

On First-Order Query Rewriting for Incomplete Database Histories

Véronique Bruyère
Institut d'Informatique
Université de Mons
Mons, Belgium

Alexandre Decan
Institut d'Informatique
Université de Mons
Mons, Belgium

Jef Wijsen
Institut d'Informatique
Université de Mons
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 selecting 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.

Keywords—database repairing; consistent query answering; temporal databases; pattern matching

I. 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.

In relational databases, primary key violations represent uncertainty. 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], [3], [4], [5]. 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	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	IBM	t_1	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>MS</td></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	MS	Ed	IBM
<u>Name</u>	Company												
Ed	IBM												
<u>Name</u>	Company												
Ed	MS												
Ed	IBM												
t_2	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>MS</td></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	MS	Ed	IBM	t_3	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>MS</td></tr></table>	<u>Name</u>	Company	Ed	MS
<u>Name</u>	Company												
Ed	MS												
Ed	IBM												
<u>Name</u>	Company												
Ed	MS												

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	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	IBM	t_1	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	IBM
<u>Name</u>	Company										
Ed	IBM										
<u>Name</u>	Company										
Ed	IBM										
t_2	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>IBM</td></tr></table>	<u>Name</u>	Company	Ed	IBM	t_3	<table><tr><th><u>Name</u></th><th>Company</th></tr><tr><td>Ed</td><td>MS</td></tr></table>	<u>Name</u>	Company	Ed	MS
<u>Name</u>	Company										
Ed	IBM										
<u>Name</u>	Company										
Ed	MS										

The query “Did MS recruit an IBM employee?” is expressed by the following formula in first-order temporal logic (FOTL):

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

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. Moreover, such rewritings may be amenable to implementation in practical database languages like (temporal extensions of) SQL.

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

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

where:

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

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 \tilde{q}_0 :

$$q_0 \equiv \tilde{q}_0 = \exists x(\varphi_{\text{IBM}}(x) \wedge \bigcirc \varphi_{\text{MS}}(x))$$

The queries $\varphi_{\text{IBM}}(x)$, $\varphi_{\{\text{IBM}, \text{MS}\}}(x)$, and $\varphi_{\text{MS}}(x)$ can be obtained using recent results on first-order query rewriting under primary keys [2], [4], [5]. In this article, the focus is on the rewriting from \tilde{q}_0 into q_0 . To study this rewriting, we use the abstraction explained in the next paragraph.

Assume a finite alphabet Σ . A *multiword* is a word over the powerset alphabet $2^\Sigma \setminus \{\emptyset\}$. For the moment, take $\Sigma = \{a, b\}$, so the powerset alphabet is the three-letter alphabet $\{\{a\}, \{b\}, \{a, b\}\}$. Every multiword W gives rise to a set of possible words (or repairs) obtained from W by selecting one representative from each set. For example, the multiword $W_0 = \{\{a\}, \{a, b\}, \{a, b\}, \{b\}\}$ has four possible words: $aaab$, $aabb$, $abab$, and $abbb$. A word w is *certain* in a multiword if it is a factor of every possible word. For example, ab is certain in W_0 , but aa is not (because aa is not a factor of $abbb$). Given a word w , let $\text{CERTAIN}(w)$ denote the set of multiwords in which w is certain. We are interested in the following problem: given a word w , is $\text{CERTAIN}(w)$ FO-definable?

The example of the preceding paragraph is our propositional abstraction of the WorksFor example introduced before. Think of a as IBM, and b as MS. The word ab corresponds to the temporal formula $\psi_0 = a \wedge \bigcirc b$, which abstracts \tilde{q}_0 . Let W be a multiword over the powerset alphabet $\{\{a\}, \{b\}, \{a, b\}\}$. Then, ψ_0 evaluates to true on every possible word of W if and only if W satisfies $\psi'_0 = \{a\} \wedge \bigcirc (\{a, b\} \text{ until } \{b\})$; the latter formula abstracts q_0 .

