

# The Uniform Measure of Simple Regular Sets of Infinite Trees<sup>1</sup>

Michał Skrzypczak<sup>a</sup>, Marcin Przybyłko<sup>a,b</sup>

<sup>a</sup>*University of Warsaw*

<sup>b</sup>*University of New Caledonia*

---

## Abstract

We consider the problem of computing the measure of a regular set of infinite binary trees. While the general case remains unsolved, we show that the measure of a language can be computed when the set is given in one of the following three formalisms: a first-order formula with no descendant relation; a Boolean combination of conjunctive queries (with descendant relation); or by a non-deterministic safety tree automaton. Additionally, in the first two cases the measure of the set is always rational, while in the third it is an algebraic number. Moreover, we provide an example of a first-order formula that uses descendant relation and defines a language of infinite trees having an irrational (but algebraic) measure.

*Keywords:* infinite trees, uniform measure, random tree, first-order logic  
*2000 MSC:* 68Q45 Formal languages and automata, 68Q87 Probability in computer science

---

## 1. Introduction

The problem of computing a measure of a set can be seen as one of the fundamental problems considered in the study of probabilistic systems. This problem has been studied mostly implicitly, as it is often one of the intermediary steps in solving stochastic games, cf. [1], in answering queries in probabilistic databases, cf. [2], or in model checking for stochastic branching processes, cf. [3].

---

<sup>1</sup>Both authors were supported by Poland's National Science Centre grant no. 2016/21/D/ST6/00491.

To us, this problem naturally arises in the study of stochastic games. Many of the games considered in the literature can be seen as instances of the stochastic version of Gale-Stewart games [4]. Such games use winning conditions expressed as sets of winning plays, i.e. a set of (in)finite words representing winning plays. Hence, computing the value of a stochastic game involves computing the measure of the winning set with respect to the probabilistic space generated by the stochastic elements of the game.

Mio [5] introduced branching games, i.e. stochastic games for which the plays are represented as (in)finite trees, rather than words. Then, the questions concerning whether a set of trees has a well-defined measure [6], and whether that measure can be computed [7] have been raised and partially answered.

*Related work.* The problem of computing the measure of an arbitrary regular set of trees has been already, explicitly or implicitly, studied. Gogacz et al. [6] prove that regular sets of trees are universally measurable. In the case of infinite trees, Chen et al. [3] show that the measure of a set accepted by deterministic automaton is computable; Michalewski and Mio [7] extend the class of sets with computable measure to the class of sets defined by the so-called game automata. In the case of finite trees, Amarilli et al. [8] show that with measures defined by fragments of the probabilistic XML, i.e. where the support of the measure consists of the trees of bounded depth, the measure is computable for arbitrary regular sets of trees. In the case of regular languages of infinite words, Staiger [9] shows that the measure of every regular language of words is computable. Note that, in all the above results, the inherent deterministic nature of the involved automata plays an important role.

The problem of computing the measure of a set of infinite trees has also been, implicitly, considered in probability games. The problem is a special case of computing the value of a stochastic game when the strategies of players are already chosen. In the case of infinite trees, Przybyłko and Skrzypczak [10] consider branching games with regular winning sets. In the case of words, for the survey of probabilistic  $\omega$ -regular games on graphs see e.g. Chatterjee and Henzinger [1].

The problem under consideration can be also seen as the problem of query evaluation in the probabilistic databases setting. For instance Amarilli et al. [11] enquire into the evaluation problem of conjunctive queries over probabilistic graphs. For an introduction to probabilistic databases see e.g. [2].

*Our contribution.* We provide algorithms that compute the measure of tree sets belonging to some restricted classes of regular languages. We show that, in the case of first-order (FO) formulae using unary predicates and child relation, the uniform measure can be computed in three-fold exponential space. We also show that in the case of Boolean combinations of conjunctive queries (BCCQ) using unary predicates, child relation, and descendant relation, the measure can be computed in exponential space. Additionally, we provide an algorithm for computing the measure of a language given by a non-deterministic safety tree automaton. As the class of languages recognisable by these automata coincides with the class of topologically closed regular languages (see Proposition 6.2), this result provides a method for languages at the basic level of topological complexity of regular tree languages. The algorithm translates the structure of a given automaton  $\mathcal{A}$  into an exponentially bigger first-order formula over the field of reals ( $\mathbb{R}$ ). Thus, decision problems about the measure of  $L(\mathcal{A})$  can be solved in doubly exponential space.

We additionally provide an example of a first-order formula over a two letter alphabet for which the defined set of trees has an irrational uniform measure. An example of a regular language with an irrational uniform measure was already presented in [7], however that language is not first-order definable.

This paper is an extended journal version of [12]. Parts of the material presented here have been included in the PhD thesis of the first author [13].

*Organization of rest of the paper.* In Section 2 we define basic notions used in this article. In Section 3 we showcase some basic properties of the uniform measure on selected examples. The computability of the measure of the regular languages defined by first-order formulae is discussed in Section 4. The computability of the measure of the regular languages defined by conjunctive queries is discussed in Section 5. In Section 6 we discuss the case of safety automata. Finally, in Section 7 we provide some computational complexity bounds. In the last section, we summarise the obtained results and propose some directions of future research.

*Acknowledgements.* The authors would like to express their gratitude to Damian Niwiński for a number of insightful comments on the topic. Also, the authors thank the anonymous referees for their careful reviews and helpful suggestions.

## 2. Preliminaries

In this section we present crucial definitions used throughout this work. We assume basic knowledge of logic, automata, and measure. For introduction to logic and automata see [14], for introduction to topology and measure see [15].

*Words and trees.* By  $\mathbb{N}$  we denote the set of natural numbers, i.e. the set  $\{0, 1, 2, \dots\}$  which, when treated as an ordinal is also denoted by  $\omega$ . An alphabet  $\Gamma$  is any non-empty finite set. A *word* is a partial function  $w: \mathbb{N} \rightarrow \Gamma$  such that the domain  $Dom(w)$  of  $w$  is  $\leq$ -closed. By  $|w|$  we denote the length of the word  $w$ , i.e. the size of its domain. By  $\varepsilon$  we denote the empty word, i.e. the unique word of length 0. If the domain of a word  $w$  is finite, then the word is called *finite*; otherwise, it is called *infinite*. Let  $n \in \mathbb{N}$  and  $\bowtie \in \{<, \leq, =\}$ , then the set of all words over an alphabet  $\Gamma$  of length  $l$  such that  $l \bowtie n$  is denoted  $\Gamma^{\bowtie n}$ , e.g.  $\{0, 1\}^{\leq 5}$  is the set of all binary words of length at most 5. Following the usual convention, the set of all finite words over an alphabet  $\Gamma$ , i.e. the set  $\Gamma^{<\omega}$ , is denoted  $\Gamma^*$  and the set of all infinite words over  $\Gamma$ , i.e. the set  $\Gamma^{=\omega}$ , is denoted  $\Gamma^\omega$ .

A word  $w$  is called a *prefix* of a word  $v$ , denoted  $w \sqsubseteq v$ , if  $Dom(w) \subseteq Dom(v)$  and for every  $i \in Dom(w)$  we have that  $w(i) = v(i)$ . By  $w \cdot v$ , or simply  $wv$ , we denote the *concatenation* of the words  $w$  and  $v$ .

A *tree* is any partial function  $t: \{\mathsf{L}, \mathsf{R}\}^* \rightarrow \Gamma$ , where the domain  $Dom(t)$  is a non-empty prefix-closed set. The elements of the set  $\{\mathsf{L}, \mathsf{R}\}$  are called *directions* (*left* and *right*, respectively) and the elements of the set  $\{\mathsf{L}, \mathsf{R}\}^*$  are called *positions*.

For a given tree  $t$ , the elements of the set  $Dom(t)$  are called *nodes* of  $t$ , or *nodes* for short. Any tree  $t$  is either *finite*, if its domain  $Dom(t)$  is finite, or *infinite*. A tree  $t$  is called a *full tree of height  $k$*  if  $Dom(t) = \{\mathsf{L}, \mathsf{R}\}^{\leq k}$ . A tree  $t$  is called a *full tree* if  $Dom(t) = \{\mathsf{L}, \mathsf{R}\}^*$ .

Let  $\Gamma$  be an alphabet. The set of all trees over the alphabet  $\Gamma$  is denoted by  $\mathcal{T}_\Gamma$ ; the set of all finite trees by  $\mathcal{T}_\Gamma^{<\omega}$ ; the set of all full trees of height  $k$  by  $\mathcal{T}_\Gamma^k$ ; the set of all full trees by  $\mathcal{T}_\Gamma^\omega$ . A tree  $t_1$  is called a *prefix* of a tree  $t_2$ , denoted  $t_1 \sqsubseteq t_2$ , if  $Dom(t_1) \subseteq Dom(t_2)$  and for every  $u \in Dom(t_1)$  we have that  $t_1(u) = t_2(u)$ . For a tree  $t$  and a node  $u \in Dom(t)$ , by  $t.u$  we denote the unique tree such that for every position  $v \in \{\mathsf{L}, \mathsf{R}\}^*$  the following holds.

$$t.u(v) \stackrel{\text{def}}{=} t(uv) \tag{1}$$

The tree  $t.u$  is called the *sub-tree of  $t$  in the node  $u$* . For a tree  $t$  and a position  $u$ , by  $\mathbb{B}_{t,u}$  we denote the set of all full trees in which  $t$  is a prefix of the sub-tree in the node  $u$ , i.e.

$$\mathbb{B}_{t,u} \stackrel{\text{def}}{=} \{t' \in \mathcal{T}_\Gamma^\omega \mid t \sqsubseteq t'.u\}, \quad (2)$$

with  $\mathbb{B}_t \stackrel{\text{def}}{=} \mathbb{B}_{t,\varepsilon}$ .

*Logic.* A tree  $t$  over an alphabet  $\Gamma$  can be seen as a relational structure  $\langle \text{Dom}(t), \varepsilon, s_L, s_R, s, \sqsubseteq, (a^t)_{a \in \Gamma} \rangle$ , where

- $\text{Dom}(t)$  is the domain of  $t$ ;
- $\varepsilon$  is the root constant;
- $s_L, s_R \subseteq \text{Dom}(t) \times \text{Dom}(t)$  are the left child relation ( $u s_L u \cdot L$ ) and the right child relation ( $u s_R u \cdot R$ ), respectively;
- $s$  is the child relation  $s_L \cup s_R$ ;
- $\sqsubseteq$  is the ancestor relation, i.e. the transitive closure of the relation  $s$ ;
- $a^t \subseteq \text{Dom}(t)$  is a subset of  $\text{Dom}(t)$ , for  $a \in \Gamma$ , and the family of sets  $(a^t)_{a \in \Gamma}$  is a partition of  $\text{Dom}(t)$ .

The partition  $(a^t)_{a \in \Gamma}$  induces the tree  $t$  in the natural way:  $t(u) = a$  if and only if  $u \in a^t$ .

The *distance* between two positions is the function  $d: \{L, R\}^* \times \{L, R\}^* \rightarrow \mathbb{N}$  defined as  $d(u, v) = |u| + |v| - 2|x|$ , where  $x$  is the longest common prefix of  $u$  and  $v$ . Equivalently, the distance between two different positions is the length of the shortest undirected path connecting the two nodes in the graph  $\langle \{L, R\}^*, s_L \cup s_R \rangle$ .

*Regular languages.* Formulae of monadic second-order logic (MSO) can quantify over nodes in trees  $\exists x, \forall x$  and over sets of nodes  $\exists X, \forall X$ . A first-order (FO) formula is an MSO formula that does not quantify over the sets of nodes. A sentence is a formula with no free variables.

We say that an MSO formula  $\varphi$  is over a signature  $\Sigma$  if  $\varphi$  is a well-formed formula built from the symbols in  $\Sigma$  together with the quantifiers and logical connectives. Let  $\Gamma$  be an alphabet, in this paper, we consider only formulae over the signatures  $\Sigma$  such that  $\Sigma \subseteq \{\varepsilon, s_L, s_R, s, \sqsubseteq\} \cup \Gamma$ .

Let  $\varphi$  be a first-order formula. We write  $t, v \models \varphi(x_1, \dots, x_k)$ , if the tree  $t$ , as a relational structure, with the valuation  $v \in \text{Dom}(t)^k$  satisfies the formula  $\varphi(x_1, \dots, x_k)$ . If  $\varphi$  is a sentence, we simply write  $t \models \varphi$ . We say that a formula  $\varphi(x_1, \dots, x_k)$  is satisfiable if there is a tree  $t$  and a tuple  $v \in \text{Dom}(t)^k$  such that  $t, v \models \varphi(x_1, \dots, x_k)$ .

Let  $\Gamma$  be an alphabet, the set defined by an MSO sentence  $\varphi$ , denoted  $L(\varphi)$ , is the set of all full trees over the alphabet  $\Gamma$  that satisfy  $\varphi$ , i.e.  $L(\varphi) \stackrel{\text{def}}{=} \{t \in \mathcal{T}_\Gamma^\omega \mid t \models \varphi\}$ . Such a set of trees is called regular. This definition of regular languages of trees is equivalent to the automata based definition, cf. e.g. [14].

*Topology and measure.* Recall that the set of all full trees over an alphabet  $\Gamma$ , denoted  $\mathcal{T}_\Gamma^\omega$ , is the set of all functions  $t: \{\mathbb{L}, \mathbb{R}\}^* \rightarrow \Gamma$ . This set can naturally be enhanced with a topology in such a way that it becomes a homeomorphic copy of the Cantor set, see Gogacz et al. [6] for more detailed definitions.

Note that the family of the sets of the form

$$\{t: \{\mathbb{L}, \mathbb{R}\}^* \rightarrow \Gamma \mid \tau \sqsubseteq t\}, \quad (3)$$

where  $\tau \in \mathcal{T}_\Gamma^{<\omega}$  is a full tree of some finite height, constitutes a *base of that topology*. A set from the basis is called a *base set*. Notice that the above set equals  $\mathbb{B}_\tau$ , see (2).

A set is *open* if it is a union, possibly empty, of some base sets; *closed* if it is the complement of an open set; *clopen* if it is both open and closed. Note that our chosen basis consists of clopen sets. Additionally, the family of clopen sets is closed under finite Boolean combinations.

The *uniform* measure  $\mu^*$  defined on the set of full trees  $\mathcal{T}_\Gamma^\omega$  is the unique complete probability Borel measure such that for every finite tree  $\tau \in \mathcal{T}_\Gamma^{<\omega}$  we have that  $\mu^*(\mathbb{B}_\tau) = |\Gamma|^{-|\text{Dom}(\tau)|}$ . In other words, this measure is such that for every node  $u \in \{\mathbb{L}, \mathbb{R}\}^*$  and label  $a \in \Gamma$ , the probability that in a random tree  $t$  the node  $u$  is labelled with the letter  $a$  is  $\frac{1}{|\Gamma|}$ , i.e.  $\mu^*(\{t \in \mathcal{T}_\Gamma^\omega \mid t(u) = a\}) = \frac{1}{|\Gamma|}$ . Notice that for any two distinct positions  $u, v$  and two letters  $a, b$  the events  $S_{u,a} = \{t \in \mathcal{T}_\Gamma^\omega \mid t(u) = a\}$  and  $S_{v,b} = \{t \in \mathcal{T}_\Gamma^\omega \mid t(v) = b\}$  are independent, i.e.

$$\mu^*(S_{u,a} \cap S_{v,b}) = \mu^*(S_{u,a}) \cdot \mu^*(S_{v,b}). \quad (4)$$

As the following theorem implies, every regular set of trees  $L$  has a well-defined uniform measure  $\mu^*(L)$ .

**Theorem 2.1** ([6]). *Every regular language  $L$  of infinite trees is universally measurable, i.e. for every complete Borel measure  $\mu$  on the set of trees, we know that  $L$  is  $\mu$ -measurable.*

Hence, the following problem is well-defined.

**Problem 2.2** (The  $\mu^*(\text{MSO})$  problem). *Is there an algorithm that given an MSO formula  $\varphi$  computes  $\mu^*(L(\varphi))$ ?*

With the  $\mu^*(\text{MSO})$  problem we associate the following decision problem.

**Problem 2.3** (The positive  $\mu^*(\text{MSO})$  problem). *Given an MSO formula  $\varphi$ , decide whether  $\mu^*(L(\varphi)) > 0$ .*

If  $\mathcal{C}$  is a class of regular languages of infinite trees, then by *the (positive)  $\mu^*(\mathcal{C})$  problem*, we understand the above where possible input languages are restricted to the class  $\mathcal{C}$ . If we restrict the class  $\mathcal{C}$  to formulae over the signature  $\Sigma$  we denote it by  $\mathcal{C}(\Sigma)$ .

The problem, in this form, was stated by Michalewski and Mio [7]. It is open in the general case, but some partial results have been obtained, see the paragraph *Related work* for details.

### 3. Simple examples

To better understand the properties of the uniform measures let us consider some simple sets of infinite trees. The presented examples not only give an insight into the behaviour of the uniform measures, but also will be used in the proofs in the following sections.

We start with a simple lemma concerning the existence of sub-trees.

