

On First-Order Query Rewriting for Incomplete Database Histories

Alexandre Decan* and Jef Wijsen

Université de Mons-Hainaut, Mons, Belgium

Abstract. Multiwords are defined as words in which single symbols can be replaced by nonempty sets of symbols. Such a set of symbols captures uncertainty about the exact symbol. Words are obtained from multiwords by selection a single symbol from every set. A pattern is *certain* in a multiword W if it occurs in every word that can be obtained from W . For a given pattern, we are interested in finding a logic formula that recognizes the multiwords in which that pattern is certain.

This problem can be seen as a special case of consistent query answering (CQA). We show how our results can be applied in CQA on database histories under primary key constraints.

1 Motivation

An *incomplete database* DB is a set of *possible* databases. Under the *certain answer* semantics, a Boolean query q evaluates to true on such an incomplete database DB if q evaluates to true on every possible database in DB ; otherwise q evaluates to false.

A database that violates primary key constraints naturally gives rise to an incomplete database: every possible database is obtained by selecting a maximal number of tuples from each relation without ever selecting two distinct tuples that agree on their primary key. In this setting, possible databases are called *repairs*, as they are obtained by “repairing” constraint violations [1].

Recently, progress has been made in computing certain answers to conjunctive queries under primary key constraints [2, 4, 9]. This article makes a first step in extending these works to database histories.

Assume a discrete, linear time scale. The following database history shows the *WorksFor* relation at four successive time points t_0, t_1, t_2, t_3 . The primary key *Name* is violated at times t_1 and t_2 .

t_0	$\frac{\text{Name Company}}{\text{Ed IBM}}$	t_1	$\frac{\text{Name Company}}{\text{Ed MS}}$ $\frac{\text{Ed IBM}}$	t_2	$\frac{\text{Name Company}}{\text{Ed MS}}$ $\frac{\text{Ed IBM}}$	t_3	$\frac{\text{Name Company}}{\text{Ed MS}}$
-------	---	-------	--	-------	--	-------	--

* Corresponding author. Author’s address: Université de Mons-Hainaut, Institut d’Informatique, Avenue du Champ de Mars 6, B-7000 Mons, Belgium; Email: alexandre.decan@umh.ac.be.

To make this history consistent, we have to delete one of the two tuples at t_1 , and one of the two tuples at t_2 . Thus, there are four repairs for this database history; one such repair is shown next.

t_0	<u>Name</u> Company	t_1	<u>Name</u> Company	t_2	<u>Name</u> Company	t_3	<u>Name</u> Company
	Ed IBM		Ed IBM		Ed IBM		Ed MS

The query “*Did MS recruit an IBM employee?*” can be expressed as the existential closure of the following formula in first-order temporal logic (FOTL):¹

$$q_0 = \text{WorksFor}(\underline{x}, \text{IBM}) \wedge \bigcirc \text{WorksFor}(\underline{x}, \text{MS})$$

The primary key arguments are underlined. It can be easily checked that q_0 evaluates to true on each of the four repairs.

We are interested in the following problem: find a FOTL query q'_0 such that for every (possibly inconsistent) database history **HIS**, q'_0 evaluates to true on **HIS** if and only if q_0 evaluates to true on every repair of **HIS**. Such a query q'_0 , if it exists, is called a *consistent FOTL-rewriting* of q_0 . The interest of consistent FOTL-rewritings should be clear: they provide certain answers on inconsistent databases, without the need for computing repairs.

It can be verified that the following query q'_0 is a consistent FOTL-rewriting of q_0 :

$$q'_0 = \varphi_{\text{IBM}} \wedge (\varphi_{\{\text{IBM}, \text{MS}\}} \text{ until } \varphi_{\text{MS}})$$

where:

$$\begin{aligned} \varphi_{\text{IBM}} &= \exists y(\text{WorksFor}(\underline{x}, y) \wedge \forall y'(\text{WorksFor}(\underline{x}, y') \rightarrow y' = \text{IBM})) \\ \varphi_{\text{MS}} &= \exists y(\text{WorksFor}(\underline{x}, y) \wedge \forall y'(\text{WorksFor}(\underline{x}, y') \rightarrow y' = \text{MS})) \\ \varphi_{\{\text{IBM}, \text{MS}\}} &= \exists y(\text{WorksFor}(\underline{x}, y) \wedge \forall y'(\text{WorksFor}(\underline{x}, y') \rightarrow y' = \text{IBM} \vee y' = \text{MS})) \end{aligned}$$

Intuitively, it is certain that some employee x moved directly from IBM to MS if in the (possibly inconsistent) database history, there is one state where x worked certainly for IBM, some later state where x certainly worked for MS, and between these two states, x was employed by either IBM or MS.