The article is organized as follows. Section II recalls some fundamental results about logics over words. Section III formalizes the problem we are interested in. Section IV gives an algorithm for deciding $\text{CERTAIN}(w)$ if w contains no variables. The algorithm immediately leads to the result that $\text{CERTAIN}(w)$ is regular, as shown in Section V. That section also provides sufficient conditions for FO-definability of $\text{CERTAIN}(w)$. Section VI deals with words that 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 . Finally, Section VII concludes the article.

II. RELATED WORK

Multiwords can be seen as classical words over some powerset alphabet. Some fundamental results on standard words will be used in this article:

- Over words, linear temporal logic has the same expressive power as FO [6], [7].
- A regular language is FO-definable if and only if it is aperiodic [8], [9].
- A language is definable in MSO (monadic second-order logic) if and only if it is regular (Büchi's theorem).

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

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

Multiwords can be seen as a syntactic generalization of *partial words*, which have been extensively studied by Blanchet-Sadri and others (a recent article is [10]). Partial words are strings over a finite alphabet that may contain one or more occurrences of the “do not know” symbol \diamond . A partial word could be encoded as a multiword, by replacing each symbol a with $\{a\}$, and \diamond with the entire alphabet. For example, over the alphabet $\{a, b, c\}$, the partial word $a\diamond b$ is encoded by the multiword $\{a\} \cdot \{a, b, c\} \cdot \{b\}$. However, the problems studied for partial words are quite different from the problems studied in this article, and semantics diverges.

III. 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 *factor* of w , denoted $w \Vdash v$, if there exist (possibly empty) words p and q such that $w = p \cdot v \cdot q$. \triangleleft

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 $2^\Sigma \setminus \{\emptyset\}$, called *powerset alphabet*.

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:

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

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

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 .

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

We have

$$\begin{aligned} \text{words}(W_0) = \{ & \text{abdabca}\underline{\text{abdabca}}\text{bcab} \\ & \text{abdabca}\underline{\text{abdabca}}\text{bdabca} \\ & \underline{\text{abdabca}}\text{bbdabca}\text{bcab} \\ & \underline{\text{abdabca}}\text{bbdabca}\text{bdabca} \}. \end{aligned}$$

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

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

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

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

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

where M_W is the structure representing W , defined as usual in FO logic over words [11, 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 unary 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

$$M_W = (\{1, 2, 3, 4\}, <, P_{\{a\}}, P_{\{b\}}, P_{\{a, b\}}),$$

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 factor among $\langle \{a\}, \{b\} \rangle$, $\langle \{a\}, \{a, b\}, \{b\} \rangle$, $\langle \{a\}, \{a, b\}, \{a, b\}, \{b\} \rangle, \dots$

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

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.

IV. 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.

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 $\text{sufpre}(u, w)$ denotes the maximal suffix of u that is a prefix of w . For a set S of words, we define $\text{sufpre}(S, w) = \{\text{sufpre}(u, w) \mid u \in S\}$.

Let S be a set of words. We write $\lfloor S \rfloor$ for the smallest set of words satisfying: $\lfloor S \rfloor \subseteq S$ and $\lfloor S \rfloor$ contains a suffix of every word in S . \triangleleft

For example, $\text{sufpre}(abcd, cde) = cd$ and $\text{sufpre}(ab, c) = \epsilon$. If $S = \{aa, ac, abc, bc, c\}$, then $\lfloor S \rfloor = \{c, aa\}$.

Lemma 1: If $\text{sufpre}(u, w) = q$, then $\text{sufpre}(u \cdot a, w) = \text{sufpre}(q \cdot a, w)$.

Proof: Straightforward. \blacksquare

We now give Lemma 2 and then illustrate its construction by an example.

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 $i \in \{1, \dots, n\}$,

$$S_i = \lfloor \text{sufpre}(S_{i-1} \cdot A_i, w) \setminus \{w\} \rfloor.$$

Let $m \in \{1, 2, \dots, n\}$. For every word $u \in \text{words}(\langle A_1, \dots, A_m \rangle)$, S_m contains a suffix of $\text{sufpre}(u, w)$ or $u \models w$.