**Lemma 3.1.** *Let  $t$  be a tree over the alphabet  $\Gamma$  and  $u \in \{\mathbf{L}, \mathbf{R}\}^*$  be a position.*

1. *If  $t$  is finite and  $L = \mathbb{B}_{t,u} = \{t' \in \mathcal{T}_\Gamma^\omega \mid t \sqsubseteq t'.u\}$  then  $\mu^*(L) = \Gamma^{-|\text{Dom}(t)|}$ .*
2. *If  $t$  is finite and  $L = \{t' \in \mathcal{T}_\Gamma^\omega \mid \exists v.(u \sqsubset v) \wedge (t \sqsubseteq t'.v)\}$  then  $\mu^*(L) = 1$ .*
3. *If  $t$  is infinite and  $L = \mathbb{B}_{t,u} = \{t' \in \mathcal{T}_\Gamma^\omega \mid t \sqsubseteq t'.u\}$  then  $\mu^*(L) = 0$ .*

*Proof.* The proof of Item 1 is straightforward. To prove Item 2, let  $L_i$  be the set  $L_i \stackrel{\text{def}}{=} \mathbb{B}_{t, u\mathbf{L}^i\mathbf{R}} = \{t' \in \mathcal{T}_\Gamma^\omega \mid t \sqsubseteq t'.(u\mathbf{L}^i\mathbf{R})\}$ . Then, for every  $j \geq 0$  we have that  $L_j \subseteq L$  and, in consequence,  $\overline{L} \subseteq \bigcap_{j \geq 0} \overline{L_j}$ . Hence, for every  $j \geq 0$  we have that

$$1 - \mu^*(L) = \mu^*(\bar{L}) \leq \mu^*\left(\bigcap_{j>i\geq 0} \bar{L}_i\right) = (1 - |\Gamma|^{-|Dom(t)|})^j,$$

where the last inequality follows from the fact that the nodes  $u_L^l$  and  $u_L^k$  are incomparable for  $k \neq l$ , thus  $L_i$  are independent and

$$\mu^*\left(\bigcap_{j>i\geq 0} \bar{L}_i\right) = \prod_{0\leq i<j} \mu^*(\bar{L}_i) = \prod_{0\leq i<j} (1 - |\Gamma|^{-|Dom(t)|}) = (1 - |\Gamma|^{-|Dom(t)|})^j.$$

Taking the limit, we conclude Item 2.

To prove Item 3, let  $t_i$  be a sequence of finite trees such that for every  $i \geq 0$  we have that  $t_i \sqsubseteq t_{i+1} \sqsubseteq t$  and  $|Dom(t_i)| < |Dom(t_{i+1})|$ . Since the sequence of sets  $\mathbb{B}_{t_i, u}$  is decreasing and its limit contains the set  $\mathbb{B}_{t, u}$ , i.e.  $\mathbb{B}_{t_i, u} \supseteq \mathbb{B}_{t_{i+1}, u} \supseteq \mathbb{B}_{t, u}$ , we have that  $\mu^*(\mathbb{B}_{t, u}) \leq \lim_{i \rightarrow +\infty} \mu^*(\mathbb{B}_{t_i, u}) = \lim_{i \rightarrow +\infty} |\Gamma|^{-|Dom(t_i)|} = 0$ .  $\square$

The above examples may suggest that the uniform measures enjoy a form of *Kolmogorov's zero-one law*: e.g. in a random tree a given finite structure exists with probability 1, whereas a given infinite structure exists with probability 0. It is not exactly the case, as can be seen by the following examples.

**Example 3.2.** Let  $\Gamma = \{a, b, c\}$ .

1. If  $L_a$  is the set of trees over the alphabet  $\Gamma$  with arbitrarily long sequences of  $a$ -labelled nodes, i.e.

$$L_a = \{t \in \mathcal{T}_\Gamma^\omega \mid \forall k \geq 0. \exists w, v \in \{L, R\}^* . ((|v| \geq k) \wedge \forall u \sqsubseteq v. t(wu) = a)\},$$

then  $\mu^*(L_a) = 1$ .

2. If  $L_{a3}$  is the language of trees over the alphabet  $\Gamma$  with an infinite  $\{a\}$ -labelled path starting at the root, i.e.

$$L_{a3} = \{t \in \mathcal{T}_\Gamma^\omega \mid \exists w \in \{L, R\}^\omega . \forall u \sqsubset w. t(u) = a\},$$

then  $\mu^*(L_{a3}) = 0$ .

3. If  $L_{a2}$  is the language of trees over the alphabet  $\{a, b\}$  with an infinite  $\{a\}$ -labelled path starting at the root, i.e.

$$L_{a2} = \{t \in \mathcal{T}_{\{a, b\}}^\omega \mid \exists w \in \{L, R\}^\omega . \forall u \sqsubset w. t(u) = a\},$$

then  $\mu^*(L_{a2}) = 0$ .

4. If  $L_{ab}$  is the language of trees over the alphabet  $\Gamma$  with an infinite  $\{a, b\}$ -labelled path starting at the root, i.e.

$$L_{ab} = \{t \in \mathcal{T}_\Gamma^\omega \mid \exists w \in \{\mathbf{L}, \mathbf{R}\}^\omega. \forall u \sqsubseteq w. t(u) \in \{a, b\}\},$$

$$\text{then } \mu^*(L_{ab}) = \frac{1}{2}.$$

*Calculating the measures.* To see Item 1, let  $t^i$  be a full tree of height  $i$  such that every node in  $\text{Dom}(t^i)$  is labelled  $a$  and let  $L^i$  be the language of trees having  $t^i$  as a prefix of its sub-tree at some node  $u$ . Then, by Item 2 of Lemma 3.1 we have that  $\mu^*(L^i) = 1$ . Moreover,  $\bigcap_{i \geq 1} L^i \subseteq L_a$  and  $L^{i+1} \subseteq L^i$ . Since every measure is monotonically continuous, we have that

$$\mu^*(L_a) \geq \mu^*\left(\bigcap_{i \geq 1} L^i\right) = \lim_{n \rightarrow +\infty} \mu^*\left(\bigcap_{n \geq i \geq 1} L^i\right) = 1.$$

Let  $\phi(t)$  stay for “in the tree  $t$  there is an infinite  $\{a\}$ -labelled path starting at the root” then Item 2 follows from the fact that the language in question is regular, thus measurable, and its measure satisfies the following equation.<sup>2</sup>

$$\begin{aligned} \mu^*(L_{a3}) &= \mu^*(\varphi(t) \wedge t(\varepsilon) \neq a) + \\ &\quad \mu^*(t(\varepsilon) = a) \cdot (\mu^*(\phi(t.\mathbf{L})) + \mu^*(\phi(t.\mathbf{R})) - \mu^*(\phi(t.\mathbf{L}) \wedge \phi(t.\mathbf{R}))) \end{aligned}$$

The equation states that the measure of  $L_{a3}$  is equal to the sum of the measures of two sets of trees. The first set consists of all trees  $t$  such that the root is not labelled  $a$  and the tree  $t$  satisfies  $\phi$ . The second set consist of all trees  $t$  such that the root is labelled  $a$ , the sub-tree at the left child of the root satisfies  $\phi$ , i.e.  $\phi(t.\mathbf{L})$ , or the sub-tree at the right child of the root satisfies  $\phi$ , i.e.  $\phi(t.\mathbf{R})$ .

Since the set of all possible trees at left child (or, at right child) of the root is the set of all trees, we get the equation

$$\mu^*(L_{a3}) = \frac{1}{3} \cdot (2\mu^*(L_{a3}) - \mu^*(L_{a3})^2) = \frac{2}{3}\mu^*(L_{a3}) - \frac{1}{3} \cdot \mu^*(L_{a3})^2$$

implying that  $\mu^*(L_{a3}) = 0$  or  $\mu^*(L_{a3}) = -1$ . Since the measure cannot be negative, we conclude that  $\mu^*(L_{a3}) = 0$ .

---

<sup>2</sup>Here, we slightly abuse the notation for the sake of readability. We write  $\mu^*(\psi)$  instead of  $\mu^*(\{t \in \mathcal{T}_\Gamma^\omega \mid \psi\})$ . Moreover, we write  $\psi(t.u)$  for the statement “the sub-tree of  $t$  in  $u$  satisfies  $\psi$ ”.

Similarly, in Item 3 we get the equation

$$\mu^*(L_{a2}) = \frac{1}{2} \cdot (2\mu^*(L_{a2}) - \mu^*(L_{a2})^2) = \mu^*(L_{a2}) - \frac{1}{2} \cdot \mu^*(L_{a2})^2 \quad (5)$$

implying that  $\mu^*(L_{a2}) = 0$ .

In Item 4, we get the equation

$$\mu^*(L_{ab}) = \frac{2}{3} \cdot (2\mu^*(L_{ab}) - \mu^*(L_{ab})^2) = \frac{4}{3}\mu^*(L_{ab}) - \frac{2}{3} \cdot \mu^*(L_{ab})^2 \quad (6)$$

implying that either  $\mu^*(L_{ab}) = \frac{1}{2}$  or  $\mu^*(L_{ab}) = 0$ . Thus we need to look at this example a bit more carefully. Consider a sequence of languages  $\{A^i\}_{i \geq 0}$ , where  $A^0 = \mathcal{T}_\Gamma^\omega$  and  $A^i$  is the language such that there is an  $\{a, b\}$ -labelled path of length  $i$  beginning at the root. Then, we claim the following.

**Claim 3.3.**  $\bigcap_{i \geq 1} A^i = L_{ab}$

*Proof.* Since an infinite path contains sub-paths of arbitrary length, we have that  $\bigcap_{i \geq 1} A^i \supseteq L_{ab}$ . For the reverse inclusion, let  $t \in \bigcap_{i \geq 1} A^i$ . Then, for every  $i \geq 0$  the tree  $t$  has an  $\{a, b\}$ -labelled path of length  $i$ . Since  $t$  is a binary tree, König's lemma assures that the tree  $t$  has an infinite  $\{a, b\}$ -labelled path. This concludes the proof of the claim.  $\square$

Now, for every  $i \geq 0$  we have that  $A^{i+1} \subseteq A^i$  and

$$\mu^*(A^{i+1}) = \frac{2}{3} \cdot (2\mu^*(A^i) - \mu^*(A^i)^2) = \frac{4}{3}\mu^*(A^i) - \frac{2}{3} \cdot \mu^*(A^i)^2. \quad (7)$$

Note that if  $\mu^*(A^i) \geq \frac{1}{2}$ , then  $\mu^*(A^{i+1}) \geq \frac{1}{2}$ . Indeed, the quadratic function  $f(x) = \frac{2}{3}(2x - x^2)$  is monotonically increasing on the interval  $[-\infty, 1]$  and we have that  $f(1) = \frac{2}{3}$  and  $f(\frac{1}{2}) = \frac{1}{2}$ . Since  $\mu^*(A^0) = 1 \geq \frac{1}{2}$ , we conclude that  $\mu^*(L_{ab}) = \frac{1}{2}$ .  $\square$

Before we proceed, observe that the languages  $L_{a3}$ ,  $L_{a2}$ , and  $L_{ab}$  can all be recognised by some non-deterministic safety automata (see Proposition 6.2 and Section 6.1). Thus, instead of computing the measures by hand, we could have invoked Theorem 6.1. It is not incidental, because the idea of inductive approximation of the measure of the set  $L_{ab}$ , expressed by Equation (7), is the cornerstone of the general construction performed in Section 6.

## 4. First-order definable languages

The ideas presented in both Lemma 3.1 and Example 3.2 allow us to compute the measures of sets of trees defined by certain first-order formulae.

### 4.1. First-order definable languages without descendant

**Theorem 4.1.** *Let  $\Gamma$  be an alphabet and let  $\varphi$  be a first-order sentence over the signature  $\{\varepsilon, s_L, s_R, s\} \cup \Gamma$ . Then the measure  $\mu^*(L(\varphi))$  is rational and computable in three-fold exponential space.*

The proof utilises the *Gaifman locality* to partition the formula into two separate sub-formulae. Intuitively, one sub-formulae describes the neighbourhood of the root while the remaining one describes the tree “far away from the root”.

*Gaifman normal form.* Let  $\mathcal{A}$  be a relational structure. The *Gaifman graph* of  $\mathcal{A}$  is the undirected graph  $G^{\mathcal{A}}$  where the set of vertices is the universe of  $\mathcal{A}$  and there is an edge between two vertices in  $G^{\mathcal{A}}$  if there is a relation  $R$  in  $\mathcal{A}$  and a tuple  $x \in R$  that contains  $u$  and  $v$ . The Gaifman distance  $d(u, v)$  between two elements  $u, v$  of the universe of  $\mathcal{A}$  is the distance between  $u$  and  $v$  in the Gaifman graph.

Notice that if we exclude the ancestor relation from the tree structure, then the Gaifman graph of a tree  $t$  is induced by the child relations only and the Gaifman distance coincides with the distance between the positions in the tree. This means, in particular, that for a fixed finite distance  $l$  and any given node  $u \in \text{Dom}(t)$  there is only finitely many positions  $v \in \text{Dom}(t)$  such that the Gaifman distance between  $u$  and  $v$  is  $l$  or less. However, if we allow the ancestor relation then any two nodes in the tree are in Gaifman distance two or less. In this section, from now on, we exclude  $\sqsubseteq$  from the signature.

For a natural number  $r \in \mathbb{N}$ , let  $[d(x, y) \leq r]$  be a first-order formula stating that the distance between  $x$  and  $y$  is at most  $r$ . This formula has two free variables  $x, y$  and its size depends on  $r$ . Similarly, the negation of that formula will be denoted  $[d(x, y) > r]$ .

We say that a first-order formula  $\varphi(x)$  is an  *$r$ -local formula around  $x$*  if the quantifiers of  $\varphi$  are restricted to the  $r$ -neighbourhood of  $x$ , i.e. if the quantifiers inside  $\varphi(x)$  (except those inside the formulae  $[d(x, y) < r]$ ) have

the form  $\forall^{\leq r}$  or  $\exists^{\leq r}$  defined as follows:

$$\begin{aligned}\exists^{\leq r} y. \psi(y) &\stackrel{\text{def}}{=} \exists y. [d(x, y) \leq r] \wedge \psi(y) \\ \forall^{\leq r} y. \psi(y) &\stackrel{\text{def}}{=} \forall y. [d(x, y) \leq r] \rightarrow \psi(y).\end{aligned}$$

We say that a first-order sentence  $\varphi$  is a *basic  $r$ -local sentence* if it is of the form

$$\varphi \stackrel{\text{def}}{=} \exists x_1, \dots, x_n \left( \bigwedge_{i=1}^n \varphi_i^r(x_i) \wedge \bigwedge_{1 \leq i < j \leq n} [d(x_i, x_j) > 2r] \right), \quad (8)$$

where  $\varphi_i^r(x)$  are  $r$ -local formulae around  $x$ .

**Theorem 4.2** (Gaifman). *Every first-order sentence is equivalent to a Boolean combination of basic  $r$ -local sentences, where  $r$  is a number depending on the size of the formula. Furthermore,  $r$  can be chosen so that  $r \leq 7^{qr(\varphi)}$ , where  $qr(\varphi)$  is the quantifier rank of  $\varphi$ .*

As proved by Heimberg et al., cf. [16], the translation to Gaifman normal form can be costly.

**Theorem 4.3** ([16]). *There is a three-fold exponential algorithm on structures of degree 3 that transforms a first-order formula into its Gaifman normal form. Moreover, there are first-order formulae for which the three-fold exponential blow-up is unavoidable.*

*Root formula.* Now we define the idea of a *root formula*, i.e. a formula that necessarily describes the neighbourhood of the root. Let  $\psi(x)$  be an  $r$ -local formula around  $x$ . We say that  $\psi(x)$  is a *root formula* if for every tree  $t \in \mathcal{T}_\Gamma^\omega$  and every position  $u \in \{\text{L}, \text{R}\}^*$  if  $t, u \models \psi(x)$  then  $d(u, \varepsilon) < r$ . Note that every unsatisfiable formula is, by the definition, a root formula.

Let  $\varphi$  be a basic  $r$ -local sentence, i.e. of the form given by (8). We say that  $\varphi_i^r$ , for  $i \in \{1, \dots, n\}$ , is a *root formula* of  $\varphi$  if  $\varphi_i^r$  is a root formula.

**Fact 4.4.** *For every satisfiable basic  $r$ -local sentence there is at most one root formula.*

In other words, only one of  $\varphi_i^r$ s can describe the  $r$ -neighbourhood of the root. As, if two such formulas  $\varphi_i^r, \varphi_j^r$  would be root formulae, then the

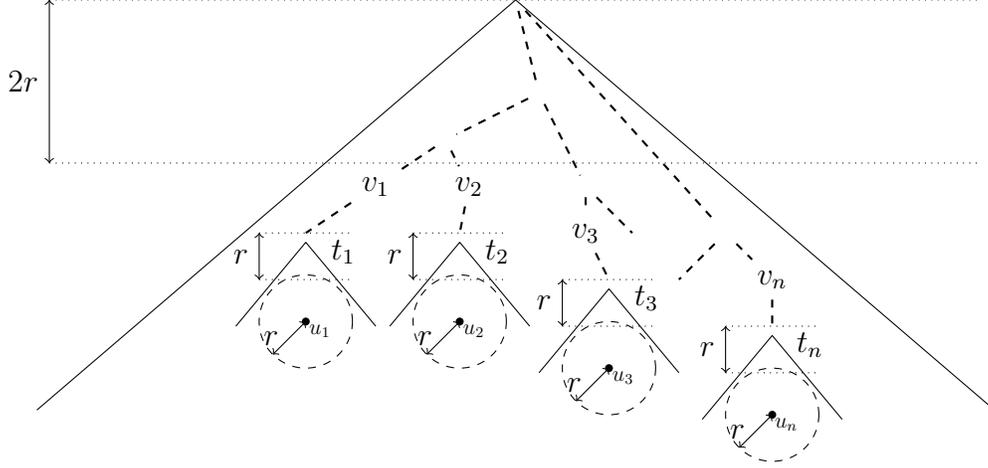


Figure 1: The tree from family  $F$  in the proof of Lemma 4.5.

variables  $x_i$  and  $x_j$  would be mapped in a distance at most  $r-1$  from the root, i.e. in a distance strictly smaller than  $2r$  from each other.