Clearly, on databases that satisfy the primary key **Name**, q_0 is equivalent to the following query \bar{q}_0 :

$$q_0 \equiv \bar{q}_0 = \varphi_{\text{IBM}} \wedge \bigcirc \varphi_{\text{MS}}$$

The queries φ_{IBM} , $\varphi_{\{\text{IBM}, \text{MS}\}}$, and φ_{MS} can be obtained by following the rewriting scheme proposed in [2] or [9]. In this article, the focus is on the rewriting from \bar{q}_0 into q'_0 . To capture this rewriting, we look at each employee’s history as a sequence; for example, Ed’s history corresponds to the following sequence:

$$W_{\text{Ed}} = \langle \varphi_{\text{IBM}}, \varphi_{\{\text{IBM}, \text{MS}\}}, \varphi_{\{\text{IBM}, \text{MS}\}}, \varphi_{\text{MS}} \rangle ,$$

¹ Since we deal only with Boolean queries, each variable is understood to be existentially quantified.

expressing that φ_{IBM} holds at time t_0 , $\varphi_{\{\text{IBM}, \text{MS}\}}$ at times t_1 and t_2 , and φ_{MS} at t_3 . Any sequence w obtained from W_{Ed} by replacing each occurrence of $\varphi_{\{\text{IBM}, \text{MS}\}}$ by either $\varphi_{\{\text{IBM}\}}$ or $\varphi_{\{\text{MS}\}}$ then corresponds to a repair. All these repair sequences w satisfy \bar{q}_0 if and only if W_{Ed} satisfies q'_0 .

We thus come to the following abstraction. Assume a finite alphabet Σ . A multiword is a word over the powerset alphabet $2^\Sigma \setminus \{\emptyset\}$. Every multiword W gives rise to a set of possible words obtained from W by selecting one representative from each set. For example, the multiword $W_0 = \langle \{a\}, \{a, b\}, \{a, b\}, \{b\} \rangle$ gives rise to four possible words: $aaab$, $aabb$, $abab$, and $abbb$. A subword w is *certain* in a multiword if it occurs in every possible word. For example, ab is certain in W_0 , but aa is not. Subwords can contain variables, which are placeholders for symbols of Σ . For example, axx is certain in W_0 if either aaa or abb is certain in W_0 . Given a subword w (possibly containing variables), let $\text{CERTAIN}(w)$ denote the set of multiwords in which w is certain. We are interested in the following problem: given a subword w , is $\text{CERTAIN}(w)$ FO-definable?

The article is organized as follows. Section 2 recalls some fundamental results about FO over words. Section 3 formalizes the problem we are interested in. Section 4 gives a straightforward algorithm for deciding $\text{CERTAIN}(w)$ if w contains no variables. The algorithm is interesting because it immediately leads to the result that $\text{CERTAIN}(w)$ is regular, as shown in Section 5. That section also provides sufficient conditions for FO-definability of $\text{CERTAIN}(w)$. Section 6 deals with subwords that contain variables. Finally, Section 7 concludes the article.

2 Recall

We are not aware of existing work on words with uncertain symbols. On the other hand, some fundamental results on standard words will be used in this article:

- Over words, linear temporal logic has the same expressive power as FO [3, 5].
- A regular language is FO-definable if and only if it is aperiodic [7, 8].

Recall that a regular language L is aperiodic if there exists an integer $k \geq 0$ such that for all words p, u, q ,

$$p \cdot u^k \cdot q \in L \text{ if and only if } p \cdot u^{k+1} \cdot q \in L.$$

3 Problem Statement

We assume a finite alphabet $\Sigma = \{a, b, c, \dots\}$ of *symbols*.

Definition 1. A word of length $n \geq 0$ is a sequence $a_1 a_2 \dots a_n$ of symbols. The length of a word w is denoted by $|w|$. The empty word has length 0 and is denoted by ϵ . The concatenation of words v and w is denoted by $v \cdot w$. The concatenation operator naturally extends to sets S, T of words: $S \cdot T = \{v \cdot w \mid v \in S, w \in T\}$. A word v is a subword of w , denoted $w \Vdash v$, if there exist (possibly empty) words p and q such that $w = p \cdot v \cdot q$.

Multiwords capture the notion of inconsistency and are defined next.

Definition 2. A multiword (over Σ) is a sequence $W = \langle A_1, \dots, A_n \rangle$ where for each $i \in \{1, \dots, n\}$, $A_i \subseteq \Sigma$ and $A_i \neq \emptyset$. Thus, a multiword can be conceived as a word (in the sense of Def. 1) relative to the alphabet $2^\Sigma \setminus \{\emptyset\}$.

For the multiword $W = \langle A_1, \dots, A_n \rangle$, we define:

$$\text{words}(W) = \{a_1 a_2 \dots a_n \mid \forall i \in \{1, \dots, n\} : a_i \in A_i\} .$$

If v is a word, then we define:

$$W \Vdash_{\text{certain}} v \text{ if and only if for every } w \in \text{words}(W), w \Vdash v .$$

We write \mathcal{M}_Σ for the set of all multiwords (over Σ).

