by Semenov [18].

Computing the intersection, union, complementation, difference and Cartesian product of sets represented by NDDs reduces to performing the corresponding operations on the languages accepted by the automata. Projection is more tricky, as the resulting automaton has to be completed in order to accept all the encodings of the vectors it recognizes [5]. Finally, since NDDs are finite-word automata, they can be determinized, as well as minimized into a canonical form.

2.2. Real number automata

Real numbers can also be encoded positionally. Let $r > 1$ be a base. An encoding w of a number $x \in \mathbb{R}$ is an infinite word $w_I \cdot \star \cdot w_F$ over $\Sigma_r \cup \{\star\}$, where $w_I \in \Sigma_r^*$ encodes the integer part $x_I \in \mathbb{Z}$ of x , and $w_F \in \Sigma_r^\omega$ its fractional part $x_F \in [0, 1]$, i.e., we have $w_F = b_1b_2b_3 \dots$ with $x_F = \sum_{i>0} b_i r^{-i}$. Note that some numbers have two distinct encodings with the same integer-part length. For example, in base 10, the number $11/2$ has the encodings $0^+ \cdot 5 \cdot \star \cdot 5 \cdot 0^\omega$ and $0^+ \cdot 5 \cdot \star \cdot 4 \cdot 9^\omega$. Such encodings are said to be *dual*. We denote by Λ_r the set of valid prefixes of base- r encodings that include a separator, i.e., $\Lambda_r = \{0, r-1\} \cdot \Sigma_r^* \cdot \star \cdot \Sigma_r^*$. For a word $w \in \Lambda_r \cdot \Sigma_r^\omega$, we denote by $[w]_r$ the real number encoded by w in base r . Similarly, for $w \in \{0, r-1\} \cdot \Sigma_r^*$, $[w]_r$ denotes the integer number encoded by w in base r , i.e., $[w]_r = [w \cdot \star \cdot 0^\omega]_r$.

Similarly to the case of integers, the base- r encoding scheme transforms a set $S \subseteq \mathbb{R}$ into a language $L(S) \subseteq \Lambda_r \cdot \Sigma_r^\omega$. A *Real Number Automaton (RNA)* is defined as a Büchi automaton that accepts the language containing all the base- r encodings of the elements of S . This representation can be generalized into *Real Vector Automata (RVA)*, suited for subsets of \mathbb{R}^n ($n > 0$) [3].

state of \mathcal{A} to q_i , omitting the separator symbol \star . Moreover, we define L_i^F as the language of all (infinite) words that can be accepted from q_i .

The language L accepted by \mathcal{A} is thus of the form $\bigcup_{i=1}^n L_i^I \cdot \star \cdot L_i^F$, where for all i , $L_i^I \subseteq \Sigma_{r_1}^*$ encodes the integer part, and $L_i^F \subseteq \Sigma_{r_1}^\omega$ the fractional part, of the encodings of numbers $x \in S$. More precisely, for every i , let $S_i^I \subseteq \mathbb{Z}$ denote the set encoded by L_i^I and let $S_i^F \subseteq [0, 1]$ denote the set encoded by $0^+ \cdot \star \cdot L_i^F$. We have $S = \bigcup_{i=1}^n (S_i^I + S_i^F)$. Note that each L_i^I is recognizable by a NDD in base r_1 and that, similarly, each language of the form $0^+ \cdot \star \cdot L_i^F$ is recognizable by a RNA.²

The definition of the sets S_i^I and S_i^F , for $i \in [1, \dots, n]$, relies on the languages L_i^I and L_i^F , whose definition depends in turn on the encoding base. We now show that those sets can also be defined independently from this base. For each $x \in \mathbb{Z}$, we define the set $S^F(x) = \{y \in [0, 1] \mid x + y \in S\}$ of its corresponding fractional parts in S . Since \mathcal{A} is in canonical form, the languages L_i^F are pairwise different, and so are the sets S_i^F . Besides, for each $i \in [1, \dots, n]$, there exists $x_i \in \mathbb{Z}$ such that $S_i^F = S^F(x_i)$. This last expression does not involve the representation base, and makes it possible to define the sets S_i^F independently from this base. Formally, we have $\{S_1^F, S_2^F, \dots, S_n^F\} = \{S^F(x) \mid x \in \mathbb{Z}\}$.

Similarly, given a set $U \subseteq [0, 1]$ of fractional parts, we define the set $S^I(U) = \{x \in \mathbb{Z} \mid \forall y \in [0, 1] : y \in U \Leftrightarrow x + y \in S\}$ of its corresponding integer parts in S . The sets S_i^I can then be defined independently from the representation base: We have, for each $i \in [1, \dots, n]$, $S_i^I = S^I(S_i^F)$.

We have thus established that the decomposition of S into $S = \bigcup_{i=1}^n (S_i^I + S_i^F)$ is independent from the representation base, and that each set S_i^I and each set S_i^F is recognizable in every base in which S is recognizable. Therefore, if S is recognizable in two relatively prime bases r_1 and r_2 , then so are S_i^I and S_i^F for every i . From Cobham's theorem, each S_i^I must then be definable in $(\mathbb{Z}, +, <)$. In order to show that S is definable in $(\mathbb{R}, \mathbb{Z}, +, <)$, it is hence sufficient to prove that each S_i^F is definable in that theory. We have thus reduced the problem of characterizing the subsets of \mathbb{R} that are simultaneously recognizable in two relatively prime bases to the same problem over the subsets of $[0, 1]$.