Note that, by the definition of satisfiability, for a tree  $t \in \mathcal{T}_r^\omega$  and a basic  $r$ -local sentence  $\varphi$  we have that  $t \models \varphi$  if and only if there is a function  $\tau: \{x_1, \dots, x_n\} \rightarrow \{\mathbb{L}, \mathbb{R}\}^*$  mapping variables  $x_1, \dots, x_n$  to nodes of  $t$  so that for every  $i \in \{1, \dots, n\}$  we have that  $t, \tau(x_i) \models \varphi_i(x_i)$  and for every pair of indices  $i \neq j$  we have that  $d(\tau(x_i), \tau(x_j)) > 2r$ .

**Lemma 4.5.** *Let  $\varphi$  be a basic  $r$ -local sentence, i.e. as in (8). If  $\varphi$  is*

- *not satisfiable, then  $\mu^*(\mathbb{L}(\varphi)) = 0$ ,*
- *satisfiable and has no root formula, then  $\mu^*(\mathbb{L}(\varphi)) = 1$ ,*
- *satisfiable and has a root formula  $\varphi^*$ ,*  
*then for every  $\tau$  that is a full tree of height  $2r+1$  we have that*

$$\mu^*(\mathbb{L}(\varphi) \cap \mathbb{B}_\tau) = \begin{cases} \mu^*(\mathbb{B}_\tau) & \text{if } \exists u \in \{\mathbb{L}, \mathbb{R}\}^{<r}. \tau, u \models \varphi^*(x); \\ 0 & \text{otherwise.} \end{cases}$$

*Proof.* If  $\varphi$  is not satisfiable then  $\mathbb{L}(\varphi) = \emptyset$  and  $\mu^*(\mathbb{L}(\varphi)) = 0$ . Therefore, let us assume that  $\varphi$  is satisfiable. By Fact 4.4 we know that there is at most one root formula in  $\varphi$ . Let  $I$  be the set of indices of non-root formulae, i.e. for  $i \in I$  we have that  $\varphi_i$  is not a root formula. Since  $\varphi$  is satisfiable, for every

$i \in I$  there are a finite tree  $t_i \in \mathcal{T}_\Gamma^{<\omega}$  and a node  $u_i \in \{\mathsf{L}, \mathsf{R}\}^*$  of length  $|u_i| > r$ , such that  $t_i, u_i \models \varphi_i(x)$  and the set  $\text{Dom}(t_i)$  contains the  $r$ -neighbourhood of  $u_i$ .

Let  $W = \{v_i\}_{i=1}^n$  be a set of  $n$   $\sqsubseteq$ -incomparable nodes such that for  $i \in I$  we have that  $|v_i| > 2r$ . Let  $F = \bigcap_{i \in I} L_i$  where  $L_i$  is the set of trees for which  $t_i$  is a prefix of a sub-tree at some node below  $v_i$ , i.e.  $L_i \stackrel{\text{def}}{=} \{t' \in \mathcal{T}_\Gamma^\omega \mid \exists u.(v_i \sqsubset u) \wedge (t_i \sqsubseteq t'.u)\}$ . By Lemma 3.1 every  $L_i$  has measure 1, thus we have that  $\mu^*(F) = 1$ . Moreover, for every tree  $t \in F$  and index  $i \in I$  there is a node  $v'_i \in \{\mathsf{L}, \mathsf{R}\}^*$  such that  $d(v_i, v'_i) > r$ ,  $v_i \sqsubseteq v'_i$  and  $t, v'_i \models \varphi_i(x)$ .

Now, if there is no root formula in  $\varphi$ , i.e.  $I = \{1 \dots, n\}$ , then  $F \subseteq \text{L}(\varphi)$ . Indeed, let  $t \in F$ , then for  $i \neq j$  we have that  $d(v'_i, v'_j) > 2r$  and we can infer that  $t \models \varphi$ . Hence, the sequence of inequalities

$$1 = \mu^*(F) \leq \mu^*(\text{L}(\varphi)) \leq 1$$

is sound and proves the second bullet of Lemma 4.5.

Consider the opposite case that there is a root formula in  $\varphi$ . Without loss of generality,  $\varphi_1$  is the root formula and  $I = \{2, \dots, n\}$ . Moreover, let  $F$  and  $v'_i$ 's be as before, let  $\tau$  be a full tree of height  $2r+1$ , as stated in the lemma, and  $t \in F \cap \mathbb{B}_\tau$  be a full tree.

If there is  $u_1 \in \{\mathsf{L}, \mathsf{R}\}^{<r}$  such that  $\tau, u_1 \models \varphi_1(x_1)$ , then we take  $v'_1 \stackrel{\text{def}}{=} u_1$ . Now, again, for  $i \neq j$  we have that  $d(v'_i, v'_j) > 2r$  and for all  $i \in I$  we have that  $t, v'_i \models \varphi_i(x_i)$ . In other words, if there is  $u_1 \in \{\mathsf{L}, \mathsf{R}\}^{<r}$  such that  $\tau, u_1 \models \varphi_1(x)$  then  $F \cap \mathbb{B}_\tau \subseteq \text{L}(\varphi) \cap \mathbb{B}_\tau$ . Moreover, since  $F$  is of measure 1, the following sequence of inequalities is sound

$$\mu^*(\mathbb{B}_\tau) = \mu^*(F \cap \mathbb{B}_\tau) \leq \mu^*(\text{L}(\varphi) \cap \mathbb{B}_\tau) \leq \mu^*(\mathbb{B}_\tau).$$

Now assume that there is no such  $u_1$ . We claim that in that case  $\text{L}(\varphi_1) \cap \mathbb{B}_\tau = \emptyset$ . Indeed, if there were a tree  $t$  in the intersection, then the definition of the root formula would provide a node  $u_1 \in \{\mathsf{L}, \mathsf{R}\}^{<r}$  making  $t, u_1 \models \varphi_1(x_1)$  true. But it would mean that  $\tau, u_1 \models \varphi_1(x_1)$  as well by the form of the quantifiers inside  $\varphi_1$ . Thus, in that case  $\mu^*(\text{L}(\varphi) \cap \mathbb{B}_\tau) = 0$ , which concludes the proof.  $\square$

Intuitively, the above lemma states that, when we consider the uniform measure and a basic  $r$ -local sentence, the behaviour of the sentence is almost surely defined by the neighbourhood of the root. This intuition can be formalised as follows.

**Lemma 4.6.** *Let  $\varphi$  be a basic  $r$ -local sentence. Then there is a sentence  $\varphi^*$  such that for every full tree  $\tau$  of height  $2r + 1$  we have that*

$$\mu^*(L(\varphi) \cap \mathbb{B}_\tau) = \mu^*(L(\varphi^*) \cap \mathbb{B}_\tau).$$

Moreover, for every full tree  $t \in \mathbb{B}_\tau$  we have that  $t \models \varphi^*$  if and only if  $\tau \models \varphi^*$ .

*Proof.* If  $\varphi$  has a root formula  $\varphi_i$ , then we take  $\varphi^* \stackrel{\text{def}}{=} \exists x. \varphi_i(x) \wedge [d(x, \varepsilon) < r]$ . If  $\varphi$  has no root formulae but is satisfiable then we take  $\varphi^* \stackrel{\text{def}}{=} \exists x. \varepsilon(x)$ . Otherwise, we take  $\varphi^* \stackrel{\text{def}}{=} \perp$ .  $\square$

The formula  $\varphi^*$  is called the *reduction* of  $\varphi$ . Before we show how to compute the reduction of a basic  $r$ -local formula, we recall a known result.

**Lemma 4.7** (Folklore). *There is an algorithm that given a first-order sentence  $\varphi$  and a finite tree  $\tau$  decides whether  $\tau \models \varphi$  in space polynomial with respect to the size of the formula  $\varphi$  and with respect to the size of the set of nodes of the tree  $\tau$ .*

*Proof.* The lemma is folklore; the property can be easily verified using an alternating polynomial time (*APT*IME) algorithm.  $\square$

Now we show how to compute the reduction.

**Lemma 4.8.** *Given a basic  $r$ -local sentence  $\varphi$  one can compute its reduction  $\varphi^*$  in space polynomial in the size of the formula and doubly exponential in the unary encoding of  $r$ .*

*Proof.* An  $r$ -local formula  $\psi(x)$  is satisfiable in some full tree if, and only if, it is satisfiable in a node  $u$  of some tree of height  $2r + 1$ , such that  $|u| < r + 1$ . Thus, to check the satisfiability of any formula  $\varphi_i$ , we need to check the trees of height at most  $2r + 1$ .

Moreover, to check whether  $\varphi_i$  is a root formula, we need to check whether  $\varphi_i$  is satisfiable and the formula  $\varphi_i(x) \wedge [d(x, \varepsilon) \geq r]$  is not satisfied in any full tree. This, again, can be checked by iterating over all trees of height at most  $2r + 3$ .

Thus, the reduction of  $\varphi$  can be computed by the algorithm `COMPUTEREDUCTION` presented in Algorithm 1. The complexity follows from Lemma 4.7.  $\square$

---

**Algorithm 1** COMPUTEREDUCTION

---

**Require:** a first-order sentence  $\varphi$  in Gaifman normal form  
   $S \leftarrow \{i \mid \varphi_i \text{ is not satisfiable}\}$   
  **if**  $|S| > 0$  **then**  
    **return**  $\perp$   
  **end if**  
   $S \leftarrow \{i \mid \varphi_i \text{ is a root formula}\}$   
  **if**  $|S| = 0$  **then**  
    **return**  $\exists x. \varepsilon(x)$   
  **else if**  $|S| = 1$  **then**  
     $i \leftarrow S.any()$   
    **return**  $\exists x. \varphi_i(x) \wedge [d(x, \varepsilon) < r]$   
  **else**  
    **return**  $\perp$   
  **end if**

---

Lemma 4.6 can be extended to Boolean combinations of basic  $r$ -local sentences by the following property of measurable sets.

**Lemma 4.9.** *Let  $M$  be a measurable space with measure  $\mu$ ,  $W$  be a  $\mu$ -measurable set, and  $\{S_i\}_{i \in I}$  be a family of  $\mu$ -measurable sets such that for every  $i \in I$  either  $\mu(W \cap S_i) = 0$  or  $\mu(W \cap S_i) = \mu(W)$ . Then, for every set  $S$  in the Boolean algebra of sets generated by  $\{S_i\}_{i \in I}$ , we have that either  $\mu(W \cap S) = 0$  or  $\mu(W \cap S) = \mu(W)$ .*

*Proof.* The proof goes by a standard inductive argument. □

Hence, by Lemma 4.5 and the above lemma, we obtain the following.

**Lemma 4.10.** *Let  $\phi$  be a Boolean combination of basic  $r$ -local sentences and  $\tau$  be a full tree of height  $2r+1$ . Then,  $\mu^*(L(\phi) \cap \mathbb{B}_\tau) = \mu^*(L(\phi^*) \cap \mathbb{B}_\tau)$ , where  $\phi^*$  is the reduction of  $\phi$ , i.e. the Boolean combination  $\phi$  with its every basic  $r$ -local sentence  $\varphi$  replaced by its reduction  $\varphi^*$ .*

Moreover,

$$\mu^*(L(\phi^*) \cap \mathbb{B}_\tau) = \begin{cases} \mu^*(\mathbb{B}_\tau) & \text{if } \tau \models \phi^*; \\ 0 & \text{otherwise.} \end{cases}$$

*Proof.* For every full tree  $\tau$  of height  $r$  we use Lemma 4.9 with  $M$  being the set of full trees  $\mathcal{T}_\Gamma^\omega$ ,  $\mu$  being the uniform measure  $\mu^*$ , and  $W$  being the set  $\mathbb{B}_\tau$ . The sets  $S_i$  are the sets of trees defined by the basic  $r$ -local sentences and  $S = L(\phi)$ . By Lemma 4.6, the assumptions of Lemma 4.9 are satisfied.  $\square$

With the above lemmas, we can finally prove Theorem 4.1.

*Proof of Theorem 4.1.* Let  $\varphi$  be a first-order sentence as in the theorem. We utilise the Gaifman locality theorem (see Theorem 4.2 on page 12) to translate the sentence  $\varphi$  into a Boolean combination  $\phi$  of basic  $r$ -local sentences. Now, let  $\phi^*$  be the reduction of  $\phi$ , as in Lemma 4.10, and let  $S = \mathcal{T}_\Gamma^{2r+1}$  be the set of all full trees of height  $h = 2r+1$ . Then,

$$\begin{aligned} \mu^*(L(\phi)) &\stackrel{1}{=} \mu^*(L(\phi) \cap (\bigcup_{\tau \in S} \mathbb{B}_\tau)) &\stackrel{2}{=} \mu^*(\bigcup_{\tau \in S} (L(\phi) \cap \mathbb{B}_\tau)) \\ &\stackrel{3}{=} \sum_{\tau \in S} \mu^*(L(\phi) \cap \mathbb{B}_\tau) &\stackrel{4}{=} \sum_{\tau \in S} \mu^*(L(\phi^*) \cap \mathbb{B}_\tau) \\ &\stackrel{5}{=} \sum_{\tau \in S \wedge \tau \models \phi^*} \mu^*(\mathbb{B}_\tau) &\stackrel{6}{=} |\{\tau \in S \mid \tau \models \phi^*\}| \cdot \frac{1}{|\Gamma|^{2h+1-1}}. \end{aligned}$$

The first equation follows from the fact that the sets in  $\{\mathbb{B}_\tau \mid \tau \in S\}$  are pairwise disjoint. The second from operations on sets and the third is a simple property of measures. The fourth follows from the first part of Lemma 4.10, while the fifth follows from the second part of this lemma. The last equation is a consequence of the fact that  $\mu^*(\mathbb{B}_\tau) = |\Gamma|^{-|\text{Dom}(\tau)|}$ .

Since  $\mu^*(L(\phi)) = \frac{|\{\tau \in S \mid \tau \models \phi\}|}{|\Gamma|^{2h+1-1}}$ , it is enough to count how many full trees of height  $h = 2r+1$  satisfy the reduction of  $\phi$ . The pseudo-code of the algorithm, called COMPUTEMEASUREFO, is presented in Algorithm 2.

The complexity upper bound comes from the fact that translating a first-order sentence  $\varphi$  into its Gaifman normal form can be done in three-fold exponential time and can produce a three-fold exponential sentence  $\phi$  in result, see [17] for details. The resulting sentence  $\phi$  is a Boolean combination of basic  $r$ -local sentences, thus, we can compute its reduction in three-fold exponential space. The function COMPUTEREDUCTION\* computes the reduction  $\phi^*$  of the Boolean combinations by replacing the sentences used in the Boolean combination with their reductions. This can be done in the required complexity, see Lemma 4.8 and note that the size of the sentence dominates the constant  $r$ . Finally, the last part of the algorithm requires us to check the sentence  $\phi^*$  against three-fold exponential number of trees of size that

is two-fold exponential in the size of the original formula. Since model checking of a first-order sentence can be done in polynomial space with respect to the size of the tree and to the size of the sentence, see Lemma 4.7, we get the upper bound.  $\square$

---

**Algorithm 2** COMPUTEMEASUREFO

---

**Require:** a first-order sentence  $\varphi$  and a positive number  $h$   
 $S \leftarrow$  the set of all full trees of height  $h$   
 $\phi \leftarrow$  COMPUTEGAIFMANFORM( $\varphi$ )  
 $\phi \leftarrow$  COMPUTEREDUCTION\*( $\phi$ )  
 $S \leftarrow \{t \in S \mid t \models \phi\}$   
**return**  $|S| \cdot |\Gamma|^{-2^{h+1}+1}$

---

The following remark follows directly from the construction.

**Remark 4.1.** *The above theorem implies that  $\mu^*(L(\varphi)) = \mu^*(L(\varphi^*))$ , where  $\varphi^*$  is the reduction of the given formula  $\varphi$ . Moreover, Lemma 4.6 implies that  $L(\varphi^*)$  is a clopen set because it is a finite union of basic sets.*

4.2. First-order definable languages with descendant

The technique used to prove Theorem 4.1 cannot be extended to formulae utilising the descendant relation because when we allow the descendant relation, the diameter of the Gaifman graph of any tree is at most two. Additionally, as presented in Proposition 4.11 below, sets of full trees defined by such formulae can have irrational measures.

**Proposition 4.11.** *There is a set of full trees over an alphabet  $\Gamma$  that is definable by a first-order formula over the signature  $\{\varepsilon, s_L, s_R, s, \sqsubseteq\} \cup \Gamma$  and the uniform measure of this set is irrational.*

*Proof.* Let  $\Gamma = \{a, b\}$ , we define a language  $L$  in the following way  $L \stackrel{\text{def}}{=} \{t \in \mathcal{T}_{\{a,b\}}^\omega \mid \text{for every path the earliest node labelled } b \text{ (if exists) is at an even depth}\}$ . We will prove that the measure  $\mu^*(L)$  is irrational, and there is a language  $L'$  definable by a first-order formula over the signature  $\{s_L, s_R, \sqsubseteq\} \cup \Gamma$  such that  $\mu^*(L') = \mu^*(L)$ . We start by computing the measure of  $L$ , then we will define  $L'$ .

Observe that the measure  $\mu^*(L)$  satisfies the following equation.

$$\mu^*(L) = \mu^*(\{t \in \mathcal{T}_{\{a,b\}}^\omega \mid t(\varepsilon)=b\}) + \mu^*(\{t \in \mathcal{T}_{\{a,b\}}^\omega \mid t(\varepsilon)=t(\mathsf{L})=t(\mathsf{R})=a\}) \cdot \mu^*(L)^4$$

The equation says that the trees in  $L$  either contain  $b$  at the root (i.e.  $t(\varepsilon) = b$ ) or contain  $a$  in the first three vertices:  $\varepsilon$ ,  $\mathsf{L}$ , and  $\mathsf{R}$ ; and the four subtrees of  $t$  under all nodes of length 2 (i.e.  $t.\mathsf{LL}$ ,  $t.\mathsf{LR}$ ,  $t.\mathsf{RL}$ , and  $t.\mathsf{RR}$ ) belong to  $L$ .