Proof: Proof by induction on m . For the induction basis, take $m = 1$ and let $u = b \in \text{words}(\langle A_1 \rangle)$ such that $b \not\models w$ (i.e. $w \neq b$). Then, S_1 contains a suffix of $\text{sufpre}(\epsilon \cdot b, w) = \text{sufpre}(b, w)$.

For the induction step, assume w.l.o.g. $S_m = \{p_1, \dots, p_k\}$ and $A_{m+1} = \{a_1, \dots, a_l\}$. Then, $S_{m+1} = \lfloor \{\text{sufpre}(p_1 \cdot a_1, w), \dots, \text{sufpre}(p_i \cdot a_j, w), \dots, \text{sufpre}(p_k \cdot a_l, w)\} \setminus \{w\} \rfloor$. 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 $u = v \cdot a_j$. We need to show that S_{m+1} contains a suffix of $\text{sufpre}(v \cdot a_j, w)$ or $v \cdot a_j \models w$. Two cases can occur:

- $v \models w$. Obviously, $v \cdot a_j \models w$.
- $v \not\models w$. By the induction hypothesis, S_m contains a suffix of $\text{sufpre}(v, w)$. We can assume w.l.o.g. $i \in \{1, \dots, k\}$ such that p_i is a suffix of $\text{sufpre}(v, w)$. We can assume word q such that $\text{sufpre}(v, w) = q \cdot p_i$. By Lemma 1, $\text{sufpre}(v \cdot a_j, w) = \text{sufpre}(q \cdot p_i \cdot a_j, w)$. If $\text{sufpre}(p_i \cdot a_j, w) = w$, then $v \cdot a_j \models w$; otherwise S_{m+1} contains a suffix of $\text{sufpre}(p_i \cdot a_j, w)$. Clearly, every suffix of $\text{sufpre}(p_i \cdot a_j, w)$ is a suffix of $\text{sufpre}(v \cdot a_j, w)$.

This concludes the proof. \blacksquare

Theorem 1 will show that the construction of Lemma 2 solves the membership problem for $\text{CERTAIN}(w)$ as follows: $W \in \text{CERTAIN}(w)$ if and only if the last element of $\langle S_0, S_1, \dots, S_n \rangle$ is the empty set. The construction is illustrated in Fig. 1 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 and the word $w = abdabcab$. The set S_8 , for example, is computed from $S_7 \cdot A_8 = \{abdabcaa, w\}$, in which $abdabcaa$ is replaced with its suffix a , and the word w is removed.

A rough time complexity analysis shows that the construction of the sequence $\langle S_0, S_1, \dots, S_n \rangle$ in Lemma 2 is in $O(|W| \cdot |\Sigma| \cdot |w|^3)$. The length of the sequence is equal to $|W| + 1$. Each S_i contains at most $|w|$ prefixes of w , and each A_i contains at most $|\Sigma|$ symbols. Each S_i is of the form

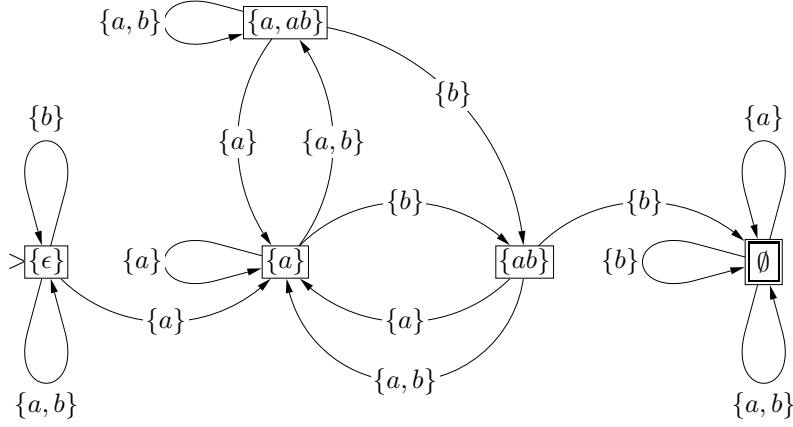


Figure 2. $\text{DFA}_{\Sigma}(abb)$ with $\Sigma = \{a, b\}$.

