

BIROn - Birkbeck Institutional Research Online

Kikot, Stanislav and Kurucz, Agi and Podolskii, Vladimir V. and Zakharyashev, Michael (2021) Deciding Boundedness of Monadic Sirups. In: PODS'21: Proceedings of the 40th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, Virtual Event, China, June 20-25, 2021.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/44865/>

Usage Guidelines:

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html>
contact lib-eprints@bbk.ac.uk.

or alternatively

Deciding Boundedness of Monadic Sirups

Stanislav Kikot

Institute for Information Transmission Problems
Moscow, Russia
staskikotx@gmail.com

Vladimir V. Podolskii

HSE University
Moscow, Russia
vpodolskii@hse.ru

Agi Kurucz

King's College London
London, UK
agi.kurucz@kcl.ac.uk

Michael Zakharyashev

Birkbeck, University of London, UK &
HSE University, Moscow, Russia
michael@dcs.bbk.ac.uk

ABSTRACT

We show that deciding boundedness (aka FO-rewritability) of monadic single rule datalog programs (sirups) is 2EXPTIME-hard, which matches the upper bound known since 1988 and finally settles a long-standing open problem. We obtain this result as a byproduct of an attempt to classify monadic ‘disjunctive sirups’—Boolean conjunctive queries q with unary and binary predicates mediated by a disjunctive rule $T(x) \vee F(x) \leftarrow A(x)$ —according to the data complexity of their evaluation. Apart from establishing that deciding FO-rewritability of disjunctive sirups with a dag-shaped q is also 2EXPTIME-hard, we make substantial progress towards obtaining a complete FO/L-hardness dichotomy of disjunctive sirups with ditree-shaped q .

CCS CONCEPTS

• **Information systems** → **Query languages**; • **Theory of computation** → **Complexity theory and logic**; **Description logics**; • **Computing methodologies** → **Knowledge representation and reasoning**.

KEYWORDS

Boundedness; monadic datalog; first-order rewritability; ontology-mediated query.

ACM Reference Format:

Stanislav Kikot, Agi Kurucz, Vladimir V. Podolskii, and Michael Zakharyashev. 2021. Deciding Boundedness of Monadic Sirups. In *Proceedings of the 40th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS '21)*, June 20–25, 2021, Virtual Event, China. ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3452021.3458332>

1 INTRODUCTION

There have been two waves in the investigation of boundedness or first-order rewritability of various types of recursive queries. The first one started in the mid 1980s, when the deductive database

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PODS '21, June 20–25, 2021, Virtual Event, China

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-8381-3/21/06...\$15.00
<https://doi.org/10.1145/3452021.3458332>

community was analysing recursion in datalog queries with the aim of optimising and parallelising their execution. One of the fundamental issues was the problem of deciding whether the depth of recursion required to evaluate a given datalog query could be bounded independently of the input data. By 2000, among other remarkable results, it had been discovered that

- boundedness of linear datalog queries with binary predicates and of ternary linear datalog queries with a single recursive rule is undecidable [26, 33];
- deciding program boundedness is 2EXPTIME-complete for monadic programs [9, 18], PSPACE-complete for linear monadic programs [18, 36], and NP-complete for linear monadic and dyadic single rule programs [37].

Interestingly, the exact complexity of deciding boundedness of monadic datalog programs with a single recursive rule, known as *sirups* since [19], has remained open so far, somewhere between NP and 2EXPTIME, to be more precise. To clarify the ‘status [of boundedness] for sirups’ is part of Open Problem 4.2.10 in [28]. According to [4], Kanellakis and Papadimitriou, who were interested in datalog programs computable in NC, and so parallelisable, ‘have investigated the case of unary sirups, and have made progress towards a complete characterization.’ Alas, that work appears to have never been completed and published.

In this paper, we finally settle the boundedness problem for monadic sirups by showing that it is 2EXPTIME-hard, which matches the upper bound for deciding boundedness of arbitrary monadic datalog programs [18] (and which should be compared with the NP–PSPACE gap between deciding boundedness of *linear* sirups and non-sirups.)

We obtained this result while surfing the second wave, which was triggered in the mid 2010s by the theory and practice of ontology-based data access (OBDA) [15, 35, 38] (recently rebranded to virtual knowledge graphs [39]). In OBDA, a typical ontology-mediated query (OMQ) takes the form $Q = (O, q)$ with a description logic (DL) ontology O and a conjunctive query (CQ) q . A fundamental problem in this setting is to decide whether a given OMQ Q is FO-rewritable, in which case finding certain answers to Q can be done by evaluating a non-recursive SQL-query using a standard RDBMS.

The ontology language *OWL 2 QL* for OBDA systems (such as Mastro¹ or Ontop²), standardised by the W3C in 2009, is based on

¹<https://www.obdasystems.com>

²<https://ontopic.biz>

DL-Lite that uniformly guarantees FO-rewritability of all OMQs with an *OWL 2 QL* ontology. Uniformly FO-rewritable tgds, aka Datalog[±] or existential rules, have also been identified; see, e.g., [17, 23, 30]. As an inevitable consequence, however, all of these ontology languages are very inexpressive.

The FO-rewritability problem for OMQs in more expressive ontology languages was attacked in [11] via a reduction to CSPs. It has been shown, among other results, that

- deciding FO-rewritability of OMQs with ontologies in expressive DLs such as \mathcal{ALC} (notational variant of multi-modal logic \mathbf{K}_n) and atomic CQs is NEXPTIME-complete [11], which becomes 2NEXPTIME-complete in the case of (non-atomic) CQs and also monadic disjunctive datalog queries [14, 20];
- any OMQ with a (Horn) \mathcal{EL} ontology and a CQ is either FO-, or linear-datalog-, or datalog-rewritable, and deciding this trichotomy is EXPTIME-complete [31, 32]; see also [6, 10] for complexity results on deciding FO-rewritability of OMQs with more expressive Horn description logic ontologies and frontier-guarded existential rules.

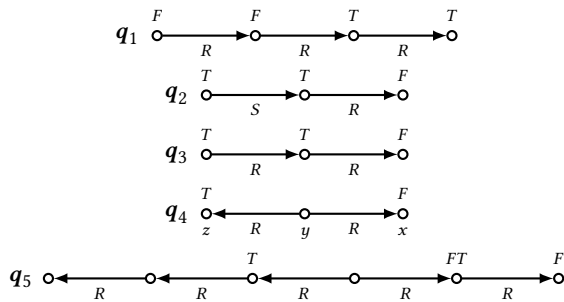
In [22], aiming to single out and classify possible causes of non-FO- or non-(linear)-datalog-rewritability of OMQs, we considered (in the DL setting) a disjunctive analogue of monadic sirups, namely, monadic disjunctive datalog programs Δ_q of the form

$$T(x) \vee F(x) \leftarrow A(x) \quad (1)$$

$$G \leftarrow q \quad (2)$$

where q is a (Boolean) CQ with unary predicates $T(x)$, $F(y)$ and arbitrary binary predicates, and G is a nullary (goal) predicate. In DL and conceptual modelling, rule (1) is known as a *covering axiom* (or constraint) $A \sqsubseteq T \sqcup F$ (as in ‘class Animal is covered by classes Male and Female’). We illustrate the zoo of ‘monadic disjunctive sirups’ by an example, where CQs are given as digraphs with labelled edges and (partially) labelled nodes.

EXAMPLE 1. Consider the CQs q_1, \dots, q_5 shown below:



For instance, in full, rule (2) in the program Δ_{q_4} looks as

$$G \leftarrow F(x), R(y, x), R(y, z), T(z).$$

Intuitively, the certain answer to the Boolean query (Δ_{q_4}, G) over a data instance \mathcal{D} (given in the form of a labelled graph) is ‘yes’ iff we can find the pattern q_4 in every graph obtained by labelling each of the A -nodes in \mathcal{D} with either T or F . As shown in [22], answering (Δ_{q_i}, G) is coNP-complete for q_1 , P-complete for q_2 , NL-complete for q_3 , L-complete for q_4 , and, in view of Example 4 below, q_5 is FO-rewritable and so in AC⁰.

Every disjunctive sirup Δ_q , in which q has a single ‘solitary’ F -node (like in q_2 – q_5), is equivalent to a monadic datalog program Π_q . For instance, Δ_{q_4} is equivalent to Π_{q_4} with three rules

$$G \leftarrow F(x), R(y, x), R(y, z), P(z)$$

$$P(x) \leftarrow T(x)$$

$$P(x) \leftarrow A(x), R(y, x), R(y, z), P(z)$$

Furthermore, for certain CQs q , boundedness of Π_q coincides with boundedness of a sirup sub-program of Π_q (see Sec. 2). In the above example, this sirup, Σ_{q_4} , comprises the last two rules of Π_{q_4} , and neither (Δ_{q_4}, G) nor (Σ_{q_4}, P) is FO-rewritable.

On the other hand, every disjunctive sirup Δ_q can be encoded as a CQ mediated by a Schema.org³ ontology. Deciding FO-rewritability of UCQs mediated by Schema.org is known to be PSPACE-hard [24].

Our first result in this paper establishes 2EXPTIME-hardness of deciding FO-rewritability in all of these cases. In Sec. 3, we show how a computation of an alternating Turing machine can be captured in terms of boundedness of the disjunctive sirup Δ_q , datalog program Π_q or its sirup sub-program Σ_q , for some CQ q . Compared to known techniques, which require multiple rules in a program or a union of multiple CQs to check properties of Turing machine computations, we achieve the same aim by means of polynomially-many small Boolean circuits that are ‘implemented’ by a *single* CQ q and check local properties of binary trees representing the expansions of Π_q .

What causes such high computational costs of recognising FO-rewritability of seemingly very primitive programs? Are there any natural classes of monadic (disjunctive) sirups whose boundedness can be checked by tractable algorithms? The 2EXPTIME-hardness proof provides three clues: first, the CQs q used in it are dags; second, each of them has two T -nodes; and, third, they contain many twin FT -nodes (as in q_5 above). In [22], we gave a complete classification of monadic disjunctive sirups Δ_q with a path CQ q and an extra disjointness constraint

$$\perp \leftarrow T(x), F(x) \quad (3)$$

(as in ‘classes Male and Female are disjoint’) according to their data complexity (AC⁰/NL/P/coNP) and rewritability type (FO/linear datalog/datalog/disjunctive datalog).

Here, in Sec. 4, we make significant progress towards a complete understanding of FO-rewritability of disjunctive sirups Δ_q with a ditree-shaped CQ q . First, we prove that twin-free CQs q as well as those that contain comparable (w.r.t. the tree order in q) solitary F - and T -nodes (like in q_1 – q_3 but not q_4 and q_5) give rise to NL-hard disjunctive sirups. In particular, this yields a tractable FO/NL-hardness dichotomy of the ditree disjunctive sirups with disjointness (3). Second, we obtain a tractable FO/L/NL-completeness trichotomy of ditree disjunctive sirups with one solitary F , one solitary T and any number of FT -twins. (This case corresponds to linear ditree sirups.) Finally, we establish an FO/L-hardness dichotomy for ditree disjunctive sirups with one solitary F and show that this dichotomy can be decided in polynomial time if the number of solitary T s in the CQs is bounded (like in our 2EXPTIME-hardness

³<https://schema.org>: ‘Many applications from Google, Microsoft, Pinterest, Yandex and others already use these vocabularies to power rich, extensible experiences’.

proof) and in exponential time otherwise. It follows that deciding FO-rewritability of such disjunctive sirups is fixed-parameter tractable if the number of solitary T s is regarded as a parameter.

The omitted proofs can be found in the full version [29].

2 PRELIMINARIES

We remind the reader (who can consult [2] for details) that a datalog program is a finite set, Π , of rules of the form

$$\forall \mathbf{x} (\gamma_0 \leftarrow \gamma_1 \wedge \dots \wedge \gamma_m) \quad (4)$$

where each γ_i is a (constant- and function-free) atom $Q(\mathbf{y})$ with $\mathbf{y} \subseteq \mathbf{x}$. As usual, we omit $\forall \mathbf{x}$ and replace \wedge with a comma. The atom γ_0 is the *head* of the rule, and $\gamma_1, \dots, \gamma_m$ comprise its *body*. The variables in the head must also occur in the body. The predicate in the head of a rule in Π is called an *IDB predicate*; non-IDB predicates in Π are *EDB predicates*. We call a rule *recursive* if its body has at least one IDB predicate; otherwise, it is an *initialisation rule*. The *arity* of Π is the maximum arity of its IDB predicates. Here, we only consider *monadic* datalog programs with at most *binary* EDBs. A *monadic sirup* is a monadic program with a single recursive rule.

A *data instance* for Π is any finite set \mathcal{D} of ground atoms with EDB predicates in Π . The set of constants in \mathcal{D} is denoted by $\text{ind}(\mathcal{D})$. For a unary IDB predicate P , a *certain answer* to the *datalog query* (Π, P) over \mathcal{D} is any $a \in \text{ind}(\mathcal{D})$ such that $\mathcal{I} \models P[a]$, for every model \mathcal{I} of Π and \mathcal{D} , or, in other words, $P(a)$ is in the closure $\Pi(\mathcal{D})$ of \mathcal{D} under the rules in Π . For a 0-ary IDB G (goal), a *certain answer* to (Π, G) over \mathcal{D} is ‘yes’ if $G \in \Pi(\mathcal{D})$, and ‘no’ otherwise.

A typical monadic datalog program, Π_q , we deal with in this paper is associated with a *conjunctive query* (CQ) q , which in our context is just a set of atoms with unary predicates F, T and arbitrary binary predicates. An atom $F(z) \in q$ is *solitary* if $T(z) \notin q$, and symmetrically for $T(z)$; a pair $T(z), F(z) \in q$ is referred to as *twins*.

For a CQ q with a single solitary $F(x)$, possibly multiple solitary $T(y_1), \dots, T(y_n)$, arbitrary twins $T(z), F(z)$ and binary atoms, the program Π_q comprises the following rules with 0-ary goal G :

$$G \leftarrow F(x), q^-, P(y_1), \dots, P(y_n) \quad (5)$$

$$P(x) \leftarrow T(x) \quad (6)$$

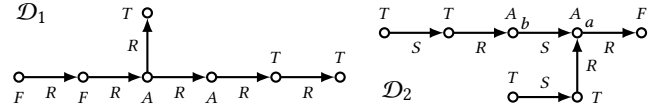
$$P(x) \leftarrow A(x), q^-, P(y_1), \dots, P(y_n) \quad (7)$$

where $q^- = q \setminus \{F(x), T(y_1), \dots, T(y_n)\}$, and A and P are fresh unary EDB and IDB predicates, respectively. One can show (see [22, 27] for details) that, for any such q , called a *1-CQ* henceforth, (Π_q, G) is equivalent to the *disjunctive datalog program* (Δ_q, G) with rules (1) and (2) in the sense that they return the same answer over any data instance \mathcal{D} . Here, as usual, a *certain answer* to (Δ_q, G) over \mathcal{D} is ‘yes’ iff $\mathcal{I} \models G$, for every model \mathcal{I} of Δ_q and \mathcal{D} .

The monadic sirups, deciding boundedness of which is proved to be 2EXPTIME-hard in Sec. 3, take the form $\Sigma_q = \{(6), (7)\}$ with a 1-CQ q and goal predicate P . Adapting a similar terminology, we refer to disjunctive datalog programs $\Delta_q = \{(1), (2)\}$ and queries (Δ_q, G) , where q may contain multiple T and F in general, as *monadic disjunctive sirups* or *d-sirups*, for short.

EXAMPLE 2. Note that recursion in d-sirups is implicit and originates in ‘proof by exhaustion’ or ‘case distinction’, which can be seen by evaluating (Δ_{q_1}, G) and (Δ_{q_2}, G) (or the corresponding (Π_{q_2}, G)), with the q_i from Example 1, over the data instances \mathcal{D}_1

and, respectively \mathcal{D}_2 below.



For instance, let \mathcal{I} be any model of Δ_{q_2} and \mathcal{D}_2 . By rule (1), each of the A -nodes a and b in \mathcal{I} is labelled by F or T . If a is an F -node, q_2 is embeddable in \mathcal{I} via the vertical R -arrow. So let a be a T -node. If b is a T -node, q_2 is embeddable in \mathcal{I} starting from a , and if b is an F -node, there is an embedding starting from b . Thus, $\mathcal{I} \models q_2$.