For example, the following multiword W_0 contains two uncertain positions with values $\{a, b\}$ and $\{c, d\}$. Curly braces are omitted at positions that are certain; for example, $\{a\}$ is written as a .

$$W_0 = \langle a, b, d, a, b, c, a, \{a, b\}, b, d, a, b, \{c, d\}, a, b, c, a, b \rangle \quad (1)$$

We have

$$\begin{aligned} \text{words}(W_0) = \{ & \underline{abdabca} \underline{abdab} \underline{cab} \underline{cab} \\ & \underline{abdabca} \underline{abdab} \underline{cab} \\ & \underline{abdabca} \underline{bbdab} \underline{cab} \underline{cab} \\ & \underline{abdabca} \underline{bbdab} \underline{cab} \} . \end{aligned}$$

The underlined positions show that $abdabcab$ is a subword of each word in $\text{words}(W_0)$. Hence, $W_0 \Vdash_{\text{certain}} abdabcab$.

For a word w , we are interested in (the complexity of) the following language:

$$\text{CERTAIN}(w) := \{W \in \mathcal{M}_\Sigma \mid W \Vdash_{\text{certain}} w\} .$$

In particular, we want to answer the question:

Given w , is $\text{CERTAIN}(w)$ FO-definable?

That is, is there a FO formula φ_w such that

$$\text{CERTAIN}(w) = \{W \in \mathcal{M}_\Sigma \mid W \models \varphi_w\} ,$$

where the logic formalism is FO over words [6, p. 124]?

For example, for the alphabet $\Sigma = \{a, b\}$, multiwords can be regarded as words with symbols taken in the alphabet $\{\{a\}, \{b\}, \{a, b\}\}$. This gives rise to three predicate symbols $P_{\{a\}}$, $P_{\{b\}}$, and $P_{\{a, b\}}$. For instance, the multiword $W = \langle \{a\}, \{a, b\}, \{a, b\}, \{b\} \rangle$ of length 4 is represented by the first-order structure $(\{1, 2, 3, 4\}, <, P_{\{a\}}, P_{\{b\}}, P_{\{a, b\}})$ where $<$ is the natural order and each P_A contains the positions in W at which A occurs: $P_{\{a\}} = \{1\}$, $P_{\{b\}} = \{4\}$, $P_{\{a, b\}} = \{2, 3\}$.

The following formula defines all the multiwords that contain a subsequence among $\langle \{a\}, \{b\} \rangle$, $\langle \{a\}, \{a, b\}, \{b\} \rangle$, $\langle \{a\}, \{a, b\}, \{a, b\}, \{b\} \rangle, \dots$

$$\exists x \exists y (x < y \wedge P_{\{a\}}(x) \wedge P_{\{b\}}(y) \wedge \forall z ((x < z \wedge z < y) \rightarrow P_{\{a, b\}}(z)))$$

It can be verified that this formula defines the language $\text{CERTAIN}(ab)$. Recall that there exists an equivalent formula in linear time logic (LTL), since over words, LTL and FO have the same expressive power.

4 Deciding $\text{CERTAIN}(w)$

Given w , we give a procedure for deciding membership of $\text{CERTAIN}(w)$. The procedure uses a construction defined in the statement of Lemma 2 and its correctness is shown in Theorem 1. The decision procedure is interesting because it can be used to show that $\text{CERTAIN}(w)$ is regular (see Section 5).

Definition 3. *Let w be a word. For all (possibly empty) words p and q , if $w = p \cdot q$, then p is a prefix of w , and q a suffix. If u is a word, then $u^{\rightarrow w}$ denotes the maximal suffix of u that is a prefix of w . For a set S of words, we define $S^{\rightarrow w} = \{u^{\rightarrow w} \mid u \in S\}$.*

Example 1. $abcd^{\rightarrow cde} = cd$ and $ab^{\rightarrow c} = \epsilon$. Clearly, $w^{\rightarrow w} = w$.

Lemma 1. *If $u^{\rightarrow w} = q$, then $(u \cdot a)^{\rightarrow w} = (q \cdot a)^{\rightarrow w}$.*

Proof. Assume $u^{\rightarrow w} = q$. We can assume a (possibly empty) word u' such that:

- $u = u' \cdot q$,
- q is a prefix of w , and
- *Maximality:* for every suffix $s \neq \epsilon$ of u' , $s \cdot q$ is not a prefix of w .

We need to show $(u' \cdot q \cdot a)^{\rightarrow w} = (q \cdot a)^{\rightarrow w}$. Obviously, $(q \cdot a)^{\rightarrow w}$ is a suffix of $(u' \cdot q \cdot a)^{\rightarrow w}$. Assume $|(u' \cdot q \cdot a)^{\rightarrow w}| > |(q \cdot a)^{\rightarrow w}|$. Then, there exists a suffix $s \neq \epsilon$ of u' such that $s \cdot q \cdot a$ is a prefix of w , contradicting the above *Maximality* condition. We conclude by contradiction $(u' \cdot q \cdot a)^{\rightarrow w} = (q \cdot a)^{\rightarrow w}$.