After substituting the appropriate values, we obtain the equation

$$\mu^*(L) = \frac{1}{2} + \frac{1}{8}\mu^*(L)^4 \tag{9}$$

which, by the *rational root theorem*, see e.g. [18] page 116, has no rational solutions.

To conclude the proof, we will describe how to define the language  $L'$ . The crux of the construction comes from the beautiful example by Potthoff, see [19, Lemma 5.1.8]. We will use the following interpretation of the lemma: one can define in first-order logic over the signature  $\{a, b, s_{\mathsf{L}}, s_{\mathsf{R}}, \sqsubseteq\}$  that a given finite tree<sup>3</sup> over the alphabet  $\{a, b\}$  satisfies the following property: every node labelled  $a$  has exactly two children and every node labelled  $b$  is a leaf on an even depth.

A discussion why the above language is in fact First-Order definable can be found in the proof of Theorem 13 on page 14 in [20]. The rough idea is that one can express in FO the notion of *zig-zags*: a pair of nodes  $u \leq v$  of a tree forms a *zig-zag* if the path between them changes direction at each step, i.e. the consecutive directions are  $\mathsf{L}, \mathsf{R}, \mathsf{L}, \mathsf{R}, \dots$  or  $\mathsf{R}, \mathsf{L}, \mathsf{R}, \mathsf{L}, \dots$ . Based on that, one can express an inductive condition, that guarantees that all leafs are at the same depth modulo 2. Finally, there is a unique leaf that is connected by a zig-zag with the root. Based on the shape of that zig-zag one can express the length of it modulo 2.

To construct  $L'$  we simply utilise the formula defining the language in the Potthoff's example to define  $L'$  by substituting:

1. the formula describing a leaf with the formula describing a first occurrence of the label  $b$  on a path:  $\varphi_{\text{leaf}}(x) \stackrel{\text{def}}{=} b(x) \wedge \forall y. (y \sqsubseteq x) \implies a(y)$ ,

---

<sup>3</sup>It is important to notice that the formula of Potthoff works over finite trees, i.e. the fact that a given tree is finite is an assumption that is not expressed by the formula itself.

2. the formula describing an internal node with the formula describing a node labelled  $a$  with no occurrences of the label  $b$  on the path:  
 $\varphi_{\text{node}}(x) \stackrel{\text{def}}{=} a(x) \wedge \forall y. (y \sqsubseteq x) \implies a(y),$

Note that the set  $L'$  agrees with  $L$  on every tree that has a label  $b$  on every infinite path from the root, because such trees are interpreted as finite trees. On the other hand, the truth value of the modified formula on trees that have an infinite path from the root with no nodes labelled  $b$ , i.e. on the set  $L_{a2}$  from Example 3.2, is of no concern to us. Indeed, as previously shown, the uniform measure of the set  $L_{a2}$  is 0.

To be more precise, for every tree  $t \in \mathcal{T}_{\{a,b\}}^\omega \setminus L_{a2}$  we have that  $t \in L \iff t \in L'$ , where  $L_{a2}$  is the language from Example 3.2. Therefore, we have that  $L \cup L_{a2} = L' \cup L_{a2}$ . Since  $\mu^*(L_{a2}) = 0$ , we have that

$$\mu^*(L) = \mu^*(L \cup L_{a2}) = \mu^*(L' \cup L_{a2}) = \mu^*(L'),$$

which concludes the proof.  $\square$

## 5. Conjunctives queries

Proposition 4.11 from the previous section implies that allowing the descendant relation in full first-order logic permits irrational values of measures. Nevertheless, we can allow use of the ancestor relation and retain both rational values and computability when we restrict the formulae to the positive existential fragment using only atomic formulae and conjunction, i.e. to the *conjunctive queries*.

Recall that introducing the ancestor/descendant relation to the tree structure causes that every two nodes in the Gaifman graph are in distance at most two from each other. Thus, for the purpose of having a relevant definition of the distance in the tree, we retain the child related notion of distance, i.e. in this section, as before, the notion of the distance is induced by the child relations only.

*Conjunctive queries.* A *conjunctive query* (CQ) over an alphabet  $\Gamma$  is a formula of first-order logic, using only conjunction and existential quantification, over unary predicates  $a(x)$ , for  $a \in \Gamma$ , the root predicate  $\varepsilon(x)$ , and binary predicates  $s_L(x, y)$ ,  $s_R(x, y)$ ,  $s(x, y)$ , and  $\sqsubseteq(x, y)$ .

An alternative way of looking at conjunctive queries is via graphs and graph homomorphisms. Intuitively, a conjunctive query can be seen as

a graph (a relational structure) in which the variables of the query constitute the vertices, the unary relations of the query label the vertices, and the binary relations form and label the edges. We call such graphs *patterns*. More formally, a pattern  $\pi$  over  $\Gamma$  is a relational structure  $\pi = \langle V, V_\varepsilon, E_L, E_R, E_s, E_{\sqsupseteq}, \lambda_\pi \rangle$ , where  $\lambda_\pi: V \rightarrow \Gamma$  is a partial labelling,  $V_\varepsilon$  is the set of root vertices, and  $G_\pi = \langle V, E_L \cup E_R \cup E_s \cup E_{\sqsupseteq} \rangle$  is a finite graph whose edges are split into left child edges  $E_L$ , right child edges  $E_R$ , child edges  $E_s$ , and ancestor edges  $E_{\sqsupseteq}$ . By  $|\pi|$  we mean the size of the underlying graph.

We say that a tree  $t = \langle \text{Dom}(t), s_L, s_R, \sqsupseteq, (a^t)_{a \in \Gamma} \rangle$  satisfies a pattern  $\pi = \langle V, V_\varepsilon, E_L, E_R, E_s, E_{\sqsupseteq}, \lambda_\pi \rangle$ , denoted  $t \models \pi$ , if there exists a homomorphism  $h: \pi \rightarrow t$ , that is a function  $h: V \rightarrow \text{Dom}(t)$  such that

1.  $h: \langle V, E_L, E_R, E_s, E_{\sqsupseteq} \rangle \rightarrow \langle \text{Dom}(t), s_L, s_R, s_L \cup s_R, \sqsupseteq \rangle$  is a homomorphism of relational structures,
2. for every  $v \in V_\varepsilon$  we have that  $h(v) = \varepsilon$ ,
3. and for every  $v \in \text{Dom}(\lambda_\pi)$  we have that  $\lambda_\pi(v) = t(h(v))$ .

Observe that, by definition, every conjunctive query can be represented as a pattern. The reverse is also true. To obtain a query introduce a variable for every vertex of the pattern, and then express every vertex label by an unary atom and every edge label by a binary atom. Since every pattern can be seen as a conjunctive query and vice versa, we will use those terms interchangeably. The class of conjunctive queries is denoted CQ, the class of formulae that are Boolean combinations of conjunctive queries is denoted BCCQ.

**Theorem 5.1.** *Let  $q$  be a conjunctive query over the signature  $\{\varepsilon, s_L, s_R, s, \sqsupseteq\} \cup \Gamma$ . Then, the uniform measure of the language<sup>4</sup>  $L(q)$  is rational and computable in exponential space.*

To prove the theorem we will modify the concept of *firm* sub-patterns, used e.g. in [21]. Intuitively, a *firm sub-pattern* is a maximal part of a conjunctive query that has to be mapped in a small neighbourhood. The overall proof strategy is similar to the first-order case: we identify those parts (in the form of *firm* sub-patterns) of the conjunctive query that have to be satisfied in the small neighbourhood of the root and those parts that can be satisfied arbitrarily far from the root. As previously, the former decide the value of the measure and the latter can be ignored.

---

<sup>4</sup>A conjunctive query is a special first-order formula, thus the set  $L(q)$  is well-defined.

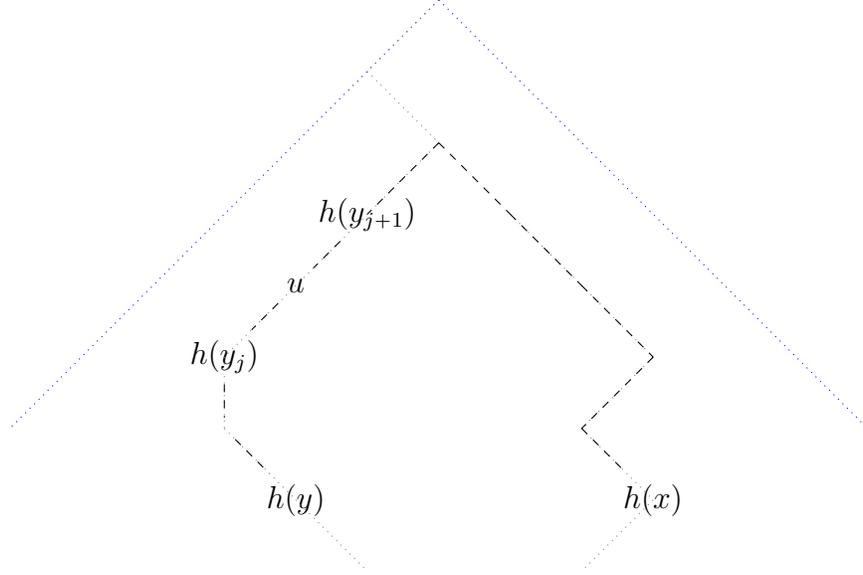


Figure 2: A possible placement of nodes in the proof of Proposition 5.2.

A sub-pattern  $\pi'$  is *firm* if it is a sub-pattern of a pattern  $\pi$  induced by vertices belonging to a maximal strongly-connected component in the *graph of connections*  $C_\pi = \langle V, E \rangle$  such that  $V$  is the set of vertices of  $G_\pi$  and  $\langle x, y \rangle \in E$  if either  $\varepsilon(x)$ ,  $xs_Ly$ ,  $ys_Lx$ ,  $xs_Ry$ ,  $ys_Rx$ ,  $xsy$ ,  $ysx$ , or  $x\sqsupseteq y$ . In particular, a pattern is firm if it has a single strongly-connected component. We say that a sub-pattern is *rooted* if it contains the predicate  $\varepsilon$ .

**Proposition 5.2.** *Let  $\pi = \langle V, V_\varepsilon, E_L, E_R, E_s, E_\sqsupseteq, \lambda_\pi \rangle$  be a firm pattern. Then, for every tree  $t$  such that  $t \models \pi$ , for every two vertices  $x, y$  in  $V$ , and for every homomorphism  $h: \pi \rightarrow t$  we have that  $d(h(x), h(y)) < |\pi|$ . Moreover, if  $\pi$  is rooted then for every vertex  $x$  we have that  $d(h(x), \varepsilon) < |\pi|$ .*

*Proof.* Let us assume otherwise and put  $n = |\pi|$ . Then, there is a tree  $t$ , a homomorphism  $h$ , and two vertices  $x, y$  such that  $t \models \pi$  and  $d(h(x), h(y)) \geq n$ . We claim that  $x$  and  $y$  cannot be in the same strongly-connected component.

Since for some  $m$  we have that  $d(h(x), h(y)) = m - 1 \geq n$ , there is a simple path connecting  $h(x)$  and  $h(y)$ . That is, there is a sequence of distinct nodes  $u_1, u_2, \dots, u_m$  such that  $u_1 = h(x)$ ,  $u_m = h(y)$  and for every  $i$ ,  $u_i$  and  $u_{i+1}$  are in a child relation, i.e.  $s(u_i, u_{i+1})$  or  $s(u_{i+1}, u_i)$ . Notice that every node in the sequence is an ancestor of one of the nodes  $h(x)$  or  $h(y)$ .

Moreover, since the path consists of  $n + 1$  nodes, there has to be a node  $u$  in the sequence such that  $u$  is not in the image of the homomorphism  $h$ . More precisely, there is a node  $u$  such that  $u = u_i$  for some  $1 \leq i \leq m$ ,  $u \notin h(\pi)$ , and one of the nodes  $h(x)$  or  $h(y)$  is a descendant of  $u$ , i.e.  $u \sqsubset h(x)$  or  $u \sqsubset h(y)$ . Without loss of generality, let us say that  $u \sqsubset h(y)$ . Or, more precisely, that  $u_L \sqsubseteq h(y)$ . See Figure 2 for a possible placement of the nodes.

If  $x$  and  $y$  were in the same strongly-connected component then there would be a path that connects  $y$  to  $x$  in the graph of connections  $C_\pi$ , i.e. a sequence of vertices  $y_1, y_2, \dots, y_k$ , for some  $k$ , such that  $y_1 = y$ ,  $y_k = x$ , and for every  $i = 1, \dots, k - 1$  there is an edge between  $y_i$  and  $y_{i+1}$  in  $C_\pi$ . In particular, this would imply that for every  $i$  we have that  $h(y_i)$  and  $h(y_{i+1})$  are  $\sqsubseteq$ -comparable. Now, there would also exist an index  $j \in \{1, \dots, k - 1\}$  such that  $h(y_{j+1}) \sqsubset u \sqsubset h(y_j)$ . Indeed, if there would be no such index, then all the vertices  $y_i$  would satisfy  $u_L \sqsubseteq h(y_i)$ , as  $y_i$  and  $y_{i+1}$  are  $\sqsubseteq$ -comparable for every index  $i$ . But this is impossible because if  $u_L \sqsubseteq h(y_i)$  for all  $i$ , then we would have that  $u_L \sqsubseteq h(y_k) = h(x)$ . Now, since  $u_L \sqsubseteq h(y)$  and  $u_L \sqsubseteq h(x)$ , then, by the definition of the distance,  $u$  would not belong to the sequence  $u_1, \dots, u_m$ . Which is a contradiction with our assumption.

Therefore, there is an index  $j$  such that  $h(y_{j+1}) \sqsubset u \sqsubset h(y_j)$ . Thus, by the definition of  $C_\pi$  we have that either  $\varepsilon(y_j)$ ,  $y_j s_L y_{j+1}$ ,  $y_{j+1} s_L y_j$ ,  $y_j s_R y_{j+1}$ ,  $y_{j+1} s_R y_j$ ,  $y_j s y_{j+1}$ ,  $y_{j+1} s y_j$ , or  $y_j \sqsubset y_{j+1}$ . Neither of child relations is possible because the distance between  $h(y_j)$  and  $h(y_{j+1})$  is at least two. Similarly, both  $y_j \sqsubset y_{j+1}$  and  $\varepsilon(y_j)$  are impossible because we have that  $u \sqsubset h(y_j)$ . Hence, there can be no such sequence  $y_1, \dots, y_k$  and we obtain a contradiction. Thus  $x$  and  $y$  cannot belong to the same strongly-connected component. This proves the first part of the lemma.

Now, if  $\pi$  is rooted then there is a vertex  $y$  such that for every homomorphism  $h$  we have that  $h(y) = \varepsilon$ . Hence, by the first part of the lemma, for all vertices  $x \in \pi$  we have that  $d(h(x), \varepsilon) = d(h(x), h(y)) < n$ .  $\square$

*Graph of firm sub-patterns.* Let  $q$  be a conjunctive query. Consider a graph  $G_q^F = \langle V, E \rangle$  where  $V$  is the set of firm sub-patterns of  $q$  and there is an edge  $\langle v_1, v_2 \rangle \in E \subseteq V \times V$  between two vertices  $v_1, v_2$  if and only if there is an  $\sqsubset$ -labelled edge between some two vertices  $w_1 \in v_1, w_2 \in v_2$ . We call this graph the *graph of firm sub-patterns* of the conjunctive query  $q$ .

**Fact 5.3.** *The directed graph  $G_q^F$  of firm sub-patterns of a conjunctive query  $q$  is acyclic and has at most one rooted firm sub-pattern. We call that sub-pattern the root sub-pattern.*

*Proof.* By the definition of the firm sub-patterns, every node with the predicate  $\varepsilon$  ends up in the same maximal strongly-connected component. The acyclicity follows directly from the fact that firm sub-patterns are the maximal strongly-connected components.  $\square$

As in the case of root formulae, the root pattern decides of the behaviour of a satisfiable conjunctive query, as expressed by the following lemma.

**Lemma 5.4.** *Let  $q$  be a conjunctive query over the signature  $\{\varepsilon, s_L, s_R, s, \sqsubseteq\} \cup \Gamma$ . Then, either*

- *$q$  is not satisfiable and  $\mu^*(L(q)) = 0$ ,*
- *$q$  is satisfiable, has no root sub-pattern, and  $\mu^*(L(q)) = 1$ ,*
- *or  $q$  is satisfiable, has a root sub-pattern  $p$ , and  $\mu^*(L(q)) = \mu^*(L(p))$ .*

*Proof.* If  $q$  is not satisfiable then  $L(q) = \emptyset$  and so  $\mu^*(L(q)) = 0$ . Let  $q$  be satisfiable, i.e. there is a tree  $t^q$  and a homomorphism  $h: G_q \rightarrow t^q$ . Let  $t^r$  be a finite tree such that  $h(G_q) \subseteq \text{Dom}(t^r)$  and let the set  $S \subseteq \mathcal{T}_\Gamma^\omega$  be the set of all trees  $t$  such that for every node  $u \in \{L, R\}^{|q|+1}$  the tree  $t^r$  is a prefix of  $t.u'$  for some node  $u'$  such that  $u \sqsubseteq u'$ . By Lemma 3.1, we have that  $\mu^*(S) = 1$ .