4. Interval boundary points

We now consider a set $S \subseteq [0, 1]$ represented by a weak deterministic RNA \mathcal{A} . We define the *interval boundary points* of S as points with specific topological properties, and establish a relation between the existence of such points and some structures in the transition graph of \mathcal{A} .

4.1. Definitions

A *neighborhood* $N_\varepsilon(x)$ of a point $x \in \mathbb{R}$, with $\varepsilon > 0$, is the set $N_\varepsilon(x) = \{y \mid |x - y| < \varepsilon\}$. A point $x \in \mathbb{R}$ is a *boundary point* of S iff all its neighborhoods contain points from S as well as from its complement \bar{S} , i.e., $\forall \varepsilon > 0 : N_\varepsilon(x) \cap S \neq \emptyset \wedge N_\varepsilon(x) \cap \bar{S} \neq \emptyset$.

A *left neighborhood* $N_\varepsilon^<(x)$ of a point $x \in \mathbb{R}$, with $\varepsilon > 0$, is the set $N_\varepsilon^<(x) = \{y \mid x - \varepsilon < y < x\}$. Similarly, a *right neighborhood* $N_\varepsilon^>(x)$ of x is defined as $N_\varepsilon^>(x) = \{y \mid x < y < x + \varepsilon\}$. A boundary point x of S is a *left interval boundary point* of S iff it admits a left neighborhood $N_\varepsilon^<(x)$ that is entirely contained in either S or \bar{S} , i.e., $\exists \varepsilon > 0 : N_\varepsilon^<(x) \subseteq S \vee N_\varepsilon^<(x) \subseteq \bar{S}$. *Right interval boundary points* are defined in the same way. A point $x \in \mathbb{R}$ is an *interval boundary point* of S iff it is a left or a right interval boundary point of S .

Each interval boundary point x of S is thus characterized by its direction (left or right), its polarity w.r.t. S (i.e., whether $x \in S$ or $x \notin S$), and the polarity of its left or right neighborhoods of sufficiently small size (i.e., whether they are subsets of S or of \bar{S}). The possible combinations define eight *types* of interval boundary points, that are illustrated in Fig. 2.

Remark that a boundary point is not necessarily an interval boundary point. For instance, in any base r , each point of $[0, 1]$ is a boundary point of the set of numbers that have an encoding ending in 0^ω , but not an interval boundary point of that set.

4.2. Recognizing interval boundary points

Recall that \mathcal{A} is a weak deterministic RNA recognizing a set $S \subseteq [0, 1]$. Let $r > 1$ be the representation base. We assume w.l.o.g. that \mathcal{A} is in canonical form (and thus that its transition relation is complete). Consider a path π of \mathcal{A} that reads an encoding w of a left interval boundary point x of S . This path eventually reaches a strongly connected component C that it does not leave. The accepting status of C corresponds to the polarity of x w.r.t. S .

We first consider the case of a component C that is neither empty nor universal. Since x is a left interval boundary point, all its sufficiently small left neighborhoods are either subsets of S or subsets of \bar{S} , depending on the type of x . Hence, from each state s of C visited infinitely many times by π , its outgoing transitions labeled by smaller digits than the one read in π must necessarily lead to either the universal or the empty strongly connected component of \mathcal{A} , i.e., those accepting respectively

² In order for such a RNA to recognize all encodings of numbers, it should also accept the words $(r_1 - 1)^+ \cdot \star \cdot (r_1 - 1)^\omega$ if $0^\omega \in L_i^F$, and $0^+ \cdot 1 \cdot \star \cdot 0^\omega$ if $(r_1 - 1)^\omega \in L_i^F$.

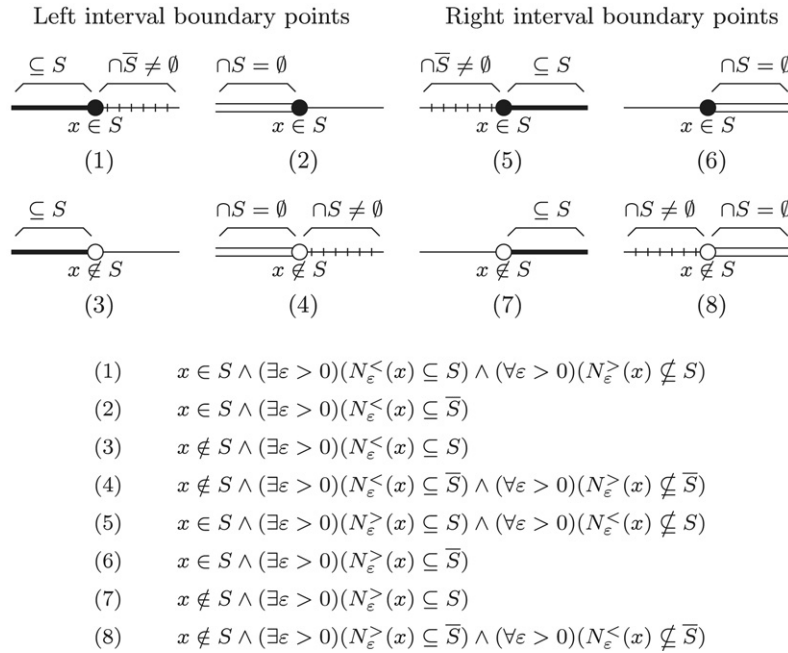


Fig. 2. Types of interval boundary points.