The construction for deciding membership of $\text{CERTAIN}(w)$ is described in the statement of Lemma 2. For a given multiword $W = \langle A_1, \dots, A_n \rangle$, we construct a sequence $\langle S_0, \dots, S_n \rangle$ such that $S_0 = \{\epsilon\}$ and for each subsequent S_i , $u \in S_i$ if and only if $u \neq w$ and $u = (p \cdot a)^{\rightarrow w}$ for some $p \in S_{i-1}$ and $a \in A_i$. By Theorem 1, $W \in \text{CERTAIN}(w)$ if and only if the last element S_n is the empty set.

The construction is illustrated next for the multiword

$$W_0 = \langle a, b, d, a, b, c, a, \{a, b\}, b, d, a, b, \{c, d\}, a, b, c, a, b \rangle$$

shown earlier (see equation (1) on page 4) and the word $w = abdabcab$.

	$S_0 = \{\epsilon\}$
$A_1 = \{a\}$	$S_1 = \{a\}$
$A_2 = \{b\}$	$S_2 = \{ab\}$
$A_3 = \{d\}$	$S_3 = \{abd\}$
$A_4 = \{a\}$	$S_4 = \{abda\}$
$A_5 = \{b\}$	$S_5 = \{abdab\}$
$A_6 = \{c\}$	$S_6 = \{abdabc\}$
$A_7 = \{a\}$	$S_7 = \{abdabca\}$
$A_8 = \{a, b\}$	$S_8 = \{a\}$
$A_9 = \{b\}$	$S_9 = \{ab\}$
$A_{10} = \{d\}$	$S_{10} = \{abd\}$
$A_{11} = \{a\}$	$S_{11} = \{abda\}$
$A_{12} = \{b\}$	$S_{12} = \{abdab\}$
$A_{13} = \{c, d\}$	$S_{13} = \{abdabc, abd\}$
$A_{14} = \{a\}$	$S_{14} = \{abdabca, abda\}$
$A_{15} = \{b\}$	$S_{15} = \{abdab\}$
$A_{16} = \{c\}$	$S_{16} = \{abdabc\}$
$A_{17} = \{a\}$	$S_{17} = \{abdabca\}$
$A_{18} = \{b\}$	$S_{18} = \{\}$

Lemma 2. Let $W = \langle A_1, \dots, A_n \rangle$ be a multiword. Let w be a nonempty word. Let $\langle S_0, S_1, \dots, S_n \rangle$ be a sequence such that $S_0 = \{\epsilon\}$ and for every $i \in \{1, \dots, n\}$, $S_i = (S_{i-1} \cdot A_i)^{\rightarrow w} \setminus \{w\}$.

Let $m \in \{1, 2, \dots, n\}$. For every word $u \in \text{words}(\langle A_1, \dots, A_m \rangle)$, either $u^{\rightarrow w} \in S_m$ or $u \Vdash w$.

Proof. Proof by induction on m .

Basis $m = 1$. Let $u = b \in \text{words}(\langle A_1 \rangle)$ such that $b \not\Vdash w$ (i.e. $w \neq b$). Then, S_1 contains $(\epsilon \cdot b)^{\rightarrow w} = b^{\rightarrow w}$.

Step. Assume w.l.o.g. $S_m = \{p_1, \dots, p_k\}$ and $A_{m+1} = \{a_1, \dots, a_l\}$. Then,

$$S_{m+1} = \{(p_1 \cdot a_1)^{\rightarrow w}, \dots, (p_i \cdot a_j)^{\rightarrow w}, \dots, (p_k \cdot a_l)^{\rightarrow w}\} \setminus \{w\} .$$

Let $u \in \text{words}(\langle A_1, \dots, A_{m+1} \rangle)$. We can assume $v \in \text{words}(\langle A_1, \dots, A_m \rangle)$ and $j \in \{1, \dots, l\}$ such that $u = v \cdot a_j$. Two cases can occur:

- $v \Vdash w$. Obviously, $u \Vdash w$.
- $v \not\Vdash w$. By the induction hypothesis, $v^{\rightarrow w} \in S_m$. We can assume w.l.o.g. $i \in \{1, \dots, k\}$ such that $p_i = v^{\rightarrow w}$. By Lemma 1, $u^{\rightarrow w} = (v \cdot a_j)^{\rightarrow w} = (p_i \cdot a_j)^{\rightarrow w}$. If $(p_i \cdot a_j)^{\rightarrow w} = w$, then $u \Vdash w$; otherwise S_{m+1} contains $u^{\rightarrow w}$.

This concludes the proof.

Lemma 3. Let $W = \langle A_1, \dots, A_n \rangle$, w , and $\langle S_0, S_1, \dots, S_n \rangle$ be as in Lemma 2.