If  $q$  has no root firm sub-pattern, then  $S \subseteq L(q)$  and we have that

$$\mu^*(L(q)) \geq \mu^*(S) = 1.$$

On the other hand, if  $q$  has a root sub-pattern  $p$  then for every tree  $t \in S$  we have that  $t \models q$  if and only if  $t \models p$ . Thus,  $L(q) \cap S = L(p) \cap S$  and since  $\mu^*(S) = 1$ , we have that

$$\mu^*(L(q)) = \mu^*(L(q) \cap S) = \mu^*(L(p) \cap S) = \mu^*(L(p)). \quad \square$$

In other words, the problem of computing the uniform measure of a set of full trees defined by a conjunctive query reduces to the following problem of counting the models of a fixed height.

**Problem 5.5** (Conjunctive queries counting).

*Input:* A conjunctive query  $q$  and a natural number  $n$  in unary.  
*Output:* Number of full trees of height  $n$  that satisfy  $q$ .

**Proposition 5.6.** *Problem 5.5 can be solved in space exponential in  $n$  and polynomial in the size of the query  $q$ .*

*Proof.* All we need is to enumerate all full trees of height  $n$  and check whether they satisfy the query. The number of such trees is exponential in  $n$  and the model checking can be done in polynomial space, with respect to both the query and the size of the tree, see Lemma 4.7.  $\square$

We infer Theorem 5.1 as an immediate consequence.

*Proof of Theorem 5.1.* In polynomial space we can check whether the query is satisfiable, see e.g. [22], and in polynomial time compute its root sub-pattern: it is folklore that one can compute all strongly-connected components of a directed graph in polynomial time.

Now, by Lemma 5.4, if the query is not satisfiable then the measure is 0; if it is satisfiable, but there is no root sub-pattern then the measure is 1. Thus, the only case left is when the query is satisfiable and has a root sub-pattern  $p$ .

If it is the case, then  $\mu^*(L(q)) = \mu^*(L(p))$ . Since  $p$  is a root sub-pattern, then by Proposition 5.2 for any tree  $t$ , any homomorphism  $h: p \rightarrow t$ , and any vertex  $v$  of the query we have that  $|h(v)| < |p|$ . Thus, for any full tree  $t \in \mathcal{T}_\Gamma^\omega$  to decide whether  $t \models p$  we only need to check the prefix of  $t$  that is of height  $|p|$ .

Since the sets  $\mathbb{B}_{t'}$ , where  $t'$  ranges over full trees of height  $|p|$ , form a partition of  $\mathcal{T}_\Gamma^\omega$ , to compute  $\mu^*(L(p))$  it is enough to iterate over all such trees and compute how many of them satisfy  $p$ .

Since every full tree of height  $|p|$  is of exponential size with respect to  $p$  and they can be iterated in exponential space, we infer that the measure  $\mu^*(L(p))$  can be computed in exponential space. To conclude the proof we notice that  $\mu^*(L(p))$  is a finite sum of the measures of sets of form  $\mathbb{B}_t$ , where  $t$  is some finite tree. Since for a finite tree  $t$  the measure of  $\mathbb{B}_t$  is rational, then so is  $\mu^*(L(p))$ .  $\square$

As in the case of first-order formulae, we can lift Theorem 5.1 to Boolean combinations of conjunctive queries.

**Corollary 5.7.** *Let  $\varphi$  be a Boolean combination of conjunctive queries. Then, the uniform measure  $\mu^*(L(\varphi))$  is rational and can be computed in exponential space.*

*Proof.* This is an immediate consequence of Lemma 4.9 and Theorem 5.1.

Indeed, by Lemma 4.9 and Theorem 5.1, patterns in  $\varphi$  can be replaced by their root patterns without the change of measure, or by the query  $\perp$  if they are unsatisfiable. The rest of the reasoning follows as in the proof of Theorem 5.1.  $\square$

We can slightly strengthen the above result if we want to only decide whether the measure is positive.

**Proposition 5.8.** *The positive  $\mu^*(BCCQ)$  problem is in NEXP.*

*Proof.* Let  $\phi$  be a Boolean combination of conjunctive queries. Let  $m$  be the maximum over the sizes of the conjunctive queries in  $\phi$ . By Lemma 5.4 and Lemma 4.9, we can translate  $\phi$  into a Boolean combination of firm, rooted conjunctive queries  $\phi^*$  such that  $\mu^*(L(\phi)) = \mu^*(L(\phi^*))$  and  $\phi^*$  is of polynomial size with respect to  $\phi$ . This can be done in exponential time and requires verifying whether the patterns in  $\phi$  are satisfiable.

Now, since every conjunctive query in  $\phi^*$  is firm and rooted, then either there is a finite tree  $t$  of depth  $2n$  such that  $t \models \phi^*$  or  $\phi^*$  is not satisfiable. If there is no such tree, then  $\mu^*(L(\phi^*)) = 0$ . If such a tree  $t$  exists, then  $\mathbb{B}_t \subseteq L(\phi^*)$  and by Lemma 3.1 we have that  $\mu^*(L(\phi^*)) > 0$ . Hence, it is enough to guess the tree  $t$  and verify that  $t \models \phi^*$ . Since  $t$  is of exponential size in  $\phi$  and model checking of a conjunctive query can be done in polynomial time, we infer the upper bound.  $\square$

## 6. Non-deterministic safety automata

In this section we provide an effective method of computing the uniform measure of those regular sets of infinite trees which are recognisable by non-deterministic safety automata. These automata, also known as automata *with no acceptance condition* can be equivalently defined as parity automata that use only the parity 0, i.e. have Rabin–Mostowski index equal to  $(0, 0)$ . The result of this section is expressed by the following theorem.

**Theorem 6.1.** *The measure  $m = \mu^*(L(\mathcal{A}))$  is an algebraic number that can be effectively computed given a safety automaton  $\mathcal{A}$ . Moreover, this computation can be done in three-fold exponential time and the decision problems about  $m$  (e.g.  $m > 0$ ,  $m > 0.5$ , etc) can be solved in two-fold exponential space.*

Before providing a formal definition of the considered class of languages, we recall a standard fact, that is often considered folklore, see [23, Section 6], [24, Section 3.2], or [25].

**Proposition 6.2.** *The following conditions are equivalent for a regular set of infinite trees  $L$ .*

1.  $L$  can be recognised by a non-deterministic safety automaton,
2.  $L$  is closed as a subset of  $\mathcal{T}_\Gamma^\omega$ .

Thus, it can be said that the procedure from Theorem 6.1 works for languages that are topologically simple.

The idea behind the procedure is as follows. If  $L$  is a regular closed set, then  $L$  is an intersection  $\bigcap_{n \in \mathbb{N}} L_n$  of a decreasing sequence of regular sets  $L_n$  that are Boolean combinations of base sets. Moreover, that sequence can be effectively computed. Intuitively, the language  $L_n$  in the above intersection describes the trees in  $L$  up to the depth  $n$ .

Having found an appropriate decomposition as above, we know that  $\mu^*(L) = \lim_{n \rightarrow \infty} \mu^*(L_n)$ , i.e. the sequence of sets  $L_n$  approximates  $L$  with respect to the measure. It turns out that to compute the measure of  $L$  it is enough to be able to track changes of the measures of the subsequent approximations. To do so, we introduce a finitely dimensional space  $\mathcal{D} \subseteq \mathbb{R}^k$  (for appropriately chosen  $k$ ) and show that there exist: an initial value  $\alpha_0 \in \mathcal{D}$  and two computable functions  $\mathcal{F}: \mathcal{D} \rightarrow \mathcal{D}$  and  $\mathcal{M}: \mathcal{D} \rightarrow \mathbb{R}$ , such that:

$$\mu^*(L_n) = \mathcal{M}(\mathcal{F}^n(\alpha_0)). \quad (10)$$

Moreover,  $\mathcal{F}$ ,  $\mathcal{M}$ , and  $\alpha_0$  can be chosen so that  $\mathcal{F}$  is continuous, monotone (w.r.t. an appropriately chosen order on  $\mathcal{D}$ ), and defined by a vector of polynomials;  $\mathcal{M}$  is continuous; and  $\alpha_0$  is the greatest element of  $\mathcal{D}$ . Therefore, repeating a variant of Kleene's Fixpoint Theorem, we infer that  $\lim_{n \rightarrow \infty} \mu^*(L_n)$  is equal to  $\mathcal{M}(\alpha_\infty)$ , where  $\alpha_\infty \in \mathcal{D}$  is the greatest fixpoint of  $\mathcal{F}$ . Now, since the function  $\mathcal{F}$  is given by a vector of polynomials, the value of  $\alpha_\infty$  can be effectively computed using Tarski's Quantifier Elimination.

*Notation.* For the sake of this section, we will use the following notation, for  $d \in \mathbb{N}$  and  $t \in \mathcal{T}_\Gamma^\omega$ :

$$t^{\leq d} \stackrel{\text{def}}{=} t \upharpoonright_{\{L, R\}^{\leq d}} \in \mathcal{T}_\Gamma^d,$$

i.e. the tree  $t^{\leq d}$  is the full tree of height  $d$  obtained by restricting  $t$  to the nodes of height at most  $d$ .

Additionally, we will consider the set  $\mathcal{T}_\Gamma^d$  for  $d \in \mathbb{N}$  as a measurable space with the uniform discrete measure  $\mu^d$ , where the probability of each finite tree  $\tau \in \mathcal{T}_\Gamma^d$  equals  $|\Gamma|^{-|\text{Dom}(\tau)|}$ . Notice that for each subset  $S \subseteq \mathcal{T}_\Gamma^d$  we have

$$\mu^d(S) = \mu^*(\{t \in \mathcal{T}_\Gamma^\omega \mid t^{\leq d} \in S\}) = \mu^*\left(\bigcup_{\tau \in S} \mathbb{B}_\tau\right). \quad (11)$$

### 6.1. Safety automata

We begin the proof by providing a formal definition of safety automata and their properties. A *safety automaton* is a tuple  $\mathcal{A} = \langle \Gamma, Q, \delta, I \rangle$  where:  $\Gamma$  is an alphabet,  $Q$  is a finite set of *states*,  $\delta \subseteq Q \times \Gamma \times Q \times Q$  is a *transition relation*, and  $I \subseteq Q$  is a set of *initial states*. Let  $t$  be a tree. A *run* of an automaton  $\mathcal{A}$  over  $t$  is a function  $\rho: \text{Dom}(t) \rightarrow Q$  such that for each  $v \in \text{Dom}(t)$  if both  $v_L$  and  $v_R$  belong to  $\text{Dom}(t)$  then

$$(\rho(v), t(v), \rho(v_L), \rho(v_R)) \in \delta.$$

Notice that if  $t$  is a finite tree of height 0, i.e.  $\text{Dom}(t) = \{\varepsilon\}$  then every function  $\rho: \{\varepsilon\} \rightarrow Q$  is a run of  $\mathcal{A}$  over  $t$ . A run  $\rho$  is *accepting* if  $\rho(\varepsilon) \in I$ . We say that  $\mathcal{A}$  *accepts* a tree  $t$  if there exists an accepting run of  $\mathcal{A}$  over  $t$ . For  $d \in \mathbb{N}$  by  $L^d(\mathcal{A}) \subseteq \mathcal{T}_\Gamma^d$  we denote the *language* of the automaton  $\mathcal{A}$ , that is the set of full trees of height  $d$  that are accepted by  $\mathcal{A}$ . Similarly,  $L(\mathcal{A}) \subseteq \mathcal{T}_\Gamma^\omega$  is the set of full trees accepted by  $\mathcal{A}$ .

Fix a safety automaton  $\mathcal{A}$  and let  $t$  be a tree over  $\Gamma$ . The *type* of  $t$ , denoted  $\text{type}(t) \subseteq Q$  is defined as the set of states  $q \in Q$  such that there exists a run  $\rho$  of  $\mathcal{A}$  over  $t$  with  $\rho(\varepsilon) = q$ . Thus,  $\mathcal{A}$  accepts  $t$  if and only if  $\text{type}(t) \cap I \neq \emptyset$ .

The following lemma is a standard application of a compactness argument.

**Lemma 6.3.** *Let  $t \in \mathcal{T}_\Gamma^\omega$  be a full tree. Then the following conditions are equivalent*

1.  $t \in L(\mathcal{A})$ ,
2. for every  $d \in \mathbb{N}$  we have  $t^{\leq d} \in L^d(\mathcal{A})$ .

*Proof.* The implication 1.  $\Rightarrow$  2. is clear from the definition. Consider the opposite implication. Let  $t$  be a full tree such that for all  $d \in \mathbb{N}$  we have  $t^{\leq d} \in L^d(\mathcal{A})$ . Assume that  $\rho_d \in \mathcal{T}_Q^d$  is a run witnessing that. Then, there exists a sub-sequence of the sequence of runs  $(\rho_d)_{d \in \mathbb{N}}$  that is point-wise convergent

in  $\mathcal{T}_Q$ . Let  $\rho_\infty$  be the limit of such a sub-sequence. Then  $\rho_\infty \in \mathcal{T}_Q^\omega$  and it is a run of  $\mathcal{A}$  over  $t$ . Moreover,  $\rho_\infty(\varepsilon) \in I$  because it is the case for each of the runs  $\rho_d$ . Therefore,  $t \in L(\mathcal{A})$ .  $\square$

We will now observe that each safety automaton  $\mathcal{A}$  can be lifted to another safety automaton (denoted  $\widehat{\mathcal{A}}$ ) that recognises the same language but its transition relation is bottom-up functional. More precisely, consider a safety automaton  $\mathcal{A}$  and let  $\widehat{\mathcal{A}} \stackrel{\text{def}}{=} \langle \Gamma, \mathcal{P}(Q), \widehat{\delta}, \widehat{I} \rangle$ , where  $\widehat{\delta}: \Gamma \times \mathcal{P}(Q) \times \mathcal{P}(Q) \rightarrow \mathcal{P}(Q)$  is defined as

$$\widehat{\delta}(a, R_L, R_R) \stackrel{\text{def}}{=} \{q \in Q \mid \exists (q, a, q_L, q_R) \in \delta. q_L \in R_L \wedge q_R \in R_R\},$$

and  $\widehat{I}$  contains those  $R \subseteq Q$  such that  $R \cap I \neq \emptyset$ . For the sake of simplicity, we will use  $\widehat{\delta}_a: \mathcal{P}(Q)^2 \rightarrow \mathcal{P}(Q)$  to denote the function  $\widehat{\delta}$  with the first argument fixed to  $a$ .

**Remark 6.1.** *Directly from the definition, we have  $L(\mathcal{A}) = L(\widehat{\mathcal{A}})$ .*

In the following constructions we will treat the set  $\mathcal{P}(Q)$  as partially ordered by the inclusion  $\subseteq$ . The following fact follows directly from the definition of  $\widehat{\delta}$ .

**Fact 6.4.** *The function  $\widehat{\delta}_a$  is monotone on both its arguments.*

## 6.2. Fixpoints

In this section we will develop tools allowing to define the previously mentioned space  $\mathcal{D}$  and the functionals  $\mathcal{F}$  and  $\mathcal{M}$ . Let us fix a safety automaton  $\mathcal{A}$ . Let  $\mathcal{D}$  be the space of probability distributions over  $\mathcal{P}(Q)$ . For each  $R \subseteq Q$  let  $\mathbb{1}_R \in \mathcal{D}$  be the distribution concentrated in  $R$ , i.e.  $\mathbb{1}_R(R) = 1$  and for each  $R' \neq R$  we have  $\mathbb{1}_R(R') = 0$ .

We will now lift the inclusion order  $\subseteq$  on  $\mathcal{P}(Q)$  to an order over  $\mathcal{D}$ . Recall that a set  $\mathcal{U} \subseteq \mathcal{P}(Q)$  is *upward-closed* if whenever  $R \in \mathcal{U}$  and  $R \subseteq R'$  then  $R' \in \mathcal{U}$  as well. Consider two probability distributions  $\alpha, \beta \in \mathcal{D}$  over  $\mathcal{P}(Q)$ . Let  $\alpha \leq \beta$  hold if and only if, for every upward-closed set  $\mathcal{U} \subseteq \mathcal{P}(Q)$ , we have the following inequality

$$\sum_{R \in \mathcal{U}} \alpha(R) \leq \sum_{R \in \mathcal{U}} \beta(R). \quad (12)$$

Intuitively, the partial order  $\leq$  over  $\mathcal{D}$  corresponds to the process of distributing a given distribution  $\beta$  over smaller sets. Such orders are known as *probabilistic powerdomains*, see e.g. [26].

Notice that  $\mathbb{1}_R \leq \mathbb{1}_S$  if and only if  $R \subseteq S$ . Also, it is easy to observe that  $\mathbb{1}_Q$  is the greatest element of  $\mathcal{D}$ .

Our aim is to simulate the process of randomly choosing letters  $a \in \Gamma$  (as done in the probabilistic spaces  $\mathcal{T}_\Gamma^d$ ) as a monotone mapping  $\mathcal{F}$  over the distributions  $\mathcal{D}$ . Fix a letter  $a \in \Gamma$ . Recall that  $\hat{\delta}_a: \mathcal{P}(Q)^2 \rightarrow \mathcal{P}(Q)$  is the transition function of the power-set automaton  $\hat{\mathcal{A}}$ . Define  $\mathcal{F}_a: \mathcal{D} \rightarrow \mathcal{D}$  as

$$\mathcal{F}_a(\alpha)(R) = \sum_{(R_L, R_R) \in \hat{\delta}_a^{-1}(R)} \alpha(R_L) \cdot \alpha(R_R). \quad (13)$$

In other words,  $\mathcal{F}_a(\alpha)$  is the distribution over  $\mathcal{P}(Q)$  given by the values of  $\hat{\delta}_a$  where its arguments are chosen independently according to the distribution  $\alpha$ . Finally, let  $\mathcal{F}: \mathcal{D} \rightarrow \mathcal{D}$  be defined as  $\mathcal{F}(\alpha) = |\Gamma|^{-1} \cdot \sum_{a \in \Gamma} \mathcal{F}_a(\alpha)$  (where the sum is taken coordinate-wise). It is easy to verify that all the functions  $\mathcal{F}_a$  and  $\mathcal{F}$  in fact produce probabilistic distributions, i.e. elements of  $\mathcal{D}$ .