the languages Σ_r^ω and \emptyset . (Recall that \mathcal{A} is assumed to be in canonical form.) It follows that, after having reached some state s in C , the path π follows the transitions within C that are labeled by the smallest possible digits, hence it eventually cycles through a loop.

In the case where C is empty or universal, a path π recognizing a left interval boundary point x of S must necessarily read a word ending with 0^ω , otherwise there would exist a left neighborhood of x entirely in S if $x \in S$, or entirely in \bar{S} if $x \notin S$.

Similar results hold for right interval boundary points, which are read by paths that eventually follow the largest possible digits in their terminal strongly connected component.

As a consequence, every base- r encoding w of an interval boundary point x of S is necessarily *ultimately periodic*, i.e., such that $w = u \cdot \star \cdot u' \cdot v^\omega$, with $u \in \{0, r-1\} \cdot \Sigma_r^*$, $u' \in \Sigma_r^*$ and $v \in \Sigma_r^+$. We then have $r^{|u'|+|v|}x = [u \cdot u' \cdot v \cdot \star \cdot v^\omega]_r$ and $r^{|u'|}x = [u \cdot u' \cdot \star \cdot v^\omega]_r$, from which we get

$$x = \frac{[u \cdot u' \cdot v]_r - [u \cdot u']_r}{r^{|u'|}(r^{|v|} - 1)}.$$

Besides, each ultimate period v of such encodings can be uniquely determined from a suitable state of \mathcal{A} associated with a direction (left or right). We thus have the following results.

Theorem 1. *Each interval boundary point of a subset of $[0, 1]$ that is recognizable by a weak deterministic RNA is a rational number.*

Theorem 2. *Let $S \subseteq [0, 1]$ be a set recognizable by a weak deterministic RNA in a base $r > 1$. The set of ultimate periods of the base- r encodings of the interval boundary points of S is finite.*

4.3. Recognizing interval boundary points in multiple bases

Consider now a set $S \subseteq [0, 1]$ that is simultaneously recognizable by weak deterministic RNA in two relatively prime bases $r_1 > 1$ and $r_2 > 1$. Let \mathcal{A}_1 and \mathcal{A}_2 denote, respectively, such RNA.

Suppose that S has infinitely many interval boundary points. From Theorem 2, there must exist some ultimate period $v \in \Sigma_{r_1}^+$ such that infinitely many interval boundary points of S have base- r_1 encodings of the form $u_i \cdot v^\omega$, with $\forall i : u_i \in 0^+ \cdot \star \cdot \Sigma_{r_1}^*$. Infinitely many of those encodings are such that u_i and v do not end with the same digit (this restriction expresses that u_i is chosen as small as possible). Furthermore, we can assume that $v \notin 0^*$ if S has infinitely many left interval boundary points, and that $v \notin (r-1)^*$ if S has infinitely many right interval boundary points. It follows that there exists such a word u_i with a length greater than the number of states in \mathcal{A}_1 . The path π of \mathcal{A}_1 that reads $u_i \cdot v^\omega$ thus decomposes into $\pi = \pi_1 \pi_2 \pi_3$, where π_2 is cyclic and reads a word $w_2 \in \Sigma_{r_1}^+$. Thus, from the definitions of the interval boundary points, there exist $w_1 \in 0^+ \cdot \star \cdot \Sigma_{r_1}^*$ and $w_2, w_3 \in \Sigma_{r_1}^*$, with $|w_2| > 0$, such that for every $k \geq 0$, the word $w_1 \cdot (w_2)^k \cdot w_3 \cdot v^\omega$ encodes an interval boundary point of S .

Each word of this form is ultimately periodic, thus it encodes in base r_1 a rational number that can also be encoded by an ultimately periodic word in base r_2 . We use the following lemma.

Lemma 3. Let $r_1 > 1$ and $r_2 > 1$ be relatively prime bases, and let $w_1 \in 0^+ \cdot \star \cdot \Sigma_{r_1}^*$, $w_2, w_3, w_4 \in \Sigma_{r_1}^*$, with $|w_2| > 0$, $|w_4| > 0$, such that the words $w_2 \cdot w_3$ and w_4 do not end with the same digit. The subset of \mathbb{Q} encoded in base r_1 by the language $w_1 \cdot (w_2)^* \cdot w_3 \cdot (w_4)^\omega$ cannot be encoded in base r_2 with only a finite number of ultimate periods.