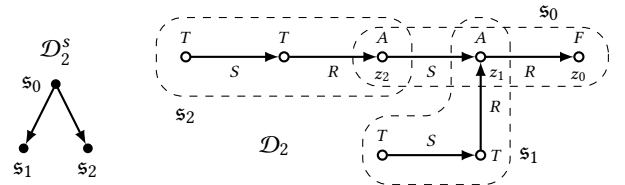
A monadic (disjunctive) datalog query (Π, Q) is *bounded* or *FO-rewritable* if there is a first-order formula $\Phi(x)$ (a sentence Φ if Q is 0-ary) such that, for any data instance \mathcal{D} , a constant $a \in \text{ind}(\mathcal{D})$ (or ‘yes’) is a certain answer to (Π, Q) over \mathcal{D} iff $\mathcal{D} \models \Phi[a]$ (respectively, $\mathcal{D} \models \Phi$), where \mathcal{D} is regarded as an FO-structure. It is known (see, e.g., [11, 20]) that in this case (Π, Q) is rewritable into a union of conjunctive queries (UCQ). It is also known [34] that FO-rewritability of datalog queries (Π, Q) can be characterised in terms of *Q-expansions*, which are defined inductively below for our special queries (Π_q, G) under the moniker ‘cactuses’.

To begin with, we set $C_G = \{F(x), q^-, T(y_1), \dots, T(y_n)\} = \{q\}$ and $\mathfrak{R}_q = \{C_G\}$. Then we take the closure of \mathfrak{R}_q under the rule **(bud)** if $T(y) \in C \in \mathfrak{R}_q$ is solitary, then we add to \mathfrak{R}_q the set of atoms obtained from C by replacing $T(y)$ with the atoms $A(x), q^-, T(y_1), \dots, T(y_n)$, in which x is renamed to y and all other variables are given *fresh* names.

The elements of the resulting (infinite if $n \geq 1$) set \mathfrak{R}_q are called *cactuses* for (Π_q, G) . We represent cactuses as labelled digraphs.

For $C \in \mathfrak{R}_q$, we refer to the copies \mathfrak{s} of (maximal subsets of) q that comprise C as *segments* and to the copy of the solitary F -node in \mathfrak{s} as its *focus*. The *skeleton* C^s of C is the ditree whose nodes are the segments \mathfrak{s} of C and edges $(\mathfrak{s}, \mathfrak{s}')$ mean that \mathfrak{s}' was attached to \mathfrak{s} by budding. The *depth* of \mathfrak{s} in C (or in C^s) is the number of edges on the branch from the root of C^s to \mathfrak{s} . The *depth* of C is the maximum depth of its segments.

EXAMPLE 3. The data instance \mathcal{D}_2 from Example 2 is (isomorphic to) a cactus from \mathfrak{R}_{q_2} obtained by applying **(bud)** to q_2 twice. The skeleton \mathcal{D}_2^s along with its three segments $\mathfrak{s}_0, \mathfrak{s}_1, \mathfrak{s}_2$ and their respective focuses z_0, z_1, z_2 are illustrated below:



In the remainder of this section, we establish a connection between boundedness of (Π_q, G) and (Σ_q, P) , for a 1-CQ q , which requires a few definitions. Every cactus $C \in \mathfrak{R}_q$ has exactly one F -node. We call it the *root-focus* of C and denote it by r . By replacing the F -label of r in C with A , we obtain a digraph C° ; the set of all such C° , for $C \in \mathfrak{R}_q$, is denoted by \mathfrak{R}_q° . The following proposition is proved by a standard induction on the derivation length:

PROPOSITION 1. For any data instance \mathcal{D} and any $a \in \text{ind}(\mathcal{D})$,

- $G \in \Pi_q(\mathcal{D})$ iff there is a homomorphism from some cactus $C \in \mathfrak{R}_q$ to \mathcal{D} ;
- $P(a) \in \Sigma_q(\mathcal{D})$ iff either $T(a) \in \mathcal{D}$ or there is a homomorphism h from some $C^\circ \in \mathfrak{R}_q^\circ$ to \mathcal{D} such that $h(r) = a$.

A 1-CQ q is called *focused* if the following condition holds:

- (foc)** for any cactuses $C, C' \in \mathfrak{R}_q$, if there is a homomorphism $h: C \rightarrow C'$, then $h(r) = r$.

The significance of this notion is shown by Example 4 below, and by the following characterisation of boundedness; cf. [34]:

PROPOSITION 2. For every focused 1-CQ q with solitary $F(x), T(y_1), \dots, T(y_n)$, the following conditions are equivalent:

- (a) (Σ_q, P) is bounded;
- (b) (Π_q, G) is bounded;
- (c) there exists $d < \omega$ such that, for every $C \in \mathfrak{R}_q$, there is a homomorphism $h: C' \rightarrow C$, for some $C' \in \mathfrak{R}_q$ of depth $\leq d$.

Conditions (b) and (c) are equivalent for every (not necessarily focused) 1-CQ q , in which case (a) is equivalent to (c) with an additional requirement that $h(r) = r$.

PROOF. (a) \Rightarrow (b) If $\Phi(x)$ is an FO-rewriting of (Σ_q, P) , then

$$\exists x, y_1, \dots, y_n, z (F(x) \wedge q' \wedge \Phi(y_1) \wedge \dots \wedge \Phi(y_n))$$

is an FO-rewriting of (Π_q, G) , where z comprises the variables in q' different from x, y_1, \dots, y_n .

(b) \Rightarrow (c) Let $\exists \mathbf{y} (q_1 \vee \dots \vee q_m)$ be a UCQ-rewriting of (Π_q, G) , where the q_i are CQs and \mathbf{y} comprises their variables. Treating the q_i as data instances, we obviously have $G \in \Pi_q(q_i)$, and so, for every $i, 1 \leq i \leq m$, there is a homomorphism from some $C_i \in \mathfrak{R}_q$ to q_i . Let d be the maximum depth of the $C_i, i = 1, \dots, m$. Consider any $C \in \mathfrak{R}_q$. Then there are homomorphisms $C_i \rightarrow q_i \rightarrow C$, for some $i, 1 \leq i \leq m$, the composition of which is the required h .

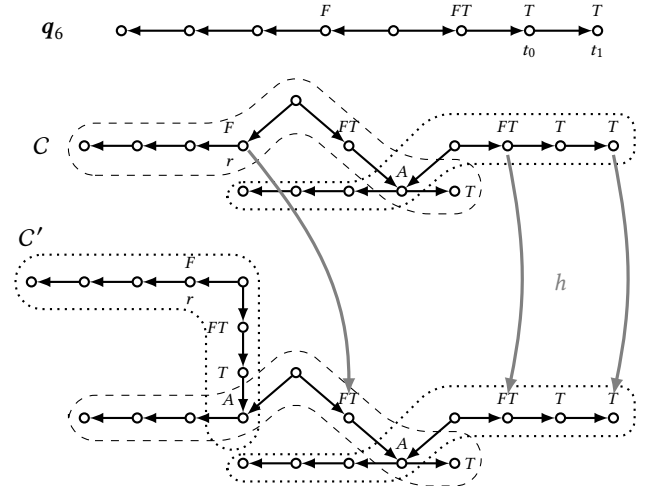
(c) \Rightarrow (a) By Prop. 1 and (c), the sentence $\exists r, \mathbf{y} (C_1 \vee \dots \vee C_m)$, where the C_i are all of the cactuses of depth $\leq d$ with root-focus r and the remaining variables \mathbf{y} , is an FO-rewriting of (Π_q, G) . We show that the formula $\Phi(r) = T(r) \vee \exists \mathbf{y} (C_1^\circ \vee \dots \vee C_m^\circ)$ is an FO-rewriting of (Σ_q, P) . Let $P(a) \in \Sigma_q(\mathcal{D})$, for some \mathcal{D} and $a \in \text{ind}(\mathcal{D})$. By Prop. 1, either $T(a) \in \mathcal{D}$, in which case $\mathcal{D} \models \Phi[a]$, or there is a homomorphism h from some $C^\circ \in \mathfrak{R}_q^\circ$ to \mathcal{D} such that $h(r) = a$. By (c), there is a homomorphism $g: C_i \rightarrow C$, for some $i \leq m$. As q is focused, $g(r) = r$, and so we can regard g as a $C_i^\circ \rightarrow C^\circ$ homomorphism. But then we obtain a homomorphism $hg: C_i^\circ \rightarrow \mathcal{D}$ with $hg(r) = a$, from which $\mathcal{D} \models \exists \mathbf{y} C_i^\circ[a]$. That $\mathcal{D} \models \Phi[a]$ implies $P(a) \in \Sigma_q(\mathcal{D})$ is trivial. \square

The next example illustrates the difference between focused and unfocused 1-CQs q as far as boundedness of (Π_q, G) and (Σ_q, P) is concerned.

EXAMPLE 4. Consider the 1-CQ q_5 from Example 1. Let C_k be the cactus obtained by applying **(bud)** k -times to $C_0 = q_5$. There are homomorphisms $h: C_1 \rightarrow C_k$, for $k \geq 2$, and so both (Π_{q_5}, G) and (Δ_{q_5}, G) are rewritable to the UCQ $C_0 \vee C_1$. For each such h , we have $h(r) = r$, so q_5 is focused and the sirup (Σ_{q_5}, P) is bounded.

Now, consider the 1-CQ q_6 below, where all of the arrows are labelled by R . It is not hard to see that, for every $C' \in \mathfrak{R}_{q_6}$ of depth ≥ 2 , there exist $C \in \mathfrak{R}_{q_6}^\circ$ of depth ≤ 1 and a homomorphism $h: C \rightarrow C'$, so (Π_{q_6}, G) and (Δ_{q_6}, G) are FO-rewritable. However,

every such h maps the root-focus F -node r to an FT -node, and so q_6 is not focused. In the picture below, C is obtained by budding at t_0 , and C' by budding first at t_1 and then at t_0 . Using Prop. 2, one can show that (Σ_{q_6}, P) is not bounded.



3 DECIDING BOUNDEDNESS OF SIRUPS

In this section, we prove the following:

THEOREM 3. The problems of deciding boundedness of monadic sirups (Σ_q, P) and monadic d -sirups (Δ_q, G) are both 2EXPTIME-hard.

Before diving into technical details, we put this theorem into the context of related work.

3.1 Related results

That deciding program boundedness of arbitrary monadic datalog queries can be done in 2EXPTIME was shown in 1988 using an automata-theoretic technique [18]. A matching lower bound for monadic queries with multiple recursive rules was finally settled in 2015 [9] using a construction from [8], which is based on the encoding of Turing machine computations from [12, 13]. For monadic sirups, the NP lower bound for the linear case [37] has remained so far the best known result (though, in view of Prop. 2 and the proof of [22, Theorem 9], it can be raised to PSPACE).

Establishing a higher lower bound for monadic sirups is difficult for two obvious reasons: monadicity and singularity. The impact of arity and the number of recursive rules on deciding boundedness of datalog programs has been studied in great detail; see [25, 33] and further references therein. For example, boundedness was shown to be undecidable first for linear datalog programs of arity 4 [21], then for those of arity 2 with multiple recursive rules [37], which were encoded in a single rule at the expense of higher arity [1]; finally, boundedness was proved to be undecidable already for linear sirups of arity 3 [33].

Intuitively, the proofs of the lower bounds mentioned above use different rules in a datalog program in order to detect and exclude different ‘defects’ in possible computations of a Turing machine. Our task in the proof of Theorem 3 will be to design such an encoding of computations that can be verified by a single CQ.

3.2 Proof idea

To achieve this, similarly to [7–9, 12, 13], we represent computations of a Turing machine by means of annotated binary trees. The design of the tree-representation of computations is such that its structure can be connected with expansions (cactuses) of a given sirup via a series of small Boolean circuits, which is the main innovation of our construction.

More precisely, we use the criterion of Prop. 2 for testing boundedness. Our aim is, given any alternating Turing machine (ATM) M deciding a language in $AEXPSPACE = 2EXPTIME$ and an input w , to construct a (dag-shaped) focused 1-CQ q of polynomial size such that the following holds:

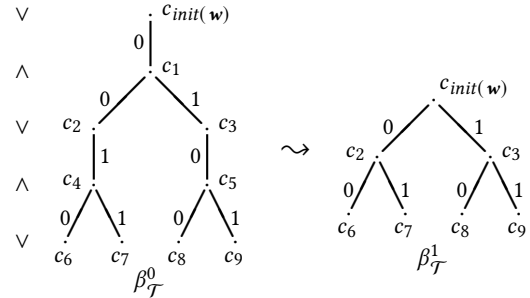
LEMMA 4. *M rejects w iff there is $K < \omega$ such that every cactus $C \in \mathfrak{R}_q$ contains a homomorphic image of some $C^- \in \mathfrak{R}_q$ of depth at most K .*

We represent both the computation space of M on w and q -cactuses by 01-trees: binary ditrees whose edges are labelled by 0 or 1, with siblings having different labels. On the one hand, we encode the computation space of M in such a way that checking whether an arbitrary 01-tree represents a rejecting computation-tree on w can be done by means of polynomially-many polynomial-size Boolean circuits (in fact, formulas). On the other hand, the 1-CQ q we associate with M and w has two solitary T -nodes, t_0 and t_1 . Thus, we can regard the skeleton C^s of any cactus $C \in \mathfrak{R}_q$ as a 01-tree, indicating which of t_0 or t_1 were budded. The 1-CQ q is assembled from gadgets implementing the Boolean circuits used for checking the above properties of computations.

3.3 Connecting computations and cactuses

3.3.1 Encoding computations by 01-trees. We assume that we are given an ATM $M = (Q, \Gamma, \delta, q_{init}, q_{accept}, q_{reject}, g)$ with states Q including $q_{init}, q_{accept}, q_{reject}$, tape alphabet Γ , transition function δ , and $g: Q \rightarrow \{\wedge, \vee\}$. For any input $w \in \Gamma^*$, a configuration of M is a triple containing information about the current state, the current position of the head, and the current content of the $2^{P(|w|)} = 2^p$ tape-cells, for some polynomial p . If its current state is q , then we call the configuration a q -configuration. The full computation space $\mathcal{T}_{M,w}$ is a finite tree whose nodes are (labelled by) configurations, with its root being the initial configuration $c_{init}(w)$ (in state q_{init} reading the leftmost symbol of w), the descendants generated by δ , and each leaf being either a q_{accept} - or a q_{reject} -configuration (a halting configuration). We assume that the depth of $\mathcal{T}_{M,w}$ is $2^{P(|w|)}$, $q_{init}, q_{accept}, q_{reject}$ are \vee -states, every non-leaf has branching 2, and \wedge - and \vee -configurations alternate on each branch. A computation-tree (of M on w) is a substructure \mathcal{T} of $\mathcal{T}_{M,w}$, which is a tree with root $c_{init}(w)$ such that every non-leaf \wedge -node (\vee -node) in \mathcal{T} has both (respectively, exactly one) of its children from $\mathcal{T}_{M,w}$ in \mathcal{T} . The tree \mathcal{T} is rejecting if it has a q_{reject} -leaf and accepting otherwise. M rejects w iff all computation-trees of M on w are rejecting, and accepts w otherwise.

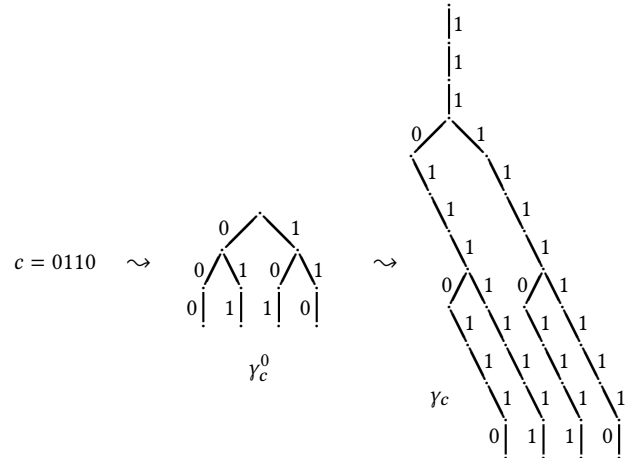
We encode a computation-tree \mathcal{T} by an infinite 01-tree $\beta_{\mathcal{T}}^+$ via a series of steps as follows. First, by our assumption on binary branching, \mathcal{T} can be considered as a (finite) 01-tree $\beta_{\mathcal{T}}^0$ (with its nodes still labelled by configurations). Next, we take the full binary ‘substructure’ $\beta_{\mathcal{T}}^1$ of the \vee -configurations in $\beta_{\mathcal{T}}^0$ as shown below:



(So the depth of $\beta_{\mathcal{T}}^1$ is 2^{p-1} .) The information about which child of each \vee -configuration is taken in $\beta_{\mathcal{T}}^0$ is provided in the encoding of the subsequent \vee -configuration. To achieve this, we fine-tune the ‘configurations-as-binary-tree-leaves’ representation of [12, 13] for our purpose. Let $d = d(|w|) > p(|w|) = p$ be a polynomial in $|w|$ such that configurations can be encoded by a 01-sequence of length 2^d . We represent each \vee -configuration c by the 01-sequence

$$\begin{array}{c} \text{state } q \quad 0 \quad \text{cell content } t_1 \quad 0 \quad \text{cell content } t_2 \quad \dots \quad 1 \quad \text{active cell } t_k \quad \dots \quad 0/1 \\ \hline \log |Q| \quad \log |\Gamma| \quad \log |\Gamma| \quad \dots \quad \log |\Gamma| \end{array}$$

where the last bit is 0 (1) iff c 's parent \wedge -configuration is a 0-child (1-child) of its parent. (By imposing some restrictions on Q and Γ , one can ensure that, given a 2^d -long 01-sequence, it is ‘easy’ to locate the first bit of each ‘cell-representation’ in it.) We encode the digits of this sequence as the leaves of a 01-tree γ_c^0 of depth $d+1$ by taking first a full binary tree of depth d , and for each of its 2^d leaves, taking a $*$ -child whenever the corresponding digit in the sequence is $*$. (Throughout, we use $*$ in 01-sequences as a wildcard for 0 or 1.) Finally, we turn γ_c^0 to a 01-tree γ_c of depth $4d+4$ by adding an incoming edge-pattern 111 above each node:



We call γ_c a c -tree (or, a configuration-tree, in general).