**Lemma 6.5.** *The function  $\mathcal{F}$  is monotone with respect to the order  $\leq$  on  $\mathcal{D}$ .*

*Proof.* Consider two distributions  $\alpha \leq \beta \in \mathcal{D}$ . Let  $\mathcal{U} \subseteq \mathcal{P}(Q)$  be an upward-closed set. Our aim is to prove the inequality (12) for  $\mathcal{F}(\alpha)$  and  $\mathcal{F}(\beta)$ . We will do it for each letter separately, what means that it is enough to show that for each  $a \in \Gamma$  we have

$$\sum_{R \in \mathcal{U}} \sum_{(R_L, R_R) \in \hat{\delta}_a^{-1}(R)} \alpha(R_L) \cdot \alpha(R_R) \leq \sum_{R \in \mathcal{U}} \sum_{(R_L, R_R) \in \hat{\delta}_a^{-1}(R)} \beta(R_L) \cdot \beta(R_R).$$

Since  $\hat{\delta}_a$  is a function, the sums on both sides of the equality take the form

$$\sum_{(R_L, R_R): \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \alpha(R_L) \cdot \alpha(R_R) \leq \sum_{(R_L, R_R): \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \beta(R_L) \cdot \beta(R_R).$$

Which is equivalent to

$$\sum_{R_L \subseteq Q} \alpha(R_L) \cdot \sum_{R_R: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \alpha(R_R) \leq \sum_{R_R \subseteq Q} \beta(R_R) \cdot \sum_{R_L: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \beta(R_L). \quad (14)$$

To obtain the last inequality, we observe that:

$$\begin{aligned}
\sum_{R_L \subseteq Q} \alpha(R_L) \cdot \sum_{R_R: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \alpha(R_R) &\stackrel{(1)}{\leq} \sum_{R_L \subseteq Q} \alpha(R_L) \cdot \sum_{R_R: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \beta(R_R) \\
&\stackrel{(2)}{=} \sum_{R_R \subseteq Q} \beta(R_R) \cdot \sum_{R_L: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \alpha(R_L) \\
&\stackrel{(3)}{\leq} \sum_{R_R \subseteq Q} \beta(R_R) \cdot \sum_{R_L: \hat{\delta}_a(R_L, R_R) \in \mathcal{U}} \beta(R_L),
\end{aligned}$$

where the consecutive (in)equalities follow from:

- (1) Fact 6.4: the family of sets  $R_R$  is upward-closed and  $\alpha \leq \beta$ ;
- (2) rearranging the sum;
- (3) the fact that the family of sets  $R_L$  is upward-closed and  $\alpha \leq \beta$  again.

Thus, we have proved (14). This concludes the proof of the lemma.  $\square$

Now we will see that, in a sense,  $\mathcal{F}$  simulates the behaviour of  $\mathcal{A}$  over  $\mathcal{T}_\Gamma^d$ . Let  $\alpha_0 = \mathbb{1}_Q$  and  $\alpha_{d+1} = \mathcal{F}(\alpha_d)$  for  $d \in \mathbb{N}$ .

**Proposition 6.6.** *For each  $d \in \mathbb{N}$  the distribution  $\alpha_d$  is the distribution of type( $\tau$ ) for a random partial tree  $\tau \in \mathcal{T}_\Gamma^d$ , i.e. for each  $R \subseteq Q$  we have*

$$\alpha_d(R) = \mu^d(\{\tau \in \mathcal{T}_\Gamma^d \mid \text{type}(\tau) = R\}).$$

*Proof.* The proof is by induction. For  $d = 0$  the claim is obvious, as  $\alpha_0 = \mathbb{1}_Q$ , and for each partial tree  $\tau \in \mathcal{T}_\Gamma^0$  we have  $\text{type}(\tau) = Q$ . Let us assume that the inductive hypothesis holds for  $d$ . Observe that

$$\mu^{d+1}(\{\tau \in \mathcal{T}_\Gamma^{d+1} \mid \text{type}(\tau) = R\}) = \sum_{a \in \Gamma} \mu^{d+1}(\{\tau \in \mathcal{T}_\Gamma^{d+1} \mid \tau(\varepsilon) = a \wedge \text{type}(\tau) = R\}).$$

Now, as a partial tree  $\tau \in \mathcal{T}_\Gamma^{d+1}$  can be seen as a pair of partial trees  $\tau_L, \tau_R \in \mathcal{T}_\Gamma^d$  merged by the letter  $\tau(\varepsilon)$ , we obtain that:

$$\begin{aligned}
&\mu^{d+1}(\{\tau \in \mathcal{T}_\Gamma^{d+1} \mid \tau(\varepsilon) = a \wedge \text{type}(\tau) = R\}) = \\
&\mu^d \times \mu^d(\{(\tau_L, \tau_R) \in \mathcal{T}_\Gamma^d \times \mathcal{T}_\Gamma^d \mid \hat{\delta}_a(\text{type}(\tau_L), \text{type}(\tau_R)) = R\}), \quad (15)
\end{aligned}$$

where the latter probability is taken in the product space. Thus, (15) equals

$$\sum_{(R_L, R_R) \in \hat{\delta}_a^{-1}(R)} \mu^d(\{\tau_L \in \mathcal{T}_\Gamma^d \mid \text{type}(\tau_L) = R_L\}) \cdot \mu^d(\{\tau_R \in \mathcal{T}_\Gamma^d \mid \text{type}(\tau_R) = R_R\}),$$

which, by the inductive assumption, equals

$$\sum_{(R_L, R_R) \in \hat{\delta}_a^{-1}(R)} \alpha_d(R_L) \cdot \alpha_d(R_R).$$

The last term is equal to  $\mathcal{F}_a(\alpha_n)(R)$  by the definition, see (13). Therefore, the inductive hypothesis holds for  $d + 1$ .  $\square$

The next three statements follow the standard way of proving Kleene's Fixpoint Theorem. However, instead of simply invoking that result, we prove these claims by hand, as it seems to be much more direct.

**Fact 6.7.** *The sequence  $(\alpha_d)_{d \in \mathbb{N}}$  is descending in the order  $\leq$  on  $\mathcal{D}$ .*

*Proof.* By monotonicity of  $\mathcal{F}$  and the fact that  $\alpha_0$  is the greatest element of  $\mathcal{D}$ .  $\square$

**Lemma 6.8.** *There exists a probabilistic distribution  $\alpha_\infty \in \mathcal{D}$  such that for each  $R \subseteq Q$  we have*

$$\alpha_\infty(R) = \lim_{d \rightarrow \infty} \alpha_d(R).$$

*Proof.* First notice that for each upward-closed set  $\mathcal{U} \subseteq \mathcal{P}(Q)$  the sequence of  $\sum_{R \in \mathcal{U}} \alpha_d(R)$  is non-increasing and bounded for  $d \rightarrow \infty$ . Thus, that sequence has a limit. Now, for each  $R \subseteq Q$  the value  $\alpha_d(R)$  can be written as a difference of sums as above for two upward-closed sets. Therefore,  $\lim_{d \rightarrow \infty} \alpha_d(R)$  exists and we can define  $\alpha_\infty(R)$  as that limit. Since  $\sum_{R \subseteq Q} \alpha_d(R)$  is always 1, the limit values also satisfy that property and thus  $\alpha_\infty$  is a probabilistic distribution.  $\square$

**Lemma 6.9.** *The distribution  $\alpha_\infty$  is the greatest fixpoint of  $\mathcal{F}$ , i.e.  $\mathcal{F}(\alpha_\infty) = \alpha_\infty$  and if  $\mathcal{F}(\alpha') = \alpha'$  for some  $\alpha' \in \mathcal{D}$  then  $\alpha' \leq \alpha_\infty$ .*

*Proof.* First notice that  $\mathcal{F}$  is continuous as a function from  $\mathbb{R}^{\mathcal{P}(Q)}$  to  $\mathbb{R}^{\mathcal{P}(Q)}$ . Therefore,

$$\mathcal{F}(\alpha_\infty) = \mathcal{F}\left(\lim_{d \rightarrow \infty} \alpha_d\right) = \lim_{d \rightarrow \infty} \mathcal{F}(\alpha_d) = \lim_{d \rightarrow \infty} \alpha_{d+1} = \alpha_\infty.$$

Now consider any fixpoint  $\alpha' \in \mathcal{D}$  of  $\mathcal{F}$ . By monotonicity of  $\mathcal{F}$  and the fact that  $\alpha_0 = \mathbb{1}_Q$ , we know that for each  $d \in \mathbb{N}$  we have  $\alpha' \leq \alpha_d$ . As the order  $\leq$  is defined by sums of values and  $\alpha_\infty$  is a point-wise limit of  $\alpha_d$ , it also holds that  $\alpha' \leq \alpha_\infty$ .  $\square$

Recall that  $\hat{I} = \{R \subseteq Q \mid R \cap I \neq \emptyset\}$  is the set of accepting states of the power-set automaton. Consider the functional  $\mathcal{M}: \mathcal{D} \rightarrow \mathbb{R}$  defined as

$$\mathcal{M}(\alpha) = \sum_{R \in \hat{I}} \alpha(R). \quad (16)$$

Clearly  $\mathcal{M}$  is monotone (by the definition of  $\leq$  and the fact that  $\hat{I}$  is upward-closed) and continuous. Thus, the following corollary holds.

**Fact 6.10.** *The sequence  $M_d \stackrel{\text{def}}{=} \mathcal{M}(\alpha_d)$  is a decreasing sequence of numbers with the limit  $M_\infty \stackrel{\text{def}}{=} \mathcal{M}(\alpha_\infty)$ .*

Moreover, Proposition 6.6 implies the following property of the values  $M_d$ .

**Remark 6.2.** *For each  $d \in \mathbb{N}$  we have  $M_d = \mu^d(\mathbb{L}^d(\mathcal{A}))$ .*

### 6.3. The measure of the language

We are now in position to provide a way of computing the measure  $\mu^*(\mathbb{L}(\mathcal{A}))$  of the considered language  $\mathbb{L}(\mathcal{A})$ .

**Proposition 6.11.** *Using the value  $M_\infty$  as defined in Fact 6.10, we have  $\mu^*(\mathbb{L}(\mathcal{A})) = M_\infty$ .*

*Proof.* Consider a depth  $d \in \mathbb{N}$  and let

$$L_d \stackrel{\text{def}}{=} \{t \in \mathcal{T}_\Gamma^\omega \mid t^{\leq d} \in \mathbb{L}^d(\mathcal{A})\}.$$

Observe that by Lemma 6.3 we know that  $\mathbb{L}(\mathcal{A}) = \bigcap_{d \in \mathbb{N}} L_d$ . Moreover, from the definition of  $\mathbb{L}^d(\mathcal{A})$  we know that for  $d \leq d'$  we have  $L_d \supseteq L_{d'}$ . Therefore,  $\mu^*(\mathbb{L}(\mathcal{A})) = \lim_{d \rightarrow \infty} \mu^*(L_d)$ . However, by the definition of the measure on  $\mathcal{T}_\Gamma^\omega$ , see (11) we know that  $\mu^*(L_d) = \mu^d(\mathbb{L}^d(\mathcal{A}))$  where the latter probability is taken in  $\mathcal{T}_\Gamma^d$ . Now, by Remark 6.2 we know that  $\mu^d(\mathbb{L}^d(\mathcal{A})) = M_d$ . Therefore, the thesis holds.  $\square$

We are now in place to prove the main theorem of this section.

*Proof of Theorem 6.1.* By Proposition 6.11 we know that  $\mu^*(L(\mathcal{A}))$  equals  $\sum_{R \in \hat{I}} \alpha_\infty(R)$ , where  $\alpha_\infty$  is the greatest fixpoint of the operator  $\mathcal{F}: \mathcal{D} \rightarrow \mathcal{D}$  (see Lemma 6.9). Since the operator  $\mathcal{F}$  is given by a vector of polynomials in  $\mathbb{R}^{\mathcal{P}(Q)}$ , one can express in first-order logic over  $\langle \mathbb{R}, +, \cdot \rangle$  that  $\alpha_\infty$  is the greatest fixpoint of  $\mathcal{F}$ . For the sake of completeness (and complexity analysis) we will provide a more precise construction here.

Fix a non-deterministic safety automaton  $\mathcal{A} = \langle \Gamma, Q, \delta, I \rangle$ . Let  $R_1, \dots, R_N$  be an enumeration of all the subsets of  $Q$ . First, consider the formulae describing: the elements of  $\mathcal{D}$  (represented as sequences of reals), upward-closed subsets of  $\mathcal{P}(Q)$  (represented as sequences of zeros and ones), and the order  $\leq$ :

$$\begin{aligned} \varphi_{\mathcal{D}}(x_1, \dots, x_N) &\stackrel{\text{def}}{=} \sum_{i=1}^N x_i = 1 \wedge \bigwedge_{i=1, \dots, N} 0 \leq x_i \leq 1, \\ \varphi_{\mathcal{P}(Q)}(l_1, \dots, l_N) &\stackrel{\text{def}}{=} \bigwedge_{i=1, \dots, N} l_i^2 = l_i \wedge \bigwedge_{R_i \subseteq R_j} l_i \leq l_j, \\ \varphi_{\leq}(\vec{x}, \vec{y}) &\stackrel{\text{def}}{=} \varphi_{\mathcal{D}}(\vec{x}) \wedge \varphi_{\mathcal{D}}(\vec{y}) \wedge \forall \vec{l}. \varphi_{\mathcal{P}(Q)}(\vec{l}) \Rightarrow \sum_{i=1}^N x_i l_i \leq \sum_{i=1}^N y_i l_i. \end{aligned}$$

As observed above, the operator  $\mathcal{F}$  is just a vector of polynomials, which means that there exists a formula  $\varphi_{\mathcal{F}}(\vec{x}, \vec{f})$  that holds for a given tuples, if  $\vec{x}, \vec{f} \in \mathcal{D}$  and  $\mathcal{F}(\vec{x}) = \vec{f}$ . Using these, one can write the following formula expressing the measure of the set  $L(\mathcal{A})$ :

$$\begin{aligned} \varphi_{\mu^*(L(\mathcal{A}))}(m) &\stackrel{\text{def}}{=} \exists \vec{x}. \varphi_{\mathcal{F}}(\vec{x}, \vec{x}) \wedge \\ &\quad (\forall \vec{y}. \varphi_{\mathcal{F}}(\vec{y}, \vec{y}) \Rightarrow \varphi_{\leq}(\vec{y}, \vec{x})) \wedge \\ &\quad \sum_{R_i \in \hat{I}} x_i = m. \end{aligned}$$

Notice that the above formula is of size polynomial in  $N$ , i.e. exponential in the number of states and transitions of  $\mathcal{A}$ . The following claim follows now directly from Lemma 6.9.

**Claim 6.12.** *The measure  $\mu^*(L(\mathcal{A}))$  is the unique real number  $m$  satisfying the formula  $\varphi_{\mu^*}(m)$ .*

Therefore, by Tarski's Quantifier Elimination [27, 28], we know that there exists a semialgebraic set [29, Chapter 2] containing a single real number  $\mu^*(L(\mathcal{A}))$ . This set can be computed in three-fold exponential time and can be used as a representation of  $\mu^*(L(\mathcal{A}))$ .

The complexity comes from results of [30]: the theory of real closed fields can be decided in deterministic exponential space. Therefore, the decision problems (i.e. positivity) about  $m = \mu^*(L(\mathcal{A}))$  can be solved in two-fold exponential space.  $\square$

As the following remark implies, the use of algebraic numbers in the above procedure is unavoidable.

**Remark 6.3.** *There exists a regular language recognisable by a non-deterministic safety automaton such that the uniform measure of the language is irrational.*

*Proof.* It is enough to notice that the language  $L$  from Proposition 4.11 is topologically closed and therefore recognisable by a non-deterministic safety automaton.  $\square$

#### 6.4. Regular languages of finite trees

In this section we discuss how the above procedure for closed regular languages can be adjusted to the case of regular languages of finite trees. The main issue in that situation is the fact that for a fixed alphabet  $\Gamma$ , the set of all finite trees  $\mathcal{T}_\Gamma^{<\omega}$  is countably infinite. Therefore, there is no “uniform” measure over  $\mathcal{T}_\Gamma^{<\omega}$ . One of the possible solutions is to consider *discounted* measures, i.e. measures where the structure of the tree is randomly generated top to bottom, and where at each stage of the process there is some *extinction probability*, i.e. the probability that we reach *an end of the structure*. Such systems, called *branching processes*, have been extensively studied, see e.g. [31]. For the study of their reachability and extinction properties see e.g. [32] and its references.

In our setting an extinction of a branching process implies generation of a finite tree. Thus, branching processes that are extinct with probability 1 generate, in the natural way, a measure on the set of finite trees.