Let $i \in \{1, \dots, n\}$. For every $p \in S_i$, there exists $u \in \text{words}(\langle A_1, \dots, A_i \rangle)$ such that $u \not\Vdash w$ and $u^{\rightarrow w} = p$.

Proof. Proof by induction on i .

Basis $i = 1$. Easy.

Step. Assume $p \in S_{i+1}$. We can assume $q \in S_i$ and $a \in A_{i+1}$ such that $p = q \cdot a$.

By the induction hypothesis, there exists $u \in \text{words}(\langle A_1, \dots, A_i \rangle)$ such that $u \not\preceq w$ and $u \rightarrow^w = q$. Then, $u \cdot a \in \text{words}(\langle A_1, \dots, A_{i+1} \rangle)$. It suffices to show that $u \cdot a \not\preceq w$ and that $(u \cdot a) \rightarrow^w = p$.

– Assume $u \cdot a \Vdash w$. Then necessarily, $w = q \cdot a$. But then $p = w$, a contradiction. We conclude by contradiction that $u \cdot a \not\preceq w$.

– By Lemma 1, $(u \cdot a) \rightarrow^w = (q \cdot a) \rightarrow^w = p \rightarrow^w$. Since p is a prefix of w , $p \rightarrow^w = p$.

This concludes the proof.

Theorem 1. Let $W = \langle A_1, \dots, A_n \rangle$, w , and $\langle S_0, S_1, \dots, S_n \rangle$ be as in Lemma 2. Then, $W \in \text{CERTAIN}(w)$ if and only if $S_n = \{\}$.

Proof. *If part.* Assume $S_n = \{\}$. Lemma 2 implies that $W \in \text{CERTAIN}(w)$.

Only-if part. Assume $S_n \neq \{\}$. Lemma 3 implies that $W \notin \text{CERTAIN}(w)$.

As a small optimization, it is not hard to verify that Theorem 1 remains valid if every S_i that contains ϵ is replaced by $\{\epsilon\}$. Formally, define:

$$\lfloor S_i \rfloor = \begin{cases} \{\epsilon\} & \text{if } \epsilon \in S_i ; \\ S_i & \text{otherwise.} \end{cases}$$

Then, the sequence $\langle S_0, \dots, S_n \rangle$ defined by $S_0 = \{\epsilon\}$ and for each $i \in \{1, \dots, n\}$, $S_i = \lfloor (S_{i-1} \cdot A_i) \rightarrow^w \setminus \{w\} \rfloor$ satisfies: $W \in \text{CERTAIN}(w)$ if and only if $S_n = \{\}$.

5 FO-definability

Assume a nonempty word w over alphabet Σ . We show that $\text{CERTAIN}(w)$ is regular and, under certain conditions, aperiodic. Theorem 1 suggests the following construction of a deterministic finite automaton that accepts $\text{CERTAIN}(w)$.

Let P' be the set of prefixes of w , and let $P = P' \setminus \{\epsilon, w\}$. We define $\text{DFA}_\Sigma(w)$ as the following deterministic finite automaton $(Q, \Omega, S_0, F, \delta)$:

1. The finite set Q of states is $2^P \cup \{\{\epsilon\}\}$. Note that $\{\epsilon\}$ is the only state containing ϵ .
2. The alphabet Ω is the ‘‘powerset alphabet’’ $2^\Sigma \setminus \{\emptyset\}$.
3. The start state is $S_0 = \{\epsilon\}$ and $F = \{\emptyset\}$ is the set of accepting states.
4. The transition function $\delta : Q \times \Omega \rightarrow Q$ is defined by:

$$\delta(S, A) = \lfloor (S \cdot A) \rightarrow^w \setminus \{w\} \rfloor .$$

For example, Fig. 1 shows $\text{DFA}_\Sigma(abb)$ for the alphabet $\Sigma = \{a, b\}$. Theorem 1 and the construction of $\text{DFA}_\Sigma(\cdot)$ immediately lead to the following result.

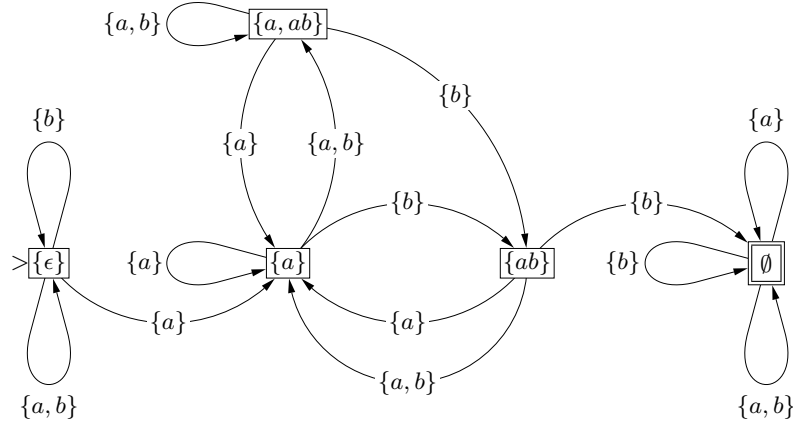


Fig. 1. $\text{DFA}_{\Sigma}(abb)$ with $\Sigma = \{a, b\}$.

