



Autonomously Computable Information Extraction

Besat Kassaie
Cheriton School of Computer Science
University of Waterloo
Waterloo, ON, Canada
bkassaie@uwaterloo.ca

Frank Wm. Tompa
Cheriton School of Computer Science
University of Waterloo
Waterloo, ON, Canada
fwtompa@uwaterloo.ca

ABSTRACT

Most optimization techniques deployed in information extraction systems assume that source documents are static. Instead, extracted relations can be considered to be materialized views defined by a language built on regular expressions. Using this perspective, we can provide an efficient verifier (using static analysis) that can be used to avoid the high cost of re-extracting information after an update. In particular, we propose an efficient mechanism to identify updates for which we can autonomously compute an extracted relation. We present experimental results that support the feasibility and practicality of this mechanism in real world extraction systems.

PVLDB Reference Format:

Besat Kassaie and Frank Wm. Tompa. Autonomously Computable Information Extraction. PVLDB, 16(10): 2431 - 2443, 2023.
doi:10.14778/3603581.3603585

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/Besatkassaie/Differential-Maintenance-Engine>.

1 INTRODUCTION

Information extraction is a crucial step in processing and understanding text. By identifying fields of interest in unstructured or semi-structured data sources, extractors provide selected data to populate relational records. For instance, an application can extract the relationships between events, dates, and venues as described in a collection of documents. *Knowledge Panels* [33] show how information extraction can improve the presentation of retrieval results in a commercial search engine. Extraction can be performed using ad-hoc text processing programs, machine learning, or reusable components offered by systems such as GATE [15] and SystemT [24]. In this paper, we address the third approach, where extractors are defined by regular expressions.

Extraction time can be a bottleneck for many applications [30, 32], and therefore efficient processing is an important consideration for information extraction. Some optimization approaches are general and can be deployed in any system for specifying extractors. For instance, Chandel et al. [5] propose an efficient algorithm for dictionary-based entity recognition, which can be used on many extraction platforms. Iperotis et al. [17] deal with many documents by keeping only “promising” documents for the extraction process.

Shen et al. [32] propose a declarative approach to information extraction using Datalog with embedded procedural predicates to express extractors. This enables the creation of execution plans for the extraction program and thus the application of cost-based optimization techniques similar to query optimization in relational databases. Similarly, Reiss et al. [30] recommend using a SQL-like declarative language, AQL (used by SystemT), to be able to exploit query optimization strategies. Considering extractions within larger applications, Jain et al. [18] treat information extraction and subsequent relational queries as an integrated system and propose optimizations that takes into account the characteristics of extractors, document retrieval methods, and join algorithms on the extracted relations.

However, none of these strategies consider that an extraction might need to be re-computed to keep extracted information synchronized with source documents as they are updated. Recognizing this situation, Chen et al. have developed an approach for incrementally updating extracted relations [6]. They do not assume that a description of the update is available, but instead compare each updated document with the previous version to find regions that have not changed. Then, based on user-provided properties of the extractor, they decide which of the extracted items from those regions can be reused.

Doleschal et al. [10] explore conditions for determining that a spanner is *split-correct*, that is, if the extracted relation can be computed by combining the extractions from sub-documents. If so, extractions from various sub-documents can be run in parallel, but additionally incremental update is applicable: re-extraction after an update can be avoided for those sub-documents that are not altered.

Freydenberger and Thompson [14] have investigated the complexity of incrementally re-evaluating spanners in the presence of updates. However, their update model assumes that a document is encoded as a fixed-length *word structure* in which (essentially) there is a special character that represents ϵ , and the only permissible update is replacing one character from $\Sigma \cup \{\epsilon\}$ by another.