For the sake of this section we will formalise that process as function from infinite trees into finite ones. Let  $\Gamma = \{a_1, \dots, a_n\}$  be an alphabet of  $n$  letters and let  $\Gamma' = \Gamma \cup \{b_1, \dots, b_n\}$ , where  $b_i$  are distinct symbols not belonging to  $\Gamma$ . A tree  $t \in \mathcal{T}_{\Gamma'}^\omega$  is called *bounded* if on each branch there is an occurrence of a symbol  $b_i$  for some  $i \in \{1, \dots, n\}$ . Notice that analogously to the language  $L_{a_2}$  from Item 3 of Example 3.2, the measure of bounded trees is 1. Now, if

a tree  $t \in \mathcal{T}_\Gamma^\omega$  is bounded, let  $f(t) \in \mathcal{T}_\Gamma^{<\omega}$  be obtained in the following way:

$$\begin{aligned} \text{Dom}(f(t)) &\stackrel{\text{def}}{=} \{u \in \text{Dom}(t) \mid \forall v \sqsubset u. t(v) \in \Gamma\}, \\ f(t)(u) &\stackrel{\text{def}}{=} \begin{cases} t(u) & \text{if } t(u) \in \Gamma, \\ a_i & \text{if } t(u) = b_i. \end{cases} \end{aligned}$$

The measure on  $\mathcal{T}_\Gamma^{<\omega}$  given by  $f \circ \mu^*$  will be denoted  $\mu^{<\omega}$ , i.e.  $\mu^{<\omega}(S) \stackrel{\text{def}}{=} \mu^*(f^{-1}(S))$  — notice that being bounded is a regular property and therefore the set  $f^{-1}(S)$  is  $\mu^*$ -measurable.

**Theorem 6.13.** *Given a regular language of finite trees  $L \subseteq \mathcal{T}_\Gamma^{<\omega}$ , one can effectively compute the value  $\mu^{<\omega}(L)$  in two-fold exponential space, and the value is algebraic.*

*Proof.* Consider the following set of full trees:

$$L' \stackrel{\text{def}}{=} \{t \in \mathcal{T}_\Gamma^\omega \mid \text{if } t \text{ is bounded then } f(t) \in L\}.$$

It is easy to observe that  $L'$  is a closed regular set of infinite trees. Moreover, given a non-deterministic tree automaton  $\mathcal{A}$  for  $L$  one can construct a non-deterministic safety automaton  $\mathcal{A}'$  for  $L'$ , and the size of  $\mathcal{A}'$  is polynomial in the size of  $\mathcal{A}$ . Since the set of bounded trees has measure 1, the  $\mu^*$ -measure of  $L'$  equals the  $\mu^{<\omega}$ -measure of  $L$ . Therefore, the result follows from Theorem 6.1.  $\square$

## 7. Complexity of computing the measure

In the previous sections, we have described the algorithms for computing the uniform measure of some classes of simple sets of infinite trees. The upper bounds that follow from the devised algorithms are expressed in Theorems 4.1, 5.1, and 6.1.

In the rest of this section, we derive lower bounds for the problem of positive measure. In the first two cases, i.e. for FO and BCCQ we will reduce directly from the problems of acceptance of Turing machines, while the lower bound in the case of non-deterministic safety automata will be more direct.

*Turing machines.* A Turing machine  $M$  is a tuple  $\langle \Gamma, Q, \delta, q_{\text{I}}, q_{\text{f}} \rangle$ , where  $\Gamma$  is an alphabet,  $Q$  is a finite set of states,  $q_{\text{I}} \in Q$  is an initial state,  $q_{\text{f}} \in Q$  is a final state, and  $\delta \subseteq Q \times \Gamma \times Q \times \Gamma \times \{-1, 0, 1\}$  is the transition relation. A run is a sequence of configurations; a configuration is a sequence of length  $2^n$  of cells; and each cell contains either:  $\text{b}$  denoting an empty cell,  $a \in \Gamma$  denoting a non-empty cell without the head, or a pair  $\langle q, a \rangle \in Q \times (\Gamma \cup \{\text{b}\})$  denoting a cell with the head of the machine over it. A run is *proper* if

- every configuration is *proper*, i.e. there is exactly one head, and after empty cells there are empty cells;
- the first configuration is *the initial configuration*, i.e. the first letter of the configuration is  $\langle q_{\text{I}}, \text{b} \rangle$ ;
- and *transitions between subsequent configurations are correct*, i.e. labels of cells without the head do not change in transition and the label of the cell with the head of the machine and the new position of the head change accordingly to the transition relation.

A run is accepting if it is proper and contains a configuration with the final state (a *final configuration*).

### 7.1. Complexity for FO without descendant

We begin by giving a lower bound in the case of first-order formulae without descendant.

**Lemma 7.1.** *The positive  $\mu^*(\text{FO})$  problem is EXPSPACE-hard.*

To prove this lemma we reduce the problem whether a given Turing machine accepts the empty input in exponential space.

*Proof.* The following problem is EXPSPACE-complete. Given a non-deterministic Turing machine  $M$  and a number  $n$  in unary, decide whether  $M$  accepts the empty input using at most  $2^n$  memory cells.

For a given Turing machine  $M$  and number  $n$  we will describe how to construct a first-order formula  $\phi$ , such that  $\mu^*(\text{L}(\phi)) > 0$  if and only if  $M$  accepts the empty input using at most  $2^n$  memory cells. The idea is that the trees that satisfy the formula  $\phi$  encode all accepting runs of  $M$  of the desired length. The size of the formula  $\phi$  will be polynomial in the value of  $n$  and the size of  $M$ . Before we proceed, recall that if  $M$  accepts the empty input

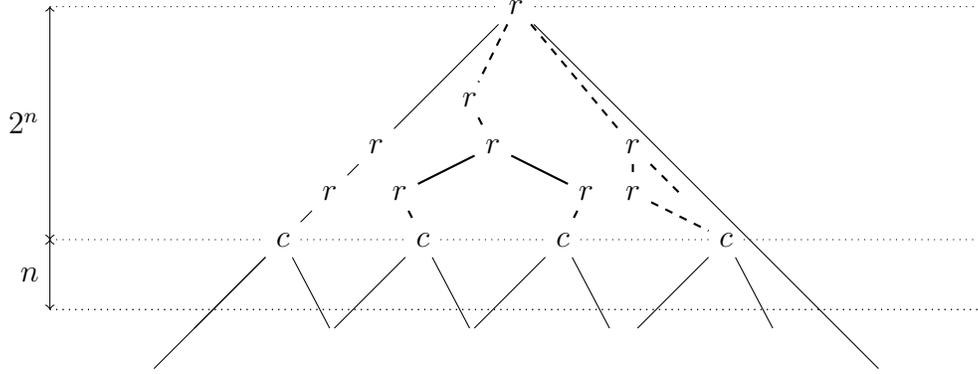


Figure 3: Encoding using first-order formulae. The bold solid lines denote some of the child relations, the dashed lines denote the induced ancestor relations. The  $c$ -labelled nodes are roots of encodings of configurations; the  $r$ -labelled nodes span the sequence of configurations. The “diamond” between second and third drawn configuration assures that they are consecutive.

using no more than  $2^n$  memory cells, then there exists an accepting run of  $M$  that uses no more than  $2^{2^n}$  steps.

A run of  $M$  will be encoded as a set of trees  $\mathbb{B}_t$ , where  $t$  is a full tree of height  $d = 2^n + n$ , such that every node up to the depth  $2^{n-1}$  is labelled with a fresh letter  $r \notin \Gamma$ , every node at depth  $2^n$  is labelled with a fresh letter  $c$ , every node between depths  $2^n + 1$  and  $2^n + n - 1$ , inclusive, is labelled with a fresh letter  $c'$ , and every node at depth  $2^n + n$  is labelled with a memory cell content. For an illustration of such a tree  $t$  see Figure 3.

Before we describe the encoding, let us recall a standard construction in which a polynomial formula can select two nodes that are in a distance that is exponential in the size of the formula.

$$\begin{aligned} \varphi_s^0(x, z) &\stackrel{\text{def}}{=} x \text{ s } z, \\ \varphi_s^{i+1}(x, z) &\stackrel{\text{def}}{=} \exists y. \forall x', y'. (x' = x \wedge y' = y) \vee (x' = y \wedge y' = z) \Rightarrow \varphi_s^i(x', y') \end{aligned} \quad (17)$$

Intuitively, the formula  $\varphi_s^n(x, z)$  states that  $x$  is an ancestor of  $z$  such that  $d(x, z) = 2^n$ . Let  $\varphi_s^{d=2^n}$  denote that formula. In a similar fashion we can define a formula  $\varphi_s^{d < n}(x, z)$  stating that  $x$  is an ancestor of  $z$  and  $d(x, z) < 2^n$ :

we simply replace the inductive construction with

$$\begin{aligned}\psi_s^0(x, z) &\stackrel{\text{def}}{=} x \text{ s } z \vee x = y, \\ \psi_s^{i+1}(x, z) &\stackrel{\text{def}}{=} \exists y. \forall x', y'. (x' = x \wedge y' = y) \vee (x' = y \wedge y' = z) \Rightarrow \psi_s^i(x', y')\end{aligned}\quad (18)$$

and write  $\varphi_s^{d < 2^n}(x, z) \stackrel{\text{def}}{=} \exists y. \psi_s^n(y, z) \wedge y s x \wedge \neg z s x$ . In the same manner we define  $\varphi_{s_L}^{d=2^n}$ ,  $\varphi_{s_L}^{d < 2^n}$  and  $\varphi_{s_R}^{d=2^n}$ ,  $\varphi_{s_R}^{d < 2^n}$ .

Now we can return to the encoding. The sub-trees of  $t$  in nodes labelled  $c$  denote configurations. Since every such sub-tree is of height  $n$ , every encoding has exactly  $2^n$  leafs encoding memory cells. Moreover, in a tree  $t$  there are exactly  $2^{2^n}$  sub-trees encoding configurations.

The only non-trivial part of the first-order formula  $\varphi_{frame}$  defining the above set of trees is the one talking about the nodes at depth  $2^n$ . This is done using  $\varphi_s^{d=2^n}(x, y)$  that expresses the fact that  $x \sqsubseteq y$  and  $d(x, y) = 2^n$ . Similarly, using  $\varphi_s^{d < 2^n}$  we can check the labels of the intermediate nodes and thus verify that a given tree has the shape as described above.

Now, all we need is to describe how to write a formula encoding an accepting run. We start with a formula  $\varphi_{proper}$  stating that every configuration is proper, i.e. there is exactly one head, and after empty cells there are empty cells. Then a formula  $\varphi_{first}$  defining initial configuration: it simply states that in the left-most configuration (defined using  $\varphi_{s_L}^{d=2^n}$ ) head is over the first memory cell and that the first memory cell is empty. Next is  $\varphi_{accepting}$  defining an accepting configuration: it simply states that there is node labelled by a letter from the set  $\{q_f\} \times (\Gamma \cup \{b\})$ .

Now we define a formula  $\varphi_{transition}$  stating that transitions between subsequent configurations are correct. To do that, we define  $\varphi_{next\_conf}(x, y)$ , stating that  $x$  and  $y$  are  $c$ -labelled roots of two consecutive configurations, and  $\varphi_{same\_cell}(x, y)$  stating that nodes  $x$  and  $y$  encode the same cell in the two consecutive configurations. More formally, the formula  $\varphi_{same\_cell}(x, y)$  states that there are two sequences of nodes  $x_0, \dots, x_n$  and  $y_0, \dots, y_n$  such that  $y_n = y, x_n = x, \varphi_{next\_conf}(x_0, y_0)$ , and for  $i > 0$   $x_i$  is a left child if and only if  $y_i$  is a left child. Given the formula  $\varphi_{same\_cell}(x, y)$ , the formula  $\varphi_{transition}$  simply states that the labels of  $x$  and  $y$  are consistent with the transition relation.

What remains is to define the formula  $\varphi_{next\_conf}(x, y)$ . It simply states that  $x$  and  $y$  are labelled  $c$  and that the path from  $x$  to  $y$  “forms a diamond”, i.e. that there are three nodes  $x_a, y_a$ , and  $z$  such that  $x_a$  is an ancestor of  $x$

reachable by  $s_R$  only,  $y_a$  is an ancestor of  $y$  reachable by  $s_L$  only, and  $x_a, y_a$  are respectively the left and the right child of  $z$ . To avoid speaking about the ancestors explicitly, we use the formulae  $\varphi_{s_L}^{d < 2^n}$  and  $\varphi_{s_R}^{d < 2^n}$ .

The resulting formula is the conjunction of the formulae  $\varphi_{frame}, \varphi_{proper}, \varphi_{first}, \varphi_{accepting}, \varphi_{transition}$ .  $\square$

## 7.2. Complexity for conjunctive queries

In the case of a conjunctive queries, we obtain exact bounds. More precisely, we observe that deciding whether the measure of a language defined by a single conjunctive query has a positive measure is *NP*-complete.

**Proposition 7.2.** *The positive  $\mu^*(CQ)$  problem is NP-complete.*

*Proof.* Let  $q$  be a conjunctive query. Then, either  $q$  is not satisfiable and  $\mu^*(q) = 0$ , or  $q$  is satisfiable and has a positive measure. That is,  $\mu^*(q) > 0$  if and only if  $q$  is satisfiable. Deciding whether a conjunctive query is satisfiable is *NP*-complete, cf. e.g. [22].  $\square$

If we allow Boolean combinations of conjunctive queries the complexity increases exponentially.

**Lemma 7.3.** *The positive  $\mu^*(BCCQ)$  problem is NEXP-hard.*

The proof is a very simple adaptation of the reduction in [21], which shows that the satisfiability of Boolean combinations of conjunctive queries with respect to recursion-free *DTDs* is *NEXP*-hard. The two differences are the lack of *DTD*, which is replaced by a Boolean combination of patterns, and the forced use of the root patterns.

*Proof.* To show the lower bound we reduce the following *NEXP*-complete problem: given a non-deterministic Turing machine  $M$  and a number  $n$  in unary, decide whether  $M$  accepts the empty input in at most  $2^n$  steps.

More precisely, for a Turing machine  $M$  and number  $n$  we will describe how to construct a Boolean combination of rooted, firm conjunctive queries  $\phi$ , such that  $\mu^*(L(\phi)) > 0$  if and only if  $M$  accepts the empty input in at most  $2^n$  steps. The idea is that the trees that satisfy the formula  $\phi$  encode all accepting runs of  $M$  of the desired length.

The construction will be similar to the one in the used in the proof of Lemma 7.1 with appropriate changes to the formulae  $\varphi_{frame}, \varphi_{proper}, \varphi_{first}, \varphi_{accepting}, \varphi_{transition}$ .

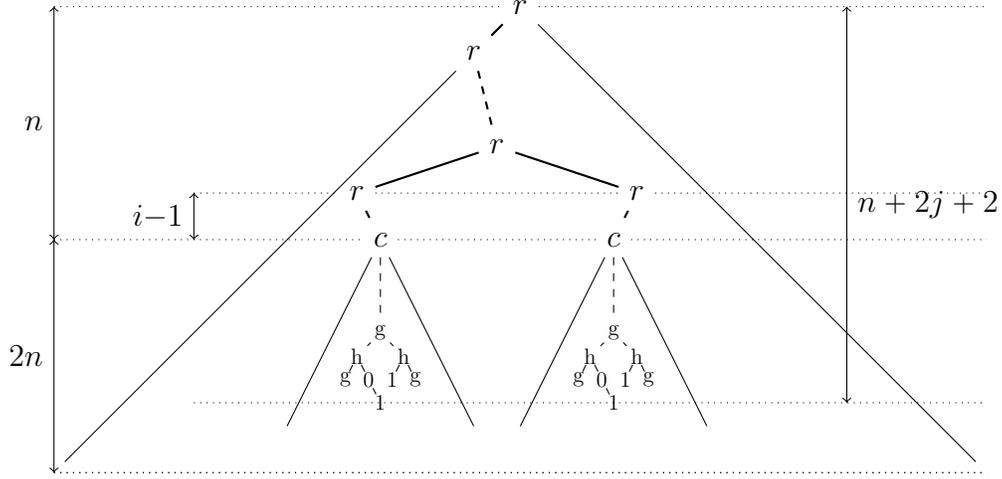


Figure 4: The structure of the encoding using conjunctive queries. The figure depicts two consecutive configurations, with an ancestor  $i$  levels above. The formula  $\varphi_{\text{same\_turn}_j}(x, y)$  chooses either both left  $h$ -labelled nodes or both right  $h$ -labelled nodes.

Since a run should accept in  $2^n$  steps, the machine will not use more than  $2^n$  memory cells, thus every configuration under consideration has length at most  $2^n$ .

A run will be encoded as follows, see Figure 4. Let  $t'$  be a full tree of height  $d + 1$ . The root of  $t'$  is labelled  $r \notin \Gamma$  and the tree  $t'.\mathbb{R}$  has every node labelled  $w \notin \Gamma$ . Now we describe the tree  $t = t'.\mathbb{L}$ . The root of  $t$  is labelled  $r$ . And every node of  $t$  at depth  $l$  such that  $1 \leq l \leq n - 1$  is labelled  $r$ . Every node of  $t$  at depth  $n$  is labelled  $c \notin \Gamma$ . The trees rooted in the  $c$ -labelled nodes encode configurations. Children of the  $c$ -labelled nodes are labelled with  $g$ . Every  $g$ -labelled node  $u$  at depth at most  $d - 3$  has  $h$ -labelled children ( $h \notin \Gamma$ ). The nodes  $u_{\text{LL}}$  and  $u_{\text{RR}}$  are labelled  $g$ ; the node  $u_{\text{LR}}$  is labelled  $0$ ; the node  $u_{\text{RL}}$  is labelled  $1$ ; and the node  $u_{\text{RR}}$  is labelled  $1$ . Nodes  $u$  at depth  $d - 2$  labelled  $g$  have similar setup with the difference that the nodes  $u_{\text{LL}}$  and  $u_{\text{RR}}$  are labelled by a letter from the alphabet  $\Gamma^t$ : they encode the memory cells. The unspecified nodes are labelled with  $w$ .