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

Figure 1. Illustration of the construction in Lemma 2.

$\lfloor R_i \rfloor$, where R_i is constructed by computing $\text{sufpre}(p \cdot a, w)$ for every $p \in S_{i-1}$ and $a \in A_i$. A naive computation of $\text{sufpre}(u, w)$ with $|u| \leq |w|$ takes time $\mathcal{O}(|w|^2)$, and the computation of $\lfloor R_i \rfloor$ is in $\mathcal{O}(|w|^3)$.

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 $u \not\models w$ and $\text{sufpre}(u, w) = p$.

Proof: Proof by induction on i . The basis of the induction, $i = 1$, is easy.

For the induction step, assume $p \in S_{i+1}$. We can assume $q \in S_i$ and $a \in A_{i+1}$ such that $w \neq p =$

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

- Assume $u \cdot a \models w$. Then necessarily, $w = q \cdot a$. But then $p = w$, a contradiction. We conclude by contradiction that $u \cdot a \not\models w$.
- By Lemma 1, $\text{sufpre}(u \cdot a, w) = \text{sufpre}(q \cdot a, w) = p$.

This concludes the proof. \blacksquare

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 $W \notin \text{CERTAIN}(w)$. \blacksquare

Intuitively, the construction of $\langle S_0, S_1, \dots, S_n \rangle$ for a word (or pattern) w and a multiword W can be thought of as executing the pattern matching algorithm of Knuth-Morris-Pratt [12] simultaneously on every word in $\text{words}(W)$. In particular, in case W is of the form $\langle \{a_1\}, \{a_2\}, \dots, \{a_n\} \rangle$, i.e. every symbol is a singleton, then $\text{words}(W)$ is the singleton $\{a_1 a_2 \dots a_n\}$ and our algorithm executes in a way that is close to the Knuth-Morris-Pratt algorithm.

The Aho-Corasick algorithm [13] extends the Knuth-Morris-Pratt algorithm from a single pattern to sets of patterns. This extension is not to be confused with the problem studied in this article, which is deciding whether a single pattern occurs in every word of a given multiword.

V. DEFINING $\text{CERTAIN}(w)$ IN FO

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 $P = P' \setminus \{w\}$. 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 $\{[S] \mid S \subseteq P\}$.
- 2) The alphabet Ω is the powerset alphabet $2^\Sigma \setminus \{\emptyset\}$.
- 3) The initial 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:

$$\delta([S], A) = [\text{sufpre}(S \cdot A, w) \setminus \{w\}].$$

For example, Fig. 2 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.

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 [8], [9]. 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. It is an open problem whether $\text{CERTAIN}(w)$ 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 1; the proof of 2 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 :

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

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

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

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

\Rightarrow Assume $P \cdot U^{k+1} \cdot Q \not\vdash_{\text{CERTAIN}} w$. 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 $m \not\vdash w$. 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 factor u_2 as follows:

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

Assume $m' \vdash w$. Since $m \not\vdash w$, any factor 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 $m' \not\vdash w$. Since $m' \in \text{words}(P \cdot U^{k+2} \cdot Q)$, $P \cdot U^{k+2} \cdot Q \not\vdash_{\text{CERTAIN}} w$.

- 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 factor $v \cdot a$ as follows:

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

Assume $m' \vdash w$. Since $m \not\vdash w$, any factor 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 $m' \not\vdash w$. It is easy to verify that $m' \in \text{words}(P \cdot U^{k+2} \cdot Q)$. It follows $P \cdot U^{k+2} \cdot Q \not\vdash_{\text{CERTAIN}} w$.

\Leftarrow Assume $P \cdot U^{k+2} \cdot Q \not\vdash_{\text{CERTAIN}} w$. 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 $m \not\vdash w$. Let m' be the word obtained from m by “omitting” the factor u_1 as follows:

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

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

- 1) $m' \not\vdash w$. Then, $P \cdot U^{k+1} \cdot Q \not\vdash_{\text{CERTAIN}} w$.
- 2) $m' \vdash w$. Since $m \not\vdash w$, any factor 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 factor 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 $m \not\Vdash w$, the factor 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 $m'' \not\Vdash w$. It is easy to verify that $m'' \in \text{words}(P \cdot U^{k+1} \cdot Q)$. Consequently, $P \cdot U^{k+1} \cdot Q \not\Vdash_{\text{certain}} w$.