In a previous short paper [21], we have noted that the problem of re-evaluating extractions after updates is analogous to the problem of maintaining materialized views [7]. We identified updates to source documents that can be reflected in the extracted relation without re-computing the extractors (i.e., irrelevant updates), and we formulated the condition of an update being *pseudo-irrelevant*, where extracted regions are merely shifted in source documents. In unpublished extended work [20], we describe verification algorithms that attempt to determine whether an update is pseudo-irrelevant with respect to an extractor. Unfortunately, the algorithms rely on normalization that causes exponential blow-up in the input size and on finding complements and intersections of

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/url(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 10 ISSN 2150-8097.
doi:10.14778/3603581.3603585

automata, which in the worst case causes exponential blow-up in time and space.

No other prior work proposes to apply static analysis to programs that specify extractors and updates in order to determine whether re-extraction can be avoided or reduced.

In this paper, we continue the investigation of pseudo-irrelevant updates, one form of autonomously computable updates, for extractors expressed as document spanners (as defined below). To that end, in this paper:

- (1) We extend the update model formalized in our earlier work to allow replacements to depend on the strings being replaced (Section 3).
- (2) We propose a new algorithm to verify that an update is pseudo-irrelevant with respect to an extractor and prove that it runs in polynomial time via three theorems (Section 5). The algorithm relies on four external tests (two for verifying the precondition and two tests for independence) that are each proven to execute in polynomial time in a proposition,¹ and it is built upon four auxiliary algorithms, each of which is proven to execute in polynomial time by an accompanying lemma.
- (3) We show experimentally that our algorithm is practical: it can be used effectively in realistic update and extraction scenarios, and it runs far faster than the time needed for re-extraction. If this approach is combined with incremental update, the overhead is relatively small, and it will often perform much faster (Section 6).

2 PRELIMINARIES

We assume that extracted views are defined using SystemT, an information extraction platform that benefits from relational database concepts to deal with text data sources [30]. SystemT processes one document at a time and populates relational tables with *spans*, directly extracted from the input document. With SystemT, users encode extractors using a SQL-like language, namely the Annotation Query Language (AQL), to manipulate tables.

The underlying principles adopted by SystemT have been formalized as *document spanners* by Fagin et al. [12]. Most of the material in this section has been introduced in that work and additional details can be found there.

2.1 Regular Expressions with Capture Variables

Given a finite alphabet Σ , a regular expression extended using variables chosen from a set V is called a *regex with capture variables* and conforms to γ in the grammar $G_S(\Sigma, V)$ as follows:

$$\begin{aligned} \gamma &:= \emptyset \mid \epsilon \mid \alpha \mid (\gamma \vee \gamma) \mid (\gamma \bullet \gamma) \mid (\gamma)^* \mid x\{\gamma\} \\ \alpha &:= \sigma \mid [\sigma, \sigma] \mid \delta \\ \delta &:= \Sigma \mid (\delta - \sigma) \end{aligned} \quad (1)$$

where σ represents any character in Σ , $[\sigma, \sigma]$ represents the disjunction of characters having their encoding between or equal to the encodings of the first and second character in the range, the terminal symbol Σ represents the disjunction of all characters in Σ , $\delta - \sigma$ represents the disjunction of all characters in δ with the

¹An additional test for constraints on the alphabet is embedded in the algorithm.

exception of σ , and (what distinguishes these expressions from conventional regular expressions) x represents any variable in V . Given a regex with capture variables r , the corresponding *regex tree* $\mathcal{T}(r)$ represents the hierarchical structure of r , in which the tree's leaves have labels \emptyset, ϵ , characters in Σ , character ranges in Σ , or the character Σ itself, and internal nodes have labels $\bullet, \vee, *, -,$ or a symbol in V . For convenience of notation, when writing a regex with capture variables, we follow common practice for regular expressions in allowing the following shorthand: omission of parentheses (relying instead on left associativity of all operations and precedence of $-$ over $*$ over \bullet over \vee) and omission of the operator \bullet .

If E is a regex with capture variables, then we denote the set of capture variables in E as $SVars(E)$. The use of a subexpression of the form $n\{g\}$ in E signifies that whenever E matches a string, the span containing a substring matched by g is to be *marked by the capture variable n* (as explained in Section 2.3). It is apparent from the grammar that capture variables can be nested.

