

EHRENFEUCHT-FRAÏSSÉ GOES AUTOMATIC FOR REAL ADDITION

FELIX KLAEDTKE

ETH Zurich, Department of Computer Science, Switzerland
E-mail address: felixkl@inf.ethz.ch

ABSTRACT. Various logical theories can be decided by automata-theoretic methods. Notable examples are Presburger arithmetic $\text{FO}(\mathbb{Z}, +, <)$ and the linear arithmetic over the reals $\text{FO}(\mathbb{R}, +, <)$, for which effective decision procedures can be built using automata. Despite the practical use of automata to decide logical theories, many research questions are still only partly answered in this area. One of these questions is the complexity of such decision procedures and the related question about the minimal size of the automata of the languages that can be described by formulas in the respective logic. In this paper, we establish a double exponential upper bound on the automata size for $\text{FO}(\mathbb{R}, +, <)$ and an exponential upper bound for the discrete order over the integers $\text{FO}(\mathbb{Z}, <)$. The proofs of these upper bounds are based on Ehrenfeucht-Fraïssé games. The application of this mathematical tool has a similar flavor as in computational complexity theory, where it can often be used to establish tight upper bounds of the decision problem for logical theories.

1. Introduction

Various logical theories admit automata-based decision procedures. The idea of using automata-theoretic methods to decide logical theories goes at least back to Büchi [7]. The elements of the domain of the logical theory are encoded by words over some alphabet in such a way that equality and the relations of the logical theory correspond to regular languages. In order to decide whether a formula is satisfiable, one constructs an automaton that precisely accepts the representatives of the elements that satisfy the formula. This automaton can be constructed by recursion over the formula structure, where standard automata constructions handle the boolean connectives and quantifiers. The satisfiability problem is thus reduced to the emptiness problem for automata.

The logical theories that admit such automata-based decision procedures are often called automatic and they have been systematically studied, e.g., in [4, 13, 14]. Prominent and practically relevant examples are the weak monadic second-order theory of one successor WS1S, Presburger arithmetic $\text{FO}(\mathbb{Z}, +, <)$, and the linear arithmetic over the reals

2000 ACM Subject Classification: F.1.1 [Computation by Abstract Devices]: Models of Computation—automata; F.4.1 [Mathematical Logic and Formal Languages]: Mathematical Logic—computational logic .

Key words and phrases: automata theory, automata-based decision procedures for logical theories, upper bounds, minimal sizes of automata, linear arithmetic over the reals, first-order equivalence, complexity.

This work was supported by the Swiss National Science Foundation (SNF).

$\text{FO}(\mathbb{R}, +, <)$, see, e.g., [5–7]. Tools like MONA [16] and LIRA [3], which have been applied to various verification problems, implement such automata-based decision procedures for logical theories such as WS1S, Presburger arithmetic, and the linear arithmetic over the reals. Furthermore, model checkers for counter systems like FAST [1, 2] use an automata-based representation of sets definable in Presburger arithmetic.

A crude complexity analysis of an automata-based decision procedure leads to a non-elementary worst-case complexity. Namely, for every quantifier alternation there is a potential exponential blow-up in the state space of the automaton. For WS1S, this worst-case scenario actually exists, since the decision problem for WS1S has a non-elementary worst-case complexity [21, 24]. However, for many other automatic logical theories, the non-elementary complexity upper bounds of automata-based decision procedures often contrasts with the known computational complexity upper bounds on the decision problems for the logical theories. Moreover, such exponential blow-ups in the state spaces of the automata are rarely observed in practice in automata-based decision procedures for Presburger arithmetic and the linear arithmetic over the reals. In fact, in many cases, one obtains a smaller automaton after eliminating a quantifier. However, only partial answers exist that explain this phenomenon.

In [15], it is shown that the size of the minimal deterministic automaton that represents a Presburger definable set is triply exponentially bounded with respect to the formula length. This upper bound is established by comparing the automata for Presburger arithmetic formulas with the formulas produced by Reddy and Loveland’s quantifier-elimination method for Presburger arithmetic [23]. The proof on the upper bound in [15] is rather tedious in the sense that several auxiliary upper bounds on the formulas that are generated by the quantifier-elimination method need to be established. These additional upper bounds depend on Reddy and Loveland’s quantifier-elimination method. With the slightly different quantifier-elimination method by Cooper [8], we obtain an upper bound on the automata size that has at least one additional exponent.

For the linear arithmetic over the reals, the approach of using quantifier-elimination methods to establish upper bounds on the automata sizes does not lead to a satisfactory result: an application of this approach establishes only a triple exponential upper bound on the automata size when using the quantifier-elimination method for the linear arithmetic over the reals described in [10]. The author is not aware of any quantifier-elimination method for the linear arithmetic over the reals that would lead to a upper bound on the automata size that is smaller than triple exponential. However, since there are decision procedures for the linear arithmetic over the reals that run in double exponential deterministic time [10], one might conjecture that the automata size is also doubly exponentially bounded.

The main result of this paper proves this conjecture. The presented proof of the double exponential upper bound is based on Ehrenfeucht-Fraïssé games (EF-games, for short from now on). It relates the states of a minimal automaton for a formula and the equivalence classes of a refinement of the equivalence relation determined by EF-games played over $(\mathbb{R}, +, <)$. This proof technique can also be used for other automatic logical theories to establish tight upper bounds on the automata sizes. As another example, we establish an exponential upper bound on the automata size for $\text{FO}(\mathbb{Z}, <)$. Note that the best known deterministic algorithms that decide $\text{FO}(\mathbb{Z}, <)$ run in exponential time [11]. In summary, the results presented in this paper shed light on the complexity of automata-based decision procedures for logical theories by identifying a relationship to EF-games.

It is worth pointing out that EF-games have already been used in similar contexts. Closely related to our work is Ladner’s work [18]. He uses EF-games to show decidability of monadic second-order theories of one successor and first fragments of it. Similar to this paper, he relates the equivalence classes determined by EF-games to automata states. However, Ladner does not focus on the automata sizes and he does not consider $\text{FO}(\mathbb{R}, +, <)$.

The use of EF-games in computational complexity theory [11] and constraint databases [22] is reminiscent of their use in this paper by partitioning the domain and connecting such a partition to the definable sets. Roughly speaking, the use EF-games for establish upper bounds on the decision problem for logical theories is as follows: The key ingredient for obtaining an upper bound for the respective logical theory is to show that the quantifiers, which can range over an infinite domain, can be relativized to a finite subset. Usually, one uses EF-games here to establish upper bounds on the sizes of such sets by analyzing the information that the formulas of a certain quantifier depth can convey. Given such a result on relativizing the quantifiers, satisfiability of a formula can be checked by an exhaustive search. The upper bounds on the sizes of the sets over which the relativized quantifiers range in turn yield upper bounds on the time and space that is needed to perform this search. For several logical theories, this use of EF-games yield tight upper bounds on the computational complexity for their decision problem.

The remainder of the paper is organized as follows. In §2, we give preliminaries. In §3, we illustrate our method by analyzing the languages that are $\text{FO}(\mathbb{Z}, <)$ -definable. In §4, we analyze the languages that are $\text{FO}(\mathbb{R}, +, <)$ -definable and establish the double exponential upper bound on the automata size. Finally, in §5, we draw conclusions. The appendix contains further proof details.

2. Preliminaries

We assume that the reader is familiar with first-order logic and automata theory over finite and infinite words. Here, we recall the needed background in these areas and fix the notation and terminology that we use in the remainder of the text.

2.1. Words and Languages

Let Σ be an alphabet. We denote the set of all finite words over Σ by Σ^* and Σ^+ denotes the set $\Sigma^* \setminus \{\varepsilon\}$, where ε is the empty word. Σ^ω is the set of all ω -words over Σ . The *concatenation* of words is written as juxtaposition. We write $|w|$ for the *length* of $w \in \Sigma^*$. We often write a word $w \in \Sigma^*$ of length $\ell \geq 0$ as $w(0) \dots w(\ell - 1)$ and an ω -word $\alpha \in \Sigma^\omega$ as $\alpha(0)\alpha(1)\alpha(2) \dots$, where $w(i)$ and $\alpha(i)$ denote the i th letter of w and α , respectively.

For a language $L \subseteq \Sigma^*$, the Nerode relation $\sim_L \subseteq \Sigma^* \times \Sigma^*$ is defined as $u \sim_L v$ iff for all $w \in \Sigma^*$, it holds that $uw \in L \Leftrightarrow vw \in L$. Analogously, for an ω -language $L \subseteq \Sigma^\omega$, we define $\sim_L \subseteq \Sigma^* \times \Sigma^*$ as $u \sim_L v$ iff for all $\gamma \in \Sigma^\omega$, it holds that $u\gamma \in L \Leftrightarrow v\gamma \in L$.

2.2. First-order Logic

The (first-order) formulas over a signature are defined as usual: they are built from variables v_0, v_1, \dots , the symbol \approx for equality, the atomic formulas over the signature, the boolean connectives \neg and \vee , and the quantifier \exists . In this paper, we only consider signatures that consist of relation symbols. The signature, its relation symbols, and the arities of its relation symbols are always clear from the context. We write $\varphi(x_1, \dots, x_r)$ when at most

the variables x_1, \dots, x_r occur free in the formula φ . The *quantifier depth* of a formula φ is recursively defined as

$$\text{qd}(\varphi) := \begin{cases} \text{qd}(\psi) & \text{if } \varphi = \neg\psi, \\ \max\{\text{qd}(\psi), \text{qd}(\psi')\} & \text{if } \varphi = \psi \vee \psi', \\ 1 + \text{qd}(\psi) & \text{if } \varphi = \exists x\psi, \text{ and} \\ 0 & \text{otherwise.} \end{cases}$$

A (first-order) structure over a signature consists of a nonempty universe U and it associates with each relation symbol in the signature a relation over U^r , where r is the arity of the relation symbol. We use \mathfrak{R} and \mathfrak{Z} to denote the structures $(\mathbb{R}, +, <)$ and $(\mathbb{Z}, <)$, respectively, where $+$ is the ternary addition relation and $<$ is the ordering relation over the reals or the integers, respectively.

Let \mathfrak{A} be a structure over some signature and with the universe A . For $a_1, \dots, a_r \in A$ and a formula $\varphi(x_1, \dots, x_r)$, we write $\mathfrak{A} \models \varphi[a_1 \dots, a_r]$ if φ is satisfied in \mathfrak{A} when the variable x_i is interpreted as a_i , for all $1 \leq i \leq r$. For the sake of brevity, we often write \bar{x} and \bar{a} instead of x_1, \dots, x_r and a_1, \dots, a_r , respectively.