Theorem 2. *Let w be a nonempty word and $W = \langle A_1, \dots, A_n \rangle$ a multiword (both relative to the alphabet Σ). Then, $W \in \text{CERTAIN}(w)$ if and only if $\text{DFA}_{\Sigma}(w)$ accepts W . Hence, $\text{CERTAIN}(w)$ is regular.*

A regular language is FO-definable if and only if it is aperiodic [7, 8]. Since our driving motivation is consistent FO-rewriting, we are interested in aperiodicity results for $\text{CERTAIN}(w)$. The following theorem shows that $\text{CERTAIN}(w)$ is aperiodic if the first (or the last) symbol of w occurs only once in w . Our proof technique explicitly relies on the condition that the first (or the last) symbol of w is unique. This uniqueness condition is not necessary for FO-definability, however; for example, $\text{CERTAIN}(aa)$ is FO-definable. We conjecture that $\text{CERTAIN}(w)$ is aperiodic for any word w .

Theorem 3.

1. *If the last symbol of w occurs only once in w , then $\text{CERTAIN}(w)$ is aperiodic.*
2. *If the first symbol of w occurs only once in w , then $\text{CERTAIN}(w)$ is aperiodic.*

Proof. We prove the first item; the proof of the second item is symmetrical. Let $w = w' \cdot a$ where a does not occur in w' . Let $k = |w|$. It suffices to show that for all multiwords P, U, Q :

$$P \cdot U^{k+1} \cdot Q \in \text{CERTAIN}(w) \text{ if and only if } P \cdot U^{k+2} \cdot Q \in \text{CERTAIN}(w).$$

That is, for all multiwords P, U, Q :

$$P \cdot U^{k+1} \cdot Q \not\ll_{\text{certain}} w \iff P \cdot U^{k+2} \cdot Q \not\ll_{\text{certain}} w .$$

Let $U = \langle A_1, \dots, A_n \rangle$. The proof is obvious if $n = 0$. Next assume $n > 0$.

$\boxed{\Rightarrow}$ Assume $P \cdot U^{k+1} \cdot Q \not\ll_{\text{certain}} w$. We can assume the existence of $p \in \text{words}(P)$, $u_1, \dots, u_{k+1} \in \text{words}(U)$, and $q \in \text{words}(Q)$ such that the word

$$m = p \cdot u_1 \cdot u_2 \cdot \dots \cdot u_{k+1} \cdot q$$

satisfies $m \not\ll w$. Two cases can occur:

1. The symbol a does not occur in $u_2 \cdot \dots \cdot u_{k+1}$. Let m' be the word obtained from m by “duplicating” the subword u_2 as follows:

$$m' = \overbrace{p \cdot u_1 \cdot u_2}^{s_1} \cdot \underbrace{u_2 \cdot \dots \cdot u_{k+1} \cdot q}_{t_2} \quad \overbrace{\phantom{u_2 \cdot \dots \cdot u_{k+1} \cdot q}}^{t_1}$$

Assume $m' \Vdash w$. Since $m \not\ll w$, any subword w of m' must “start in” s_1 and “end in” t_1 . Since $|w| = k$, w must actually end in t_2 . But then a must occur in $u_2 \cdot \dots \cdot u_{k+1}$, a contradiction. We conclude by contradiction that $m' \not\ll w$. Since $m' \in \text{words}(P \cdot U^{k+2} \cdot Q)$, $P \cdot U^{k+2} \cdot Q \not\ll_{\text{certain}} w$.

2. a occurs in $u_2 \cdot \dots \cdot u_{k+1}$. Then, we can assume a word v with $|v \cdot a| = |U|$ and words o and r such that:

$$m = p \cdot o \cdot v \cdot a \cdot r \cdot q$$

Let m' be the word obtained from m by “duplicating” the subword $v \cdot a$ as follows:

$$m' = \overbrace{p \cdot o \cdot v \cdot a}^{s_3} \cdot \underbrace{v \cdot a \cdot r \cdot q}_{t_3}$$

Assume $m' \Vdash w$. Since $m \not\ll w$, any subword w of m' must start in s_3 and end in t_3 . But then w contains two occurrences of a , a contradiction. We conclude by contradiction that $m' \not\ll w$. It is easy to verify that $m' \in \text{words}(P \cdot U^{k+2} \cdot Q)$. It follows $P \cdot U^{k+2} \cdot Q \not\ll_{\text{certain}} w$.