Example 2.1. Let Σ be the set of Latin alphanumeric, punctuation, and space characters (the last represented by $_$). $\gamma_{fullDate}$ is a regex with capture variables:²

$$\gamma_{fullDate} = \Sigma^* \boxed{\text{mdate}=\text{}} F\{Y\{\gamma_d Y_d \gamma_d Y_d\} \boxed{-} M\{\gamma_d Y_d\} \boxed{-} D\{\gamma_d Y_d\}\} \Sigma^*$$

where $\gamma_d = [\boxed{0}, \boxed{9}]$ and $SVars(E) = \{F, Y, M, D\}$. For this example, we say that Y, M, D are *nested variables* and that F is *exposed*. We extend conventional set notation to write $Y \subset F$ and $\boxed{-} \in F$.

2.2 Document Spans

A document D is a finite string over some alphabet: $D \in \Sigma^*$ (Figure 1). A *span* of document D , denoted $[i, j)$ ($1 \leq i \leq j \leq |D| + 1$), specifies the start and end offsets of a substring in D , which is in turn denoted $D_{[i, j)}$, and extends from offset i through offset $j - 1$.

Example 2.2. In the document presented in Figure 1, $D_{[56, 70)}$ represents the substring James.F. Allen.

$[i, i)$ denotes an empty span at offset i . Spans $s_1 = [i_1, j_1)$ and $s_2 = [i_2, j_2)$ are identical if and only if $i_1 = i_2$ and $j_1 = j_2$.

A substantial portion of our work is based on investigating various relationships between spans. Allen has defined a set of 13 possible relationships between non-empty intervals [1]. These can be extended to capture the same basic relationships among spans (including empty spans) as summarized in Table 1. All possible relationships among spans can be described by disjunctions of these basic relationships; for example, “ X overlaps Y ” ($\Gamma_{X \cap Y}$) can be expressed as the disjunction of the last nine basic relationships.³ “ X overlaps but is not equal to Y ” ($\Gamma_{X \cap \neq Y}$) can be similarly expressed as the disjunction of the fifth through the twelfth basic relationships.

²Throughout this paper, characters in Σ appearing in a formula are represented like this to distinguish them from the regular expression's meta-characters.

³The definition of overlapping spans given by Fagin et al. [12] is asymmetric for empty spans, i.e., given a span $[i, j)$ the empty span at $[i, i)$ is considered overlapping with $[i, j)$ while the empty span at $[j, j)$ is considered disjoint from $[i, j)$. We treat both as overlapping.

```

< a r t i c l e _ k e y = " c a c m / A l l e n 8 3 " _ m d a t e = " 2 0 1 1 - 0 6 - 0 7 " > < a u
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50
t h o r > J a m e s _ F . _ A l l e n < / a u t h o r > < t i t l e > M a i n t a i n i n g _ K n o
51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100
w l e d g e _ a b o u t _ T e m p o r a l _ I n t e r v a l s . < / t i t l e > < / a r t i c l e >
101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150

```

Figure 1: A sample input document D for our running example.

Table 1: Allen’s interval relationships extended to spans.