Next, we take the full binary 01-tree $\beta_{\mathcal{T}}^1$ above (whose nodes are labelled by \vee -configurations), and turn it to a 01-tree $\beta_{\mathcal{T}}$ (now without node labels) as follows. We add an incoming edge-pattern 0010 above the root, stick the root of a c -tree to each node labelled by some c , and add an outgoing edge-pattern 001 below each node before branching; see Fig. 1. Note that $\beta_{\mathcal{T}}$ is of depth $e = e(|w|)$, for some exponential function e . For any configuration c , if the c -tree γ_c is a substructure of $\beta_{\mathcal{T}}$, then we call the root node of γ_c a main node (of c) and say that it represents c in $\beta_{\mathcal{T}}$; see \bullet -nodes in Fig. 1.

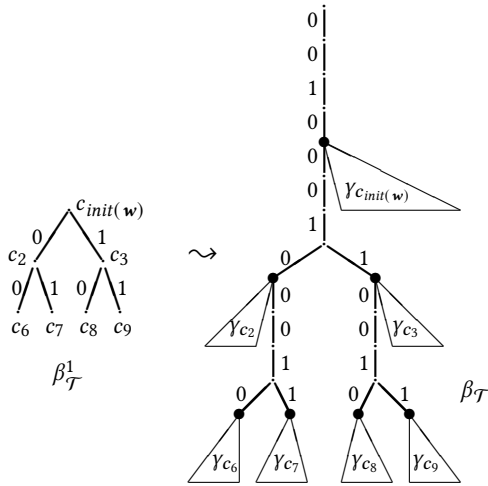
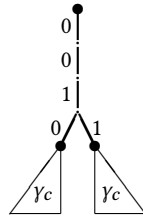


Figure 1: The 01-tree $\beta_{\mathcal{T}}$.

We also consider an infinite ‘version’ of $\beta_{\mathcal{T}}$. We obtain the infinite 01-tree $\beta_{\mathcal{T}}^+$ from $\beta_{\mathcal{T}}$ by repeatedly sticking the following pattern to the main node of each halting configuration c :



In other words, we assume δ to be such that after reaching a halting configuration c , c is repeated forever on every branch of $\mathcal{T}_{M,w}$.

An infinite 01-tree β is *ideal* if it can be constructed by starting with $\beta_{\mathcal{T}_0}^+$ for some computation-tree \mathcal{T}_0 , and then by repeatedly attaching to each of the leaves (that must be leaves of some configuration-tree) the root of some $\beta_{\mathcal{T}}^+$, where each \mathcal{T} can be any computation-tree.

Observe that each branch of an ideal tree is infinite. We are interested in finite ‘middle-bits’ of ideal trees. We call a subtree of an ideal tree having a main node (of not necessarily $c_{init(w)}$) as root a *desired tree*. Given some $M < \omega$ and a 01-tree β , by an M -cut of β we mean the 01-tree obtained by cutting all longer than M branches in β at depth M . The pretty baroque design above ensures that there is a polynomial list of polynomially detectable properties that identify desired trees up to isomorphism (see Claim 4.1 below). In the next subsection, we discuss these properties.

3.3.2 Characterising exponential computations polynomially. We investigate certain polynomial neighbourhoods of nodes in 01-trees, and collect a polynomial list of their properties that fully characterise those situations that can occur in a desired tree. For each of the properties P below, in Sec. 3.4 we describe in detail how to give a small Boolean circuit φ_P having specific *input-types* such that, when evaluated at a node \mathbf{a} of some 01-tree β , P fails at \mathbf{a} iff there is some 01-sequence \mathbf{b} such that \mathbf{b} is gathered from the neighbourhood of \mathbf{a} in β according to the input-types of φ_P and $\varphi_P[\mathbf{b}] = 1$.

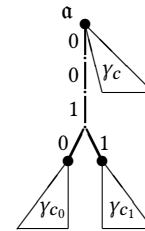
Given $n < \omega$ and a node \mathbf{a} of depth $\leq n$ in a 01-tree β , for any $k \leq n$, we denote by $P_{\mathbf{a}}^k$ the k -long suffix of the path ending at \mathbf{a} in β . To begin with, observe that every path in a desired tree that is longer than $4d + 6$ must contain a main node, and main nodes can be identified by the property ‘the path leading to the node ends with a 001*-pattern’. So, given a node \mathbf{a} in a 01-tree β , we say that \mathbf{a} is *good* in β , if either the depth of \mathbf{a} in β is $< 4d + 11$, or $P_{\mathbf{a}}^{4d+11}$ contains a 001*-pattern; see Sec. 3.4.1.

Next, we describe proper branching-patterns in a desired tree. It is easy to see that if the path leading to a node \mathbf{a} does contain a 001*-pattern, then there exist unique k, ℓ and w such that $4 \leq k \leq 4d + 11$, $P_{\mathbf{a}}^k = 001*(111*)^\ell w$, and either $\ell \leq d$ and w is a prefix of 001, or $\ell < d$ and w is a prefix of 111. Moreover, ℓ and w characterise the children of \mathbf{a} . We call \mathbf{a} *properly branching* in β if the following conditions (pb1)–(pb4) hold:

- (pb1) if either w is empty and $\ell = 0$, or $w = 001$, or $w = 111$ and $\ell < d - 1$, then \mathbf{a} has two children;
- (pb2) if either w is empty and $0 < \ell < d$, or $w = 1$, or $w = 11$, or $w = 00$, then \mathbf{a} has no 0-child;
- (pb3) if either w is empty and $\ell = d$, or $w = 0$, then \mathbf{a} has no 1-child;
- (pb4) if $w = 111$ and $\ell = d - 1$, then \mathbf{a} has only one child;

see Sec. 3.4.2. Note that leaves are never properly branching.

Next, we ensure that the ‘building-block’ computation-trees in an ideal tree are properly represented in a 01-tree β (provided that all of its nodes are properly branching). First, after each leaf of a configuration-tree, the representation of a proper computation-tree from $c_{init(w)}$ should start. In such a case, the main node \mathbf{a} of $c_{init(w)}$ can be identified by an incoming path ending with a $111*001*$ -pattern. Then detecting whether the c -tree with root \mathbf{a} is not a $c_{init(w)}$ -tree requires checking polynomial information. We call \mathbf{a} *properly initialising* in β if whenever the depth of \mathbf{a} in β is ≥ 8 , $P_{\mathbf{a}}^8$ is of the form $111*001*$, and \mathbf{a} is the root of a c -tree, then $c = c_{init(w)}$; see Sec. 3.4.4. Second, the computation steps described by δ should be properly represented. We call \mathbf{a} *properly computing* in β if, whenever the following pattern is present at \mathbf{a} in β



then the triple (c, c_0, c_1) of v -configurations ‘matches’ the transition function δ of M . In order to detect that this is not the case, one needs to check polynomial information ‘around’ \mathbf{a} in β ; see Sec. 3.4.3.

We call \mathbf{a} *correct* in β if \mathbf{a} is good, properly branching, properly initialising and properly computing in β . Otherwise, \mathbf{a} is called *incorrect* in β . Now it is straightforward to show that the collected properties of \mathbf{a} -neighbourhoods characterise desired trees:

CLAIM 4.1. *For any $M < \omega$, any 01-tree β and any node \mathbf{a} with $P_{\mathbf{a}}^4 = 001*$, the M -cut $\beta_{\mathbf{a}}^M$ of the subtree of β with root \mathbf{a} is isomorphic to the M -cut of a desired tree iff every node of depth $< M$ in $\beta_{\mathbf{a}}^M$ is correct in $\beta_{\mathbf{a}}^M$.*

We also need to detect the presence of nodes representing q_{reject} -configurations in computation-trees; see Sec. 3.4.5.

3.3.3 Cactus homomorphisms vs rejecting computations. As our 1-CQ q will have one solitary F -node and two solitary T -nodes t_0 and t_1 , there are four possible kinds of non-root segments in any cactus $C \in \mathfrak{R}_q$ denoted $q_{TT}^-, q_{AT}^-, q_{TA}^-$ and q_{AA}^- . For example, q_{TT}^- is obtained by replacing the F -label of the solitary F -node in q by A ; in cactuses different from q , leaf segments are of this form. The segment q_{TA}^- is obtained by replacing both the F -label of the solitary F -node and the T -label of the t_1 -node by A . As q itself does not contain A , if $h: C \rightarrow C'$ is a homomorphism, for some $C, C' \in \mathfrak{R}_q$, then the focus of every non-root segment s in C (labelled by A) is mapped by h to the focus of some non-root segment s' in C' . Our q will also satisfy **(foc)**: for every homomorphism $h: C \rightarrow C'$ between cactuses $C, C' \in \mathfrak{R}_q$, h maps the only solitary F -node in C (the focus of its root segment) to the only solitary F -node in C' . So we say that h maps a segment s into a segment s' if h maps the focus of s to the focus of s' .

Now the proof of Theorem 3 can be completed as follows: Using Claim 4.1, we prove in Appendix A that to obtain Lemma 4 it suffices to construct a 1-CQ q such that **(foc)** holds and, for every $C \in \mathfrak{R}_q$,

- (leaf)** there is a homomorphism $h: q_{TT}^- \rightarrow C$ mapping q_{TT}^- into some non-leaf segment s of C iff either s is incorrect or s represents a q_{reject} -configuration in the skeleton C^s of C ;
- (branch)** if h maps q_{TT}^- into a non-leaf segment s that is not properly branching in C^s due to violating **(pb1)**, but s is correct in C^s according to the other properties, then
 - $h(t_0) = t_0$ and $h(t_1) \neq t_0$, if $s = q_{TA}^-$;
 - $h(t_1) = t_1$ and $h(t_0) \neq t_1$, if $s = q_{AT}^-$.

After defining the focused 1-CQ q in Secs. 3.5.1–3.5.3, we show in Sec. 3.5.4 that, for every $C \in \mathfrak{R}_q$, **(leaf)** and **(branch)** are satisfied, completing the proof of Lemma 4.

3.4 Boolean formulas

We describe polynomially-many polynomial-size Boolean circuits (in fact, Boolean formulas) that test the (failure) of the properties of a node a in a 01-tree β , given in Sec. 3.3.2. For each such formula, we also define some *input-types*, describing where around the tested node a the input 01-sequence for the formula should be ‘gathered’ from. In defining the input-types we use the following terminology: for $n < \omega$, the *n-long uppath* (of a in β) is the reverse of the n -long suffix of the path ending at a in β ; while an *n-long downpath* is the n -long prefix of some path starting at a in β .

3.4.1 Checking goodness. One can clearly define a Boolean formula $\text{GOOD}(x_1, \dots, x_{4d+11})$ such that, for any $4d+11$ -long 01-sequence b , $\text{GOOD}[b] = 1$ iff b does not contain the reverse of a 001^* -pattern. The input should be gathered from the $4d+11$ -long uppath.

3.4.2 Checking proper branching-patterns. For each of conditions **(pb1)**–**(pb4)** in Sec. 3.3.2, we have a different family of formulas.

(pb1) For every k with $4 \leq k \leq 4d+11$, we define a Boolean formula $\text{MUSTBRANCH}^k(x_1, \dots, x_k)$ such that, for any k -long 01-sequence $b = (b_1, \dots, b_k)$, $\text{MUSTBRANCH}^k[b] = 1$ iff b is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where either w is empty and $\ell = 0$, or $w = 001$, or $w = 111$ and $\ell < d-1$. For example, if $k = 4$

then we have

$$\text{MUSTBRANCH}^4[b] = 1 \quad \text{iff } b \text{ is the reverse of } 001^*.$$

The input should be gathered from the k -long uppath.

(pb2) For every k with $4 \leq k \leq 4d+11$, we define a Boolean formula $\text{NOBRANCH}_0^k(x_1, \dots, x_{k+1})$ such that, for any $k+1$ -long 01-sequence $b = (b_1, \dots, b_{k+1})$, $\text{NOBRANCH}_0^k[b] = 1$ iff $b_{k+1} = 0$ and (b_1, \dots, b_k) is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where either w is empty and $0 < \ell < d$, or $w = 1$, or $w = 11$, or $w = 00$. The input for (x_1, \dots, x_k) should be gathered from the k -long uppath, and for x_{k+1} from a 1-long downpath.

(pb3) For every k with $4 \leq k \leq 4d+11$, we define a polynomial size Boolean formula $\text{NOBRANCH}_1^k(x_1, \dots, x_{k+1})$ such that, for any $k+1$ -long 01-sequence $b = (b_1, \dots, b_{k+1})$, $\text{NOBRANCH}_1^k[b] = 1$ iff $b_{k+1} = 1$ and (b_1, \dots, b_k) is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where either w is empty and $\ell = d$, or $w = 0$. The input for (x_1, \dots, x_k) should be gathered from the k -long uppath, and for x_{k+1} from a 1-long downpath.

(pb4) For every k with $4 \leq k \leq 4d+11$, we define a polynomial size Boolean formula $\text{NOBRANCH}^k(x_1, \dots, x_{k+2})$ such that, for any $k+2$ -long 01-sequence $b = (b_1, \dots, b_{k+2})$, $\text{NOBRANCH}^k[b] = 1$ iff $b_{k+1} \neq b_{k+2}$ and (b_1, \dots, b_k) is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where $w = 111$ and $\ell = d-1$. The input for (x_1, \dots, x_k) should be gathered from the k -long uppath, and for each of x_{k+1} and x_{k+2} from a 1-long downpath.

3.4.3 Checking proper computation steps. This is an adaptation of the technique of [12, 13] to our representation. Suppose that

- n_Q is such that n_Q -long 01-sequences are in one-to-one correspondence with the states in Q ,
- n_Γ is such that $(n_\Gamma - 1)$ -long 01-sequences are in one-to-one correspondence with the symbols in Γ ,
- $n_Q + 2^P \cdot n_\Gamma + 1 = 2^d$

(see the picture in Sec. 3.3.1 on representing configurations with 2^d -long 01-sequences).

First, we define a Boolean formula $\text{HEAD}(x_1, \dots, x_{4(d+1)})$ such that, for any $4(d+1)$ -long 01-sequence b , $\text{HEAD}[b] = 1$ iff b describes a path in a desired tree starting at a main node with 1, and ending at the first bit of the representation of some cell-content of some configuration c (that is, it is the $(n_Q + i \cdot n_\Gamma + 1)$ th bit of the 01-sequence representing c , for some $i < 2^P$). Similarly, for $* = 0, 1$, we define a Boolean formula $\text{HEAD}^*(x_1, \dots, x_{4(d+1)+4})$ such that, for any $4(d+1)+4$ -long 01-sequence b , $\text{HEAD}^*[b] = 1$ iff b describes a path in a desired tree starting at a main node with 001^*1 , and ending at the first bit of the representation of some cell-content of some configuration.

Next, we define a Boolean formula

$$\text{SAMECELL}(x_1, \dots, x_{4(d+1)}, y_1, \dots, y_{4(d+1)+4}, z_1, \dots, z_{4(d+1)+4})$$

such that, for any $4(d+1)$ -long 01-sequence b and $4(d+1)+4$ -long 01-sequences b^0, b^1 , we have $\text{SAMECELL}[b, b^0, b^1] = 1$ iff $\text{HEAD}[b] = 1$, $\text{HEAD}^*[b^*] = 1$, and the three paths b, b^0, b^1 end at the same cell of a configuration and its two children-configurations. (The number i of this cell, for some $i < 2^P$, can be identified from b .)

Next, we define a Boolean formula $\text{STATE}(x_1, \dots, x_{4(d+1) \cdot n_Q})$ such that, for any $4(d+1)$ -long 01-sequences $\mathbf{b}_1, \dots, \mathbf{b}_{n_Q}$,

$$\text{STATE}[\mathbf{b}_1, \dots, \mathbf{b}_{n_Q}] = 1$$

iff for every $j \leq n_Q$, \mathbf{b}_j describes a $4(d+1)$ -long path in a desired tree starting at a main node with 1, and ending at the j th bit of the representation of some configuration c . (So, whenever $\mathbf{b}_j = (b_j^1, \dots, b_j^{4(d+1)})$ for each $j \leq n_Q$, then $(b_1^{4(d+1)}, \dots, b_{n_Q}^{4(d+1)})$ encodes a state in Q .) Similarly, for $*$ = 0, 1, we define a Boolean formula $\text{STATE}^*(x_1, \dots, x_{4(d+1)+4 \cdot n_Q})$ such that, for any $4(d+1) + 4$ -long 01-sequences $\mathbf{b}_1, \dots, \mathbf{b}_{n_Q}$, $\text{STATE}^*[\mathbf{b}_1, \dots, \mathbf{b}_{n_Q}] = 1$ iff for every $j \leq n_Q$, \mathbf{b}_j describes a path in a desired tree starting at a main node with 001*1, and ending at the j th bit of the representation of some configuration.