Proof. For every $k \geq 0$, we define $x_k = [w_1 \cdot (w_2)^k \cdot w_3 \cdot (w_4)^\omega]_{r_1}$. The prefix w_1 can be decomposed into $w_1 = w'_1 \cdot \star \cdot w''_1$, with $w'_1 \in \{0, r_1 - 1\} \cdot \Sigma_{r_1}^*$ and $w''_1 \in \Sigma_{r_1}^*$. We have for every $k > 0$,

$$\begin{aligned} r_1^{|w''_1|+k|w_2|+|w_3|+|w_4|} x_k &= [w'_1 \cdot w''_1 \cdot (w_2)^k \cdot w_3 \cdot w_4 \cdot \star \cdot (w_4)^\omega]_{r_1}, \\ r_1^{|w''_1|+k|w_2|+|w_3|} x_k &= [w'_1 \cdot w''_1 \cdot (w_2)^k \cdot w_3 \cdot \star \cdot (w_4)^\omega]_{r_1}, \end{aligned}$$

hence

$$r_1^{|w''_1|+k|w_2|+|w_3|} (r_1^{|w_4|} - 1) x_k = [w'_1 \cdot w''_1 \cdot (w_2)^k \cdot w_3 \cdot w_4]_{r_1} - [w'_1 \cdot w''_1 \cdot (w_2)^k \cdot w_3]_{r_1},$$

which gives

$$x_k = \frac{y_k}{r_1^{|w''_1|+k|w_2|+|w_3|} (r_1^{|w_4|} - 1)}, \quad (1)$$

with $y_k = (r_1^{|w_4|} - 1)[w'_1 \cdot w''_1 \cdot w_2^k \cdot w_3]_{r_1} + [0 \cdot w_4]_{r_1}$. Remark that y_k is an integer, but cannot be a multiple of r_1 . Indeed, we have $y_k \bmod r_1 = ([0 \cdot w_4]_{r_1} - [w'_1 \cdot w''_1 \cdot w_2^k \cdot w_3]_{r_1}) \bmod r_1$, which is non-zero thanks to the hypothesis on the last digits of $w_2 \cdot w_3$ and w_4 .

We now develop the expression of y_k . For every $k > 0$, we have

$$\begin{aligned} y_k &= (r_1^{|w_4|} - 1) \left(r_1^{k|w_2|+|w_3|} [w'_1 \cdot w''_1]_{r_1} + r_1^{|w_3|} \frac{r_1^{k|w_2|} - 1}{r_1^{|w_2|} - 1} [0 \cdot w_2]_{r_1} + [0 \cdot w_3]_{r_1} \right) + [0 \cdot w_4]_{r_1} \\ &= \frac{z_k}{r_1^{|w_2|} - 1}, \end{aligned}$$

with

$$\begin{aligned} z_k &= a r_1^{k|w_2|} + b, \\ a &= r_1^{|w_3|} (r_1^{|w_4|} - 1) ((r_1^{|w_2|} - 1) [w'_1 \cdot w''_1]_{r_1} + [0 \cdot w_2]_{r_1}), \quad \text{and} \\ b &= -r_1^{|w_3|} (r_1^{|w_4|} - 1) [0 \cdot w_2]_{r_1} + (r_1^{|w_2|} - 1) (r_1^{|w_4|} - 1) [0 \cdot w_3]_{r_1} + (r_1^{|w_2|} - 1) [0 \cdot w_4]_{r_1}. \end{aligned}$$

Substituting in (1), we get

$$x_k = \frac{z_k}{r_1^{|w''_1|+k|w_2|+|w_3|} (r_1^{|w_2|} - 1) (r_1^{|w_4|} - 1)}. \quad (2)$$

Since $z_k = (r_1^{|w_2|} - 1) y_k$ and $y_k \bmod r_1 \neq 0$, we have $z_k \bmod r_1 \neq 0$, hence $b \neq 0$. Consider a prime factor f of r_1 , and define l as the greatest integer such that f^l divides b . For every $k > l$, we have $z_k \bmod f^l = 0$ and $z_k \bmod f^{l+1} = b \bmod f^{l+1} \neq 0$. It follows that the reduced rational expression of x_k , i.e., $x_k = n_k/d_k$ with $n_k, d_k \in \mathbb{Z}$, $d_k > 0$ and $\gcd(n_k, d_k) = 1$, is such that f^{k-l} divides d_k for every $k > l$. Indeed, the numerator of (2) is not divisible by f^{l+1} whereas its denominator is divisible by f^{k+1} .

Assume now, by contradiction, that the set $\{x_k \mid k \geq 0\}$ can be represented in base r_2 using only a finite number of ultimate periods. Then, there exists an ultimate period $v \in \Sigma_{r_2}^+$ such that for infinitely many values of k , we have

$$x_k = [u'_k \cdot \star \cdot u''_k \cdot v^\omega]_{r_2},$$

with $u'_k \in \{0, r_2 - 1\} \cdot \Sigma_{r_2}^*$ and $u''_k \in \Sigma_{r_2}^*$. We thus have, for these values of k ,

$$\begin{aligned} r_2^{|u''_k|+|v|} x_k &= [u'_k \cdot u''_k \cdot v \cdot \star \cdot v^\omega]_{r_2}, \\ r_2^{|u''_k|} x_k &= [u'_k \cdot u''_k \cdot \star \cdot v^\omega]_{r_2}, \end{aligned}$$

hence

$$x_k = \frac{[u'_k \cdot u''_k \cdot v]_{r_2} - [u'_k \cdot u''_k]_{r_2}}{r_2^{|u''_k|} (r_2^{|v|} - 1)}.$$

Since $(r_2^{|v|} - 1)$ is bounded, and r_2 is relatively prime with r_1 by hypothesis, the denominator of this expression can only be divisible by a bounded number of powers of f , which contradicts our previous result. \square

Together with [Theorem 2](#), this lemma contradicts our assumption that S has infinitely many interval boundary points. We thus have the following theorem.