1	$\Gamma_{(X<Y)}$	X precedes Y	$\left. \begin{array}{l} \Sigma^* X \{ \Sigma^* \} \Sigma^+ Y \{ \Sigma^* \} \Sigma^* \\ \Sigma^* X \{ \Sigma^+ \} Y \{ \Sigma^+ \} \Sigma^* \end{array} \right\}$	7	$\Gamma_{(XdY)}$	X during Y	$\left. \begin{array}{l} \Sigma^* Y \{ \Sigma^+ X \{ \Sigma^* \} \Sigma^+ \} \Sigma^* \\ \Sigma^* Y \{ X \{ \Sigma^* \} \Sigma^+ \} \Sigma^* \end{array} \right\}$
2	$\Gamma_{(Y>X)}$	Y is preceded by X		8	$\Gamma_{(YdX)}$	Y contains X	
3	$\Gamma_{(XmY)}$	X meets Y	$\left. \begin{array}{l} \Sigma^* X \{ \Sigma^+ \} Y \{ \Sigma^+ \} \Sigma^* \\ \Sigma^* (X \vdash) \Sigma^+ (Y \vdash) \Sigma^+ (\vdash X) \Sigma^+ (\vdash Y) \Sigma^{*a} \end{array} \right\}$	9	$\Gamma_{(XsY)}$	X starts Y	$\left. \begin{array}{l} \Sigma^* Y \{ \Sigma^+ X \{ \Sigma^* \} \Sigma^+ \} \Sigma^* \\ \Sigma^* Y \{ \Sigma^+ X \{ \Sigma^* \} \} \Sigma^* \end{array} \right\}$
4	$\Gamma_{(YmiX)}$	Y is met by X		10	$\Gamma_{(YsiX)}$	Y is started by X	
5	$\Gamma_{(XoY)}$	X overhangs Y	$\left. \begin{array}{l} \Sigma^* (X \vdash) \Sigma^+ (Y \vdash) \Sigma^+ (\vdash X) \Sigma^+ (\vdash Y) \Sigma^{*a} \\ \Sigma^* X \{ Y \{ \Sigma^* \} \} \Sigma^* \vee \Sigma^* X \{ \epsilon \} Y \{ \epsilon \} \Sigma^* \end{array} \right\}$	11	$\Gamma_{(XfY)}$	X finishes Y	$\left. \begin{array}{l} \Sigma^* Y \{ \Sigma^+ X \{ \Sigma^* \} \} \Sigma^* \\ \Sigma^* Y \{ \Sigma^+ X \{ \Sigma^* \} \} \Sigma^* \end{array} \right\}$
6	$\Gamma_{(YoiX)}$	Y is overhung by X		12	$\Gamma_{(YfiX)}$	Y is finished by X	
	13	$\Gamma_{(X=Y)}$	X is equal to Y				

^aThis abuse of notation represents a spanner (as described in Section 2.3) matching an automaton with operators that open and close the variables X and Y as indicated.

2.3 Extractors Expressed by Document Spanners

Determining membership in a language defined by a regex with capture variables E can be accomplished by executing a corresponding *vset-automaton* $\mathcal{A}(E)$. Given E and $V = SVars(E)$, $\mathcal{A}(E)$ is a non-deterministic finite state automaton augmented with a designated set (initially empty) and two operators for each variable $x \in V$, namely $x \vdash$ (“open x ”) and $\vdash x$ (“close x ”). Besides including standard character transitions, $\mathcal{A}(E)$ also includes *operation transitions* that, instead of consuming a character from the input string, insert the variable x into the designated set if the transition label is $x \vdash$ and remove x from the designated set if the label is $\vdash x$. A document D is accepted by $\mathcal{A}(E)$ if, after scanning the whole input, we end up in a final state and the designated set is empty. A *matching* of E against document D is an accepting run in $\mathcal{A}(E)$, where for each variable $x \in SVars(E)$, the spans marked by x each begin with the offset in D when x is inserted into the designated set and end with the offset in D when x is removed from that set.

If E is a regex with capture variables, it specifies a *document spanner*, denoted as $\llbracket E \rrbracket$, with $SVars(\llbracket E \rrbracket) = SVars(E)$. Applying a document spanner to a document D produces a *span relation*, i.e., a relation that contains spans of D . Thus, $\llbracket E \rrbracket$ is a function mapping strings over Σ^* to $\mathcal{S}^{|SVars(E)|}$ where \mathcal{S} is the set of all spans of D . To ensure that the span relation is in first-normal form with no null values, we restrict our attention to a specific class of document spanners, namely *functional document spanners*, that mark exactly one span for each variable for all accepting runs, regardless of the input document D . In particular, for a given document D , the spanner specified by E produces a span relation $\llbracket E \rrbracket(D)$ in which there is one column for each variable from $SVars(E)$ appearing in E , each row corresponds to a matching of E against D when the variables are ignored, and the value in a row for the column corresponding to $x \in SVars(E)$ is the span marked by x .