$\boxed{\Leftarrow}$ Assume $P \cdot U^{k+2} \cdot Q \not\ll_{\text{certain}} w$. We can assume $p \in \text{words}(P)$, $u_1, \dots, u_{k+2} \in \text{words}(U)$, and $q \in \text{words}(Q)$ such that the word

$$m = p \cdot u_1 \cdot \dots \cdot u_{k+1} \cdot u_{k+2} \cdot q$$

satisfies $m \not\ll w$. Let m' be the word obtained from m by “omitting” the subword u_1 as follows:

$$m' = \overbrace{p}^{s_4} \cdot \underbrace{u_2 \cdot \dots \cdot u_k \cdot u_{k+1} \cdot u_{k+2} \cdot q}_{t_5} \quad \overbrace{\phantom{u_2 \cdot \dots \cdot u_k \cdot u_{k+1} \cdot u_{k+2} \cdot q}}^{t_4}$$

Clearly, $m' \in \text{words}(P \cdot U^{k+1} \cdot Q)$. Two cases can occur:

1. $m' \not\ll w$. Then, $P \cdot U^{k+1} \cdot Q \not\ll_{\text{certain}} w$.

2. $m' \Vdash w$. Since $m \not\ll w$, any subword w of m' must start in s_4 and end in t_4 . Since $|w| = k$, w must actually end in t_5 . Then a must occur in $u_2 \cdot \dots \cdot u_k$. We can assume a word v with $|v| = |U|$ and words o and r such that:

$$m = p \cdot u_1 \cdot o \cdot a \cdot v \cdot r \cdot q$$

Let m'' be the word obtained from m by “omitting” the subword v as follows:

$$m'' = \overbrace{p \cdot u_1 \cdot o \cdot a}^{s_6} \cdot \overbrace{r \cdot q}^{t_6}$$

Assume $m'' \Vdash w$. Since $m \not\ll w$, the subword w of m'' must start in s_6 and end in t_6 . But then w contains two occurrences of a , a contradiction. We conclude by contradiction that $m'' \not\ll w$. It is easy to verify that $m'' \in \text{words}(P \cdot U^{k+1} \cdot Q)$. Consequently, $P \cdot U^{k+1} \cdot Q \not\ll_{\text{certain}} w$.

This concludes the proof.

Since aperiodic regular languages are expressible in FO [7, 8], we have the following result.

Corollary 1. *If the symbol a does not occur in w , then $\text{CERTAIN}(a \cdot w)$ and $\text{CERTAIN}(w \cdot a)$ are FO-definable.*

6 Words With Variables

We define v-words as words in which variables can replace symbols. These variables are placeholders for symbols. For example, the v-word xax represents any word of length 3 in which the first and the last symbol are equal, and the second symbol is a .

We assume a countable set $\mathbf{vars} = \{x, y, z, x_1, y_1, z_1, \dots\}$ of *variables* disjoint from the alphabet Σ .

Definition 4. *A v-word is a sequence $s_1 s_2 \dots s_n$ where for each $i \in \{1, \dots, n\}$, $s_i \in \Sigma \cup \mathbf{vars}$. The notions of length and concatenation (see Def. 1) naturally extend to v-words.*

The relation \Vdash is extended as follows. A valuation is a total function $\theta : \mathbf{vars} \cup \Sigma \rightarrow \Sigma$ that is the identity on Σ . For the v-word $v = s_1 s_2 \dots s_n$, we define $\theta(v) = \theta(s_1) \cdot \theta(s_2) \cdot \dots \cdot \theta(s_n)$. Note that if v contains no variables (i.e. v is a word), then $\theta(v) = v$. If w is a word and v a v-word, then we define:

$$w \Vdash v \text{ if and only if for some valuation } \theta, \theta(v) \text{ is a subword of } w.$$

The relation \Vdash_{certain} directly carries over to v-words: if W is a multiword and v a v-word, then $W \Vdash_{\text{certain}} v$ if and only if for every $w \in \text{words}(W)$, $w \Vdash v$.

The problem statement in Section 3 carries over to v-words. For example, $\text{CERTAIN}(xx)$ contains some multiword W if each $w \in \text{words}(W)$ contains two successive occurrences of the same symbol. Assume a binary alphabet $\Sigma =$

$\{a, b\}$. Clearly, $\text{CERTAIN}(aa) \cup \text{CERTAIN}(bb) \subseteq \text{CERTAIN}(xx)$, and the inclusion is strict: for $W_1 = \langle a, \{a, b\}, b \rangle$, we have $W_1 \notin \text{CERTAIN}(aa) \cup \text{CERTAIN}(bb)$ and $W_1 \in \text{CERTAIN}(xx)$.

Theorem 4 shows that $\text{CERTAIN}(xx)$ is not FO-definable.

Theorem 4. *Let Σ be a finite alphabet with $|\Sigma| \geq 2$. Then,*

1. $\text{CERTAIN}(xx)$ is in \mathbf{P} ;
2. $\text{CERTAIN}(xx)$ is not FO-definable.