Theorem 4. *If a set $S \subseteq [0, 1]$ is simultaneously recognizable by weak deterministic RNA in two relatively prime bases, then it has finitely many interval boundary points.*

We therefore call a set that satisfies the hypotheses of [Theorem 4](#) a *finite-boundary set*.

5. Finite-boundary sets

Our goal is now to characterize the structure of the transition graph of RNA that recognize finite-boundary sets. We start by establishing some properties that hold for all weak deterministic RNA, and then focus on the specific case of finite-boundary sets.

5.1. Properties of weak deterministic RNA

Let \mathcal{A} be a weak deterministic RNA, assumed to be in canonical form, recognizing a subset of $[0, 1]$ in a base $r > 1$. Consider a strongly connected component C of \mathcal{A} , from which each outgoing transition leads to either the universal or the empty strongly connected component, i.e., those accepting respectively the languages Σ_r^ω and \emptyset .

Lemma 5. *Let π be a minimal (resp. maximal) infinite path within C , i.e., a path that follows from each visited state the transition of C labeled by the smallest (resp. largest) possible digit. The destination states of all outgoing transitions from states visited by π , and labeled by a smaller (resp. larger) digit than the one read in π , are identical.*

Proof. We first study the case of two transitions t_1 and t_2 originating from the same state s visited by π , that are respectively labeled by digits d_1, d_2 smaller than the digit d read from s in π . Among the digits that satisfy this condition, one can always find consecutive values, hence it is sufficient to consider the case where $d_2 = d_1 + 1$.

Let σ be a finite path that reaches s from the initial state of \mathcal{A} . By appending to σ suffixes that read $d_1 \cdot (r - 1)^\omega$ and $d_2 \cdot 0^\omega$, one obtains paths that recognize dual encodings of the same number, hence these paths must be either both accepting or both non-accepting. Therefore, t_1 and t_2 share the same destination.

Consider now transitions t_1 and t_2 from distinct states s_1 and s_2 visited by π , labeled by smaller digits than those – respectively denoted d_1 and d_2 – read in π . We can assume w.l.o.g. that s_1 and s_2 are consecutive among the states visited by π that have such outgoing transitions. In other words, the subpath of π that links s_1 to s_2 is labeled by a word of the form $d_1 \cdot 0^k$, with $d_1 > 0$ and $k \geq 0$. The digit d_2 is thus the first non-zero digit read by π after d_1 .

Let σ' be a finite path that reaches s_1 from the initial state of \mathcal{A} . Appending to σ' suffixes that read $(d_1 - 1) \cdot (r - 1)^\omega$ and $d_1 \cdot 0^\omega$ yields paths that read dual encodings of the same number, hence these paths must be either both accepting or both non-accepting. The destinations of the transitions that leave C from s_1 and s_2 must thus be identical.

The case of maximal paths is handled in the same way. \square

Lemma 6. *There exists a state $s \in C$ from which either the outgoing transition labeled by 0, or the one labeled by $r - 1$, leads to the empty strongly connected component if C is accepting, and to the universal one if C is non-accepting.*

Proof. Consider first the case of an accepting component C . By contradiction, suppose that from each state $s \in C$, the destinations of the transitions labeled by 0 and by $r - 1$ are either states of C , or belong to the universal strongly connected component.

Since \mathcal{A} is in canonical form and the outgoing transitions from C lead, by hypothesis, to the universal or the empty strongly connected component, there exists a transition from a state $s' \in C$ labeled by a digit $c \in \Sigma_r$, ending in the empty component. W.l.o.g., we assume that c is the smallest possible digit that satisfies this condition. By the contradiction hypothesis, we have $c > 0$. Hence, for every $c' < c$, the destination of the transition originating from s' and labeled by c' is either a state of C , or the universal strongly connected component. We choose $c' = c - 1$, and define w as the label of a path from the initial state of \mathcal{A} to s' . The word $w \cdot c \cdot 0^\omega$ is not accepted, since reading c from s leads to the empty component. On the other hand, the word $w \cdot (c - 1) \cdot (r - 1)^\omega$ is accepted by \mathcal{A} . Indeed, reading c from s leads to either the universal component, or to a state of C from which, by hypothesis, $(r - 1)^\omega$ must be accepted. We thus have a contradiction with the fact that \mathcal{A} accepts all the encodings of the numbers it recognizes.

The proof is similar for the case of a non-accepting component C . \square

The following result now expresses a constraint on the trivial strongly connected components of the fractional part of \mathcal{A} (i.e., the part of \mathcal{A} reached after reading one occurrence of the symbol \star).