Notice that for  $d = 3n$ , we have  $2^n$  configurations, nodes labelled  $c$ , with  $2^n$  memory cells each, and that all those requirements can be easily forced by a polynomial in size Boolean combination of patterns. Indeed, the above frame can be expressed by a Boolean combination  $\varphi_{\text{frame}}$ , in which the used patterns forbid ill-labelled children and enforce good labelling of the roots.

Thus, all that is left to define is a Boolean combination of patterns that will enforce that a tree as above encodes an accepting run. We will do this by stating that the configurations are proper ( $\varphi_{proper}$ ); that there is a configuration with an accepting state  $\varphi_{accepting}$ ; that the first configuration is the initial configuration ( $\varphi_{first}$ ); and that each transition between two consecutive configurations is consistent with the transition relation  $\delta$  ( $\varphi_{transition}$ ). The first two requirements are easily expressible, it is also easy to write a formula saying that every configuration is proper, i.e. that there is only one head and after blanks there are blanks.

Now, we describe how to express that the transitions are proper. We start with family of formulae  $\varphi_{next\_conf_i}(x, y)$ ,  $i \in \{1, 2, \dots, n\}$ , where  $\varphi_{next\_conf_i}(x, y)$  holds for two nodes that are roots of two consecutive configurations with a common ancestor  $i$  levels above. The formula  $\varphi_{next\_conf_i}(x, y)$  says that  $x$  and  $y$  are  $c$ -labelled and there are two sequences  $x_1, \dots, x_{i-1}, x$  and  $y_1, \dots, y_{i-1}, y$ , and a node  $z$  such that  $x_1$  is the left child of  $z$ ,  $y_1$  is the right child of  $z$ ,  $x_1, \dots, x_{i-1}, x$  are right-child connected  $y_1, \dots, y_{i-1}, y$  are left-child connected, they “form a diamond” as on Figure 4. In a similar way, we can define  $\varphi_{next\_cell_i}(x, y)$  choosing nodes that are two consecutive memory cells in the same configuration.

Now, we will describe how to construct a family of formulae  $\varphi_{same\_cell_i}(x, y)$ ,  $j \in \{1, 2, \dots, n\}$ , that chooses the same memory cell in two consecutive  $\varphi_{next\_conf_i}$  configurations. We start by observing that if  $x_c$  and  $y_c$  are the  $c$ -labelled nodes of the appropriate configurations then  $x$  and  $y$  are the same memory cell if at each level of the path from  $x_c$  to  $x$  we choose the same turns as on the appropriate levels of the path from  $y_c$  to  $y$ . In other words, if  $x_c, x_1, \dots, x$  is child connected path from  $x_c$  to  $x$  and  $y_c, y_1, \dots, y$  is child connected path from  $y_c$  to  $y$ , then  $x_i$  is the left child of its parent if and only if  $y_i$  is the left child of its parent. This can be achieved by a family of formulae  $\varphi_{same\_turn_j}(x, y)$ .

The formula  $\varphi_{same\_turn_j}(x, y)$  states that  $x$  and  $y$  are at depth  $2j + n$ , are  $h$ -labelled, and are both left or both right children. We achieve the last one by stating that there are two sequences of nodes  $x_1, \dots, x_{n+2j+1}$  and  $y_1, \dots, y_{n+2j+1}$  such that for  $j = 1, \dots, n + 2j + 1$  we have that  $x_j s x_{j+1}, y_j s y_{j+1}$ ,  $x$  is an ancestor of  $x_{n+2j+1}$ ,  $y$  is an ancestor of  $y_{n+2j+1}$ ,  $x_1 = y_1$  and both  $y_{n+2j+1}, x_{n+2j+1}$  are labelled with 1. Note, that  $\varphi_{same\_turn_i}(x, y)$  if and only if  $x$  and  $y$  are on the same depth and both are left or both are right children.

Finally, we define  $\varphi_{same\_cell_i}(x, y)$  by stating that  $x$  and  $y$  are at depth  $d$  and the paths from their respective  $c$ -labelled ancestors take the same turns.

With  $\varphi_{\text{same\_cell}_i}(x, y)$  defined, for every transition  $\sigma$  that is not in  $\delta$  we can write a pattern  $\varphi_\sigma$  stating that a transition between two configurations is incorrect. The pattern states that for some  $i, j, k$  there are two nodes  $x, y$  such that  $\varphi_{\text{same\_cell}_i}(x, y)$  and there also are two nodes  $y', y''$  such that  $\varphi_{\text{next\_cell}_j}(y, y')$ ,  $\varphi_{\text{next\_cell}_k}(y', y'')$  and the labels of  $x, y, y', y''$  encode  $\sigma$ . Hence,  $\varphi_{\text{transition}}$  is simply the conjunction of the negations of the formulae  $\varphi_\sigma$ .

The final formula is the conjunction of the formulae  $\varphi_{\text{frame}}, \varphi_{\text{proper}}, \varphi_{\text{first}}, \varphi_{\text{accepting}}, \varphi_{\text{transition}}$ , modified so that the used patterns are firm. We do that by quantifying every variable by formula  $\varphi_{\text{limited\_depth}_d}(x)$  stating that there is a node  $y$  at depth  $d + 1$  such that  $x$  is an ancestor of  $y$ .  $\square$

Since the lower bound from Lemma 7.3 matches the upper bound obtained in Proposition 5.8 we infer the following.

**Theorem 7.4.** *The positive  $\mu^*(\text{BCCQ})$  problem is NEXP-complete.*

*Counting complexity.* We will now provide a short explanation how the above results can be read in terms of counting complexity classes. However, as the main scope of the article is expressed in terms of classical decision procedures, the provided explanation is brief, only indicating how this approach can be taken.

A careful inspection of the proof of hardness, i.e. the proof of Lemma 7.3, reveals that the described reduction in fact reduces the problem of computing the number of runs of a non-deterministic Turing machine running in exponential time to the problem of computing the standard measure of a language defined by a Boolean combination of conjunctive queries. The crucial observation is that under the above defined reduction, every accepting run of the machine defines an unique prefix of fixed depth, say  $k$ . Moreover, each such run corresponds to exactly the set of trees that conform to that prefix. Since for every run the size of the corresponding prefix is  $2^k$ , there exists a computable constant  $c$  such that the number of accepting runs, say  $a$ , and the measure, say  $\mu$ , are related by the equation  $c \cdot a = \mu$ . Thus, it shows that the quantitative  $\mu^*(\text{BCCQ})$  problem is  $\#EXP$ -hard.

On the other hand, a careful analysis of the proofs of Theorem 5.1 and Corollary 5.7 allow us to claim that computing the measure belongs to the class  $\#EXP$ . Let  $k$  be the size of the input query. A simple algorithm guesses a full binary tree of height  $k$  that is the prefix witnessing the satisfiability of the root patterns and then verifies it in deterministic exponential time. Since

the size of every such prefix is exponential in  $k$  and the prefixes define a partition into sets of equal measure, instead of counting the prefixes, it is enough to count the accepting runs.

Thus, we state the following.

**Proposition 7.5.** *The problem of computing the measure of BCCQ (i.e. Boolean combinations of conjunctive queries) is  $\sharp$ EXP-complete.*

Notice that the arguments above use crucially the idea that in the case of conjunctive queries the measure is determined by the close neighbourhood of the root. This allows us to simply enumerate all the prefixes of some fixed depth and count the positive occurrences. A similar observation can be made in the case of first-order formulae not using descendant relation, but most likely does not extend to the case of first-order formulae using descendant nor to the case of safety automata. Since both classes admit languages of irrational uniform measure and every prefix of fixed depth is of rational measure, simple counting is not enough.

### 7.3. Complexity for safety automata

In this section we provide a lower bound for the problem of computing the measure of a language given by a non-deterministic safety automaton. In that case, contrary to the cases of FO and BCCQ, the given model is asymmetric — the considered model is not closed under complement. Therefore, the problems of positive measure and full measure (i.e. whether  $\mu^*(L) = 1$ ) are not obviously inter-reducible. Therefore, to obtain a better lower bound we will focus on the later problem.

**Proposition 7.6.** *The problem whether  $\mu^*(L(\mathcal{A})) = 1$  for a non-deterministic safety automaton  $\mathcal{A}$  is EXP-hard.*

*Proof.* The result follows directly from EXP-hardness of the problem of *universality* of non-deterministic automata over finite trees (i.e. whether  $L(\mathcal{A}) = \mathcal{T}_\Gamma^{<\omega}$ ), see e.g. [33, Theorem 1.7.7], and the reduction provided in Theorem 6.13.  $\square$

## 8. Conclusions and future work

We have shown that there exists an algorithm that, given a first-order sentence  $\varphi$  over the signature  $\Gamma \cup \{\varepsilon, s_R, s_L, s\}$ , computes  $\mu^*(L(\varphi))$  in three-fold

exponential space. We also have shown that there exists an algorithm that, given a Boolean combination of conjunctive queries  $\varphi$  over the signature  $\Gamma \cup \{\varepsilon, s_R, s_L, s, \sqsubseteq\}$ , computes the uniform measure  $\mu^*(L(\varphi))$  in exponential space. Both these algorithms base on the fact that the measure of the considered language coincides with the measure of an effectively computable clopen regular set of trees (see Remark 4.1 and the proof of Theorem 5.1). Such languages are always of computable rational measure and thus the results follow. As witnessed by Proposition 4.11 this technique cannot work for the full first-order logic, because it can define languages of irrational measures.

Additionally to the above results, we provide an algorithm for computing the measure of a language given by a non-deterministic safety automaton. The approach in that case must be different, as witnessed by Remark 6.3. The provided algorithm reduces the problem to a first-order formula over the reals and then invokes the celebrated Tarski's quantifier elimination procedure. This result can be adjusted to the case of regular languages of finite trees, as discussed in Section 6.4.

We have also studied the computational complexity of the problems of positive measure. The upper and lower bounds match in the case of conjunctive queries (*NP*-complete) and Boolean combinations of conjunctive queries (*NEXP*-complete). The bounds in the case of first-order logic without descendant do not match: the problem is *EXPSPACE*-hard while the provided algorithm runs in three-fold exponential space. Similarly, the bounds for non-deterministic safety automata do not match. We think that establishing the exact bounds of the problems poses an interesting direction of future research.

The substantial gap in the case of first-order formulae stems from two reasons. The first reason is applicable not only in the case of first-order formulae and can be stated as follows: the provided algorithms solve the computational quantitative problem, i.e. *compute the measure*, while the provided reductions utilise decision versions of the qualitative problem, i.e. *decide whether the measure is 0* or *decide whether the measure is 1*. This is especially evident in the case of Boolean combinations of conjunctive queries, where the problem of deciding positive measure is *NEXP*-complete, computing the measure is  $\sharp$ *EXP*-complete, and we provide an algorithm that computes the measure in exponential space. To further support this observation, note that in the cases where we have obtained matching complexity bounds we have provided specific algorithms to solve the positive measure problems.

The second reason is more algorithm specific. The proposed algorithm

uses the Gaifman locality to compute a new formula in Gaifman normal form. This results in necessarily three-fold exponential blow-up in the size of the formula, see [17]. On the other hand, the transformation does not alter the bound on the height of the counted prefixes, as it remains exponential in the size of the original formula. Notice that our reduction exploits the blow-up in the prefix size, but not in the size of the new formula. We conjecture that the actual complexity of the problem should match the lower bound, as we think that it is possible to circumvent the blow-up in the formula size by performing the translation on-line and computing only the necessary parts of the normal form on demand.

The results of this paper are expressed in terms of trees with binary branching. The results directly extend to trees of other fixed branching, i.e. ternary trees. The situation is a bit more subtle with  $\omega$ -branching trees: although they can be naturally interpreted into binary trees, the interpretation is not definable in all the considered formalisms (e.g. in FO with successor). Even worse, it seems that in the case of unbounded finitely branching trees there is no natural definition of an uniform measure: intuitively, every measure on such trees has to be effectively discounted, in the sense that having more children has to be less probable than having less. Due to these difficulties, we decided to focus on the binary formulations. If there is a need of translating them to some other model, one should carefully check whether it is possible to interpret that model within binary branching trees in a way definable in the considered logic.

Note that the considered measure respects a form of a 0–1-law. By Lemma 3.1, if  $t$  is a finite tree then with probability 1 it appears as a sub-tree in a random tree. It would be interesting to extend the enquiry to measures that do not possess such a property. Such measures can be expressed, for example, by graphs or by branching boards, cf. [10].

Obviously, the most interesting problem is to find an algorithm that can compute the uniform measure of an arbitrary regular language of infinite trees. While we know that languages with irrational measures exist, we conjecture that for any regular language of trees  $L$  the uniform measure  $\mu^*(L)$  is algebraic.

## References

- [1] K. Chatterjee, T. A. Henzinger, A survey of stochastic  $\omega$ -regular games, *J. Comput. Syst. Sci.* 78 (2) (2012) 394–413.

- [2] D. Suciu, D. Olteanu, C. Ré, C. Koch, Probabilistic Databases, Synthesis Lectures on Data Management, Morgan & Claypool Publishers, 2011.
- [3] T. Chen, K. Dräger, S. Kiefer, Model checking stochastic branching processes, in: MFCS, 2012, pp. 271–282.
- [4] D. Gale, F. M. Stewart, Infinite games with perfect information, in: Contributions to the theory of games, Vol. 2 of Annals of Mathematics Studies, no. 28, Princeton University Press, 1953, pp. 245–266.
- [5] M. Mio, Probabilistic modal  $\mu$ -calculus with independent product, Logical Methods in Computer Science 8 (4) (2012) 1–36.
- [6] T. Gogacz, H. Michalewski, M. Mio, M. Skrzypczak, Measure properties of regular sets of trees, Inf. Comput. 256 (2017) 108–130.
- [7] H. Michalewski, M. Mio, On the problem of computing the probability of regular sets of trees, in: FSTTCS, 2015, pp. 489–502.
- [8] A. Amarilli, P. Bourhis, P. Senellart, Provenance circuits for trees and treelike instances, in: ICALP, 2015, pp. 56–68.
- [9] L. Staiger, The Hausdorff measure of regular  $\omega$ -languages is computable, Bulletin of the EATCS 66 (1998) 178–182.
- [10] M. Przybyłko, M. Skrzypczak, On the complexity of branching games with regular conditions, in: MFCS, Vol. 58, 2016, pp. 78:1–78:14.
- [11] A. Amarilli, M. Monet, P. Senellart, Conjunctive queries on probabilistic graphs: Combined complexity, in: PODS, 2017, pp. 217–232.
- [12] M. Przybyłko, Tree games with regular objectives, in: GandALF, 2014, pp. 231–244.
- [13] M. Przybyłko, Stochastic games and their complexity, Ph.D. thesis, University of Warsaw (2019).
- [14] W. Thomas, Languages, automata, and logic, in: Handbook of Formal Languages, Springer, 1996, pp. 389–455.

- [15] A. Kechris, *Classical descriptive set theory*, Springer-Verlag, New York, 1995.
- [16] L. Heimberg, D. Kuske, N. Schweikardt, An optimal Gaifman normal form construction for structures of bounded degree, in: *LICS*, 2013, pp. 63–72.
- [17] A. Dawar, M. Grohe, S. Kreutzer, N. Schweikardt, Model theory makes formulas large, in: *ICALP*, 2007, pp. 913–924.
- [18] L. Bunt, P. Jones, J. Bedient, *The Historical Roots of Elementary Mathematics*, Dover Books Explaining Science, Dover Publications, 1988.
- [19] A. Potthoff, *Logische Klassifizierung regulärer Baumsprachen*, Ph.D. thesis, Universität Kiel (1994).
- [20] M. Bojańczyk, Tree-walking automata, in: *LATA*, 2008, pp. 1–2.
- [21] F. Murlak, M. Ogiński, M. Przybyłko, Between tree patterns and conjunctive queries: Is there tractability beyond acyclicity?, in: *MFCS*, 2012, pp. 705–717.
- [22] H. Bjørklund, W. Martens, T. Schwentick, Conjunctive query containment over trees, *Journal of Computer and System Sciences* 77 (3) (2011) 450–472.
- [23] I. Walukiewicz, Deciding low levels of tree-automata hierarchy, *Electr. Notes Theor. Comput. Sci.* 67 (2002) 61–75.
- [24] F. Cavallari, *Regular tree languages in the first two levels of the Borel hierarchy*, Ph.D. thesis, HEC Lausanne (UNIL) and University of Turin (2018).
- [25] D. Janin, G. Lenzi, On the logical definability of topologically closed recognizable languages of infinite trees, *Computers and Artificial Intelligence* 21 (3) (2002) 1–19.
- [26] C. Jones, G. Plotkin, Probabilistic powerdomain of evaluations, in: *LICS*, 1989, pp. 186–195.
- [27] A. Tarski, *A Decision Method for Elementary Algebra and Geometry*, University of California Press, 1951.

- [28] G. E. Collins, Hauptvortrag: Quantifier elimination for real closed fields by cylindrical algebraic decomposition, in: Automata Theory and Formal Languages, 1975, pp. 134–183.
- [29] J. Bochnak, M. Coste, M.-F. Roy, Real Algebraic Geometry, Vol. 36 of A Series of Modern Surveys in Mathematics, Springer-Verlag Berlin Heidelberg, 1998.
- [30] M. Ben-Or, D. Kozen, J. Reif, The complexity of elementary algebra and geometry, Journal of Computer and System Sciences 32 (2) (1986) 251–264.
- [31] T. E. Harris, The Theory of Branching Processes, Springer–Verlag Berlin Heidelberg, 1963.
- [32] K. Etessami, M. Yannakakis, Recursive Markov decision processes and recursive stochastic games, J. ACM 62 (2) (2015) 11:1–11:69.
- [33] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, M. Tommasi, Tree automata techniques and applications, Available on: <http://www.grappa.univ-lille3.fr/tata>, release October, 12th 2007 (2007).