Proof. For the first item, let $W = \langle A_1, \dots, A_n \rangle$. We construct, in polynomial time, a sequence $\langle B_1, \dots, B_n \rangle$ where $B_1 = A_1$ and for each $k \in \{2, \dots, n\}$,

$$B_k = \begin{cases} \{\} & \text{if } B_{k-1} = \{\} \\ A_k \setminus B_{k-1} & \text{if } |B_{k-1}| = 1 \\ A_k & \text{otherwise} \end{cases}$$

It is easy to show, by induction on increasing j , that for each $j \in \{1, \dots, n\}$,

$$B_j = \{a \in A_j \mid \langle A_1, \dots, A_{j-1}, \{a\} \rangle \not\in_{\text{CERTAIN}} xx\} .$$

Then, $W \in \text{CERTAIN}(xx)$ if and only if $B_n = \{\}$.

For the second item, for every $i \geq 0$, let

$$W_i = \langle \{a\}, \overbrace{\{a, b\}, \dots, \{a, b\}}^{i \text{ times}}, \{b\} \rangle .$$

It is easy to verify that $W_i \in \text{CERTAIN}(xx)$ if and only if i is odd. The latter condition is not FO-definable.

7 Conclusion

Every multiword defines a set of possible words. A subword w is “certain” in a multiword W if it occurs in every possible word of W . Given w , a significant question is whether $\text{CERTAIN}(w)$, i.e. the set of multiwords in which w is certain, is FO-definable. We showed the following results:

1. $\text{CERTAIN}(w)$ is regular if w is variable-free.
2. $\text{CERTAIN}(w)$ is FO-definable if w is variable-free and the first or the last symbol of w occurs only once in w .
3. $\text{CERTAIN}(xx)$ is not FO-definable.

These results are useful in consistent first-order rewriting under primary keys in historical databases. For example, our results imply that the following FOTL query q_1 (“Did some employee remain in the same company for two successive time points?”) has no consistent first-order rewriting:

$$q_1 = \exists x \exists y (\text{WorksFor}(x, y) \wedge \circ \text{WorksFor}(x, y)) .$$

On the other hand, the following FOTL query q_2 (“Did MS recruit an IBM employee who had worked for IBM for more than k time instants?”) has a consistent first-order rewriting:

$$q_2 = \exists x(\text{WorksFor}(\underline{x}, \text{IBM}) \wedge \circ^1 \text{WorksFor}(\underline{x}, \text{IBM}) \\ \wedge \dots \\ \wedge \circ^k \text{WorksFor}(\underline{x}, \text{IBM}) \wedge \circ^{k+1} \text{WorksFor}(\underline{x}, \text{MS})) ,$$

where $\circ^i \varphi = \overbrace{\circ \circ \dots \circ}^{i \text{ times}} \varphi$. Moreover, there exists an effective procedure for constructing a consistent first-order rewriting of q_2 . This follows from Theorem 3 and the fact that the translation from an aperiodic regular language into FO is effective. So far, we have not studied algorithmic simplifications that may apply in our framework.

The results presented in this article are preliminary. In particular, we are currently addressing the following tasks:

1. Show the conjecture that $\text{CERTAIN}(w)$ is FO-definable if w is variable-free.
2. Show the conjecture that $\text{CERTAIN}(w)$ is regular for any v -word w .
3. Find a syntactic characterization of the v -words w for which $\text{CERTAIN}(w)$ is FO-definable.

References

1. M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In *PODS*, pages 68–79. ACM Press, 1999.
2. A. Fuxman and R. J. Miller. First-order query rewriting for inconsistent databases. *J. Comput. Syst. Sci.*, 73(4):610–635, 2007.
3. D. M. Gabbay, A. Pnueli, S. Shelah, and J. Stavi. On the temporal basis of fairness. In *POPL*, pages 163–173, 1980.
4. L. Grieco, D. Lembo, R. Rosati, and M. Ruzzi. Consistent query answering under key and exclusion dependencies: algorithms and experiments. In O. Herzog, H.-J. Schek, N. Fuhr, A. Chowdhury, and W. Teiken, editors, *CIKM*, pages 792–799. ACM, 2005.
5. J. A. W. Kamp. *Tense Logic and The Theory of Linear Order*. PhD thesis, University of California, Los Angeles, 1968.
6. L. Libkin. *Elements of Finite Model Theory*. Springer, 2004.
7. R. McNaughton and S. Papert. *Counter-free Automata*. MIT Press, Cambridge, MA, 1971.
8. M. P. Schützenberger. On finite monoids having only trivial subgroups. *Information and Control*, 8(2):190–194, 1965.
9. J. Wijsen. On the consistent rewriting of conjunctive queries under primary key constraints. In M. Arenas and M. I. Schwartzbach, editors, *DBPL*, volume 4797 of *Lecture Notes in Computer Science*, pages 112–126. Springer, 2007.