Example 2.3. Let γ_d be as defined in Example 2.1. Applying the spanner represented by $Y_{partialDate} = \Sigma^* F \{ M \{ \gamma_d \gamma_d \} \ominus D \{ \gamma_d \gamma_d \} \} \Sigma^*$ to the document in Figure 1 results in the span relation in Figure 2.

F	M	D
[38, 43)	[38, 40)	[41, 43)
[41, 46)	[41, 43)	[44, 46)

Figure 2: The extracted relation $\llbracket Y_{partialDate} \rrbracket(D)$, where D is depicted in Figure 1.

Definition 2.4. Throughout this paper, a functional document spanner used for the purpose of information extraction is called an *extractor*, the regex with capture variables defining it is called an *extraction formula*, and the span relation produced for a document is called an *extracted relation*.

Let S , S_1 , and S_2 be extractors where the last two are union-compatible (i.e., $SVars(S_1) = SVars(S_2)$); $X \subseteq SVars(S)$; and $x, y \in SVars(S)$. An algebra over spanners can be defined with operators:

- (1) union: $S_1 \cup S_2$ having variables $SVars(S_1 \cup S_2) = SVars(S_1)$,
- (2) projection: $\pi_X(S)$ having variables X ,
- (3) natural join: $S \bowtie S_1$ having variables $SVars(S) \cup SVars(S_1)$, and
- (4) binary string selection: $\zeta_{x,y}^{\bar{x}} S$ having variables $SVars(S)$.

The set of *core spanners* (corresponding to the core of SystemT’s AQL) includes extractors specified by any extraction formula (so-called *primitive extractors*) together with all extractors in the closure of core spanners under this algebra. Applying a core spanner to any document D is equivalent to applying each included primitive spanner to D and then applying the corresponding relational operators to the extracted relations. The set of core spanners define the extractors and updates subject to analysis in this paper.

2.4 Efficient Construction of Extractors

Given one or more spanners as input, we use various algebraic operations defined over spanners to verify properties of those spanners statically. Specifically, we convert the inputs to *eVset-automata*, a variant of *vset-automata* with the same expressivity, as proposed by Morciano [26]. By preserving three properties of *eVset-automata*

(namely, *well-behaved*, *pruned*, and *operation-closed*), Morciano is able to construct eVset-automata in polynomial time for simulating the application of projection, union, and join over spanners. His thesis also shows that converting a regex with capture variables to an automaton and checking for emptiness can be done in polynomial time. It is trivial to show that renaming variables can also be accomplished in polynomial time with eVset-automata.

2.5 Model Architecture

In this work, we hypothesize systems that include a document database \mathbb{D} and a set of core spanners $\{\mathbb{E}_1, \dots, \mathbb{E}_e\}$ specifying extractors that run over \mathbb{D} . The union of span relations produced by \mathbb{E}_k against the document database is stored in a relation \mathbb{T}_k that includes an additional column to associate a document identifier with the spans for the corresponding extracted relation. These tables serve as materialized views of the document database.

3 DOCUMENT UPDATE MODEL

Updates can add documents to or delete documents from the database \mathbb{D} , or they can change documents already in \mathbb{D} . In this paper, we concentrate on the latter form of update, where substring *replacement*, *deletion*, and *insertion* are basic update operations that we wish to support. A change to the text is typically preceded by some browsing activities or search operations to locate update positions in a target document. We use an extractor as our search mechanism and to specify which strings to replace.

Definition 3.1. An *update expression* is denoted as $Repl(g, U)$, where g is an update formula with one exposed capture variable that determines which spans contain strings to be replaced and U is a string that specifies the replacement value.

Given an update expression $Repl(g, U)$, the string U may include named back-references to g , similar to the mechanisms defined in several programming languages: $U \in ((\Sigma - \$) \cup (\$(SVars(g))))^*$.