Next, we define a Boolean formula $\text{CELL}(x_1, \dots, x_{4(d+1) \cdot n_\Gamma})$ such that, for any $4(d+1)$ -long 01-sequences $\mathbf{b}_1, \dots, \mathbf{b}_{n_\Gamma}$,

$$\text{CELL}[\mathbf{b}_1, \dots, \mathbf{b}_{n_\Gamma}] = 1$$

iff there is $i < 2^P$ such that, for every $j \leq n_\Gamma$, \mathbf{b}_j describes a path in a desired tree starting at a main node with 1, and ending at the j th bit of the representation of the i th cell's content in some configuration. (So, $\text{HEAD}[\mathbf{b}_1] = 1$, and if $\mathbf{b}_j = (b_j^1, \dots, b_j^{4(d+1)})$ for each $j \leq n_\Gamma$, then $(b_2^{4(d+1)}, \dots, b_{n_\Gamma}^{4(d+1)})$ encodes a symbol in Γ .) Similarly, for $*$ = 0, 1, we define a polynomial size Boolean formula $\text{CELL}^*(x_1, \dots, x_{4(d+1)+4 \cdot n_\Gamma})$ such that, for any $4(d+1) + 4$ -long 01-sequences $\mathbf{b}_1, \dots, \mathbf{b}_{n_\Gamma}$, $\text{CELL}^*[\mathbf{b}_1, \dots, \mathbf{b}_{n_\Gamma}] = 1$ iff there is $i < 2^P$ such that, for every $j \leq n_\Gamma$, \mathbf{b}_j describes a path in a desired tree starting at a main node with 001*1, and ending at the j th bit of the representation of the i th cell's content in some configuration (In particular, $\text{HEAD}^*[\mathbf{b}_1] = 1$.)

Next, for $z \in \{0, 1\}$, we take the following tuples of variables:

- $\mathbf{s} = (s_1, \dots, s_{4(d+1) \cdot n_Q})$, which is to be gathered from n_Q -many $4(d+1)$ -long downpaths (representing the \vee -state in c);
- $\mathbf{v} = (v_1, \dots, v_{4(d+1) \cdot n_\Gamma})$, which needs to be gathered from n_Γ -many $4(d+1)$ -long downpaths (representing the active cell's content in c);
- $\mathbf{s}^0 = (s_1^0, \dots, s_{4(d+1)+4 \cdot n_Q}^0)$, $\mathbf{s}^1 = (s_1^1, \dots, s_{4(d+1)+4 \cdot n_Q}^1)$, each of which needs to be gathered from n_Q -many $4(d+1)+4$ -long downpaths (representing the \vee -states in c_0, c_1);
- $\mathbf{t}^z = (t_1^z, \dots, t_{4(d+1) \cdot n_\Gamma}^z)$, $\mathbf{t}^{z0} = (t_1^{z0}, \dots, t_{4(d+1)+4 \cdot n_\Gamma}^{z0})$, and $\mathbf{t}^{z1} = (t_1^{z1}, \dots, t_{4(d+1)+4 \cdot n_\Gamma}^{z1})$, for $z \in \{0, 1\}$, where \mathbf{t}^z is to be gathered from n_Γ -many $4(d+1)$ -long downpaths, and each of \mathbf{t}^{z0} and \mathbf{t}^{z1} is to be gathered from n_Γ -many $4(d+1)+4$ -long downpaths ($\mathbf{t}^z, \mathbf{t}^{z0}, \mathbf{t}^{z1}$ represent the i th cell's contents in c, c_0, c_1 , for some $i < 2^P$, when the z \wedge -child of c is taken in $\mathcal{T}_{M,x}$);
- $\mathbf{z}^0 = (z_1^0, \dots, z_{4(d+1)+4}^0)$ and $\mathbf{z}^1 = (z_1^1, \dots, z_{4(d+1)+4}^1)$, each of which needs to be gathered from a $4(d+1) + 4$ -long downpath (representing the respective bits identifying the parent \wedge -configuration of c_0 and c_1).

For $z \in \{0, 1\}$, we can define a Boolean formula STEP^z such that, for any 01-sequence $\mathbf{b} = (\mathbf{s}, \mathbf{v}, \mathbf{s}^0, \mathbf{s}^1, \mathbf{t}, \mathbf{t}^0, \mathbf{t}^1, \mathbf{z}^0, \mathbf{z}^1)$, $\text{STEP}^z[\mathbf{b}] = 1$ iff $\mathbf{z}^0 = 001011 \dots 1z$, $\mathbf{z}^1 = 001111 \dots 1z$, $\text{STATE}[\mathbf{s}] = 1$, $\text{CELL}[\mathbf{v}] = 1$,

$\text{STATE}^0[\mathbf{s}^0] = 1$, $\text{STATE}^1[\mathbf{s}^1] = 1$, $\text{CELL}[\mathbf{t}] = 1$, $\text{CELL}^0[\mathbf{t}^0] = 1$, $\text{CELL}^1[\mathbf{t}^1] = 1$, and

$$\text{SAMECELL}[t_1, \dots, t_{4(d+1)}, t_1^0, \dots, t_{4(d+1)+4}^0, t_1^1, \dots, t_{4(d+1)+4}^1] = 1,$$

but the information provided by \mathbf{b} is inconsistent with the transition function δ in the sense that when the z \wedge -child of c is chosen as the common parent of c_0 and c_1 in the computation-tree \mathcal{T} , the content-triple of the i th cells of c, c_0, c_1 is wrong, where i is identified from the input in $(t_1, \dots, t_{4(d+1)}, t_1^0, \dots, t_{4(d+1)+4}^0, t_1^1, \dots, t_{4(d+1)+4}^1)$.

Finally, we define $\text{STEP}(\mathbf{s}, \mathbf{v}, \mathbf{s}^0, \mathbf{s}^1, \mathbf{t}^0, \mathbf{t}^{00}, \mathbf{t}^{01}, \mathbf{t}^1, \mathbf{t}^{10}, \mathbf{t}^{11}, \mathbf{z}^0, \mathbf{z}^1)$ as the disjunction

$$\begin{aligned} &\text{STEP}^0(\mathbf{s}, \mathbf{v}, \mathbf{s}^0, \mathbf{s}^1, \mathbf{t}^0, \mathbf{t}^{00}, \mathbf{t}^{01}, \mathbf{z}^0, \mathbf{z}^1) \vee \\ &\text{STEP}^1(\mathbf{s}, \mathbf{v}, \mathbf{s}^0, \mathbf{s}^1, \mathbf{t}^1, \mathbf{t}^{10}, \mathbf{t}^{11}, \mathbf{z}^0, \mathbf{z}^1). \end{aligned}$$

It is not hard to see that there is \mathbf{b} with $\text{STEP}[\mathbf{b}] = 1$ iff the information about the configuration-triple (c, c_0, c_1) encoded in \mathbf{b} is inconsistent with δ .

3.4.4 Checking proper initialisation. We take the following tuples of variables:

- $\mathbf{y} = (y_1, \dots, y_8)$, which is to be gathered from the 8-long uppath (representing the last 8-bits of the path leading to the main node of a configuration c);
- $\mathbf{s} = (s_1, \dots, s_{4(d+1) \cdot n_Q})$, to be gathered from n_Q -many $4(d+1)$ -long downpaths (representing the state in c);
- $\mathbf{w}^j = (w_1^j, \dots, w_{4(d+1) \cdot n_\Gamma}^j)$, for $j < |\mathbf{w}|$, each of which needs to be gathered from n_Γ -many $4(d+1)$ -long downpaths (representing the contents of the first $|\mathbf{w}|$ -many cells in c);
- $\mathbf{t} = (t_1, \dots, t_{4(d+1) \cdot n_\Gamma})$, which needs to be gathered from n_Γ -many $4(d+1)$ -long downpaths (representing the contents of some cell in c).

Then we can define a Boolean formula INIT such that, for any 01-sequence $\mathbf{b} = (\mathbf{y}, \mathbf{s}, \mathbf{w}^1, \dots, \mathbf{w}^{|\mathbf{w}|}, \mathbf{t})$, $\text{INIT}[\mathbf{b}] = 1$ iff \mathbf{y} is the reverse of some pattern 111*001*, $\text{STATE}[\mathbf{s}] = 1$, for all $1 \leq j \leq |\mathbf{x}|$, $\text{CELL}[\mathbf{w}^j] = 1$ and $(w_1^j, \dots, w_{4(d+1)}^j)$ ends at the $(n_Q + (j-1) \cdot n_\Gamma + 1)$ th bit of the 01-sequence representing configurations, $\text{CELL}[\mathbf{t}] = 1$, but the information provided by \mathbf{b} is inconsistent with c being $c_{\text{init}(\mathbf{w})}$ (at the cell identified by the prefix $(t_1, \dots, t_{4(d+1)})$ of \mathbf{t}).

3.4.5 Representing q_{reject} -configurations. This can clearly be done by a formula $\text{REJECT}(s_1, \dots, s_{4(d+1) \cdot n_Q})$ for which $\text{REJECT}[\mathbf{s}] = 1$ iff $\text{STATE}[\mathbf{s}] = 1$ and the sequence

$$(s_{4(d+1)}, s_{4(d+1)-2}, \dots, s_{4(d+1) \cdot n_Q})$$

encodes q_{reject} . The input should be gathered from n_Q -many $4(d+1)$ -long downpaths.

3.5 Query design

The dag-shaped 1-CQ q having one solitary F -node, two solitary T -nodes and many FT -twins will be such that properties (**foe**), (**leaf**) and (**branch**) given in Sec. 3.3.3 hold, for all $C, C' \in \mathfrak{R}_q$.

3.5.1 Overall query structure. To simplify notation, in our pictures we omit the R -labels from R -arrows, and use extra labels (different from F, T) on nodes, say B on a , as a shorthand for a B -arrow (a, a') to a *fresh* node a' . Letters other than upper case italics (greek, lower case italics and bold) are used as pointers and are not part of q .

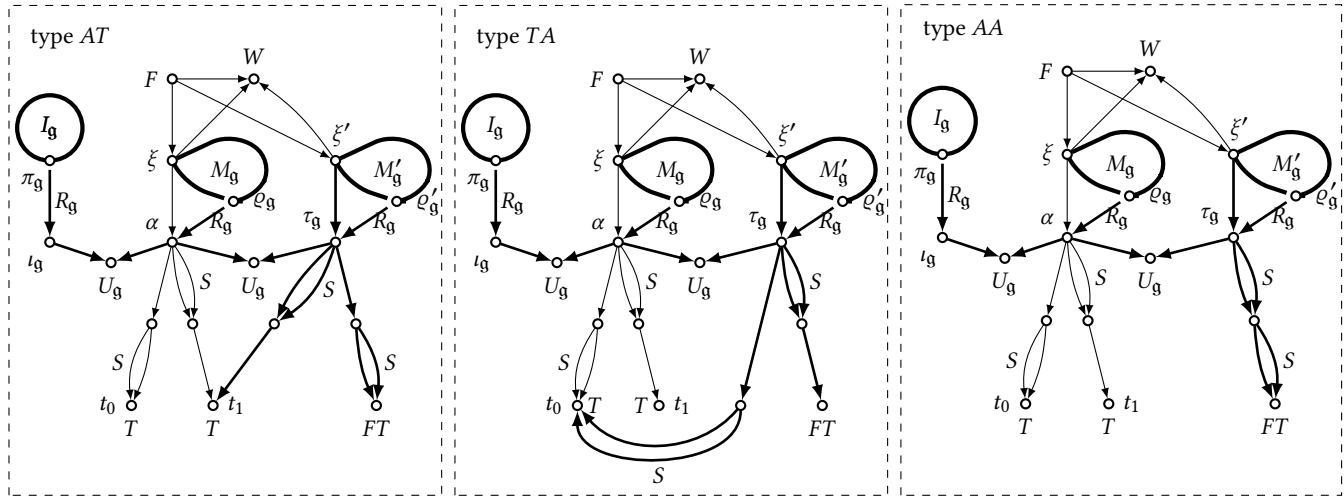
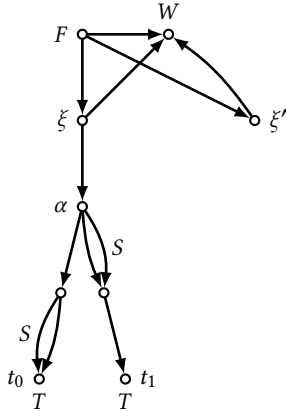


Figure 2: Frames of type AT, TA and AA.

The 1-CQ q has the following simple *base block* containing all of the solitary F - and T -nodes of q :



Wired to the base in q are *gadgets* g that *implement* the Boolean formulas φ_g defined in Sec. 3.4. Each gadget g has four components: two isomorphic copies of its *main block* M_g and M'_g , an *input block* I_g and a *frame*. The frame wires the gadget to the base and can be of one of the three *types* AT, TA and AA, which are shown in Fig. 2 (with the base block being indicated in each case by thin lines). We say that g is of type Z if its frame is of type Z . The frame of g has a few distinguished nodes: π_g and ι_g via which I_g is R_g -wired to the base block, ϱ_g via which M_g is R_g -wired to the base block, ϱ'_g and τ_g via which M'_g is R_g -wired to the base block (where the edge-labelling R_g is also unique for gadget g), the single FT -twin of g (none of I_g , M_g and M'_g contains any FT -twins), and two nodes labelled by U_g .

It is easy to see that q satisfies **(foc)**: its F -node has successors, while none of the FT -nodes does. Further, we observe that if $h: q_{TT}^- \rightarrow C$ is a homomorphism mapping q_{TT}^- into some non-leaf segment s , then there must be a gadget g such that $h(\alpha) = \tau_g$, for the α -node in q_{TT}^- and the τ_g -node in s . Then, because of the U_g -nodes, $h(\iota_g) = \alpha$ must hold, for the ι_g -node in q_{TT}^- and the α -node in s . Therefore, $h(\pi_g) = \varrho_g$ and the I_g -block of q_{TT}^- must be mapped

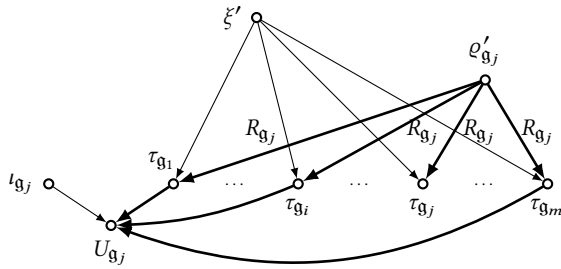
by h to the M_g -block of s , forcing the input to interact with the formula.

Given a gadget g and a homomorphism $h: q_{TT}^- \rightarrow C$ mapping q_{TT}^- into a non-leaf segment s of some cactus C , we say that g is *triggered* by h at s if $h(\iota_g) = \alpha$, for the ι_g -node in q_{TT}^- and the α -node in s . We say that g is *triggered* at s if there is a homomorphism h triggering g at s . Observe that if s is of the form q_Z^- , for some $Z \in \{AT, TA, AA\}$, and g is triggered at s , then g is either of type AA or of type Z .

Each gadget g in q ‘implements’ some Boolean formula $\varphi_g(\mathbf{y})$ checking some property of desired trees at node s in the skeleton 01 -tree C^s of the cactus C . So the input values for the variables in \mathbf{y} are ‘collected’ from an environment of s in C^s . This collection process is regulated by the input-types of each φ_g . We have the following gadgets in q , each implementing some formula described in Sec. 3.4:

- (g1) a type AA gadget implementing GOOD;
- (g2) for every k with $4 \leq k \leq 4d + 11$ and every $Z \in \{AT, TA\}$, a type Z gadget implementing MUSTBRANCH^k ;
- (g3) for every k with $4 \leq k \leq 4d + 11$ and every $*$ in $\{0, 1\}$, a type AA-gadget implementing NOBRANCH^k_* ;
- (g4) for every k with $4 \leq k \leq 4d + 11$, a type AA-gadget implementing NOBRANCH^k ;
- (g5) a type AA gadget implementing STEP;
- (g6) a type AA gadget implementing INIT;
- (g7) a type AA gadget implementing REJECT.