Let $m, r \in \mathbb{N}$, $\bar{a} \in A^r$ and $\bar{b} \in A^r$. We write $\bar{a} \equiv_m^r \bar{b}$ if for all formulas $\varphi(x_1, \dots, x_r)$ with $\text{qd}(\varphi) \leq m$, it holds that $\mathfrak{A} \models \varphi[\bar{a}] \Leftrightarrow \mathfrak{A} \models \varphi[\bar{b}]$. Note that the relation \equiv_m^r partitions the elements of A^r . The equivalence classes of \equiv_m^r can be game-theoretically characterized by so-called Ehrenfeucht-Fraïssé games. For details on these games, see, for instance, [9]. Instead of working directly with \equiv_m^r , we work with refinements of it, since the reasoning about a well-chosen refinement of \equiv_m^r simplifies matters. In particular, it might be difficult for \equiv_m^r to directly establish an upper bound on the index of \equiv_m^r , to identify elements $\bar{a}, \bar{b} \in A^r$ that are in the same equivalence class, and to find a representative of an equivalence class.

2.3. Representation of Sets Definable in Real Addition

Boigelot, Jodogne, and Wolper have shown in [5] that every first-order definable set $X \subseteq \mathbb{R}^r$ in \mathfrak{R} determines an ω -language L that is in the Borel class $F_\sigma \cap G_\delta$. In other words, L can be accepted by a so-called weak deterministic Büchi automaton. In fact, Boigelot, Jodogne, and Wolper have established in [5] a stronger result. First, they have proved the result for an extension of \mathfrak{R} with the additional predicate \mathbb{Z} . Second, for a formula $\varphi(x_1, \dots, x_r)$ over this extended structure, they have shown how to effectively construct a weak deterministic Büchi automaton that represents the set $\{\bar{a} \in \mathbb{R}^r : \mathfrak{R} \models \varphi[\bar{a}]\}$.

We recall the representation of subsets of \mathbb{R}^r by ω -languages from [5]. In the remainder of the text, let $\varrho > 1$ and $\Sigma := \{0, \dots, \varrho - 1\}$ be fixed. ϱ is called the *base*. Let $r \geq 1$.

- (a) \mathbf{V}_r denotes the set of all ω -words over the alphabet $\Sigma^r \cup \{\star\}$ of the form $v \star \gamma$, where $v \in (\Sigma^r)^+$ and $\gamma \in (\Sigma^r)^\omega$.
- (b) Let $v \star \gamma$ be an ω -word in \mathbf{V}_r with $v(0) = (v_1, \dots, v_r)$. The ω -word $v \star \gamma$ represents the vector of real numbers with r components

$$\langle v \star \gamma \rangle := -\varrho^{|v|-1} \cdot \begin{pmatrix} b_1 \\ \vdots \\ b_r \end{pmatrix} + \sum_{0 < i < |v|} \varrho^{|v|-i-1} \cdot v(i) + \sum_{i \geq 0} \varrho^{-i-1} \cdot \gamma(i),$$

where $b_i := \lceil \frac{v_i}{\varrho} \rceil$, for $1 \leq i \leq r$. Observe that $b_i = 0$ if $v_i = 0$, and $b_i = 1$, otherwise. Here, scalar multiplication is as usual and vector addition is componentwise. Note that we do not distinguish between vectors and tuples.

- (c) For a formula $\varphi(x_1, \dots, x_r)$, we define $L(\varphi) := \{\alpha \in \mathbf{V}_r : \mathfrak{R} \models \varphi[\langle \alpha \rangle]\}$.

Note that the encoding $v \star \gamma \in \mathbf{V}_1$ of a real number is based on the ϱ 's complement representation. The symbol \star plays the role of a decimal point, separating the integer part v from the fractional part γ . Furthermore, the first letter determines whether a “track” represents a number that is greater than or equal to 0, or a number that is less than or equal to 0. Note that the ω -words $0 \star 0^\omega$ and $(\varrho - 1) \star (\varrho - 1)^\omega$ both represent the number 0, where b^ω denotes the infinite repetition of the letter $b \in \Sigma$.

We overload the notation $\langle \cdot \rangle$ by using it also for finite nonempty prefixes in \mathbf{V}_r . For $v \in (\Sigma^r)^+$ and $v' \in (\Sigma^r)^*$, we write $\langle v \rangle$ and $\langle v \star v' \rangle$ for $\langle v \star \bar{0}^\omega \rangle$ and $\langle v \star v' \bar{0}^\omega \rangle$, respectively, where $\bar{0}$ denotes the vector $(0, \dots, 0) \in \Sigma^r$.

3. Automata Upper Bound for the Ordering over the Integers

Before looking at the ω -languages that can be described by the first-order logic over \mathfrak{R} , we look at a simpler case. Namely, we investigate the languages that can be described by formulas over \mathfrak{Z} . We establish an exponential upper bound on the automata size for these languages. The purpose of investigating this simpler case first is twofold. First, it introduces the main concepts, which we also use in §4 for the ω -languages definable in the first-order logic over \mathfrak{R} . Second, it demonstrates the generality of the approach. The results in this section illustrate the relationship between the equivalence classes of a refinement of the equivalence relation \equiv_m^r and the equivalence classes of the Nerode relation of a language described by a formula $\varphi(x_1, \dots, x_r)$ over \mathfrak{Z} with $\text{qd}(\varphi) \leq m$.

Throughout this section, formulas are over \mathfrak{Z} 's signature, and m and r range over the natural numbers. We start with some definitions. For a formula $\varphi(x_1, \dots, x_r)$, we define the language

$$K(\varphi) := \{v \in (\Sigma^r)^+ : \mathfrak{Z} \models \varphi[\langle v \rangle]\}.$$

We partition \mathbb{Z}^r by the equivalence relation E_m^r that is defined as

$$\bar{a} E_m^r \bar{b} \quad \text{iff} \quad \begin{array}{l} \text{sign}(a_i - a_j - c) = \text{sign}(b_i - b_j - c), \\ \text{for all } c, i, j \in \mathbb{N} \text{ with } c \leq m \text{ and } 1 \leq i, j \leq r, \end{array}$$

where $\bar{a}, \bar{b} \in \mathbb{Z}^r$, and $\text{sign}(x) := 0$ if $x < 0$ and $\text{sign}(x) := 1$, otherwise, for $x \in \mathbb{R}$. Intuitively speaking, $\bar{a}, \bar{b} \in \mathbb{Z}^r$ are in the same equivalence class of E_m^r if the distances between their components are equal up to the threshold m .

Before we launch into the proof of establishing an upper bound on the size of the minimal deterministic automaton for a formula $\varphi(x_1, \dots, x_r)$, we give an outline: (i) We show that $E_{2\text{qd}(\varphi)}^r$ refines $\equiv_{\text{qd}(\varphi)}^r$. (ii) We establish an upper bound on the index of $E_{2\text{qd}(\varphi)}^r$. (iii) We show that $E_{2\text{qd}(\varphi)}^r$ has a congruence property with respect to word concatenation. (iv) By using (i) and (iii), we show that $E_{2\text{qd}(\varphi)}^r$ determines an equivalence relation on $(\Sigma^r)^+$ that refines the Nerode relation $\sim_{K(\varphi)}$. Finally, from (ii) we derive an upper bound on the index of $\sim_{K(\varphi)}$. Note that the equivalence classes of $\sim_{K(\varphi)}$ can be viewed as the states of the minimal deterministic finite automaton that accepts $K(\varphi)$. The properties (i) to (iv) correspond to the Lemmas 3.1 to 3.4, respectively, which are given below.

Lemma 3.1. *The equivalence relation E_{2m}^r refines the equivalence relation \equiv_m^r . That means, $\bar{a} E_{2m}^r \bar{b}$ implies $\bar{a} \equiv_m^r \bar{b}$, for all $\bar{a}, \bar{b} \in \mathbb{Z}^r$.*

To prove Lemma 3.1, we apply a standard technique from model theory. First, we show that the family $(E_n^s)_{s,n \in \mathbb{N}}$ of equivalence relations has the following property:

$$\text{If } \bar{a} E_{2m+1}^r \bar{b} \text{ then for every } a' \in \mathbb{Z}, \text{ there is some } b' \in \mathbb{Z} \text{ such that } (\bar{a}, a') E_{2m}^{r+1} (\bar{b}, b'). \quad (1)$$

Properties of this kind are often called back-and-forth properties in the literature. Note that E_{2m+1}^r is symmetric. Second, we complete the proof by an induction over m , where we use the property (1) in the induction step. The proof is given in the appendix.

Lemma 3.2. *The index of E_m^r is at most $r! \cdot (m+1)^r$.*

Proof. There are at most $r!$ many possibilities to order the r elements increasingly. If in such an ordering the distance between the i th element x and the $(i+1)$ st element y is greater than or equal to m , we have that $\text{sign}(y-x-c) = 1$, for all $c \in \mathbb{N}$ with $c \leq m$. We obtain that the index is at most $r! \cdot (m+1)^r$. ■

Lemma 3.3. *Let $u, v \in (\Sigma^r)^+$. If $\langle u \rangle E_m^r \langle v \rangle$ then $\langle uw \rangle E_m^r \langle vw \rangle$, for all $w \in (\Sigma^r)^*$.*

Proof. Let $n := |w|$, $\bar{a} := (a_1, \dots, a_r) := \langle u \rangle$, $\bar{b} := (b_1, \dots, b_r) := \langle v \rangle$, and $\bar{d} := (d_1, \dots, d_r) := \langle \bar{0}w \rangle$. We have that $\langle uw \rangle = \varrho^n \bar{a} + \bar{d}$ and $\langle vw \rangle = \varrho^n \bar{b} + \bar{d}$. Furthermore, it holds that $d_i < \varrho^n$, for all $i \in \{1, \dots, r\}$. Let $i, j, c \in \mathbb{N}$ with $1 \leq i, j \leq r$ and $c \leq m$. We have to show that

$$\text{sign}(\varrho^n(a_i - a_j) + d_i - d_j - c) = \text{sign}(\varrho^n(b_i - b_j) + d_i - d_j - c). \quad (2)$$

Case $a_i - a_j = 0$. We have that $\text{sign}(a_i - a_j) = 1 = \text{sign}(a_j - a_i)$. From the assumption $\bar{a} E_m^r \bar{b}$, it follows that $\text{sign}(a_i - a_j) = \text{sign}(b_i - b_j)$ and $\text{sign}(a_j - a_i) = \text{sign}(b_j - b_i)$, and hence, $b_i - b_j = 0$. Obviously, the equality (2) holds.