Example 3.2. Consider $Repl(\gamma_{DOI}, \text{"DOI: $(C)"})$, where

$$\gamma_{DOI} = \Sigma^* \langle ee \rangle (F \{ \text{https://doi.org/} C \{ (\Sigma - \langle \rangle)^* \} \} \vee F \{ \text{doi:} C \{ (\Sigma - \langle \rangle)^* \} \} \langle ee \rangle \Sigma^*$$

F is the exposed variable, indicating that strings enclosed by $\langle ee \rangle$... $\langle /ee \rangle$ and starting with a DOI specifier should be replaced. The second argument for the update expression indicates that each replacement should consist of the string "DOI:" followed by whatever matches the capture variable C followed by another quotation mark.

3.1 Update Formulas

An *update formula* is a functional extraction formula with one exposed variable (v) and conforms to $\bar{\gamma}$ in the grammar $G_U(\Sigma, V, v)$:

$$\bar{\gamma} := (\bar{\gamma} \vee \bar{\gamma}) \mid (\gamma' \bullet \bar{\gamma}) \mid (\bar{\gamma} \bullet \gamma') \mid v\{\gamma''\} \quad (2)$$

where γ' is a variable-free regular expression, γ'' is a regex with capture variables V , and the exposed variable v is an additional capture variable (the *update variable*) within which all other capture variables are nested. Given an update formula g , the update variable is denoted as $UVar(g)$. The extraction formulas presented in Examples 2.1 and 2.3 conform to $\bar{\gamma}$ where F is the update variable.

Definition 3.3. An update formula is in *normalized form* if it is written as $\bigvee_{i=1}^k g_i$ where each g_i is a formula conforming to $\hat{\gamma}$:

$$\hat{\gamma} := \gamma' \bullet v\{\gamma''\} \bullet \gamma'$$

where each γ' is again a variable-free (possibly empty) regular expression and γ'' is again a regex with capture variables V .

To normalize an update formula, all disjunctions that have the update variable v in their disjuncts⁴ can be “pulled up” over concatenations in the corresponding extended regex tree to create separate disjuncts at the outermost level of the formula.

Example 3.4. The expression γ_{DOI} in Example 3.2 conforms to grammar (2), and its corresponding normalized form is:

$$\gamma_{DOI} = \Sigma^* \langle ee \rangle F \{ \text{https://doi.org/} C \{ (\Sigma - \langle \rangle)^* \} \} \langle /ee \rangle \Sigma^* \vee \Sigma^* \langle ee \rangle F \{ \text{doi:} C \{ (\Sigma - \langle \rangle)^* \} \} \langle ee \rangle \Sigma^*$$

LEMMA 3.5. Given any update formula g , Algorithm 1 normalizes it by producing the list of disjuncts $\Delta(g, UVar(g))$ in time that is polynomial in $|g|$, where $|x|$ denotes the length of x .

PROOF. The proof of correctness is by induction on the height of the expression tree, which can be constructed and traversed (using recursive descent) in linear time in the input length. For each occurrence of $UVar(g)$, a new expression is added to the list, which makes the output size $O(m * n)$, where $n = |g|$ and $m < n$ is the number of occurrences of $UVar(g)$ in g . □

3.2 Update Semantics

The functional document spanner that is represented by an update formula g maps every document D to a span relation, which we call the *update relation* and denote as $\llbracket g \rrbracket(D)$. When the spanner is used for updating a document D , substrings of D associated with the spans in the update relation are simultaneously replaced by new values specified by U .

More precisely, given a document D and applying $Repl(g, U)$ to D produces a new document $Repl(g, U)(D)$ that is identical to D except that each substring $s_i \in D$ corresponding to a span marked by $UVar(g)$ is replaced by the string U_i where, for any $v \in SVars(g)$, U_i is the same as U except that every occurrence of a substring $\$(v)$ in U is replaced by the string in s_i corresponding to the span marked by v .