Now suppose g_1, \dots, g_m are all of the gadgets in q . We want to ensure that when a gadget g_j is triggered by some h , then the other gadgets are not triggered (that is, the ι_{g_i} -node for every $i \neq j$ can be mapped by h to itself). So, in addition to the above, for every j , we connect the ι_{g_j} -node via an U_{g_j} -labelled node to the τ_{g_i} -node, for all $i \neq j$. We also want to ensure that when g_j is triggered by some h , then the M_{g_i} -block can be h -mapped to the M'_{g_i} -block for every i (not just for j). So we not only R_{g_j} -connect ϱ'_{g_j} with τ_{g_i} , but also add R_{g_j} -arrows connecting ϱ'_{g_j} with all of the τ_{g_i} :



The proof of the following claim is provided in Appendix B:

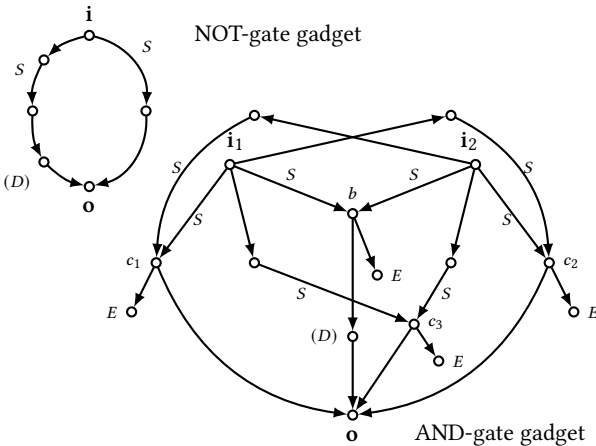
CLAIM 4.2. A gadget g in q is triggered at s iff there is b_g^s such that b_g^s is gathered from 'around' s in C^s according to the input-types for φ_g and $\varphi_g[b_g^s] = 1$.

Claim 4.2 will be used in Sec. 3.5.4 to show that every $C \in \mathcal{R}_q$ satisfies properties (leaf) and (branch) given in Sec. 3.3.3.

3.5.2 Main blocks in gadgets. Here we give a uniform description of the main block M_g of each gadget g in q . Apart from the label W , which is uniform through the gadgets, for each particular g , there are a few additional labels on some nodes in M_g , M'_g and I_g . These are always specific to g , but we omit indicating this to simplify notation.

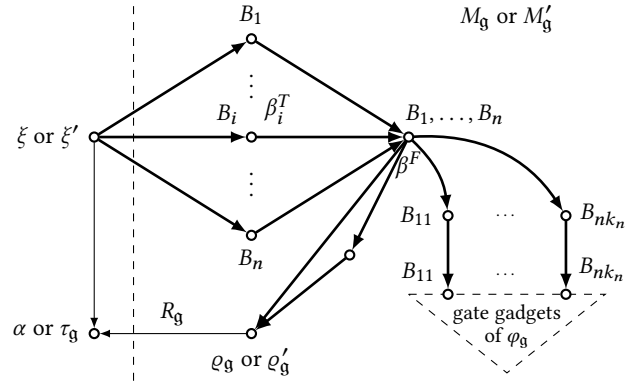
A Boolean formula $\varphi_g(\mathbf{y})$ is regarded as a ditree whose vertices are called gates. Leaf gates are labelled by the variables from the list $\mathbf{y} = (y_1, \dots, y_n)$, with each y_i labelling k_i -many leaves of $\varphi_g(\mathbf{y})$. Each non-leaf g is either an AND-gate (having 2 children) or a NOT-gate (having 1 child), with the outgoing edge(s) leading to the input(s) of g . Given an assignment b of 0 or 1 to the input-variables y_i , we compute the value of each gate in φ_g under b as usual in Boolean logic.

We encode the gate-structure of $\varphi_g(\mathbf{y})$ by the M_g -block (and also by its copy M'_g) as follows. With each non-leaf gate g in φ_g we associate a fresh copy of its gadget shown below (where D in brackets means that D is only present when the gate in question is the root of φ_g):

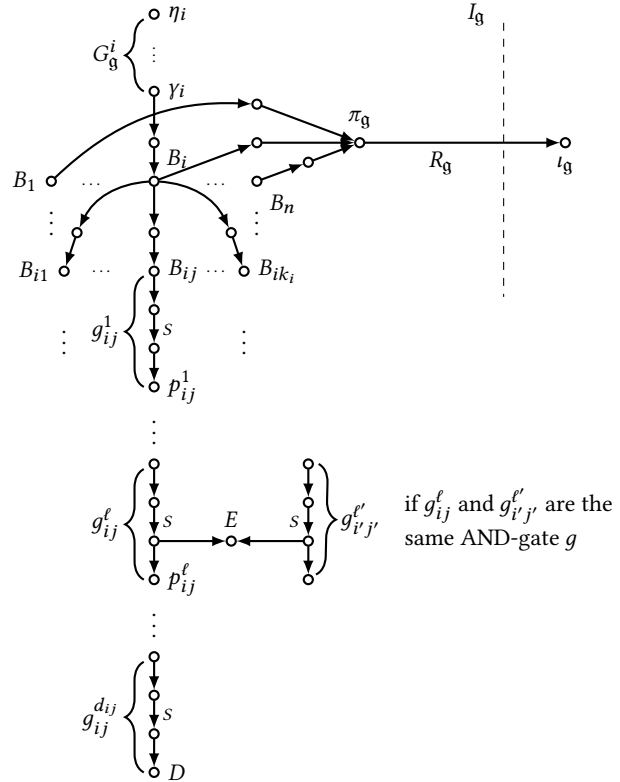


Each branch of φ_g is characterised by a pair (i, j) such that the leaf node of the branch is labelled by the j th copy y_i^j of the variable y_i , for some i, j with $1 \leq i \leq n$ and $1 \leq j \leq k_i$. For each pair (i, j) , we introduce a label B_{ij} . Suppose that g_1 and g_2 are the inputs of an

AND-gate g . Then, for each $m = 1, 2$, if g_m is a non-leaf gate, then we merge node o of the g_m -gadget with node i_m of the g -gadget; and if g_m is labelled by y_i^j , we merge node i_m of the g -gadget with the lower B_{ij} -node in M_g . We proceed similarly with NOT-gates. The picture below shows how M_g (and its copy M'_g) looks like (where, apart from the B_{ij} , we also label some nodes with B_i , for $1 \leq i \leq n$):



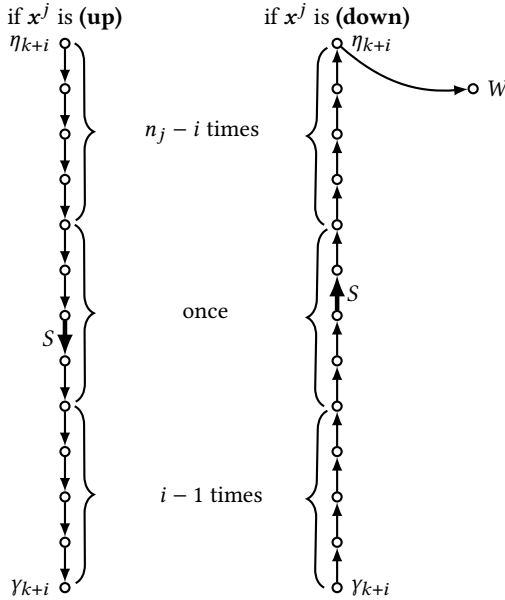
3.5.3 Input blocks in gadgets. Given a Boolean formula $\varphi_g(\mathbf{y})$ with $\mathbf{y} = (y_1, \dots, y_n)$, its input block I_g consists of a uniformly describable part (depending on φ_g and n) and a gathering block G_g^i , for each i with $1 \leq i \leq n$ (depending on the input-types of $\varphi_g(\mathbf{y})$). For each branch of φ_g characterised by (i, j) , let $g_{ij}^1, \dots, g_{ij}^{d_{ij}}$ be the sequence of non-leaf gates from leaf to root on the branch with leaf y_i^j . The structure of the input block I_g is shown below:



Finally, we describe the gathering blocks G_g^i in I_g . The Boolean formula in g takes the form $\varphi_g(\mathbf{y}) = \varphi_g(x^1, \dots, x^m)$ where each tuple $x^j = (x_{n_1}^j, \dots, x_{n_j}^j)$ of variables can be of two input-types:

- (**up**) either x^j is gathered from the (unique) n_j -long *uppath* (the reverse of the suffix of the path ending at \mathfrak{s} in $C^{\mathfrak{s}}$);
- (**down**) or x^j is gathered from an n_j -long *downpath* (the prefix of a path starting at \mathfrak{s} in $C^{\mathfrak{s}}$).

So suppose $y_{k+1}, \dots, y_{k+n_j}$ are among the variables of φ_g such that $(y_{k+1}, \dots, y_{k+n_j}) = x^j$ for some j . Then, for each i with $1 \leq i \leq n_j$, G_g^{k+i} is shown below:



In case x^j is like in (**down**), the W -node (of the base block) is a common successor of the η_{k+i} -nodes, for every $i = 1, \dots, n_j$, which ensures that the input bits for $y_{k+1}, \dots, y_{k+n_j}$ are all gathered from the same n_j -long downpath; see Appendix B for an example.

3.5.4 Proving that q satisfies (leaf) and (branch). Suppose C is a cactus in \mathfrak{R}_q and \mathfrak{s} is a non-leaf segment in the skeleton $C^{\mathfrak{s}}$ of C . Then \mathfrak{s} is of the form $q_{Z_{\mathfrak{s}}}^-$ for some $Z_{\mathfrak{s}} \in \{AT, TA, AA\}$.

(**leaf**) (\Rightarrow) Suppose $h: q_{TT}^- \rightarrow C$ is a homomorphism mapping q_{TT}^- into \mathfrak{s} . Then there is a gadget g that is triggered by h at \mathfrak{s} (that is, the $h(t_g) = \alpha$ for the t_g -node in q_{TT}^- and the α -node in \mathfrak{s}). By Claim 4.2, there is $b_g^{\mathfrak{s}}$ such that $b_g^{\mathfrak{s}}$ is gathered from ‘around’ \mathfrak{s} in $C^{\mathfrak{s}}$ according to the input-types for φ_g and $\varphi_g[b_g^{\mathfrak{s}}] = 1$. Now we have a case distinction, depending on the gadget g , as listed in Sec. 3.5.1. Each gadget implements a formula whose input-types and behaviour are described in Sec. 3.4:

(**g1**) g is the type AA gadget implementing GOOD (cf. Sec. 3.4.1). By the input-types of GOOD, $b_g^{\mathfrak{s}}$ is the $4d + 11$ -long uppath, and it does not contain the reverse of a 001^* -pattern. Thus, \mathfrak{s} is not good in $C^{\mathfrak{s}}$.

(**g2**) There exist some k with $4 \leq k \leq 4d + 11$ and $Z \in \{AT, TA\}$ such that g is the type Z gadget implementing MUSTBRANCH^k (cf. Sec. 3.4.2). By the input-types of MUSTBRANCH^k , $b_g^{\mathfrak{s}}$ is the k -long uppath, and it is the reverse of a sequence of the form $001^*(111^*)^\ell w$,

where either w is empty and $\ell = 0$, or $w = 001$, or $w = 111$ and $\ell < d - 1$. On the other hand, $Z_{\mathfrak{s}} = Z$ must hold, and so $C^{\mathfrak{s}}$ is not branching at \mathfrak{s} . As branching at \mathfrak{s} is required in condition (**pb1**) of being properly branching, it follows that \mathfrak{s} is not properly branching, and so it is incorrect in $C^{\mathfrak{s}}$.

(**g3**) There exist some k with $4 \leq k \leq 4d + 11$ and $* \in \{0, 1\}$ such that g is the type AA gadget implementing NOBRANCH_*^k (cf. Sec. 3.4.2). Suppose, say, that $* = 0$ (the case when $* = 1$ is similar). By the input-types of NOBRANCH_*^k , $b_g^{\mathfrak{s}} = (e^{\mathfrak{s}}, b^{\mathfrak{s}})$, where $e^{\mathfrak{s}}$ is the k -long uppath and $b^{\mathfrak{s}}$ is a 1-long downpath. Also, $b^{\mathfrak{s}} = 0$ and $e^{\mathfrak{s}}$ is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where either w is empty and $0 < \ell < d$, or $w = 1$, or $w = 11$, or $w = 00$. As having a 0-child is forbidden in condition (**pb2**) of being properly branching, it follows that \mathfrak{s} is not properly branching, and so it is incorrect in $C^{\mathfrak{s}}$.

(**g4**) There exists some k with $4 \leq k \leq 4d + 11$ such that g is the type AA gadget implementing NOBRANCH^k (cf. Sec. 3.4.2). By the input-types of NOBRANCH^k , $b_g^{\mathfrak{s}} = (e^{\mathfrak{s}}, b_1^{\mathfrak{s}}, b_2^{\mathfrak{s}})$, where $e^{\mathfrak{s}}$ is the k -long uppath and each of $b_1^{\mathfrak{s}}$ and $b_2^{\mathfrak{s}}$ is a 1-long downpath. Also, $b_1^{\mathfrak{s}} \neq b_2^{\mathfrak{s}}$, and $e^{\mathfrak{s}}$ is the reverse of a sequence of the form $001^*(111^*)^\ell w$, where $w = 111$ and $\ell = d - 1$. As having two children is forbidden in condition (**pb4**) of being properly branching, it follows that \mathfrak{s} is not properly branching, and so \mathfrak{s} is incorrect in $C^{\mathfrak{s}}$.

(**g5**) g is the type AA gadget implementing STEP. By the input-types of STEP, $b_g^{\mathfrak{s}}$ should have gathered data about some \vee -configuration c and its two ‘subsequent’ \vee -configurations c_0, c_1 . As explained in Sec. 3.4.3, (c, c_0, c_1) is inconsistent with δ , and so \mathfrak{s} is not properly computing in $C^{\mathfrak{s}}$. Thus, it is incorrect in $C^{\mathfrak{s}}$.

(**g6**) g is the type AA gadget implementing INIT. By the input-types of INIT, $b_g^{\mathfrak{s}}$ should have gathered data about the 8-long uppath and some \vee -configuration c . As explained in Sec. 3.4.4, the part of $b_g^{\mathfrak{s}}$ gathered from the 8-long uppath is the reverse of some pattern 111^*001^* , but $c \neq c_{\text{init}(w)}$. Thus, \mathfrak{s} is not properly initialising in $C^{\mathfrak{s}}$, and so it is incorrect in $C^{\mathfrak{s}}$.

(**g7**) g is the type AA gadget implementing REJECT. As explained in Sec. 3.4.5, by the input-types of REJECT, $b_g^{\mathfrak{s}}$ should have gathered data about some state q , and $q = q_{\text{reject}}$ must hold. Therefore, \mathfrak{s} represents a q_{reject} -configuration in $C^{\mathfrak{s}}$, as required.

(**leaf**) (\Leftarrow) Again, we have cases (**g1**)–(**g7**). In each case, we have a formula φ_g for which some input $b_g^{\mathfrak{s}}$ can be gathered from ‘around’ \mathfrak{s} in $C^{\mathfrak{s}}$ according to its input-types and for which $\varphi_g[b_g^{\mathfrak{s}}] = 1$. So by Claim 4.2, there is a homomorphism $h: q_{TT}^- \rightarrow C$ mapping q_{TT}^- into \mathfrak{s} and triggering g at \mathfrak{s} .

(**branch**) Suppose there is a homomorphism $h: q_{TT}^- \rightarrow C$ mapping q_{TT}^- into some non-leaf segment \mathfrak{s} . Then some gadget g is triggered by h at \mathfrak{s} . By our assumption on \mathfrak{s} and by the proof of the (\Rightarrow) direction of (**leaf**) above, it follows that g can only be the Z type gadget implementing MUSTBRANCH^k , where $Z \in \{AT, TA\}$ is such that $\mathfrak{s} = q_{Z}^-$. An inspection of Fig. 2 shows that $h(t_0) = t_0$ and $h(t_1) \neq t_0$ whenever $Z = TA$, and $h(t_1) = t_1$ and $h(t_0) \neq t_1$ whenever $Z = AT$. Therefore, (**branch**) holds.

This completes the proof that q satisfies (**leaf**) and (**branch**).

3.6 OMQs with Schema.org and DL-Lite_{bool}

Schema.org, founded by Google, Microsoft, Yahoo and Yandex and developed by an open community process, comprises a set of rules

$P(x) \leftarrow Q(x)$, for unary or binary predicates P and Q , together with domain and range constraints such as

$$T(x) \vee F(x) \leftarrow S(x, y) \quad (8)$$

$$T(y) \vee F(y) \leftarrow R(x, y) \quad (9)$$

For example, according to the Schema.org ontology, the range of the binary relation $\text{musicBy}(x, y)$ is covered by the union of $\text{MusicGroup}(y)$ and $\text{Person}(y)$. In the syntax of description logic *DL-Lite_{bool}* [5], rules (8) and (9) are written as

$$\exists S \sqsubseteq T \sqcup F \quad \text{and} \quad \exists R^- \sqsubseteq T \sqcup F$$

Given any d-sirup (Δ_q, G) , denote by Δ'_q the ‘Schema.org ontology’ obtained by replacing (1) in Δ_q with rule (9), for a fresh R .