Case $b_i - b_j = 0$. This case is symmetric to the case $a_i - a_j = 0$ above.

Case $a_i - a_j \neq 0$ and $b_i - b_j \neq 0$. For showing (2), it suffices to show the equality

$$\text{sign}\left(a_i - a_j + \frac{d_i - d_j - c}{\varrho^n}\right) = \text{sign}\left(b_i - b_j + \frac{d_i - d_j - c}{\varrho^n}\right). \quad (3)$$

- If $m = 0$, we have that $c = 0$ and thus $\left|\frac{d_i - d_j - c}{\varrho^n}\right| \leq \frac{|d_i - d_j|}{\varrho^n} \leq \frac{\varrho^n - 1}{\varrho^n} < 1$. Since $a_i - a_j \neq 0$ and $b_i - b_j \neq 0$ and by the assumption $\bar{a} E_0^r \bar{b}$, we conclude that the equality (3) holds.
- If $m > 0$, we have that $\left|\frac{d_i - d_j - c}{\varrho^n}\right| \leq \frac{|d_i - d_j| + |c|}{\varrho^n} \leq \frac{\varrho^n - 1 + |c|}{\varrho^n} \leq \frac{m(\varrho^n - 1) + m}{\varrho^n} = m$. The equality (3) follows from the assumption $\bar{a} E_m^r \bar{b}$. ■

Lemma 3.4. *Let φ be a formula with at most r free variables and with quantifier depth at most m . If $\langle u \rangle E_{2m}^r \langle v \rangle$ then $u \sim_{K(\varphi)} v$, for all $u, v \in (\Sigma^r)^+$.*

Proof. We prove the lemma by contraposition. Assume that $u \not\sim_{K(\varphi)} v$, i.e., there is a word $w \in \Sigma^*$ such that $uw \in K(\varphi) \not\equiv vw \in K(\varphi)$. It follows that $\langle uw \rangle \not\equiv_m^r \langle vw \rangle$. By Lemma 3.1, we conclude that $\langle uw \rangle E_{2m}^r \langle vw \rangle$ does not hold. By Lemma 3.3, we obtain that $\langle u \rangle E_{2m}^r \langle v \rangle$ does not hold. ■

Theorem 3.5. *Let φ be a formula. The index of $\sim_{K(\varphi)}$ is at most $1 + 2^{n^2}$, where n is the length of the formula φ , i.e., φ consists of n symbols.*

Proof. Let r be the number of free variables of φ and $m := \text{qd}(\varphi)$. Note that $n \geq r + m + 1$. Without loss of generality, we assume that $r > 0$. By Lemma 3.2, we have that the index of E_{2m}^r is at most $r! \cdot (2^m + 1)^r \leq 2^{r^2 + rm + r} \leq 2^{rn} \leq 2^{n^2}$. From Lemma 3.4, it follows that $\sim_{K(\varphi)}$ partitions $(\Sigma^r)^+$ in at most 2^{n^2} equivalence classes. Note that the empty word can be in an equivalence class that is distinct from all the others. ■

4. Automata Upper Bound for Real Addition

In this section, we establish an upper bound on the automata size for the first-order logic over \mathfrak{R} . The structure of the proof is similar to the proof in the previous section §3. However, it is more involved. In §4.1, we define a family $(F_n^s)_{s,n \in \mathbb{N}}$ of equivalence relations. In §4.1 and §4.2, we show that $(F_n^s)_{s,n \in \mathbb{N}}$ has similar properties as the family $(E_n^s)_{s,n \in \mathbb{N}}$ defined in §3. Namely, (1) we show that each $F_{2^{m+2}}^r$ refines \equiv_m^r and (2) we establish a relationship between the equivalence classes of the congruence relations determined by the definable ω -languages and equivalence classes of refinements of the equivalence relations $(F_n^s)_{s,n \in \mathbb{N}}$. Finally, in §4.3, we derive the double exponential upper bound on the size of a minimal Büchi automaton that accepts the ω -language of a formula of the first-order logic over \mathfrak{R} .

In the following, formulas are always over \mathfrak{R} 's signature, and r and m range over the natural numbers.

4.1. Partitioning the Reals by First-order Formulas

The results, which we use later, and their presentation of this subsection are based on Chapter 22 of Kozen's book [17]. Since subtle modifications are made, we provide proofs in the full version of the paper. At the end of this subsection, we comment on these modifications and their implications.

An *integer affine function* of arity r is a function $f : \mathbb{R}^r \rightarrow \mathbb{R}$ defined by a linear polynomial with integer coefficients, i.e., there are $c_0, \dots, c_r \in \mathbb{Z}$ such that for all $x_1, \dots, x_r \in \mathbb{R}$, it holds that $f(x_1, \dots, x_r) = c_0 + \sum_{1 \leq i \leq r} c_i x_i$. For such a function, f^* denotes the function with $f^*(x_1, \dots, x_r) = \sum_{1 \leq i \leq r} c_i x_i$, for all $x_1, \dots, x_r \in \mathbb{R}^r$. We define $\|f\| := \max\{0, |c_1|, \dots, |c_r|\}$. Let A^r be the set of all integer affine functions of arity r and

$$B_m^r := \{f \in A^r : \|f\| \leq m \text{ and } |f(\bar{0})| \leq rm\}.$$

Definition 4.1. We partition \mathbb{R}^r by the equivalence relation F_m^r that is defined as

$$\bar{a} F_m^r \bar{b} \quad \text{iff} \quad \text{for all } f \in B_m^r, \text{ sign}(f(\bar{a})) = \text{sign}(f(\bar{b})),$$

where $\bar{a}, \bar{b} \in \mathbb{R}^r$.

Note that F_m^r decomposes \mathbb{R}^r into cells. Each such cell is described by a conjunction of linear inequations, where the absolute values of the coefficients of the inequations are bounded. Moreover, we remark that the technique that we present in the following by connecting such partitions to first-order logic and Ehrenfeucht-Fraïssé games is reminiscent of techniques in computational complexity (see [11]) and constraint databases (see [22]). A novel insight is that these partitions are also connected to the relation $\sim_{L(\varphi)}$ for a formula φ . We start with some properties about the family $(F_n^s)_{s,n \in \mathbb{N}}$ of equivalence relations. The proofs are given in the appendix.

Lemma 4.2. *Let $\bar{a}, \bar{b} \in \mathbb{R}^r$. If $\bar{a} F_{4m^2}^r \bar{b}$ then for all $a' \in \mathbb{R}$, there is some $b' \in \mathbb{R}$ such that $(\bar{a}, a') F_m^{r+1} (\bar{b}, b')$.*

Similar to Lemma 3.1, we obtain the following lemma by using Lemma 4.2.

Lemma 4.3. *For all $\bar{a}, \bar{b} \in \mathbb{R}^r$, it holds that if $\bar{a} F_{2^{m+2}}^r \bar{b}$ then $\bar{a} \equiv_m^r \bar{b}$.*

The following two lemmas show how to obtain a set $R \subseteq \mathbb{R}^r$ such that each equivalence class of F_m^r has at least one representative in R . Let σ be an equivalence class of F_m^r and

let σ' be an equivalence class of F_n^{r+1} , where $n \in \mathbb{N}$. We say that σ' is *consistent* with σ if $(\sigma \times \mathbb{R}) \cap \sigma' \neq \emptyset$.

Lemma 4.4. *For each equivalence class σ of F_m^1 , we have that*

$$\sigma \cap \left\{ \frac{d}{c} : c, d \in \mathbb{Z} \text{ with } c \neq 0 \text{ and } |c|, |d| \leq 2m^2 \right\} \neq \emptyset.$$

Lemma 4.5. *Let $r > 1$, $\bar{a} \in \mathbb{R}^r$, where σ is the equivalence class of \bar{a} with respect to $F_{2m^2}^r$. For every equivalence class σ' of F_m^{r+1} that is consistent with σ , we have that*

$$\sigma' \cap \left\{ \left(\bar{a}, \frac{f(\bar{a})+d}{c} \right) : f \in B_2^r \text{ and } c, d \in \mathbb{Z} \setminus \{0\} \text{ with } |c| \leq 2, \text{ and } |d| < 2 \right\} \neq \emptyset$$

if $m = 1$, and, for $m \neq 1$, we have that

$$\sigma' \cap \left\{ \left(\bar{a}, \frac{f(\bar{a})}{c} \right) : f \in B_{2m^2}^r \text{ and } c \in \mathbb{Z} \setminus \{0\} \text{ with } |c| \leq 2m^2 \right\} \neq \emptyset.$$

Remark 4.6. Before we proceed to establish the upper bound on the size of the minimal automata representation for the set defined by a formula φ , we point out the differences between the family $(F_n^s)_{s,n \in \mathbb{N}}$ of equivalence relations and the family of equivalence relations defined in Kozen's book [17] in Chapter 22.

Kozen is not using the function $\text{sign} : \mathbb{R} \rightarrow \{0, 1\}$ but the standard signum function $\text{sgn} : \mathbb{R} \rightarrow \{-1, 0, 1\}$ that is defined as $\text{sgn}(x) := -1$ if $x < 0$, $\text{sgn}(x) := 1$ if $x > 0$, and $\text{sgn}(0) := 0$. In Kozen's book, two elements $\bar{a}, \bar{b} \in \mathbb{R}^r$ are related iff $\text{sgn}(f(\bar{a})) = \text{sgn}(f(\bar{b}))$, for all integer affine function $f \in A^r$ with $\|f\| \leq m$ and $|f(\bar{0})| \leq m$.

There are two differences to our definition. First, we use the function sign instead of the function sgn . This difference is actually irrelevant. Using sign instead of sgn in the definition in Kozen's book would not change the equivalence relations. However, we found the reasoning in the proofs when using the function sign slightly simpler. Second and more relevant, we require $|f(\bar{0})| \leq rm$ instead of $|f(\bar{0})| \leq m$. The proofs of the Lemmas 4.2 to 4.5 follow the lines of the proofs of the corresponding lemmas in Kozen's book. However, there are subtle differences, e.g., in Lemma 4.5, we have the special case for $m = 1$, which is not needed in the corresponding lemma in Kozen's book.

An immediate consequence of only requiring this weaker restriction on the functions $f \in A^r$ is that the equivalence relation F_m^r refines the corresponding equivalence relation as defined in Kozen's book. The purpose for having finer equivalence relations is the following: For a formula $\varphi(x_1, \dots, x_r)$, we show in §4.2 that the equivalence classes of $\sim_{L(\varphi)}$ are related to the equivalence classes of a certain relation in the family $(F_n^s)_{s,n \in \mathbb{N}}$. Without the weaker requirement we were not able to establish a similar relationship. The problem can be pinpointed to Lemma 4.8, which is crucial in relating the equivalence relations. The corresponding statement of Lemma 4.8 would not be correct when using the equivalence relations as defined in Kozen's book.