Note that if U is the empty string, then the update results in the deletion of the substrings corresponding to spans marked by $UVar(g)$; otherwise, wherever an empty span $[i, i)$ is marked by $UVar(g)$, the replacement, in effect, inserts a string before the i^{th} character (or at the end of the string if $i = |D| + 1$).

An update yields a specific functional mapping between the original and updated documents: $ReplSpan(g, U)([i, j]) \rightarrow [i', j']$ which is a mapping from spans to spans. This is properly defined when $[i, j)$ is disjoint from all spans marked by $UVar(g)$ or when $[i, j)$ is the complete span marked by $UVar(g)$ (Figure 3). If $[i, j)$ is disjoint from all spans marked by $UVar(g)$ when updating D , then $D'_{[i', j']} = D_{[i, j)}$.

⁴Because the formulas are functional, if a capture variable appears in one disjunct, it must appear in all alternative disjuncts.

Algorithm 1: Normalize g with respect to v .

Input: extraction formula g with exposed variable v
Output: list of disjuncts $\Delta(g, v)$
Precondition: g is functional

```

1  $\mathcal{T} \leftarrow \text{ExpressionTree}(g)$ ;
2 return  $\text{normalize}(\mathcal{T}, v)$ 

3 Function  $\text{normalize}(\mathcal{T}: \text{expression tree}, v: \text{variable}): \text{list}$ 
   Precondition:  $v$  is in  $\mathcal{T}$ 
4    $\Delta(g, v) \leftarrow \text{list}()$ ;
5   if  $\mathcal{T}.\text{root} == v$  then
6      $\Delta(g, v).\text{add}(\text{toRegExp}(\mathcal{T}))$ 
7   end
8   if  $\mathcal{T}.\text{root} == \vee$  then
9     /* normalize subtrees */
10     $\Delta(g, v).\text{add}(\text{normalize}(\mathcal{T}.\text{left}, v))$ ;
11     $\Delta(g, v).\text{add}(\text{normalize}(\mathcal{T}.\text{right}, v))$ 
12  end
13  if  $\mathcal{T}.\text{root} == \bullet$  then
14    /* if  $v$  occurs in right (left) subtree, normalize right (left)
15     subtree and then concatenate every expression in  $\Delta(g, v)$  with
16     left (right) expression */
17    if  $v$  in  $\text{subtree}(\mathcal{T}.\text{right})$  then
18       $\Delta(g, v).\text{add}(\text{normalize}(\mathcal{T}.\text{right}, v))$ ;
19       $\text{toRegExp}(\mathcal{T}.\text{left}) \bullet \Delta(g, v)$ 
20    else
21       $\Delta(g, v).\text{add}(\text{normalize}(\mathcal{T}.\text{left}, v))$ ;
22       $\Delta(g, v) \bullet \text{toRegExp}(\mathcal{T}.\text{right})$ 
23    end
24  end
25  return  $\Delta(g, v)$ 
26 end

```

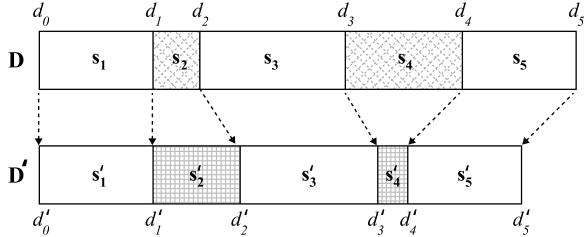


Figure 3: A sample document D and its updated peer D' . Filled areas in D are marked by the update variable. ReplSpan maps $[d_i, d_{i+1})$ to $[d'_i, d'_{i+1})$.

Definition 3.6. An update spanner specified by g is said to be *well-defined* if the following two properties hold for every document and each pair of accepting runs: (1) if the pair of spans for $UVar(g)$ are not equal, they must not overlap, and (2) if the pair of spans for $UVar(g)$ are equal, then for each variable appearing in U , the pair of spans marked for that variable must be equal.