Lemma 7. *From any trivial strongly connected component of the fractional part of \mathcal{A} , there must exist a reachable strongly connected component that is neither empty, trivial, nor universal.*

Proof. The proof is by contradiction. Let $\{s\}$ be a trivial strongly connected component of the fractional part of \mathcal{A} . Assume that all paths from s eventually reach the universal or the empty strongly connected component, after passing only through trivial components. As a consequence, the language accepted from s is of the form $L \cdot \Sigma_r^\omega$, where $L \subset \Sigma_r^*$ is finite. We can require w.l.o.g. that all words in L share the same length l . Note that L cannot be empty or equal to Σ_r^l , since s does not belong to the empty or universal components.

Each word in Σ_r^l can be seen as the base- r encoding of an integer in the interval $[0, r^l - 1]$. Since L is neither empty nor universal, there exist two words $w_1, w_2 \in \Sigma_r^l$ that do not both belong to L or to $\Sigma_r^l \setminus L$, and that encode two consecutive integers n and $n + 1$. Then, $u \cdot w_2 \cdot 0^\omega$ and $u \cdot w_1 \cdot (r - 1)^\omega$ encode the same number in base r , where u is the label of an arbitrary path from the initial state of \mathcal{A} to s . This contradicts the fact that \mathcal{A} accepts all the encodings of the numbers it recognizes. \square

5.2. Properties of RNA recognizing finite-boundary sets

Theorem 8. *Let \mathcal{A} be a weak deterministic RNA, supposed to be in canonical form, recognizing a finite-boundary set $S \subseteq [0, 1]$. Each non-trivial, non-empty and non-universal strongly connected component of the fractional part of \mathcal{A} takes the form of a single cycle. Moreover, from each such component, the only reachable strongly connected components besides itself are the empty or the universal ones.*

Proof. Let C be a non-empty and non-universal strongly connected component of the fractional part of \mathcal{A} , from which the only reachable components (besides itself) are the empty and the universal one. From Lemma 7, C cannot be trivial. By Lemma 6, there exists a state $s \in C$ and a digit $d \in \{0, r - 1\}$ such that the outgoing transition from s labeled by d leads to the empty component if C is accepting, and to the universal one otherwise.

Consider first the case where $d = 0$. The path π from s that stays within C and follows the transitions with the smallest possible digits is cyclic. From the definition of the interval boundary points of S , the label of π corresponds to a suffix of the encoding of a left interval boundary point of S . If on the other hand $d = r - 1$, then the path π from s that stays within C and follows the transitions with the largest possible digits is cyclic as well, and recognizes a suffix of the encoding of a right interval boundary point of S .

In both cases, if C contains other cycles, or if C is reachable from other non-trivial strongly connected components in the fractional part of \mathcal{A} , then π can be prefixed by infinitely many reachable paths from an entry state of the fractional part of \mathcal{A} to s . This contradicts the fact that S has only finitely many interval boundary points. \square

This result characterizes quite precisely the shape of the fractional part of a weak deterministic RNA recognizing a finite-boundary set: Its transition graph is first composed of a bottom layer of strongly connected components containing only the universal and the empty one, and then a (possibly empty) layer of single-cycle components leading to the bottom layer. Thanks to Lemma 5, the transitions that leave a single-cycle component with a smaller (or larger) digit all lead to the same empty or universal component (which may differ for the smaller and larger cases). Thus, each single-cycle component can simply be characterized by its label and the polarity of its smaller and greater alternatives. Finally, the two layers of non-trivial strongly connected components can be reached through an acyclic structure of trivial components, such that from each of them, there is at least one outgoing path leading to a single-cycle component. This structure is illustrated in Fig. 3.

As a consequence, we are now able to describe the language accepted by such a RNA.

Theorem 9. *Let \mathcal{A} be a weak deterministic RNA recognizing a finite-boundary set $S \subseteq [0, 1]$ encoded in a base $r > 1$. For every word $v \in \Sigma_r^+$, let $L_<(v)$ (resp. $L_>(v)$) denote the language of all the infinite words over Σ_r that are lexicographically smaller (resp. larger) than v^ω . The language $L(\mathcal{A})$ accepted by \mathcal{A} can be expressed as*

$$L(\mathcal{A}) = \bigcup_i L' \cdot w_i \cdot \Sigma_r^\omega \cup \bigcup_i L' \cdot w'_i \cdot (v_i)^\omega \cup \bigcup_i L' \cdot w''_i \cdot L_<(v'_i) \cup \bigcup_i L' \cdot w'''_i \cdot L_>(v''_i) \cup L_0 \cup L_1,$$

where each union is finite (and possibly empty), $\forall i : w_i, w'_i, w''_i, w'''_i, v_i, v'_i, v''_i \in \Sigma_r^*$ with $|v_i| > 0, |v'_i| > 0, |v''_i| > 0$, $L' = 0^+ \cdot \star$, L_0 is either empty or equal to $(r - 1)^+ \cdot \star \cdot (r - 1)^\omega$, and L_1 is either empty or equal to $0^+ \cdot 1 \cdot \star \cdot 0^\omega$.