4.2. Relationship to Languages

In this subsection, we establish a relationship between the equivalence relation $F_{2^{2m+2+1}}^r$ and the congruence relation $\sim_{L(\varphi)}$, where $\varphi(x_1, \dots, x_r)$ is a formula with $\text{qd}(\varphi) \leq m$. Namely, we show that $F_{2^{2m+2+1}}^r$ determines a refinement of the congruence relation $\sim_{L(\varphi)}$.

We start with a technical lemma. Its proof is straightforward and we therefore omit it. In the following, we will use it without explicitly referring to it.

Lemma 4.7. *For $f \in A^r$, $u \in (\Sigma^r)^+$, $u' \in (\Sigma^r)^*$, and $\gamma \in (\Sigma^r)^\omega$, the following facts hold:*

- (1) $f(\langle uu' \rangle) = f(\bar{0}) + \varrho^{|u'|} f^*(\langle u \rangle) + f^*(\langle \bar{0}u' \rangle)$, and
(2) $f(\langle u \star u' \gamma \rangle) = f(\bar{0}) + f^*(\langle u \star u' \rangle) + \varrho^{-|u'|} f^*(\langle \bar{0} \star \gamma \rangle)$.

The next two lemmas show that the equivalence relations in the family $(F_n^s)_{s,n \in \mathbb{N}}$ have congruence properties on words with respect to concatenation and show how their equivalence classes relate to the equivalence classes of the congruence relation $\sim_{L(\varphi)}$. We want to point out a subtle technical detail, which is reflected in the (b)-parts of the lemmas. The words $u \star u'$ and $u \star u' \bar{0}$ represent the same vector of real numbers, i.e., $\langle u \star u' \rangle = \langle u \star u' \bar{0} \rangle$. Therefore, $u \star u'$ and $u \star u' \bar{0}$ represent the same equivalence class in F_m^r . However, $u \star u'$ and $u \star u' \bar{0}$ might not be in the same equivalence class with respect to $\sim_{L(\varphi)}$. Observe that appending an ω -word $\gamma \in (\Sigma^r)^\omega$ to $u \star u'$ and $u \star u' \bar{0}$ may yield representations of different vectors of real numbers, i.e., $\langle u \star u' \gamma \rangle \neq \langle u \star u' \bar{0} \gamma \rangle$. In particular, $u \star u' \gamma$ and $u \star u' \bar{0} \gamma$ may represent different equivalence classes in F_m^r .

Lemma 4.8. *For all $u, v \in (\Sigma^r)^+$ and $u', v' \in (\Sigma^r)^*$, the following two facts hold:*

- (a) *If $\langle u \rangle F_m^r \langle v \rangle$ then for all $w \in (\Sigma^r)^*$, $\langle uw \rangle F_m^r \langle vw \rangle$.*
(b) *If $\langle u \star u' \rangle F_{2m}^r \langle v \star v' \rangle$ and $|u'| \geq |v'|$ then for all $\gamma \in (\Sigma^r)^\omega$, $\langle u \star u' \gamma \rangle F_m^r \langle v \star v' \bar{0}^k \gamma \rangle$ with $k = \min\{|u'| - |v'|\} \cup \{k \in \mathbb{Z} : \varrho^k \geq rm\}$.*

Proof. For $r = 0$, there is nothing to prove. In the following, we assume that $r > 0$.

(a) We prove (a) by contraposition. Assume that for some $w \in (\Sigma^r)^*$, it is not the case that $\langle uw \rangle F_m^r \langle vw \rangle$, i.e., there is some $f \in B_m^r$ with $\text{sign}(f(\langle uw \rangle)) \neq \text{sign}(f(\langle vw \rangle))$. Without loss of generality, we assume that $f(\langle uw \rangle) < 0$ and hence $f(\langle vw \rangle) \geq 0$. The other cases can be reduced to this case by using the function $g \in B_m^r$ with $g(\bar{x}) = -f(\bar{x})$, for all $\bar{x} \in \mathbb{R}^r$.

We have that $\varrho^{|w|} f^*(\langle u \rangle) + f(\langle \bar{0}w \rangle) < 0$ and $\varrho^{|w|} f^*(\langle v \rangle) + f(\langle \bar{0}w \rangle) \geq 0$. Obviously, it must hold that $f^*(\langle u \rangle) \neq f^*(\langle v \rangle)$. If $\text{sign}(f^*(\langle u \rangle)) \neq \text{sign}(f^*(\langle v \rangle))$ then $\langle u \rangle F_m^r \langle v \rangle$ does not hold and we are done. So, assume that $\text{sign}(f^*(\langle u \rangle)) = \text{sign}(f^*(\langle v \rangle))$. If $|f^*(\langle u \rangle)| \leq rm$ or $|f^*(\langle v \rangle)| \leq rm$ then we are also done by choosing an appropriate function $g \in B_m^r$ with $\text{sign}(g(\langle u \rangle)) \neq \text{sign}(g(\langle v \rangle))$. So, assume that $|f^*(\langle u \rangle)|, |f^*(\langle v \rangle)| > rm$. Note that $|f^*(\langle \bar{0}w \rangle)| \leq (\varrho^{|w|} - 1)rm$.

– If $f^*(\langle v \rangle) < -rm$, we obtain a contradiction to the assumption $f(\langle vw \rangle) \geq 0$, since

$$\begin{aligned} \varrho^{|w|} f^*(\langle v \rangle) + f(\langle \bar{0}w \rangle) &= \varrho^{|w|} f^*(\langle v \rangle) + f^*(\langle \bar{0}w \rangle) + f(\bar{0}) \\ &< -\varrho^{|w|} rm + (\varrho^{|w|} - 1)rm + rm \leq 0. \end{aligned}$$

– If $f^*(\langle v \rangle) > rm$, we conclude that $f^*(\langle u \rangle) > rm$. Analogously, as in the above case, we obtain a contradiction to the assumption $f(\langle uw \rangle) < 0$.

(b) Let f be an arbitrary function in B_m^r and $\gamma \in (\Sigma^r)^\omega$. We have to show that $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle))$. Since $B_m^r \subseteq B_{2m}^r$, it follows from the assumption $\langle u \star u' \rangle F_{2m}^r \langle v \star v' \rangle$ that $\text{sign}(f(\langle u \star u' \rangle)) = \text{sign}(f(\langle v \star v' \rangle))$. That means, either (1) $f(\langle u \star u' \rangle), f(\langle v \star v' \rangle) < 0$ or (2) $f(\langle u \star u' \rangle), f(\langle v \star v' \rangle) \geq 0$ holds. Since the case (1) can be reduced to the case (2) by considering the function $g(\bar{x}) = -f(\bar{x})$, for all $\bar{x} \in \mathbb{R}^r$, we restrict ourselves to (2).

For the sake of readability, we use the abbreviations $a := f^*(\langle u \star u' \rangle)$, $b := f^*(\langle v \star v' \rangle)$, and $c := f^*(\langle \bar{0} \star \gamma \rangle)$. Note that

$$f(\langle u \star u' \gamma \rangle) = f(\bar{0}) + a + c\varrho^{-|u'|} \quad \text{and} \quad f(\langle v \star v' \bar{0}^k \gamma \rangle) = f(\bar{0}) + b + c\varrho^{-|v'|-k}. \quad (4)$$

If $c \geq 0$ then $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle)) = 1$. In the following, assume $c < 0$.

Case $a \neq b$. With the assumption $\langle u \star u' \rangle F_{2m}^r \langle v \star v' \rangle$ we conclude that $a, b > 2rm$. Note that $|f^*(\langle \bar{0} \star \alpha \rangle)| \leq rm$, for all $\alpha \in (\Sigma^r)^\omega$. It follows that

$$f^*(\langle u \star u' \gamma \rangle) = a + c\rho^{-|u'|} > 2rm - rm \geq rm.$$

The reasoning for $f^*(\langle v \star v' \bar{0}^k \gamma \rangle) > rm$ is similar. Since $|f(\bar{0})| \leq rm$, we have that $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle)) = 1$.

Case $a = b$. For $k = |u'| - |v'|$, it immediately follows from the equalities in (4) that $f(\langle u \star u' \gamma \rangle) = f(\langle v \star v' \bar{0}^k \gamma \rangle)$, and hence $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle))$. For $a = b = -f(\bar{0})$, it is also straightforward to see from the two equalities in (4) that $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle))$. For the rest of the proof, assume $k = \min\{k \in \mathbb{Z} : \rho^k \geq rm\}$ and $b \neq -f(\bar{0})$. Moreover, for $|c| \cdot \rho^{-|u'|} > f(\bar{0}) + a$, it follows directly from the equalities (4) that $\text{sign}(f(\langle u \star u' \gamma \rangle)) = \text{sign}(f(\langle v \star v' \bar{0}^k \gamma \rangle)) = 0$. So, we also assume that $|c| \cdot \rho^{-|u'|} \leq f(\bar{0}) + a$. Observe that $f(\langle u \star u' \gamma \rangle) \geq 0$. Furthermore, observe that $f(\bar{0}) + b \geq \rho^{-|v'|}$. We have that $f(\langle v \star v' \bar{0}^k \gamma \rangle) \geq \frac{1}{\rho^{|v'|}} + \frac{c}{\rho^{|v'|+k}} = \frac{\rho^k - |c|}{\rho^{|v'|+k}} \geq \frac{rm - rm}{\rho^{|v'|+k}} \geq 0$. ■

Lemma 4.9. *Let $\varphi(x_1, \dots, x_r)$ be a formula with $\text{qd}(\varphi) \leq m$. For all $u, v \in (\Sigma^r)^+$ and $u', v' \in (\Sigma^r)^*$, the following two facts hold:*

- (a) *If $\langle u \rangle F_{2^{2m+2+1}}^r \langle v \rangle$ then $u \sim_{L(\varphi)} v$.*
- (b) *If $\langle u \star u' \rangle F_{2^{2m+2+1}}^r \langle v \star v' \rangle$ and $|u'| \geq |v'|$ then $u \star u' \sim_{L(\varphi)} v \star v' \bar{0}^k$ with $k = \min\{|u'| - |v'| \} \cup \{k \in \mathbb{Z} : \rho^k \geq rm\}$.*

Proof. We only show (a). The proof for (b) is analogous and we omit it. From Lemma 4.8(a), it follows that $\langle uw \rangle F_{2^{2m+2+1}}^r \langle vw \rangle$, for all $w \in (\Sigma^r)^*$. With Lemma 4.8(b), we obtain that $\langle uw \star \gamma \rangle F_{2^{2m+2}}^r \langle vw \star \gamma \rangle$, for all $w \in (\Sigma^r)^*$ and $\gamma \in (\Sigma^r)^\omega$. By Lemma 4.3, we conclude that $\langle uw \star \gamma \rangle \equiv_m^r \langle vw \star \gamma \rangle$, for all $w \in (\Sigma^r)^*$ and $\gamma \in (\Sigma^r)^\omega$. In particular, we have that $uw \star \gamma \in L(\varphi) \Leftrightarrow vw \star \gamma \in L(\varphi)$, for all $w \in (\Sigma^r)^*$ and $\gamma \in (\Sigma^r)^\omega$. From this it follows that $u \sim_{L(\varphi)} v$, since for any ω -word α not in V_r , we have that $u\alpha, v\alpha \notin L(\varphi)$. ■

4.3. Upper Bounds

We establish an upper bound on the index of F_m^r , from which we then derive an upper bound on the automata size. We start with a simple lemma.