This concludes the proof. \blacksquare

Since aperiodic regular languages are expressible in FO [8], [9], 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.

VI. 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 xxx 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 $\text{vars} = \{x, y, z, x_1, y_1, z_1, \dots\}$ 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 $i \in \{1, \dots, n\}$, $s_i \in \Sigma \cup \text{vars}$. 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 $\theta : \text{vars} \cup \Sigma \rightarrow \Sigma$ 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:

$w \Vdash v$ iff for some valuation θ , $\theta(v)$ is a factor 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 $w \in \text{words}(W)$, $w \Vdash v$. \triangleleft

The problem statement in Section III 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 2 stated that $\text{CERTAIN}(w)$ is regular if w is variable-free. We now show that $\text{CERTAIN}(w)$ remains regular if w contains variables. Notice from the example in the previous paragraph that an automaton for $\text{CERTAIN}(xx)$ is *not* simply the union of the automata for $\text{CERTAIN}(aa)$ and $\text{CERTAIN}(bb)$. Rather than taking an automaton approach as in Section V, we will show that for any v-word v , $\text{CERTAIN}(v)$ can be defined by a formula in

monadic second-order logic (MSO). Regularity then follows by Büchi's theorem.

We first look at MSO formulas defining $\text{CERTAIN}(w)$ if w is variable-free. We know from Theorem 2 and Büchi's theorem that such formula exists (there exists even an $\exists\text{MSO}$ definition since $\text{MSO} = \exists\text{MSO}$ over words). Table I shows an MSO formula (actually a $\forall\text{MSO}$ formula) defining $\text{CERTAIN}(ab)$ when the alphabet is $\Sigma = \{a, b\}$. That is, for any multiword W , $W \Vdash_{\text{certain}} ab$ if and only if W satisfies the MSO formula of Table I. The monadic second-order variables X_a and X_b encode all possible words of W : if $X_a(y)$ is true, then a is chosen at position y , and if $X_b(y)$ is true, then b is chosen at position y . The premise states that for each position in multiword W , either a or b has to be chosen; a can only be chosen at position y if the multiword W contains $\{a\}$ or $\{a, b\}$ at position y ; b can only be chosen at position y if the multiword W contains $\{b\}$ or $\{a, b\}$ at y . Notice that the premise depends on the alphabet only. The consequence states that two successive positions contain a and b in that order; the successor relation $S(z_1, z_2)$ means that position z_2 is immediately after z_1 .

After the example of Table I, it is now easy to construct an MSO formula that defines $\text{CERTAIN}(w)$ for any word w over a fixed alphabet Σ . Moreover, it is an easy exercise to generalize this construction to v-words, as illustrated by Table II for the v-word xx . So we have the following result.

Theorem 4: For every v-word v , $\text{CERTAIN}(v)$ is regular.

We are particularly interested in conditions on v that guarantee aperiodicity of $\text{CERTAIN}(v)$, or equivalently, FO-definability of $\text{CERTAIN}(v)$. For words without variables, such condition was settled by Theorem 3, and it is an open question whether there exists a variable-free word w such that $\text{CERTAIN}(w)$ is not FO-definable. On the other hand, the following theorem shows the existence of a v-word v such that $\text{CERTAIN}(v)$ is not FO-definable.