Proof. The accepting paths of \mathcal{A} end up in either a cyclic component, or the universal one. There are finitely many cyclic components, and each of them can only be reached by finitely many prefixes. Hence, the accepting paths that end up in a cyclic component recognize a language of the form $\bigcup_i L' \cdot w'_i \cdot (v_i)^\omega$.

We now consider the paths that end up in the universal component. Before reaching this component, such paths can visit either only trivial components, or trivial components followed by one single cyclic component. The paths in the former category recognize a language of the form $\bigcup_i L' \cdot w_i \cdot \Sigma_r^\omega$. For the latter case, we know that after leaving a cyclic component labeled by $v \in \Sigma_r^+$, the accepting paths read either the words that are smaller than v^ω , or those that are larger than v^ω , or a combination of both. These paths thus recognize a language of the form $\bigcup_i L' \cdot w''_i \cdot L_<(v'_i) \cup \bigcup_i L' \cdot w'''_i \cdot L_>(v''_i)$.

Finally, the terms L_0 and L_1 are introduced in order to deal with the dual encodings of 0 and 1, since \mathcal{A} must accept all the encodings of the elements of S . \square

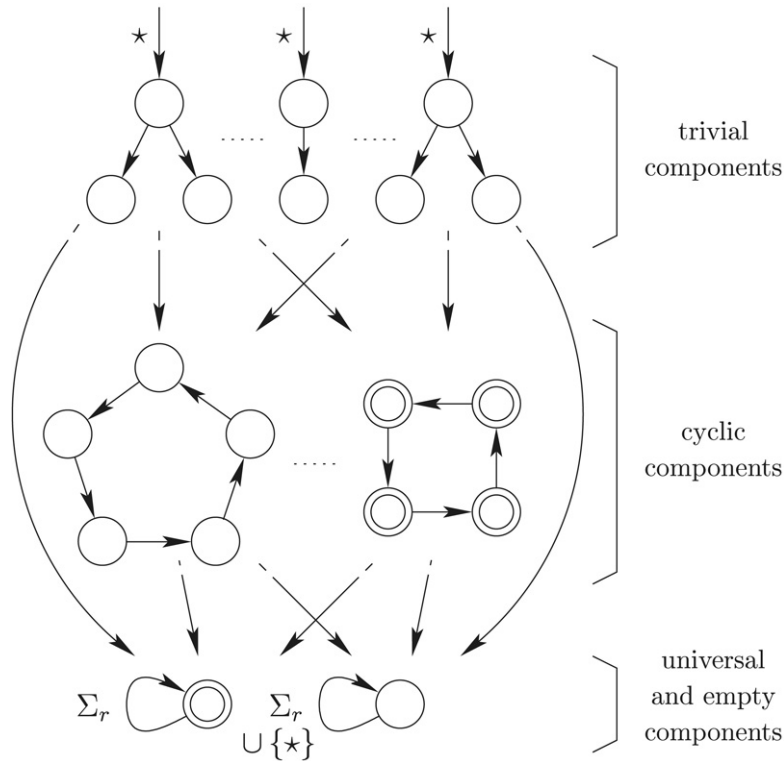


Fig. 3. Fractional part of a RNA representing a finite-boundary set.

In the expression given by Theorem 9, each term of the union encodes a subset of $[0, 1]$ that is definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$: $L' \cdot w_i \cdot \Sigma_r^\omega$ defines an interval $[a, b]$, with $a, b \in \mathbb{Q}$, the terms $L' \cdot w'_i \cdot (v_i)^\omega$, L_0 and L_1 correspond to single rational numbers $c \in \mathbb{Q}$, and the languages $L' \cdot w''_i \cdot L_{<}(v'_i)$ and $L' \cdot w'''_i \cdot L_{>}(v'_i)$ correspond respectively to intervals $[a, b[$ and $]a, b]$, with $a, b \in \mathbb{Q}$. This shows that the set $S \subseteq [0, 1]$ recognized by \mathcal{A} is definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. Combining this result with Theorem 4, as well as the reduction discussed in Section 3, we get our main result:

Theorem 10. *If a set $S \subseteq \mathbb{R}$ is simultaneously recognizable by weak deterministic RNA in two relatively prime bases, then it is definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$.*

Corollary 11. *A set $S \subseteq \mathbb{R}$ is recognizable by weak deterministic RNA in every base $r > 1$ iff it is definable in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$.*

6. Conclusions and future work

The main contribution of this work is to show that the subsets of \mathbb{R} that can be recognized by weak deterministic RNA in all integer bases $r > 1$ are exactly those that are definable in the first-order additive theory of the real and integer numbers $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. Our central result is actually stronger, stating that recognizability in two relatively prime bases r_1 and r_2 is sufficient for forcing definability in $\langle \mathbb{R}, \mathbb{Z}, +, < \rangle$. Using the same argument as in the proof of Lemma 3, this result can directly be extended to bases r_1 and r_2 that do not share the same set of prime factors. This differs slightly from the statement of Cobham's original theorem, which considers instead bases that are multiplicatively independent, i.e., that cannot be expressed as integer powers of the same integer [11,8]. Unfortunately, our approach does not easily generalize to multiplicatively independent bases, since Lemma 3 then becomes invalid. This issue will be addressed in a forthcoming article.