Lemma 4.10. *The cardinality of B_m^r is at most $(2rm + 1)(2m + 1)^r$.*

Using the Lemmas 4.4, 4.5, and 4.10, we establish an upper bound on the index of F_m^r and thus also on the index of \equiv_m^r .

Lemma 4.11. *The index of F_m^r is at most $\max\{1, m^{2^{3+r}} \cdot 2^{2^{3+r}}\}$.*

Proof. For $r = 0$ or $m = 0$, the index of F_m^r is 1. In the following, we assume that $r > 0$ and $m > 0$. To establish an upper bound on the index of F_m^r , we consider the sequences $(e_i)_{i \in \mathbb{N}}$ and $(c_i)_{i \in \mathbb{N}}$, which are recursively defined by $e_0 := c_0 := 0$, $e_i := 2e_{i-1} + 2 + 2i$, and $c_i := c_{i-1} + e_{i-1} + 2 + i \log_2(8i)$, for $i > 0$. In the appendix, we show that $e_i, c_i \leq 2^{3+i}$, for all $i \in \mathbb{N}$.

To show the claimed upper bound on the index of F_m^r , it suffices to show that the index of F_m^r is at most $2^{e_r} \cdot m^{e_r}$.

- Assume $r = 1$. From Lemma 4.4, it follows that any equivalence class of F_m^1 has a representative in the set $\{\frac{d}{c} : c, d \in \mathbb{Z} \text{ with } c \neq 0 \text{ and } |c|, |d| \leq 2m^2\}$. The cardinality of that set is at most $(5m^2)(4m^2) = 20m^4 \leq 2^{c_1} \cdot m^{e_1}$.
- Assume $r > 1$ and $F_{2m^2}^{r-1}$ has at most $2^{c_{r-1}} \cdot (2m^2)^{e_{r-1}} \leq 2^{c_{r-1}+e_{r-1}} \cdot m^{2e_{r-1}}$ equivalence classes. Let σ be an equivalence class of $F_{2m^2}^{r-1}$ and $\bar{a} \in \sigma$.

For $m = 1$, it follows from Lemma 4.5 that all equivalence classes of F_m^r that are consistent with σ have a representative in the set $\{(\bar{a}, \frac{f(\bar{a})+d}{c}) : f \in B_2^{r-1}, c, d \in \mathbb{Z} \setminus \{0\} \text{ with } |c| \leq 2 \text{ and } |d| < 2\}$. With Lemma 4.10, we obtain the upper bound

$$(4(r-1)+1) \cdot 5^{r-1} \cdot 4 \cdot 2 = (32r-24) \cdot 5^{r-1} \leq 32r \cdot 8^{r-1} \leq 4 \cdot 8^r \cdot r^r = 2^{2+r \log_2(8r)}$$

on the cardinality of that set. Similarly, for $m > 1$, it follows from Lemma 4.5 that all equivalence classes of F_m^r that are consistent with σ have a representative in the set $\{(\bar{a}, \frac{f(\bar{a})}{c}) : f \in B_{2m^2}^{r-1} \text{ and } c \in \mathbb{Z} \setminus \{0\} \text{ with } |c| \leq 2m^2\}$. An upper bound on the cardinality of that set is

$$(4(r-1)m^2+1)(4m^2+1)^{r-1}4m^2 \leq (5rm^2)^r 4m^2 \leq 2^{2+3r+r \log_2 r} m^{2r+2} \leq 2^{2+r \log_2(8r)} m^{2r+2}.$$

For $m = 1$ or $m > 1$, an upper bound on the number of equivalence classes of F_m^r is

$$(2^{c_{r-1}+e_{r-1}} m^{2e_{r-1}}) \cdot (2^{2+r \log_2(8r)} m^{2r+2}) = 2^{c_r} \cdot m^{e_r}. \quad \blacksquare$$

Theorem 4.12. *Let φ be a formula. The index of $\sim_{L(\varphi)}$ is at most $2^{2^{8+n}}$, where n is the length of the formula φ , i.e., the number of symbols of φ .*

Proof. Let r be the number of free variables in φ and $m := \text{qd}(\varphi)$. We use $F_{2^{2m+2+1}}^r$ to define a refinement R of $\sim_{L(\varphi)}$. First, the singleton $\{\varepsilon\}$ is an equivalence class of R . Second, the set of words with at least two occurrences of the letter \star is another equivalence class of R . The equivalence class of a word $v \in (\Sigma^r)^+$ of R is $\{u \in (\Sigma^r)^+ : \langle v \rangle F_{2^{2m+2+1}}^r \langle u \rangle\}$.

It remains to define the equivalence classes of R on $F := \{v \star v' : v \in (\Sigma^r)^+ \text{ and } v' \in (\Sigma^r)^*\}$. For $v \star v' \in F$, let $S := \{u \star u' \in F : \langle v \star v' \rangle F_{2^{2m+2+1}}^r \langle u \star u' \rangle\}$. R chops S into equivalence classes, assuming $|v'| \leq |u'|$, for all $u \star u' \in S$:

- For $k \in \{0, \dots, \lceil \log_\rho r 2^{2^{m+2}+1} \rceil - 1\}$, the equivalence class of $v \star v' \bar{0}^k$ of R is $\{u \star u' \in S : |u'| = |v'| + k\}$.
- For $k = \lceil \log_\rho r 2^{2^{m+2}+1} \rceil$, the equivalence class of $v \star v' \bar{0}^k$ of R is $\{u \star u' \in S : |u'| \geq |v'| + k\}$.

Note that any word $u \star u' \in S$ relates to exactly one word $v \star v' \bar{0}^k$.

With Lemma 4.9 at hand, it is easy to see that R refines $\sim_{L(\varphi)}$. It remains to prove an upper bound on the index of R . Note that $n \geq m + r \geq 1$. By Lemma 4.11, an upper bound on the index of $F_{2^{2m+2+1}}^r$ is

$$(2^{2^{m+2}+1})^{2^{3+r}} \cdot 2^{2^{3+r}} = 2^{2^{m+2} \cdot 2^{3+r} + 2^{3+r} + 2^{3+r}} = 2^{2^{5+r+m} + 2^{4+r}} \leq 2^{2^{6+n}}.$$

Hence, R partitions $(\Sigma^r)^+$ into at most $2^{2^{6+n}}$ equivalence classes and F is partitioned into at most $2^{2^{6+n}} \cdot \lceil \log_\rho r 2^{2^{m+2}+1} \rceil \geq 2^{2^{6+n}} \cdot r(2^{m+2} + 1) \geq 2^{2^{6+n}} \cdot 2^{3+n} \geq 2^{2^{7+n}}$ equivalence classes. From this, we derive the upper bound $2^{2^{8+n}}$ on R 's index. \blacksquare

Remark 4.13. Since for any formula φ , $L(\varphi)$ is an ω -language in the Borel class $F_\sigma \cap G_\delta$ [5], we can—similar to deterministic finite automata—view the equivalence classes of $\sim_{L(\varphi)}$ as the states of a minimal deterministic Büchi automaton that accepts $L(\varphi)$. For further details, see [20] and [19]. Thus, Theorem 4.12 establishes a double exponential upper bound with respect to the formula length on the size of the minimal number of states of any Büchi automaton that accepts $L(\varphi)$.

Remark 4.14. The double exponential upper bound on the automata size is tight, i.e., there is a family of formulas $(\varphi_n)_{n \in \mathbb{N}}$ such that for each $n \in \mathbb{N}$, the length of φ_n is linear in n and the index of $\sim_{L(\varphi_n)}$ is double exponential in n . Details are in the appendix. An analogous result with a similar proof has already been shown in [15] for Presburger arithmetic.

5. Conclusion

This paper presented a new method to reason about the sizes of automata that represent first-order definable sets of automatic structures. The method consists of identifying a relationship between the states of a minimal deterministic automaton for a formula and the equivalence classes of a refinement of the equivalence relation determined by Ehrenfeucht-Fraïssé games. We applied the presented method to establish tight upper bounds on the minimal sizes of automata that represent sets definable in $\text{FO}(\mathbb{Z}, <)$ and $\text{FO}(\mathbb{R}, +, <)$. For $\text{FO}(\mathbb{R}, +, <)$, previously proposed techniques based on quantifier-elimination methods [15] failed to establish a double exponential upper bound on the automata size. We hope that the new insights will eventually lead to more efficient automata constructions that can be used to decide $\text{FO}(\mathbb{R}, +, <)$ more efficiently.

As future work, we want to investigate how, and to what extent, the upper bounds on the automata sizes depend on how elements of a structure are encoded as words. The word encoding of integers and reals that we have used in this paper is based on the ϱ 's complement representation, for some $\varrho \in \mathbb{N}$ with $\varrho \geq 2$. There are various other word encodings of numbers so that, e.g., $\text{FO}(\mathbb{Z}, <)$ admits an automata-based decision procedure. For a study on the impact of encodings in automatic structures, see, e.g., [14]. We also plan to apply the presented technique to establish further upper bounds on automata sizes for other automatic structures and use it to simplify the proofs of previously established upper bounds. For instance, for Presburger arithmetic, we expect that we can use a family of equivalence relations that is similar to the one used in this paper for $\text{FO}(\mathbb{R}, +, <)$. However, we have to adjust the bounds on the coefficients and take the definable divisibility relations into account.

Acknowledgments. The author thanks David Basin, Cas Cremers, Matthias Schmalz, and the anonymous reviewers for their comments on earlier versions of the paper.