Theorem 5: Let Σ be a finite alphabet with $|\Sigma| \geq 2$. Then,

- 1) $\text{CERTAIN}(xx)$ is in **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 $j \in \{1, \dots, n\}$,

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

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

$$\forall X_a \forall X_b \left(\left(\begin{array}{l} \forall y (X_a(y) \vee X_b(y)) \\ \wedge \neg \exists y (X_a(y) \wedge X_b(y)) \\ \wedge \forall y (X_a(y) \rightarrow P_{\{a\}}(y) \vee P_{\{a,b\}}(y)) \\ \wedge \forall y (X_b(y) \rightarrow P_{\{b\}}(y) \vee P_{\{a,b\}}(y)) \end{array} \right) \rightarrow \exists z_1 \exists z_2 \left(\begin{array}{l} S(z_1, z_2) \\ \wedge X_a(z_1) \\ \wedge X_b(z_2) \end{array} \right) \right)$$

Table I
MSO DEFINITION OF CERTAIN(ab) FOR ALPHABET $\Sigma = \{a, b\}$.

$$\forall X_a \forall X_b \left(\left(\begin{array}{l} \forall y (X_a(y) \vee X_b(y)) \\ \wedge \neg \exists y (X_a(y) \wedge X_b(y)) \\ \wedge \forall y (X_a(y) \rightarrow P_{\{a\}}(y) \vee P_{\{a,b\}}(y)) \\ \wedge \forall y (X_b(y) \rightarrow P_{\{b\}}(y) \vee P_{\{a,b\}}(y)) \end{array} \right) \rightarrow \exists z_1 \exists z_2 \left(\begin{array}{l} S(z_1, z_2) \\ \wedge X_a(z_1) \leftrightarrow X_a(z_2) \\ \wedge X_b(z_1) \leftrightarrow X_b(z_2) \end{array} \right) \right)$$

Table II
MSO DEFINITION OF CERTAIN(xx) FOR ALPHABET $\Sigma = \{a, b\}$.

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. ■

VII. CONCLUSION

Every multiword defines a set of possible words. A word or v-word 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.
- 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:

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

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 following tasks remain for future research:

- 1) Show the conjecture that $\text{CERTAIN}(w)$ is FO-definable if w is variable-free.
- 2) 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*. ACM Press, 1999, pp. 68–79.
- [2] A. Fuxman and R. J. Miller, “First-order query rewriting for inconsistent databases,” *J. Comput. Syst. Sci.*, vol. 73, no. 4, pp. 610–635, 2007.
- [3] L. Grieco, D. Lembo, R. Rosati, and M. Ruzzi, “Consistent query answering under key and exclusion dependencies: algorithms and experiments,” in *CIKM*, O. Herzog, H.-J. Schek, N. Fuhr, A. Chowdhury, and W. Teiken, Eds. ACM, 2005, pp. 792–799.
- [4] J. Wijsen, “Consistent query answering under primary keys: a characterization of tractable queries,” in *ICDT*, ser. ACM International Conference Proceeding Series, R. Fagin, Ed., vol. 361. ACM, 2009, pp. 42–52.
- [5] —, “On the consistent rewriting of conjunctive queries under primary key constraints,” *Inf. Syst.*, 2009.
- [6] D. M. Gabbay, A. Pnueli, S. Shelah, and J. Stavi, “On the temporal basis of fairness,” in *POPL*, 1980, pp. 163–173.
- [7] J. A. W. Kamp, “Tense logic and the theory of linear order,” Ph.D. dissertation, University of California, Los Angeles, 1968.
- [8] R. McNaughton and S. Papert, *Counter-free Automata*. Cambridge, MA: MIT Press, 1971.
- [9] M. P. Schützenberger, “On finite monoids having only trivial subgroups,” *Information and Control*, vol. 8, no. 2, pp. 190–194, 1965.

- [10] F. Blanchet-Sadri, R. Mercas, and G. Scott, "A generalization of Thue freeness for partial words," *Theor. Comput. Sci.*, vol. 410, no. 8-10, pp. 793–800, 2009.
- [11] L. Libkin, *Elements of Finite Model Theory*. Springer, 2004.
- [12] D. E. Knuth, J. H. M. Jr., and V. R. Pratt, "Fast pattern matching in strings," *SIAM J. Comput.*, vol. 6, no. 2, pp. 323–350, 1977.
- [13] A. V. Aho and M. J. Corasick, "Efficient string matching: An aid to bibliographic search," *Commun. ACM*, vol. 18, no. 6, pp. 333–340, 1975.