PROPOSITION 5. *A d-sirup (Δ_q, G) is FO-rewritable iff (Δ'_q, G) is FO-rewritable.*

PROOF. (\Rightarrow) Suppose Φ is a UCQ-rewriting of (Δ_q, G) and Φ' result of replacing every $A(y)$ in Φ with $\exists x R(x, y)$. We claim that Φ' is an FO-rewriting of (Δ'_q, G) . Indeed, suppose \mathcal{D}' is any data instance for (Δ'_q, G) . Without loss of generality we may assume that it does not contain atoms $A(a)$. Let \mathcal{D} be the result of adding $A(b)$ to \mathcal{D}' whenever $R(a, b) \in \mathcal{D}'$. Then $\Delta_q, \mathcal{D} \models G$ iff $\Delta'_q, \mathcal{D}' \models G$, and also $\mathcal{D} \models \Phi$ iff $\mathcal{D}' \models \Phi'$, from which $\Delta'_q, \mathcal{D}' \models G$ iff $\mathcal{D}' \models \Phi'$.

(\Leftarrow) Suppose Φ' is a UCQ-rewriting of (Δ'_q, G) and Φ is the result of replacing every $R(x, y)$ in Φ' with $A(y)$. Let \mathcal{D} be a data instance for (Δ_q, G) . Without loss of generality we may assume that it does not contain atoms of the form $R(a, b)$. Let \mathcal{D}' be the result of adding $R(a, b)$, for a fresh a , to \mathcal{D} whenever $A(b) \in \mathcal{D}$. Then $\Delta_q, \mathcal{D} \models G$ iff $\Delta'_q, \mathcal{D}' \models G$, and also $\mathcal{D} \models \Phi$ iff $\mathcal{D}' \models \Phi'$, from which $\Delta'_q, \mathcal{D}' \models G$ iff $\mathcal{D}' \models \Phi'$. \square

As a consequence of Theorem 3 and Proposition 5, we obtain the following theorem, which is an improvement on [24, Theorem 11] showing PSPACE-hardness of deciding FO-rewritability of UCQs mediated by Schema.org ontologies.

THEOREM 6. *Deciding FO-rewritability of CQs mediated by a Schema.org or DL-Lite_{bool} ontology is 2EXP-TIME-hard.*

4 MONADIC D-SIRUPS WITH A DITREE CQ

The high lower bound obtained in the previous section can be regarded as a formal confirmation of the empirical fact that finding transparent syntactic, let alone practical criteria of FO-rewritability for sufficiently general classes of monadic (d-)sirups is a notoriously difficult problem. The only positive results in this direction we know of are the syntactic NC/P dichotomy of binary *chain* sirups [4] (see also [3]) and the complete AC⁰/NL/P/coNP tetrachotomy of monadic *path* d-sirups without twins [22].

The CQs used in the proof of Theorem 3 were directed acyclic graphs with one solitary F -node, two solitary T -nodes, and multiple FT -twins. The question we try to answer in this section is whether the restriction of the set of CQs admitted in d-sirups to those that are *rooted directed trees* as graphs (*ditree* CQs, for short) makes deciding FO-rewritability of d-sirups (Δ_q, G) easier, having in mind a complete syntactic classification of such d-sirups as an ultimate (possibly unrealistic) aim. Note for starters that, by Example 1, the

data complexity of evaluating d-sirups with a ditree CQ ranges from AC⁰ to L, NL, P, and coNP.

A CQ q is *minimal* if there is no $q \rightarrow q'$ homomorphism, for any proper subCQ q' of q . As well-known, checking minimality of tree-shaped CQs can be done in polynomial time; see, e.g., [16]. We denote the root node of q by r and write $x \leq_q y$ to say that there is a (directed) path from x to y in q , and $x <_q y$ if $x \leq_q y$ and $x \neq y$. A pair (x, y) is *<_q-comparable* if either $x \leq_q y$ or $y \leq_q x$, otherwise (x, y) is *<_q-incomparable*. If $x \leq_q y$ then $\delta_q(x, y)$ is the number of edges between x and y . The *distance* between any x and y is $\partial_q(x, y) = \delta_q(\inf_q(x, y), x) + \delta_q(\inf_q(x, y), y)$, where $\inf_q(x, y)$ is the unique node such that $\inf_q(x, y) \leq_q x$, $\inf_q(x, y) \leq_q y$ and $z \leq_q \inf_q(x, y)$ whenever $z \leq_q x$ and $z \leq_q y$. The subscript q in $\leq_q, <_q, \delta_q, \inf_q$ and ∂_q will be dropped if understood.

If t is a solitary T -node and f is solitary F -node, we call (t, f) a *solitary pair*. We say that a solitary pair (t, f) is of *minimal distance*, if $\partial(t, f) \geq \partial(t, f)$ for any solitary pair (t, f) . A *<-incomparable* solitary pair (t, f) is called *symmetric* if the CQ obtained by removing the labels F, T from f, t and cutting the branches below them is symmetric with respect to r (see q_4 in Example 1). A ditree q is *quasi-symmetric* if it has no *<-comparable* solitary pairs, and every solitary pair (t, f) of minimal distance is symmetric.

As follows from [22] (where F and T are interchangeable),

- (a) if q has no solitary F , then (Δ_q, G) is FO-rewritable;
- (b) if q has one solitary F , then (Δ_q, G) is datalog-rewritable (and so in P for data complexity);
- (c) if q has one solitary F and one solitary T , then (Δ_q, G) is linear-datalog-rewritable (and so in NL);
- (d) if q has one solitary F , one solitary T and is quasi-symmetric, then (Δ_q, G) is symmetric-linear-datalog-rewritable (and so in L).

The following result identifies a large and tractable class of d-sirups with a ditree CQ whose evaluation is NL-hard:

THEOREM 7. *Suppose q is a minimal ditree CQ with at least one solitary F , at least one solitary T and such that either*

- (i) *there is a <-comparable solitary pair (t, f) or*
- (ii) *q is not quasi-symmetric and has no FT-twins.*

Then evaluating the d-sirup (Δ_q, G) is NL-hard.

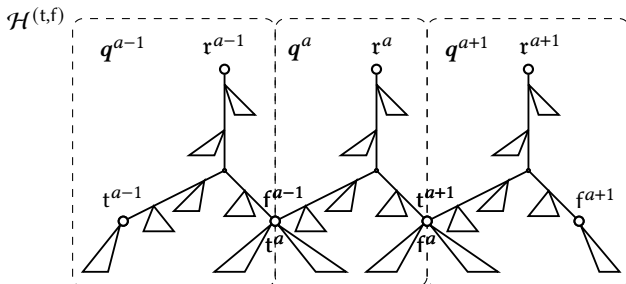
PROOF. The proof is by reduction of the NL-complete reachability problem for dags. Given a dag $G = (V, E)$ with nodes $s, t \in V$, we construct a data instance \mathcal{D}_G as follows. We pick a solitary pair (t, f) such that, in case (i), (t, f) is *<-comparable* and there is no solitary T - or F -node between t and f ; and, in case (ii), (t, f) is of minimal distance, *<-incomparable*, and not symmetric. Then, in both cases, we replace each $e = (u, v) \in E$ by a fresh copy q^e of q in which t^e is renamed to u with $T(u)$ replaced by $A(u)$, and f^e is renamed to v with $F(v)$ replaced by $A(v)$. The dag \mathcal{D}_G comprises the q^e , for $e \in E$, as well as $T(s)$ and $F(t)$. We show that $s \rightarrow_G t$ iff the answer to (Δ_q, G) over \mathcal{D}_G is ‘yes’.

(\Rightarrow) If $s = v_0, \dots, v_n = t$ is a path in G with $e_i = (v_i, v_{i+1}) \in E$, for $i < n$, then for any model \mathcal{I} of Δ_q and \mathcal{D}_G , there is some $i < n$ such that $\mathcal{I} \models T(v_i)$ and $\mathcal{I} \models F(v_{i+1})$, and so the identity map from q to its copy q^{e_i} is a $q \rightarrow \mathcal{I}$ homomorphism.

(\Leftarrow) If $s \rightarrow_G t$, we define a model \mathcal{I} of Δ_q and \mathcal{D}_G by labelling with T the A -nodes in \mathcal{D}_G that (as nodes of G) are reachable from s

(via a directed path in G) and with F the remaining ones. We claim that if one of (i) or (ii) holds, then there is no homomorphism from \mathbf{q} to \mathcal{I} , and so the answer to $(\Delta_{\mathbf{q}}, G)$ over \mathcal{D}_G is 'no'. Indeed, take any map h from \mathbf{q} to \mathcal{I} , and consider the substructure $\mathcal{H}^{(t,f)}$ of \mathcal{D}_G comprising those copies $\mathbf{q}^{e_1}, \dots, \mathbf{q}^{e_n}$ of \mathbf{q} that have a non-empty intersection with $h(\mathbf{q})$. To simplify notation, we set $\mathbf{q}^j = \mathbf{q}^{e_j}$. Then \mathcal{I} can be regarded as a model of $\mathcal{H}^{(t,f)}$. The (quite arduous case-distinction) proof in [29] shows that h cannot be a homomorphism from \mathbf{q} to \mathcal{I} .

Here, we only sketch the proof for case (ii) when \mathbf{q} is not quasi-symmetric, and we may also assume that \mathbf{q} has no \prec -comparable solitary pairs. Suppose $h: \mathbf{q} \rightarrow \mathcal{I}$ is a homomorphism, and let a be such that $h(\mathbf{r}) \in \mathbf{q}^a$. As (t, f) is \prec -incomparable, $\mathcal{H}^{(t,f)}$ consists of (at most) three copies $\mathbf{q}^a, \mathbf{q}^{a-1}$ and \mathbf{q}^{a+1} of \mathbf{q} , and looks as shown in the picture below:



As (t, f) is not symmetric, \mathcal{I} is such that the 'contacts' between the \mathbf{q} -copies are either both in $F^{\mathcal{I}}$ or both in $T^{\mathcal{I}}$.

The following 'structural' claim (tracking the possible locations of $h(f)$ and $h(t)$) is proved in the full version [29] (it is also used in the proof of Theorem 11 below):

CLAIM 7.1. *Suppose (t, f) is \prec -incomparable and of minimal distance (though \mathbf{q} might contain FT-twins). If $m = \inf_{\mathbf{q}}(t, f)$ then $h(m)$ is in \mathbf{q}^a , and one of the following holds:*

- (1) $m^a \prec_{\mathbf{q}^a} h(m) \prec_{\mathbf{q}^a} t^a$, $h(t)$ is in \mathbf{q}^{a-1} with $f^{a-1} \prec_{\mathbf{q}^{a-1}} h(t)$, and $h(f) = t^a$;
- (2) $m^a \prec_{\mathbf{q}^a} h(m) \prec_{\mathbf{q}^a} f^a$, $h(f)$ is in \mathbf{q}^{a+1} with $t^{a+1} \prec_{\mathbf{q}^{a+1}} h(f)$, and $h(t) = f^a$;
- (3) $h(m) = m^a$, $h(f) = f^a$, and $h(t)$ is in \mathbf{q}^a with $h(t) \prec_{\mathbf{q}^a} f^a$;
- (4) $h(m) = m^a$, $h(t) = t^a$, and $h(f)$ is in \mathbf{q}^a with $h(f) \prec_{\mathbf{q}^a} t^a$.

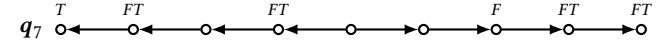
However, if \mathbf{q} contains neither FT-twins nor \prec -comparable solitary pairs, none of (1)–(4) in Claim 7.1 can happen. \square

Denote by $\Delta_{\mathbf{q}}^+$ the d-sirup $(\Delta_{\mathbf{q}}, G)$ extended by an extra rule $\perp \leftarrow T(x), F(x)$ saying that the predicates F and T are disjoint, and so $\Delta_{\mathbf{q}}^+$ with \mathbf{q} containing an FT-twin is inconsistent. As shown in [22], $(\Delta_{\mathbf{q}}, G)$ is L-hard when \mathbf{q} has at least one solitary F and at least one solitary T but no FT-twins. So we have:

COROLLARY 8. *Every d-sirup $(\Delta_{\mathbf{q}}^+, G)$ with a ditree \mathbf{q} is either FO-rewritable (if \mathbf{q} contains FT-twins), or L-hard (if \mathbf{q} is quasi-symmetric without FT-twins), or NL-hard (otherwise).*

The non-quasi-symmetric CQs \mathbf{q} that are outside the scope of Theorem 7 are those that have FT-twins and only contain \prec -incomparable solitary pairs. That Theorem 7 does not hold for such CQs is demonstrated by \mathbf{q}_5 in Example 1 (cf. Claim 7.1 (3)), \mathbf{q}_6 in Example 4 (cf. Claim 7.1 (4)), and $\mathbf{q}_7, \mathbf{q}_8$ below (cf. Claim 7.1 (1)), for

all of which $(\Delta_{\mathbf{q}}, G)$ is FO-rewritable. As before, the omitted labels on the arrows are all R .



Our next result, used in tandem with Theorem 7, gives an FO/L-hardness dichotomy for d-sirups $(\Delta_{\mathbf{q}}, G)$ with a ditree 1-CQ \mathbf{q} (having a single solitary F). If such a \mathbf{q} has k -many solitary T -nodes, each of which is \prec -incomparable with the F -node, we call it a Λ -CQ of span k .

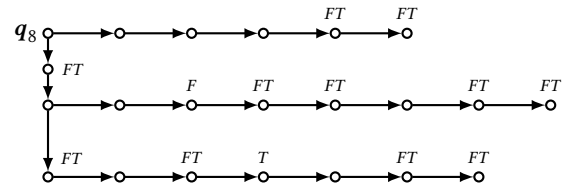
THEOREM 9. (i) *For any Λ -CQ \mathbf{q} , either the d-sirup $(\Delta_{\mathbf{q}}, G)$ is FO-rewritable or evaluating it is L-hard.*

(ii) *For Λ -CQs of span k , deciding this FO/L-dichotomy can be done in time $p(|\mathbf{q}|)2^{p'(k)}$, for some polynomials p and p' . Thus, deciding FO-rewritability of d-sirups with a Λ -CQ is fixed-parameter tractable, if the Λ -CQ's span is regarded as a parameter.*

PROOF. Let \mathbf{q} be a Λ -CQ with solitary T -nodes $T(y_1), \dots, T(y_k)$. By Prop. 2, $(\Delta_{\mathbf{q}}, G)$ is FO-rewritable iff there exists $d < \omega$ such that any cactus (for \mathbf{q}) contains a homomorphic image of some cactus of depth $\leq d$. The neighbourhood of a segment \mathfrak{s} in a cactus C consists of \mathfrak{s} itself and those segments that, in the skeleton C^s , are the children, parent and siblings of \mathfrak{s} —at most $2k + 1$ segments in total. Since \mathbf{q} is a ditree, in which the F -node is \prec -incomparable with any T -node, the following holds for any cactuses C, C' (for \mathbf{q}):

CLAIM 9.1. *Suppose $h: C \rightarrow C'$ is a homomorphism that maps the root of a segment \mathfrak{s} in C to a node in a segment \mathfrak{s}' in C' . Then the nodes in \mathfrak{s} are mapped by h to nodes in the neighbourhood of \mathfrak{s}' .*

EXAMPLE 5. Consider the Λ -CQ \mathbf{q}_8 of span 1 below. We invite the reader to verify that there is a homomorphism $h: C_2 \rightarrow C_i$, for $i \geq 3$ (where C_i is obtained by i -many applications of (bud) to $C_0 = \mathbf{q}_8$) such that the h -image of the leaf segment in C_2 intersects three segments in C_i . It is not hard to see that $(\Delta_{\mathbf{q}}, G)$ is FO-rewritable to $\exists z (C_0 \vee C_1 \vee C_2)$.



In any skeleton C^s , we label by $i \in \{1, \dots, k\}$ every edge that results from budding $T(y_i)$. The neighbourhood of any segment \mathfrak{s} is given by the triple $t = (P, i_{\mathfrak{s}}, C)$, where $P \subseteq \{1, \dots, k\}$ comprises the labels on the edges from the parent \mathfrak{s}' of \mathfrak{s} , $i_{\mathfrak{s}}$ is the label on $(\mathfrak{s}', \mathfrak{s})$, and $C \subseteq \{1, \dots, k\}$ are the labels on the edges to \mathfrak{s} 's children. If \mathfrak{s} is the root of C^s , $P = \emptyset$ and we set $i_{\mathfrak{s}} = 0$; if \mathfrak{s} is a leaf, $C = \emptyset$. We refer to \mathfrak{s} as the central segment of t , and to t as the type of \mathfrak{s} ; we call it a root type if \mathfrak{s} is the root of C^s , and a leaf type if \mathfrak{s} is a leaf. A cactus C is acyclic if none of the branches in C^s has two nodes of the same type.