References

- [1] S. BARDIN, A. FINKEL, J. LEROUX, AND L. PETRUCCI, *FAST: Fast acceleration of symbolic transition systems*, in Proc. of the 15th Int. Conf. on Computer Aided Verification (CAV), vol. 2725 of Lect. Notes Comput. Sci., Springer, 2003, pp. 118–121.
- [2] S. BARDIN, J. LEROUX, AND G. POINT, *FAST extended release*, in Proc. of the 18th Int. Conf. on Computer Aided Verification (CAV), vol. 4144 of Lect. Notes Comput. Sci., Springer, 2006, pp. 63–66.

- [3] B. BECKER, C. DAX, J. EISINGER, AND F. KLAEDTKE, *LIRA: Handling constraints of linear arithmetics over the integers and the reals*, in Proc. of the 19th Int. Conf. on Computer Aided Verification (CAV), vol. 4590 of Lect. Notes Comput. Sci., Springer, 2007, pp. 307–310.
- [4] A. BLUMENSATH AND E. GRÄDEL, *Finite presentations of infinite structures: Automata and interpretations*, Theory Comput. Syst., 37 (2004), pp. 641–674.
- [5] B. BOIGELOT, S. JODOGNE, AND P. WOLPER, *An effective decision procedure for linear arithmetic over the integers and reals*, ACM Trans. Comput. Log., 6 (2005), pp. 614–633.
- [6] B. BOIGELOT AND P. WOLPER, *Representing arithmetic constraints with finite automata: an overview*, in Proc. of the 18th Int. Conf. on Logic Programming (ICLP), vol. 2401 of Lect. Notes Comput. Sci., Springer, 2002, pp. 1–19.
- [7] J. BÜCHI, *Weak second-order arithmetic and finite automata*, Z. Math. Logik Grundlagen Math., 6 (1960), pp. 66–92.
- [8] D. C. COOPER, *Theorem proving in arithmetic without multiplication*, in Proc. of the 7th Annual Machine Intelligence Workshop, B. Meltzer and D. Michie, eds., Edinburgh University Press, 1972, pp. 91–100.
- [9] H.-D. EBBINGHAUS, J. FLUM, AND W. THOMAS, *Mathematical Logic*, Springer, 2nd ed., 1994.
- [10] J. FERRANTE AND C. RACKOFF, *A decision procedure for the first order theory of real addition with order*, SIAM J. Comput., 4 (1975), pp. 69–76.
- [11] J. FERRANTE AND C. W. RACKOFF, *The Computational Complexity of Logical Theories*, vol. 718 of Lect. Notes Math., Springer, 1979.
- [12] M. J. FISCHER AND M. O. RABIN, *Super-exponential complexity of Presburger arithmetic*, in Symp. on Applied Mathematics, vol. VII of SIAM-AMS Proceedings, 1974, pp. 27–41.
- [13] B. KHOUSSAINOV AND A. NERODE, *Automatic presentations of structures*, in Proc. of the Int. Workshop on Logical and Computational Complexity (LCC), vol. 960 of Lect. Notes Comput. Sci., Springer, 1995, pp. 367–392.
- [14] B. KHOUSSAINOV, S. RUBIN, AND F. STEPHAN, *Definability and regularity in automatic structures*, in Proc. of the 21st Annual Symp. on Theoretical Aspects of Computer Science (STACS), vol. 2996 of Lect. Notes Comput. Sci., Springer, 2004, pp. 440–451.
- [15] F. KLAEDTKE, *On the automata size for Presburger arithmetic*, in Proc. of the 19th Annual IEEE Symp. on Logic in Computer Science (LICS), IEEE Computer Society Press, 2004, pp. 110–119. Accepted for publication in ACM Trans. Comput. Log..
- [16] N. KLARLUND, A. MØLLER, AND M. I. SCHWARTZBACH, *MONA implementation secrets*, Int. J. Found. Comput. Sci., 13 (2002), pp. 571–586.
- [17] D. KOZEN, *Theory of Computation*, Springer, 2006.
- [18] R. E. LADNER, *Application of model theoretic games to discrete linear orders and finite automata*, Inform. and Control, 33 (1977), pp. 281–303.
- [19] C. LÖDING, *Efficient minimization of deterministic weak ω -automata*, Inform. Process. Lett., 79 (2001), pp. 105–109.
- [20] O. MALER AND L. STAIGER, *On syntactic congruences for omega-languages*, Theoret. Comput. Sci., 181 (1997), pp. 93–112.
- [21] A. R. MEYER, *Weak monadic second-order theory of successor is not elementary-recursive*, in Logic Colloquium, vol. 453 of Lect. Notes Math., Springer, 1975, pp. 132–154.
- [22] J. PAREDAENS, J. VAN DEN BUSSCHE, AND D. VAN GUCHT, *First-order queries on finite structures over the reals*, SIAM J. Comput., 27 (1998), pp. 1747–1763.
- [23] C. REDDY AND D. W. LOVELAND, *Presburger arithmetic with bounded quantifier alternation*, in Proc. of the 10th Annual ACM Symp. on Theory of Computing (STOC), ACM Press, 1978, pp. 320–325.
- [24] L. STOCKMEYER, *The complexity of decision problems in automata theory and logic*, PhD thesis, Department of Electrical Engineering, MIT, Boston, MA, USA, 1974.

Appendix A. Proof Details

A.1. Proof of Lemma 3.1

We first show the following back-and-forth property about the family $(E_n^s)_{s,n \in \mathbb{N}}$ of equivalence relations:

If $\bar{a} E_{2^{m+1}}^r \bar{b}$ then for every $a' \in \mathbb{Z}$, there is some $b' \in \mathbb{Z}$ such that $(\bar{a}, a') E_{2^m}^{r+1} (\bar{b}, b')$. (1)

Assume $\bar{a} E_{2^{m+1}}^r \bar{b}$. For $a' \in \mathbb{Z}$, let $k \in \{1, \dots, r\}$ be an index so that $d := |a_k - a'|$ is minimal. We choose

$$b' := \begin{cases} b_k + \min\{d, 2^m\} & \text{if } a_k \leq a', \\ b_k - \min\{d, 2^m\} & \text{otherwise.} \end{cases}$$

For $d = 0$, the proof is straightforward, since $a' = a_k$ and also $b' = b_k$ by definition. Assume $d > 0$. By definition, we have that $a' < a_k \Leftrightarrow b' < b_k$. Without loss of generality, we assume that $a' < a_k$ and $b' < b_k$. The other case is symmetric. It remains to show that $\text{sign}(a_i - a' - c) = \text{sign}(b_i - b' - c)$ and $\text{sign}(a' - a_i - c) = \text{sign}(b' - b_i - c)$, for all $c, i \in \mathbb{N}$ with $c \leq 2^m$ and $1 \leq i \leq r$. We show that $\text{sign}(a' - a_i - c) = \text{sign}(b' - b_i - c)$. The case $\text{sign}(a_i - a' - c) = \text{sign}(b_i - b' - c)$ can be proved analogously and we omit it.

Case $a' - a_i < 0$. We have that $a_k \leq a_i$, otherwise k is not minimal. From the assumption $\bar{a} E_{2^{m+1}}^r \bar{b}$, it follows that $b_k \leq b_i$. By definition, we have that $b' - b_i = b_k - b_i - \min\{d, 2^m\} < 0$. So, it holds that $a' - a_i - c < 0$ and $b' - b_i - c < 0$.

Case $a' - a_i \geq 0$. If $d \leq 2^m$, we have that

$$a' - a_i = a_k - d - a_i = a_k - a_i - d \quad \text{and} \quad b' - b_i = b_k - d - b_i = b_k - b_i - d.$$

From the assumption $\bar{a} E_{2^{m+1}}^r \bar{b}$, it follows that $\text{sign}(a' - a_i - c) = \text{sign}(b' - b_i - c)$. Note that $c + d \leq 2^{m+1}$.

If $d > 2^m$ then $a' - a_i > 2^m$ because of the choice of k . It follows that $\text{sign}(a' - a_i - c) = 1$. From the choice of k , we conclude that $a_k - a_i \geq 2^{m+1}$. From the assumption $\bar{a} E_{2^{m+1}}^r \bar{b}$, it follows that $b_k - b_i \geq 2^{m+1}$, and thus,

$$b' - b_i = b_k - 2^m - b_i = b_k - b_i - 2^m \geq 2^{m+1} - 2^m = 2^m.$$

That means, we have that $\text{sign}(b' - b_i - c) = 1$. This completes the proof of the back-and-forth property (1).

Now, we prove the lemma by induction over $m \geq 0$. For $m = 0$, it is straightforward to see that if $\bar{a} E_1^r \bar{b}$ then \bar{a} and \bar{b} satisfy the same formulas φ with $\text{qd}(\varphi) = 0$. Note that an atomic formula is of the form $x < y$ or $x \approx y$. So, we have that $\bar{a} \equiv_0^r \bar{b}$.

Assume that the claim is true for $m \geq 0$. Furthermore, assume that $\bar{a} E_{2^{m+1}}^r \bar{b}$. We have to show that $\bar{a} \equiv_{m+1}^r \bar{b}$. Since every formula with at most quantifier depth $m + 1$ is logically equivalent to a boolean combination of formulas of the form $\exists y \varphi$ with $\text{qd}(\varphi) \leq m$, it suffices to show that $\exists \bar{a} \models \exists y \varphi[\bar{a}] \Leftrightarrow \exists \bar{b} \models \exists y \varphi[\bar{b}]$, where φ is a formula with $\text{qd}(\varphi) \leq m$. By symmetry, we only need to prove the direction from left to right, i.e., $\exists \bar{a} \models \exists y \varphi[\bar{a}]$ implies $\exists \bar{b} \models \exists y \varphi[\bar{b}]$. If $\exists \bar{a} \models \varphi[\bar{a}, a']$, for some $a' \in \mathbb{Z}$, then by the property (1), there is some $b' \in \mathbb{Z}$ such that $(\bar{a}, a') E_{2^m}^{r+1} (\bar{b}, b')$. By the induction hypothesis, we have that $(\bar{a}, a') \equiv_m^{r+1} (\bar{b}, b')$. We conclude that $\exists \bar{b} \models \varphi[\bar{b}, b']$ and hence, $\exists \bar{b} \models \exists y \varphi[\bar{b}]$.

A.2. Proof of Lemma 4.2

For $m = 0$, there is nothing to prove. For $r = 0$, we have to show that for every $a' \in \mathbb{R}$ there is some $b' \in \mathbb{R}$ such that $a' F_m^1 b'$. Obviously, $b' := a'$ works. In the following, we assume that $r, m > 0$.