Another contribution is a detailed characterization of the transition graph of weak deterministic RNA that represent subsets of \mathbb{R} defined in first-order additive arithmetic. This characterization could be turned into efficient data structures for handling such RNA. In particular, since their fractional parts recognize a finite union of interval and individual rational values, an efficient representation might be based on symbolic data structures such as BDDs [9] for handling large but finite enumerations. Another possible application is the extraction of formulas from automata-based representations of sets [13,14].

Finally, another goal will be to extend our results to sets in higher dimensions, i.e., to generalize Semenov's theorem [18] to automata over real vectors.

References

- [1] J. Avigad, Y. Yin, Quantifier elimination for the reals with a predicate for the powers of two, Theoretical Computer Science 370 (2007) 48–59.
- [2] B. Boigelot, Symbolic methods for exploring infinite state spaces, Ph.D. Thesis, Université de Liège, 1998.

- [3] B. Boigelot, L. Bronne, S. Rassart, An improved reachability analysis method for strongly linear hybrid systems, in: Proc. 9th CAV, in: Lecture Notes in Computer Science, vol. 1254, Springer, Haifa, June 1997, pp. 167–177.
- [4] B. Boigelot, S. Jodogne, P. Wolper, An effective decision procedure for linear arithmetic over the integers and reals, *ACM Transactions on Computational Logic* 6 (3) (2005) 614–633.
- [5] B. Boigelot, L. Latour, Counting the solutions of Presburger equations without enumerating them, *Theoretical Computer Science* 313 (2004) 17–29.
- [6] B. Boigelot, S. Rassart, P. Wolper, On the expressiveness of real and integer arithmetic automata, in: Proc. 25th ICALP, in: Lecture Notes in Computer Science, vol. 1443, Springer, Aalborg, July 1998, pp. 152–163.
- [7] J. Brusten, Etude des propriétés des RVA, Graduate Thesis, Université de Liège, May 2006.
- [8] V. Bruyère, G. Hansel, C. Michaux, R. Villemaire, Logic and p -recognizable sets of integers, *Bulletin of the Belgian Mathematical Society* 1 (2) (1994) 191–238.
- [9] R.E. Bryant, Symbolic Boolean manipulation with ordered binary decision diagrams, *ACM Computing Surveys* 24 (3) (1992) 293–318.
- [10] J.R. Büchi, On a decision method in restricted second order arithmetic, in: Proc. International Congress on Logic, Methodology and Philosophy of Science, Stanford University Press, Stanford, 1962, pp. 1–12.
- [11] A. Cobham, On the base-dependence of sets of numbers recognizable by finite automata, *Mathematical Systems Theory* 3 (1969) 186–192.
- [12] O. Kupferman, M.Y. Vardi, Complementation constructions for nondeterministic automata on infinite words, in: Proc. 11th TACAS, in: Lecture Notes in Computer Science, vol. 3440, Springer, Edinburgh, April 2005, pp. 206–221.
- [13] L. Latour, Presburger arithmetic: From automata to formulas, Ph.D. Thesis, Université de Liège, 2005.
- [14] J. Leroux, A polynomial time Presburger criterion and synthesis for number decision diagrams, in: Proc. 20th LICS, IEEE Computer Society, Chicago, June 2005, pp. 147–156.
- [15] C. Löding, Efficient minimization of deterministic weak ω -automata, *Information Processing Letters* 79 (3) (2001) 105–109.
- [16] M. Presburger, Über die Vollständigkeit eines gewissen Systems der Arithmetik ganzer Zahlen, in welchem die Addition als einzige Operation hervortritt, in: *Comptes Rendus du Premier Congrès des Mathématiciens des Pays Slaves*, Warsaw, 1929, pp. 92–101.
- [17] S. Safra, On the complexity of ω -automata, in: Proc. 29th Symposium on Foundations of Computer Science, IEEE Computer Society, October 1988, pp. 319–327.
- [18] A.L. Semenov, Presburger-ness of predicates regular in two number systems, *Siberian Mathematical Journal* 18 (1977) 289–299.
- [19] L. van den Dries, The field of reals with a predicate for the powers of two, *Manuscripta Mathematica* 54 (1985) 187–195.
- [20] R. Villemaire, The theory of $\langle \mathbb{N}, +, V_k, V_l \rangle$ is undecidable, *Theoretical Computer Science* 106 (2) (1992) 337–349.
- [21] T. Wilke, Locally threshold testable languages of infinite words, in: Proc. 10th STACS, in: Lecture Notes in Computer Science, vol. 665, Springer, Würzburg, 1993, pp. 607–616.
- [22] P. Wolper, B. Boigelot, An automata-theoretic approach to Presburger arithmetic constraints, in: Proc. 2nd SAS, in: Lecture Notes in Computer Science, vol. 983, Springer, Glasgow, September 1995.
- [23] P. Wolper, B. Boigelot, Verifying systems with infinite but regular state spaces, in: Proc. 10th CAV, in: Lecture Notes in Computer Science, vol. 1427, Springer, Vancouver, June 1998, pp. 88–97.