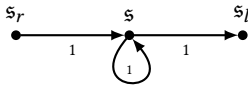
Let \mathfrak{G} be the digraph whose nodes are all possible types and there is an edge (t, t') labelled by $i \in \{1, \dots, k\}$ iff some skeleton C^s has an edge $(\mathfrak{s}, \mathfrak{s}')$ labelled by i with \mathfrak{s} being of type t and \mathfrak{s}' of type t' . Let χ_C be the canonical homomorphism of C^s to \mathfrak{G} (mapping the segments of C^s to their types). For a subgraph \mathfrak{H} of \mathfrak{G} denote by \mathfrak{H}^s

the result of replacing the types in \mathfrak{H} with their central segments and glueing them at A -nodes as indicated by the types and edges in \mathfrak{H} , mimicking (**bud**). We call this operation the $\bar{\cdot}$ -closure of \mathfrak{H} .

A node v of type (P, i, C) in a subgraph \mathfrak{H} of \mathfrak{G} is *realisable in \mathfrak{H}* if v has exactly one outgoing edge labelled by j in \mathfrak{H} , for each $j \in C$. We call \mathfrak{H} *realisable* if it has exactly one *source* (a node without incoming edges) of root type and all nodes in \mathfrak{H} are realisable.

A *periodic structure* is a triple $\mathfrak{P} = (B, P, E)$ satisfying the following conditions. First, we take some realisable subgraph \mathfrak{H} of \mathfrak{G} and define the *pre-periodic part* B to be the subgraph of \mathfrak{H} induced by those nodes v in \mathfrak{H} , for which there are no arbitrarily long paths from the source to v . The *periodic part* P is induced by the remaining nodes in \mathfrak{H} . Finally, the *post-periodic part* E comprises a set R of nodes in P , intersecting any directed cycle in P , and a family of *acyclic* subgraphs \mathfrak{H}_v , $v \in R$, of \mathfrak{G} with unique source v and such that all of the \mathfrak{H}_v 's nodes are realisable in \mathfrak{H}_v , for any $v \in R$. Denote by $\bar{\mathfrak{P}} = (\bar{B}, \bar{P}, \bar{E})$ the triple obtained by taking the $\bar{\cdot}$ -closure of the components B , P and E in \mathfrak{P} .

To illustrate, for $k = 1$, in the only periodic structure with non-empty P shown below, B comprises the root segment s_r , P the segment s (with two A -nodes), and E the leaf segment s_l . There are also three ‘degenerate’ periodic structures with empty P and E .



The *acyclic version* of a rooted digraph G is constructed as follows. We consider each path π starting in the root and ending at the first repeating node v on π with the last edge (u, v) . For all such v and (u, v) , we add to G a fresh node v' and replace (u, v) by (u, v') .

The proof of the following criterion can be found in the full version [29]:

CLAIM 9.2. *The d-sirup (Δ_q, G) is FO-rewritable iff, for any periodic structure $\mathfrak{P} = (B, P, E)$ with $P \neq \emptyset$, one of the following holds:*

- (h1) *there is a homomorphism from some cactus to the $\bar{\cdot}$ -closure of the acyclic version of $B \cup P$;*
- (h2) *there is a homomorphism from the root segment of some cactus to \bar{P} ;*
- (h3) *there is a homomorphism from the root segment of one of the \mathfrak{H}_v to \bar{E} .*

On the other hand, we have the following claim, which is proved in [29] and establishes an FO/L-hardness dichotomy of d-sirups with a Λ -CQ:

CLAIM 9.3. *If none of conditions (h1)–(h3) holds, then evaluating (Δ_q, G) is L-hard.*

In the full version [29], we show that checking the criterion of Claim 9.2 for Λ -CQs of span k can be done in time $p(|q|)2^{p'(k)}$, for some polynomials p and p' . \square

As a consequence of Theorems 7 and 9, we obtain the dichotomy:

COROLLARY 10. *Any d-sirup (Δ_q, G) with a ditree 1-CQ q is either FO-rewritable or L-hard. Deciding this dichotomy, parameterised by the number of solitary T -nodes in CQs, is fixed-parameter tractable.*

This result is in sharp contrast to 2EXPTIME-completeness of deciding FO-rewritability of d-sirups (Δ_q, G) with a dag 1-CQ q

having two solitary T -nodes. We hope that, using the techniques of [22, 31], this dichotomy can be extended to a complete FO/L/NL/P-tetrachotomy of all d-sirups with a ditree 1-CQ. As a first step, we obtain the following trichotomy:

THEOREM 11. *For any a ditree CQ q with one solitary F and one solitary T , (Δ_q, G) is either FO-rewritable, or L-complete, or NL-complete. Deciding this trichotomy can be done in polynomial time.*

PROOF. We use the results of [22] listed as items (c) and (d) on page 12. Let t and f be the solitary T - and F -nodes in q . If (t, f) is $<$ -comparable then (Δ_q, G) is NL-complete by (c) and Theorem 7 (i). If q is quasi-symmetric, then (Δ_q, G) is in L by (d); L-hardness is shown in [29] by a reduction of graph reachability (using a construction that is similar to the one in the proof of Theorem 7).

Otherwise, we consider two models \mathcal{I} over the structure $\mathcal{H}^{(t,f)}$ (defined in the proof-sketch of Theorem 7 (ii)): one has both ‘contacts’ in $F^{\mathcal{I}}$, the other in $T^{\mathcal{I}}$. We check whether there exists a homomorphism from q to either of these models: If neither, then (Δ_q, G) is NL-hard by the the proof of Theorem 7 (ii). If at least one of them is possible, then (Δ_q, G) is FO-rewritable by Prop. 2 (as one can use the $q \rightarrow \mathcal{I}$ homomorphism to define homomorphisms from some depth ≤ 2 cactus to any larger cactus). Details can be found in the full paper [29]. \square

5 CONCLUSIONS

In this paper, we settled the long-standing open problem on the complexity of deciding boundedness of monadic single rule datalog programs. Namely, we proved this problem to be 2EXPTIME-complete—that is, as hard as deciding program boundedness of arbitrary monadic datalog programs [18]. The main innovation of our proof is that we look at the computations of ATMs and the expansions of sirups through the lens of Boolean circuits and show how these circuits can be ‘implemented’ in dag-shaped CQs to verify the correctness of computations encoded by the expansions.

We obtained this result while trying to classify a somewhat different type of basic recursive programs called monadic disjunctive sirups. The disjunctive rule $F(x) \vee T(x) \leftarrow A(x)$ can make answering a Boolean CQ it mediates in the d-sirup range between AC^0 and $coNP$. Deciding FO-rewritability of monadic d-sirups (as well as of Schema.org and *DL-Lite_{bool}* ontology-mediated queries) was shown to be between 2EXPTIME and 2NEXPTIME, and so a complete classification of monadic d-sirups according to their data complexity can be as illusory as the classification of monadic sirups, which has been challenging the datalog community since the 1980s.

On the other hand, this paper shows that d-sirups with ditree CQs are less impenetrable, and we believe a complete classification is possible, though it could be quite tricky and laborious. This problem as well as pinpointing the exact complexity of deciding FO-rewritability of monadic d-sirups (2EXPTIME vs 2NEXPTIME) are left for future work.

ACKNOWLEDGMENTS

This work was supported by the UK EPSRC grant EP/S032282, HSE University Basic Research Program, and Russian Science Foundation 20-11-20203 (Section 4). Thanks are due to the anonymous reviewers for their comments and constructive suggestions.

REFERENCES

- [1] Serge Abiteboul. 1989. Boundedness is Undecidable for Datalog Programs with a Single Recursive Rule. *Inf. Process. Lett.* 32, 6 (1989), 281–287. [https://doi.org/10.1016/0020-0190\(89\)90019-7](https://doi.org/10.1016/0020-0190(89)90019-7)
- [2] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of Databases*. Addison-Wesley. <http://webdam.inria.fr/Alice/>
- [3] Foto N. Afrati, Manolis Gergatsoulis, and Francesca Toni. 2003. Linearisability on datalog programs. *Theor. Comput. Sci.* 308, 1-3 (2003), 199–226. [https://doi.org/10.1016/S0304-3975\(02\)00730-2](https://doi.org/10.1016/S0304-3975(02)00730-2)
- [4] Foto N. Afrati and Christos H. Papadimitriou. 1993. The Parallel Complexity of Simple Logic Programs. *J. ACM* 40, 4 (1993), 891–916. <https://doi.org/10.1145/153724.153752>
- [5] Alessandro Artale, Diego Calvanese, Roman Kontchakov, and Michael Zakharyashev. 2009. The DL-Lite Family and Relations. *J. Artif. Intell. Res.* 36 (2009), 1–69. <https://doi.org/10.1613/jair.2820>
- [6] Pablo Barceló, Gerald Berger, Carsten Lutz, and Andreas Pieris. 2018. First-Order Rewritability of Frontier-Guarded Ontology-Mediated Queries. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, Jérôme Lang (Ed.). ijcai.org, 1707–1713. <https://doi.org/10.24963/ijcai.2018/236>
- [7] Michael Benedikt, Pierre Bourhis, Georg Gottlob, and Pierre Senellart. 2020. Monadic Datalog, Tree Validity, and Limited Access Containment. *ACM Trans. Comput. Log.* 21, 1 (2020), 6:1–6:45.
- [8] Michael Benedikt, Pierre Bourhis, and Pierre Senellart. 2012. Monadic Datalog Containment. In *Automata, Languages, and Programming - 39th International Colloquium, ICALP 2012, Warwick, UK, July 9-13, 2012, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 7392)*, Artur Czumaj, Kurt Mehlhorn, Andrew M. Pitts, and Roger Wattenhofer (Eds.). Springer, 79–91. https://doi.org/10.1007/978-3-642-31585-5_11
- [9] Michael Benedikt, Balder ten Cate, Thomas Colcombet, and Michael Vanden Boom. 2015. The Complexity of Boundedness for Guarded Logics. In *30th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2015, Kyoto, Japan, July 6-10, 2015*. IEEE Computer Society, 293–304. <https://doi.org/10.1109/LICS.2015.36>
- [10] Meghyn Bienvenu, Peter Hansen, Carsten Lutz, and Frank Wolter. 2016. First Order-Rewritability and Containment of Conjunctive Queries in Horn Description Logics. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, Subbarao Kambhampati (Ed.). IJCAI/AAAI Press, 965–971. <http://www.ijcai.org/Abstract/16/141>
- [11] Meghyn Bienvenu, Balder ten Cate, Carsten Lutz, and Frank Wolter. 2014. Ontology-Based Data Access: A Study through Disjunctive Datalog, CSP, and MMSNP. *ACM Trans. Database Syst.* 39, 4 (2014), 33:1–33:44. <https://doi.org/10.1145/2661643>
- [12] Henrik Björklund, Wim Martens, and Thomas Schwentick. 2008. Optimizing Conjunctive Queries over Trees Using Schema Information. In *Mathematical Foundations of Computer Science 2008, 33rd International Symposium, MFCS 2008, Torun, Poland, August 25-29, 2008, Proceedings (Lecture Notes in Computer Science, Vol. 5162)*, Edward Ochmanski and Jerzy Tyszkiewicz (Eds.). Springer, 132–143. https://doi.org/10.1007/978-3-540-85238-4_10
- [13] Henrik Björklund, Wim Martens, and Thomas Schwentick. 2018. Conjunctive query containment over trees using schema information. *Acta Informatica* 55, 1 (2018), 17–56. <https://doi.org/10.1007/s00236-016-0282-1>
- [14] Pierre Bourhis and Carsten Lutz. 2016. Containment in Monadic Disjunctive Datalog, MMSNP, and Expressive Description Logics. In *Principles of Knowledge Representation and Reasoning: Proceedings of the Fifteenth International Conference, KR 2016, Cape Town, South Africa, April 25-29, 2016*, Chitta Baral, James P. Delgrande, and Frank Wolter (Eds.). AAAI Press, 207–216. <http://www.aaai.org/ocs/index.php/KR/KR16/paper/view/12847>
- [15] Diego Calvanese, Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini, and Riccardo Rosati. 2007. Tractable Reasoning and Efficient Query Answering in Description Logics: The DL-Lite Family. *J. Autom. Reason.* 39, 3 (2007), 385–429. <https://doi.org/10.1007/s10817-007-9078-x>
- [16] Chandra Chekuri and Anand Rajaraman. 2000. Conjunctive query containment revisited. *Theor. Comput. Sci.* 239, 2 (2000), 211–229. [https://doi.org/10.1016/S0304-3975\(99\)00220-0](https://doi.org/10.1016/S0304-3975(99)00220-0)
- [17] C. Civili and R. Rosati. 2012. A Broad Class of First-Order Rewritable Tuple-Generating Dependencies. In *Proc. of the 2nd Int. Datalog 2.0 Workshop (Lecture Notes in Computer Science, Vol. 7494)*. Springer, 68–80.
- [18] Stavros S. Cosmadakis, Haim Gaifman, Paris C. Kanellakis, and Moshe Y. Vardi. 1988. Decidable Optimization Problems for Database Logic Programs (Preliminary Report). In *STOC*. 477–490.
- [19] Stavros S. Cosmadakis and Paris C. Kanellakis. 1986. Parallel Evaluation of Recursive Rule Queries. In *Proceedings of the Fifth ACM SIGACT-SIGMOD Symposium on Principles of Database Systems, March 24-26, 1986, Cambridge, Massachusetts, USA*, Avi Silberschatz (Ed.). ACM, 280–293. <https://doi.org/10.1145/6012.15421>
- [20] Cristina Feier, Antti Kuusisto, and Carsten Lutz. 2019. Rewritability in Monadic Disjunctive Datalog, MMSNP, and Expressive Description Logics. *Logical Methods in Computer Science* 15, 2 (2019). [https://doi.org/10.23638/LMCS-15\(2:15\)2019](https://doi.org/10.23638/LMCS-15(2:15)2019)
- [21] Haim Gaifman, Harry G. Mairson, Yehoshua Sagiv, and Moshe Y. Vardi. 1987. Undecidable Optimization Problems for Database Logic Programs. In *Proceedings of the Symposium on Logic in Computer Science (LICS '87), Ithaca, New York, USA, June 22-25, 1987*. IEEE Computer Society, 106–115.
- [22] Olga Gerasimova, Stanislav Kikot, Agi Kurucz, Vladimir V. Podolskii, and Michael Zakharyashev. 2020. A Data Complexity and Rewritability Tetrachotomy of Ontology-Mediated Queries with a Covering Axiom. In *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning, KR 2020, Rhodes, Greece, September 12-18, 2020*, Diego Calvanese, Esra Erdem, and Michael Thielscher (Eds.). 403–413. <https://doi.org/10.24963/kr.2020/41>
- [23] Georg Gottlob, Giorgio Orsi, and Andreas Pieris. 2014. Query Rewriting and Optimization for Ontological Databases. *ACM Trans. Database Syst.* 39, 3 (2014), 25:1–25:46. <https://doi.org/10.1145/2638546>
- [24] André Hernich, Carsten Lutz, Ana Ozaki, and Frank Wolter. 2015. Schema Joint as a Description Logic. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, Qiang Yang and Michael J. Wooldridge (Eds.). AAAI Press, 3048–3054. <http://ijcai.org/Abstract/15/430>
- [25] Gerd G. Hillebrand, Paris C. Kanellakis, Harry G. Mairson, and Moshe Y. Vardi. 1991. Tools for Datalog Boundedness. In *Proceedings of the Tenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, May 29-31, 1991, Denver, Colorado, USA*, Daniel J. Rosenkrantz (Ed.). ACM Press, 1–12. <https://doi.org/10.1145/113413.113414>
- [26] Gerd G. Hillebrand, Paris C. Kanellakis, Harry G. Mairson, and Moshe Y. Vardi. 1995. Undecidable Boundedness Problems for Datalog Programs. *J. Log. Program.* 25, 2 (1995), 163–190. [https://doi.org/10.1016/0743-1066\(95\)00051-K](https://doi.org/10.1016/0743-1066(95)00051-K)
- [27] Mark Kaminski, Yavor Nenov, and Bernardo Cuenca Grau. 2016. Datalog rewritability of Disjunctive Datalog programs and non-Horn ontologies. *Artif. Intell.* 236 (2016), 90–118. <https://doi.org/10.1016/j.artint.2016.03.006>
- [28] Paris C. Kanellakis. 1990. Elements of Relational Database Theory. In *Handbook of Theoretical Computer Science, Volume B: Formal Models and Semantics*, Jan van Leeuwen (Ed.). Elsevier and MIT Press, 1073–1156. <https://doi.org/10.1016/b978-0-444-88074-1.50022-6>
- [29] Stanislav Kikot, Agi Kurucz, Vladimir Podolskii, and Michael Zakharyashev. 2021. Deciding boundedness of monadic sirups. (2021). Full version available at <https://www.dcs.bbk.ac.uk/~michael/PODS-21.pdf>.
- [30] Mélanie König, Michel Leclère, Marie-Laure Mugnier, and Michaël Thomazo. 2015. Sound, complete and minimal UCQ-rewriting for existential rules. *Semantic Web* 6, 5 (2015), 451–475. <https://doi.org/10.3233/SW-140153>
- [31] Carsten Lutz and Leif Sabellek. 2017. Ontology-Mediated Querying with the Description Logic EL: Trichotomy and Linear Datalog Rewritability. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, Carles Sierra (Ed.). ijcai.org, 1181–1187. <https://doi.org/10.24963/ijcai.2017/164>
- [32] Carsten Lutz and Leif Sabellek. 2019. A Complete Classification of the Complexity and Rewritability of Ontology-Mediated Queries based on the Description Logic EL. *CoRR* abs/1904.12533 (2019). [arXiv:1904.12533](http://arxiv.org/abs/1904.12533) <http://arxiv.org/abs/1904.12533>
- [33] Jerzy Marcinkowski. 1999. Achilles, Turtle, and Undecidable Boundedness Problems for Small DATALOG Programs. *SIAM J. Comput.* 29, 1 (1999), 231–257. <https://doi.org/10.1137/S0097539797322140>
- [34] Jeffrey F. Naughton. 1986. Data Independent Recursion in Deductive Databases. In *Proceedings of the Fifth ACM SIGACT-SIGMOD Symposium on Principles of Database Systems, March 24-26, 1986, Cambridge, Massachusetts, USA*, Avi Silberschatz (Ed.). ACM, 267–279. <https://doi.org/10.1145/6012.15420>
- [35] Antonella Poggi, Domenico Lembo, Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Riccardo Rosati. 2008. Linking Data to Ontologies. *J. Data Semant.* 10 (2008), 133–173. https://doi.org/10.1007/978-3-540-77688-8_5
- [36] Ron van der Meyden. 2000. Predicate Boundedness of Linear Monadic Datalog is in PSPACE. *Int. J. Found. Comput. Sci.* 11, 4 (2000), 591–612. <https://doi.org/10.1142/S0129054100000351>
- [37] Moshe Y. Vardi. 1988. Decidability and Undecidability Results for Boundedness of Linear Recursive Queries. In *Proceedings of the Seventh ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, March 21-23, 1988, Austin, Texas, USA*, Chris Edmondson-Yurkanan and Mihalis Yannakakis (Eds.). ACM, 341–351. <https://doi.org/10.1145/308386.308470>
- [38] Guohui Xiao, Roman Kontchakov, Domenico Lembo, Antonella Poggi, Riccardo Rosati, and Michael Zakharyashev. 2018. Ontology-Based Data Access: A Survey. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, Jérôme Lang (Ed.). ijcai.org, 5511–5519. <https://doi.org/10.24963/ijcai.2018/777>
- [39] Guohui Xiao, Linfang Ding, Benjamin Cogrel, and Diego Calvanese. 2019. Virtual Knowledge Graphs: An Overview of Systems and Use Cases. *Data Intell.* 1, 3 (2019), 201–223. https://doi.org/10.1162/dint_a_00011