Assume that $\bar{a} F_{4m^2}^r \bar{b}$ and let $a' \in \mathbb{R}$. We need to show the existence of some $b' \in \mathbb{R}$ such that

$$\text{sign}(f(\bar{a}) + ca' + p) = \text{sign}(f(\bar{b}) + cb' + p), \quad (5)$$

for all $f \in B_m^r$ and $c, p \in \mathbb{Z}$ with $|c|, |p| \leq m$.

Case $c = 0$. The equality (5) simplifies to $\text{sign}(f(\bar{a}) + p) = \text{sign}(f(\bar{b}) + p)$. We have to show that the latter equality holds. Let $h(\bar{x}) := f(\bar{x}) + p$, for all $\bar{x} \in \mathbb{R}^r$. We have that $h \in B_{4m^2}^r$, since $\|h\| \leq m$ and $|h(\bar{0})| \leq rm + m \leq 2rm \leq 4rm^2$. By the assumption $\bar{a} F_{4m^2}^r \bar{b}$, we conclude that $\text{sign}(f(\bar{a}) + p) = \text{sign}(f(\bar{b}) + p)$.

Case $c \neq 0$. The equality (5) can be rewritten to

$$a' < -\frac{f(\bar{a})+p}{c} \quad \text{iff} \quad b' < -\frac{f(\bar{b})+p}{c}.$$

To show the existence of b' , it suffices to show that the numbers $\frac{f(\bar{a})+p}{c}$, for all $f \in B_m^r$ and $c, p \in \mathbb{Z}$ with $c \neq 0$ and $|c|, |p| \leq m$ lie in the same order on the real line as the corresponding numbers $\frac{f(\bar{b})+p}{c}$. This is the case iff for all $f, g \in B_m^r$ and $c, d, p, q \in \mathbb{Z}$ with $c, d \neq 0$ and $|c|, |d|, |p|, |q| \leq m$, we have that

$$\frac{f(\bar{a})+p}{c} < \frac{g(\bar{a})+q}{d} \quad \text{iff} \quad \frac{f(\bar{b})+p}{c} < \frac{g(\bar{b})+q}{d},$$

or in other words,

$$\text{sign}(df(\bar{a}) + dp - cg(\bar{a}) - cq) = \text{sign}(df(\bar{b}) + dp - cg(\bar{b}) - cq). \quad (6)$$

For the function $h(\bar{x}) := df(\bar{x}) + dp - cg(\bar{x}) - cq$, for all $\bar{x} \in \mathbb{R}^r$, we have that $\|h\| \leq 2m^2$ and $|h(\bar{0})| \leq 2rm^2 + 2m^2 \leq 4rm^2$, i.e., $h \in B_{4m^2}^r$. The equality (6) holds because of the assumption $\bar{a} F_{4m^2}^r \bar{b}$.

A.3. Proof of Lemma 4.3

It suffices to prove that $F_{2^3 \cdot 2^{m-2}}^r$ refines \equiv_m^r , for every $m, r \in \mathbb{N}$, since $2^{3 \cdot 2^m - 2} \leq 2^{4 \cdot 2^m} = 2^{2^{m+2}}$. We prove the claim by induction over $m \in \mathbb{N}$. The base case $m = 0$ is straightforward. Note that $2^{3 \cdot 2^0 - 2} = 2$ and recall that the atomic formulas are of the form $x \approx y$, $x < y$, and $x + y \approx z$. For the step case, we assume that the claim is true for some $m \geq 0$, i.e., $F_{2^3 \cdot 2^{m-2}}^r$ refines \equiv_m^r , where $r \in \mathbb{N}$. We have to show that if $\bar{a} F_{2^3 \cdot 2^{m+1-2}}^r \bar{b}$ then $\bar{a} \equiv_{m+1}^r \bar{b}$. In the following, we assume that $\bar{a} F_{2^3 \cdot 2^{m+1-2}}^r \bar{b}$.

Since every formula with at most quantifier depth $m + 1$ is logically equivalent to a boolean combination of formulas of the form $\exists y \varphi$ with $\text{qd}(\varphi) \leq m$, it suffices to show that $\mathfrak{R} \models \exists y \varphi[\bar{a}] \Leftrightarrow \mathfrak{R} \models \exists y \varphi[\bar{b}]$, where φ is a formula with $\text{qd}(\varphi) \leq m$. By symmetry, we only need to prove the direction from left to right, i.e., $\mathfrak{R} \models \exists y \varphi[\bar{a}]$ implies $\mathfrak{R} \models \exists y \varphi[\bar{b}]$. If $\mathfrak{R} \models \varphi[\bar{a}, a']$, for some $a' \in \mathbb{R}$, then by Lemma 4.2, there is some $b' \in \mathbb{R}$ such that $(\bar{a}, a') F_{2^3 \cdot 2^{m-2}}^{r+1} (\bar{b}, b')$. Note that $4 \cdot (2^{3 \cdot 2^m - 2})^2 = 4 \cdot 2^{3 \cdot 2^{m+1} - 4} = 2^{3 \cdot 2^{m+1} - 2}$. By the induction hypothesis, we have that $(\bar{a}, a') \equiv_{m+1}^r (\bar{b}, b')$. We conclude that $\mathfrak{R} \models \varphi[\bar{b}, b']$ and hence, $\mathfrak{R} \models \exists x \varphi[\bar{b}]$.

A.4. Proof of Lemma 4.4

We prove that the set $\{\frac{d}{c} : c, d \in \mathbb{Z} \text{ with } c \neq 0 \text{ and } |c|, |d| \leq 2m^2\}$ contains for each equivalence class in F_m^1 a representative. For $m = 0$ this is obvious. In the remainder of the proof, we assume that $m > 0$.

Let $a, b \in \mathbb{R}$. We have that $a F_m^1 b$ iff $\text{sign}(f(a)) = \text{sign}(f(b))$, for all $f \in A_m^1$, i.e., $f(a) < 0 \Leftrightarrow f(b) < 0$. This is equivalent to $ca < d \Leftrightarrow cb < d$, for all $c, d \in \mathbb{Z}$ with $c \neq 0$ and $|c|, |d| \leq m$. Note that for $c = 0$ the equivalence is trivially satisfied no matter of $a, b \in \mathbb{R}$. Since $c \neq 0$, the equivalence can be rewritten to $a < \frac{d}{c} \Leftrightarrow b < \frac{d}{c}$. We conclude that each equivalence class of F_m^1 contains an element $a \in \mathbb{R}$ that satisfies one of the following properties:

- (i) a is the rational number $\frac{d}{c}$, where $c, d \in \mathbb{Z}$ with $c \neq 0$ and $|c|, |d| \leq m$.
- (ii) a is a rational number in an interval strictly between $\frac{d}{c}$ and $\frac{d'}{c'}$, where $c, c', d, d' \in \mathbb{Z}$ with $c, c' \neq 0$ and $|c|, |c'|, |d|, |d'| \leq m$.
- (iii) a is a rational number strictly less than $-m$.
- (iv) a is a rational number strictly greater than m .

That means, we can choose the representatives a of the equivalence classes of F_m^1 as follows. For (i), we take $a := \frac{d}{c}$. For (ii), we take the midpoint of the interval: $a := \frac{1}{2}(\frac{d}{c} + \frac{d'}{c'}) = \frac{c'd + cd'}{2cd}$. Note that $|2cd| \leq 2m^2$ and $|c'd + cd'| \leq 2m^2$. For (iii), we take $a := -m - 1$. For (iv), we take $a := m + 2$. Note that $|-m - 1| = |m + 1| \leq -2m \leq 2m^2$. Recall the assumption $m > 0$.

A.5. Proof of Lemma 4.5

The claim when $m = 0$ is trivial. So, assume $m > 0$. Let $a', b' \in \mathbb{R}$. We have that $(\bar{a}, a') F_m^{r+1} (\bar{a}, b')$ iff for all $f \in B_m^r$ and $c, c' \in \mathbb{Z}$ with $c' \neq 0$ and $|c|, |c'| \leq m$, we have that

$$\text{sign}(c + f(\bar{a}) + c'a') = \text{sign}(c + f(\bar{a}) + c'b').$$

This equality is equivalent to

$$a' < -\frac{c+f(\bar{a})}{c'} \quad \text{iff} \quad b' < -\frac{c+f(\bar{a})}{c'}.$$

That means, all equivalence classes of F_m^{r+1} that are consistent with σ can be represented by a pair (\bar{a}, a') , where a' is a rational number that satisfies one of the following properties:

- (i) a' is the rational number $\frac{c+f(\bar{a})}{c'}$, where $f \in B_m^r$, and $c, c' \in \mathbb{Z}$ with $c' \neq 0$ and $|c|, |c'| \leq m$.
- (ii) a' is a rational number in an interval strictly between the adjacent rational numbers $\frac{c+f(\bar{a})}{c'}$ and $\frac{d+g(\bar{a})}{d'}$, where $f, g \in B_m^r$, and $c, c', d, d' \in \mathbb{Z}$ with $c', d' \neq 0$ and $|c|, |c'|, |d|, |d'| \leq m$.
- (iii) a' is a rational number strictly less than the smallest rational number of the form $\frac{c+f(\bar{a})}{c'}$, where $f \in B_m^r$, $c, c' \in \mathbb{Z}$ with $c' \neq 0$ and $|c|, |c'| \leq m$.
- (iv) a' is a rational number strictly greater than the largest rational number of the form $\frac{c+f(\bar{a})}{c'}$, where $f \in B_m^r$, $c, c' \in \mathbb{Z}$ with $c' \neq 0$ and $|c|, |c'| \leq m$.

We first consider the unproblematic cases. For (i), we choose $a' := \frac{f(\bar{a})+c}{c'}$. Note that $|f(\bar{0}) + c| \leq rm + m \leq 2rm^2$. For (iii), we choose $a' := \frac{c+f(\bar{a})}{c'} - 1 = \frac{c+f(\bar{a})-c'}{c'}$. For (iv), we choose $a' := \frac{c+f(\bar{a})}{c'} + 1 = \frac{c+f(\bar{a})+c'}{c'}$. Note that $|c+f(\bar{0}) \pm c'| \leq m+rm+m = rm+2m \leq 2rm^2$.