of **(leaf)**, there is a homomorphism from q_{TT}^- to C mapping q_{TT}^- into s .

So in both cases (i) and (ii), we have shown that there is a segment s of depth $\leq K$ in C^s such that **(correct)** holds. However, s is not necessarily in the branch \mathcal{B} . Let s_m be the last ancestor of s in \mathcal{B} , and list the segments $s = s_0, s_1, \dots, s_m$ on the path leading upwards from s to s_m . Let C' be obtained from C by cutting at s_i every branch of C^s going through s_i other than the one going to s , for every $i \leq m$. (In particular, \mathcal{B} is cut at s_m which is of depth $\leq K$.) Let s_i^* denote the segment corresponding to s_i in C' . Then s_0^* is a leaf in C' , and so $s_0^* = q_{TT}^-$. Also, for each $i > 0$,

- either $s_i^* = s_i$
- or $s_i = q_{AA}^-$ and s_i^* is either q_{AT}^- or q_{TA}^- .

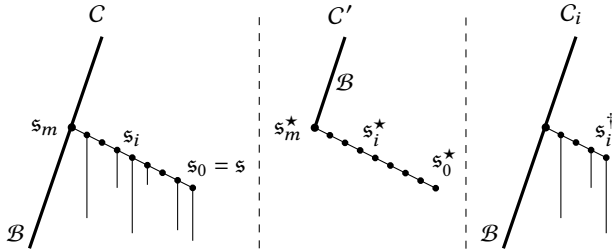
We claim that, for every $i \leq m$, there is some homomorphism $h_i: s_i^* \rightarrow C$ mapping s_i^* into s_i and such that

if s_{i-1}^* is the j -child of s_i^* , for $j = 0, 1$, then

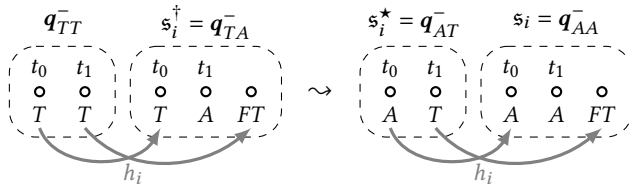
$$h_i \text{ maps the } t_j\text{-node of } s_i^* \text{ to the } t_j\text{-node of } s_i. \quad (10)$$

This will be enough for building a homomorphism from C' to C : we take these h_i on each s_i^* , and the isomorphism on any other segment.

Indeed, if $i = 0$ then the h_0 in **(correct)** is suitable. If $i > 0$ and $s_i^* = s_i$, then the isomorphism is suitable for h_i . So suppose that $s_i^* \neq s_i$ (so $s_i = q_{AA}^-$). We consider the case when $s_i^* = q_{AT}^-$, that is, s_{i-1}^* is a 0-child of s_i^* (the case when $s_i^* = q_{TA}^-$ is similar). Let C_i be obtained from C by cutting at s_i the branch leading to s . Let s_i^\dagger denote the segment corresponding to s_i in C_i , that is, $s_i^\dagger = q_{TA}^-$.



By **(correct)**, s_i is correct in C^s , and so s_i is properly branching in C^s . Thus, s_i^\dagger is incorrect in C_i^s because it violates condition **(pb1)** in Sec. 3.3.2. On the other hand, s_i^* is correct in C_i^s in all the other aspects (this is because apart from s_i and some of its descendants, every other segment is the same in both cactuses C and C_i). Therefore, by the (\Leftarrow) direction of **(leaf)**, there is a homomorphism $h_i: q_{TT}^- \rightarrow C_i$ mapping q_{TT}^- into s_i^\dagger . Also, by **(branch)**, the same h_i is a homomorphism from s_i^* to C , mapping s_i^* to s_i and such that (10) holds:



So in any case we showed that there exists a $C' \rightarrow C$ homomorphism, for some subcactus C' of C where branch \mathcal{B} is cut at some depth $\leq K$. If C' still has branches longer than K , we repeat

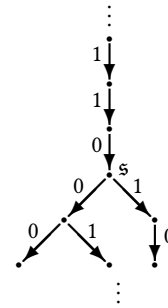
the above process for a long branch in C' to obtain a $C'' \rightarrow C'$ homomorphism for some C'' , and so on. At the end, we obtain a cactus C^- of depth $\leq K$ homomorphically mapping into C , which completes the proof of Lemma 4.

B PROOF OF CLAIM 4.2

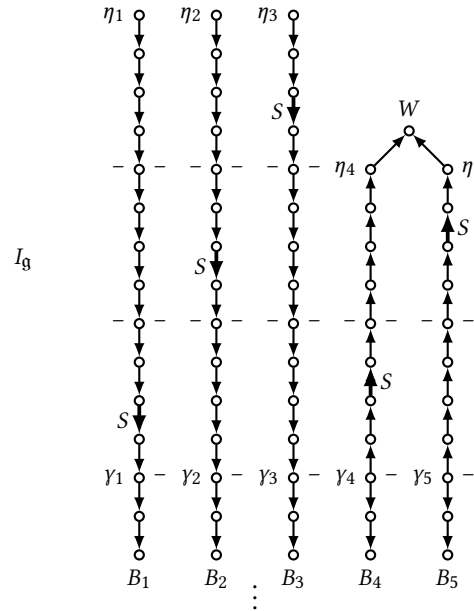
(\Rightarrow) Suppose that, for some C and s , a gadget g implementing a formula $\varphi_g(y_1, \dots, y_n)$ is triggered at s . Then there is a $h: q_{TT}^- \rightarrow C$ homomorphism mapping the I_g -block in q_{TT}^- to the M_g -block in s . In particular, $h(t_g) = \alpha$, and so $h(\pi_g) = \varrho_g$. Thus, for every $i \leq n$, the B_i -node in I_g must also be mapped to one of the two B_i -nodes in the M_g -block of s (either β_i^T or β_i^F). However, which of these two B_i -nodes is the image depends on the truth-value b_i^s of the gathered input $b_g^s = (b_1^s, \dots, b_n^s)$ on the variable y_i . We claim that

- (i) if $b_i^s = 0$, then the B_i -node in I_g is mapped by h to β_i^F ;
- (ii) if $b_i^s = 1$, then the B_i -node in I_g is mapped by h to β_i^T .

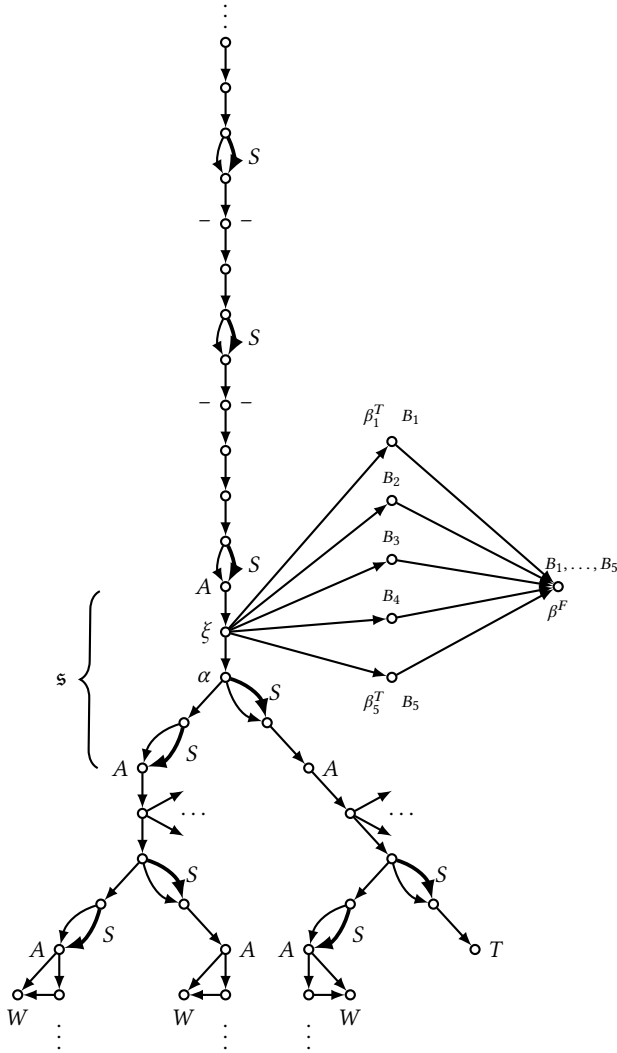
Instead of proving (i) and (ii), here we give an illustrative example. Suppose $\varphi_g(y_1, \dots, y_5)$ is such that (y_1, y_2, y_3) should be gathered from the 3-long uppath, and (y_4, y_5) from a 2-long downpath. Suppose the 'environment' of s in C^s looks like this:



Then if h is a homomorphism triggering g at s , then the possible inputs b_g^s that can be gathered are 01100, 01101, or 01110, because h should map the pattern



to the pattern shown below:



(We are also using that the parts of gadgets that are not depicted above do not contain W -nodes, so the h -image cannot ‘stray’ there when taking a downpath.)

It remains to see how h maps the remaining part of the I_g -block into the M_g -block of \mathfrak{s} . We claim that for every non-leaf gate g in φ_g , if g_{ij}^ℓ is an occurrence of g on some branch, then the end-node p_{ij}^ℓ of the RSR -pattern corresponding to g_{ij}^ℓ in I_g is mapped in such a way that

- (iii) $h(p_{ij}^\ell)$ is the \mathfrak{o} -node of the gadget for g , whenever the value of g under $\mathbf{b}_g^\mathfrak{s}$ is 0;
- (iv) $h(p_{ij}^\ell)$ is the (D) -node of the gadget for g , whenever the value of g under $\mathbf{b}_g^\mathfrak{s}$ is 1.

We prove this by induction on the tree-structure of φ_g , going from leaves to root. Take some gate g , and let g_{ij}^ℓ be an occurrence of g .

First, suppose that g is an AND-gate. There are many cases, depending on the truth-values of g and its two inputs g_1 and g_2

under $\mathbf{b}_g^\mathfrak{s}$, and also on whether each of the g_i is a leaf gate or not. We consider just two cases, the other ones are similar.

- Suppose that the value of g under $\mathbf{b}_g^\mathfrak{s}$ is 0, $\ell = 1$ (and so g_1 is a leaf labelled by y_i), and $b_i^\mathfrak{s} = 1$. Suppose that g_2 is also a leaf gate, and so g_2 has value 0 under $\mathbf{b}_g^\mathfrak{s}$. Let $g_{i'j'}^1$ be an occurrence of g_2 . By (ii), the B_{ij} -node in I_g is mapped by h to the upper B_{ij} -node in the M_g -block of \mathfrak{s} . So the first R -edge of the RSR -pattern corresponding to g_{ij}^1 is mapped to the R -edge connecting the two B_{ij} -nodes. Thus, the S -edge of the RSR -pattern corresponding to g_{ij}^1 must be mapped to an S -edge starting at the i_1 -node of the g -gadget. Similarly, by (i), the $B_{i'j'}$ -node in I_g is mapped by h to the lower $B_{i'j'}$ -node in the M_g -block of \mathfrak{s} . So the S -edge of the RSR -pattern corresponding to $g_{i'j'}^1$ must be mapped to an S -edge following an R -edge starting at the i_2 -node of the g -gadget. As h preserves E , the end-nodes of these two S -edges in the g -gadget must coincide, and so it must be node c_1 . So $h(p_{ij}^1)$ is the \mathfrak{o} -node of the g -gadget.
- Suppose that the value of g under $\mathbf{b}_g^\mathfrak{s}$ is 1, and both of its inputs are non-leaf gates having value 1 under $\mathbf{b}_g^\mathfrak{s}$. Suppose $g_{ij}^{\ell-1}$ is an occurrence of g_1 and $g_{i'j'}^{\ell-1}$ is an occurrence of g_2 . By the IH, $h(p_{ij}^{\ell-1})$ is the (D) -node of the gadget for g_1 , and $h(p_{i'j'}^{\ell-1})$ is the (D) -node of the gadget for g_2 . Then the S -edges of the RSR -patterns corresponding to $g_{ij}^{\ell-1}$ and $g_{i'j'}^{\ell-1}$ must be mapped, respectively, to S -edges starting at the i_1 - and i_2 -nodes of the g -gadget. As h preserves E , the end-nodes of these two S -edges in the g -gadget must coincide, and so it must be node b . So $h(p_{ij}^\ell)$ is the (D) -node of the g -gadget, as required.

The case when g is a NOT-gate can be handled similarly, thereby completing the proof of (iii) and (iv). As h preserves D , it follows that $\varphi_g[\mathbf{b}_g^\mathfrak{s}] = 1$.

(\Leftarrow) If there is $\mathbf{b}_g^\mathfrak{s}$ such that $\mathbf{b}_g^\mathfrak{s}$ is gathered from ‘around’ \mathfrak{s} in $C^\mathfrak{s}$ according to the input-types for φ_g and $\varphi_g[\mathbf{b}_g^\mathfrak{s}] = 1$, then we define a function $h: \mathbf{q}_{TT}^- \rightarrow C$ by taking

- $h(\alpha) = \tau_g$ for the τ_g -node of \mathfrak{s} ,
- $h(i_g) = \alpha$ for the α -node of \mathfrak{s} ,

and mapping

- the I_g -block to the M_g -block of \mathfrak{s} following the structure of $\mathbf{b}_g^\mathfrak{s}$ and φ_g as described above,
- the I_{g_i} -block of every gadget g_i different from g to the I_{g_i} -block of \mathfrak{s} .
- the M_{g_i} -block of every gadget g_i to the M'_{g_i} -block of \mathfrak{s} , and
- the M'_{g_i} -block of every gadget g_i also to the M'_{g_i} -block of \mathfrak{s} .

Using the interaction-regulating mechanism between different gadgets described in Sec. 3.5.1, it is easy to see that h is a homomorphism, and g is triggered by h at \mathfrak{s} .