For (ii), we choose the midpoint of the interval, i.e.,

$$a' := \frac{1}{2} \left(\frac{c+f(\bar{a})}{c'} + \frac{d+g(\bar{a})}{d'} \right) = \frac{d'c+d'f(\bar{a})+c'd+c'g(\bar{a})}{2c'd'}.$$

Let $h(\bar{x}) := d'c + d'f(\bar{x}) + c'd + c'g(\bar{x})$, for $\bar{x} \in \mathbb{R}^r$. We have that $\|h^*\| \leq 2m^2$. Whenever $|h(\bar{0})| \leq 2rm^2$, it follows that $h \in B_{2m^2}^r$. We can guarantee this if $m > 1$, since $|h(\bar{0})| \leq rm + rm + m^2 + m^2 = 2rm + m^2$. If $m > 1$, it follows from a straightforward induction over m that $2rm + m^2 \leq 2rm^2$. For $m = 1$, it can be the case that $|h(\bar{0})| = 2r + 1$. In such a case, let $h' \in B_{2m^2}^r$ be the function with $h'(\bar{x}) := h(\bar{x}) \pm d$, for $\bar{x} \in \mathbb{R}^r$, where $d \in \{-1, 1\}$ is chosen in such a way that the offset of h' is either $2r$ or $-2r$.

A.6. Proof Details of Lemma 4.11

To complete the proof of Lemma 4.11, we need to show that $e_i \leq 2^{i+3}$ and $c_i \leq 2^{i+3}$, for all $i \in \mathbb{N}$. We establish these upper bounds by showing the following facts:

- (1) For all $i \in \{0, \dots, 8\}$, we have that $e_i, c_i \leq 2^{i+3}$.
- (2) For all $i \geq 9$, we have that $e_i \leq e'_i$, where $e'_0 := 1$ and $e'_j := 2e'_{j-1} + 3j$ when $i > 0$.
- (3) For all $i \geq 9$, we have that $c_i \leq c'_i$, where $c'_0 := 1$ and $c'_j := c'_{j-1} + e'_{j-1} + 2j^2$ when $j > 0$.
- (4) For all $i \geq 9$, we have that $e'_i, c'_i \leq 2^{i+3}$.

For convenience, we give the first 11 elements of the sequences in the following table. For readability, we approximate $\log_2(8i)$ conservatively by $\lceil \log_2(8i) \rceil$ in the sequence $(c_i)_{i \in \mathbb{N}}$.

i	0	1	2	3	4	5	6	7	8	9	10
2^{i+3}	8	16	32	64	128	256	512	1024	2048	4096	8192
e_i	0	4	14	36	82	176	366	748	1514	3048	6118
e'_i	1	5	16	41	94	203	424	869	1762	3551	7132
c_i	0	5	19	50	108	222	436	846	1644	3223	6343
c'_i	1	4	17	51	124	268	543	1065	2062	3986	7737

The fact (1) is straightforward by checking the inequalities $e_i \leq 2^{i+3}$ and $c_i \leq 2^{i+3}$ by hand, for all $i \in \{0, \dots, 8\}$.

We show the fact (2) by induction over i . The base case $i = 9$ is obvious. For the step case, assume that the inequality $e_i \leq e'_i$ holds for some $i \geq 9$. We have that

$$e_{i+1} = 2e_i + 2 + 2(i+1) \leq 2e'_i + 2 + 2(i+1) = 2e'_i + 4 + 2i \leq 2e'_i + 3i = e'_{i+1}.$$

We show the fact (3) by induction over i . The base case $i = 9$ is straightforward. For the step case, assume that the inequality $c_i \leq c'_i$ holds for some $i \geq 9$. We have that

$$\begin{aligned} c_{i+1} &= c_i + e_i + 2 + (i+1) \log_2(8(i+1)) \leq c'_i + e'_i + 2 + (i+1) \log_2(8(i+1)) \\ &\leq c'_i + e'_i + 2 + (i+1)(i+1) \leq c'_i + e'_i + 2(i+1)^2 = c'_{i+1}. \end{aligned}$$

For the fact (4), we first show that $e'_{i-1} \leq 2^{i+2} - 2i^2$, for all $i \geq 9$. We prove this inequality by induction over i . The base case $i = 9$ is straightforward. For the step case, assume the inequality $e'_{i-1} \leq 2^{i+2} - 2i^2$ holds, for some $i \geq 9$. We have that

$$e'_i = 2e'_{i-1} + 3i \leq 2(2^{i+2} - 2i^2) + 3i = 2^{i+3} - (4i^2 - 3i) \leq 2^{i+3} - 2(i+1)^2.$$

The last inequality holds because $4i^2 - 3i \geq 2(i+1)^2$, for all $i \geq 4$.

Now, assume that $i \geq 9$. We have that

$$e'_i = 2e'_{i-1} + 3i \leq 2 \cdot 2^{2+i} - 2i^2 + 3i \leq 2^{3+i}.$$

We show the other inequality $c'_i \leq 2^{3+i}$ by induction over i . The base case for $i = 9$ is straightforward. For the step case, assume the inequality $c'_i \leq 2^{3+i}$ holds, for some $i \geq 9$. We have that

$$c'_{i+1} = c'_i + e'_i + 2(i+1)^2 \leq 2^{3+i} + 2^{i+3} - 2(i+1)^2 + 2(i+1)^2 = 2^{4+i}.$$

A.7. Worst-case Example for Real Addition

We provide a worst-case example that shows that the established double exponential upper bound on the automata size for $\text{FO}(\mathbb{R}, +, <)$ is tight. We use the formulas $M_n(x, y, z)$ defined by Fischer and Rabin in [12], for $n \geq 0$. For $a, b, c \in \mathbb{R}$, it holds that

$$\mathfrak{R} \models M_n[a, b, c] \quad \text{iff} \quad ab = c \text{ and } a \in \mathbb{N} \text{ with } a < 2^{2^n}.$$

For $n \geq 0$, the length of the formula M_n is linear in n . Note that for $x \in \mathbb{R}$, we have that $\mathfrak{R} \models M_n[x, 0, 0]$ iff $x \in \mathbb{N}$ and $x < 2^{2^n}$.

We state the following lemma, which is proved in [15].

Lemma A.1. *Let $\ell \in \mathbb{N}$ with $\ell \geq 1$. For all $z \in \mathbb{N}$ with $\varrho^{\ell-1} \leq z \leq \varrho^\ell - 2$, there are $x, y, z' \in \mathbb{N}$ with $x, y, z' < \varrho^\ell$ such that $xy = \varrho^\ell z + z'$.*

The proof for the lower bound on the automata size is based on the following lemma about the set

$$\text{MULT}_m := \{(a, b, c) \in \mathbb{N}^3 : ab = c \text{ and } a, b < \varrho^m\},$$

for $m \in \mathbb{N}$.

Lemma A.2. *Let $m \in \mathbb{N}$ and let $S \subseteq \mathbb{R}^3$ be the graph of a partial function from \mathbb{N}^2 to \mathbb{N} with $\text{MULT}_m \subseteq S$. If S is definable in the first-order logic over \mathfrak{R} then the index of the congruence relation of the ω -language $\{u \in \mathbb{V}_3 : \langle u \rangle \in S\}$ is at least ϱ^m .*

Proof. For $m = 0$, the claim is trivial since the index of \sim_L of every ω -language L is at least 1. In the following, we assume that $m > 0$ and let $T := \{u \in \mathbb{V}_3 : \langle u \rangle \in S\}$.

Let K be the set of words of the form $(0, 0, 0)(0, 0, b_{m-1}) \dots (0, 0, b_0) \in (\Sigma^3)^*$ with $b_{m-1} \neq 0$ and if $b_i = \varrho - 1$, for all $1 \leq i < m$, then $b_0 \leq \varrho - 2$. Let $w \in K$ and let z be the integer that is encoded by the third track of w . It holds that

$$\varrho^{m-1} \leq z \leq \varrho^m - 2.$$

From Lemma A.1, it follows that there are $x, y, z' \in \mathbb{N}$ with $x, y, z' < \varrho^m$ such that

$$xy = \varrho^m z + z'.$$

We conclude that for every prefix u of a word in K , there is a word $v \in (\Sigma^3)^*$ such that $\langle uv \rangle \in \text{MULT}_m$.

Now, let L be the set of all prefixes of K . Let $u, u' \in L \setminus \{\varepsilon\}$ with $u \neq u'$. Moreover, let $v \in (\Sigma^3)^*$ with $\langle uv \rangle \in \text{MULT}_m$. The first and second tracks of uv and $u'v$ encode both the pair (x, y) . The third tracks of uv and $u'v$ are different. It follows that $\langle u'v \rangle \notin \text{MULT}_m$. Since $\text{MULT}_m \subseteq S$ and S is the graph of a partial function, we have that $u \not\sim_T u'$. We conclude that every word in L is in a distinct equivalence class of \sim_T .

In the following, we determine the cardinality of L . For $0 \leq i \leq m+1$, let $L_i := \{w \in L : |w| = i\}$. We have that $L_0 = \{\varepsilon\}$, $L_1 = \{(0, 0, 0)\}$, $L_2 = \{(0, 0, 0)b : b \in \Sigma^3 \setminus \{\bar{0}\}\}$, $L_i = \{wb : w \in L_{i-1} \text{ and } b \in \Sigma^3\}$, for $3 \leq i \leq m$, and $L_{m+1} = K$. It holds that

$$\begin{aligned} |L| &= |L_0| + |L_1| + |L_2| + |L_3| + \cdots + |L_m| + |L_{m+1}| \\ &= 1 + 1 + (\varrho - 1) + (\varrho - 1)\varrho + \cdots + (\varrho - 1)\varrho^{m-2} + (\varrho - 1)\varrho^{m-1} - 2 \\ &= \varrho^m - 1. \end{aligned}$$

We conclude that the index of \sim_L is at least ϱ^m : for every word in L there is a distinct equivalence class and there is one equivalence class for the ω -words not on V_3 . ■

Theorem A.3. *Let $n \in \mathbb{N}$. The size of the index the congruence relation $\sim_{L(\varphi)}$ is at least $2^{\lfloor \frac{2^n}{\log_2 \varrho} \rfloor}$, where $\varphi(x, y, z) := M_{n+1}[x, y, z] \wedge \exists u(u + u \approx u \wedge M_{n+1}(y, u, u))$.*

Proof. First, note that $S := \{(a, b, c) \in \mathbb{N}^3 : \mathfrak{R} \models M_{n+1}[a, b, c] \text{ and } \mathfrak{R} \models M_{n+1}[b, 0, 0]\}$ is the graph of a partial function from \mathbb{N}^2 to \mathbb{N} . Let $m := \lfloor \frac{2^n}{\log_2 \varrho} \rfloor$. It holds that $\text{MULT}_m \subseteq S$ since

$$(\varrho^m - 1)^2 < \varrho^{2m} = 2^{2 \lfloor \frac{2^n}{\log_2 \varrho} \rfloor \log_2 \varrho} \leq 2^{2^{n+1}}.$$

The claim follows directly from Lemma A.2. ■