PROPOSITION 3.7. *Whether or not an update spanner specified by g is well-defined can be verified in polynomial time.*

PROOF. The proof has two similar parts, one for each property. First we show that detecting whether there can be overlaps in spans that will be changed can be done by constructing a specific spanner and testing for emptiness. Then we show that detecting whether a single updated span might have ambiguously marked subspans can also be done by constructing a specific spanner and testing for emptiness. All operations take polynomial time.

Given an update formula g and symbols $X \notin SVars(g)$ and $Y \notin SVars(g)$, we create the spanner

$$\text{conflicts}(g) = \pi_X(\rho_{UVar(g) \rightarrow X}(\llbracket g \rrbracket)) \bowtie \llbracket \Gamma_{(X \cap Y)} \rrbracket \bowtie \pi_Y(\rho_{UVar(g) \rightarrow Y}(\llbracket g \rrbracket))$$

Recall that $\Gamma_{(X \cap Y)}$ is the disjunction of the fifth through the twelfth basic relationships in Table 1; that is, we are testing whether spans in $\pi_{UVar(g)}(\llbracket g \rrbracket(D))$ could include two unequal spans that cover identical sub-spans. If $\text{conflicts}(g) = \emptyset$, then g satisfies the first condition. The operators utilized to construct $\text{conflicts}(g)$ use time that is at most quadratic in the size of inputs (Section 2.4).

Similarly, given an update formula g , symbols $X \notin SVars(g)$ and $Y \notin SVars(g)$, and $Z \in SVars(g) \setminus \{UVar(g)\}$, we create the spanner

$$\text{ambig}(g, Z) = \pi_{\{UVar(g), X\}}(\rho_{Z \rightarrow X}(\llbracket g \rrbracket)) \bowtie \llbracket \Gamma_{(X \neq Y)} \rrbracket \bowtie \pi_{\{UVar(g), Y\}}(\rho_{Z \rightarrow Y}(\llbracket g \rrbracket))$$

where $\Gamma_{(X \neq Y)}$ is the disjunction of the first 12 basic relationships in Table 1; that is, spans in $\llbracket g \rrbracket(D)$ could include two rows that match on $UVar(g)$ but do not match on Z . If $\text{ambig}(g, Z) = \emptyset$ for all $Z \in SVars(g) \setminus \{UVar(g)\}$, then g satisfies the second condition. Clearly the time complexity to check this property is also polynomial in the size of the input. \square

Example 3.8. The spanner presented in Example 2.3 is not well-defined, since $\text{conflicts}(Y_{\text{partialDate}}) \neq \emptyset$ (M in one match for F can overlap D in another match).

Example 3.9. The spanner represented by the following extraction formula is easily verified to be not well-defined: the context ensures that it is free of conflicts, but it is ambiguous:

$$Y'_{\text{date}} = \Sigma^* \langle \text{now} \rangle F\{Y\{y_d^*\}M\{y_d^*\}D\{y_d^*\}\langle \text{now} \rangle \Sigma^*$$

4 AUTONOMOUSLY COMPUTABLE UPDATES

Given an update expression $\text{Repl}(g, U)$ and an extractor $\llbracket E \rrbracket$, we wish to determine, for all potential input documents, whether the extracted materialized view can be kept consistent with the updated source documents without running the extractor after updating the documents in the database. This problem is similar to filtering out irrelevant updates or applying updates autonomously to relational materialized views [4, 16, 21].

Definition 4.1. An update expression $\text{Repl}(g, U)$ is *irrelevant* with respect to an extractor $\llbracket E \rrbracket$ if for every input document, applying $\llbracket E \rrbracket$ to $\text{Repl}(g, U)(D)$ produces an extracted relation that is identical to applying $\llbracket E \rrbracket$ to D . That is, if $D' = \text{Repl}(g, U)(D)$, then $\llbracket E \rrbracket(D') = \llbracket E \rrbracket(D)$.

If an update expression is not irrelevant with respect to an extractor, it may still be that the modification to the extracted relation can be computed without re-running the extractor.

Definition 4.2. An update expression $Repl(g, U)$ is *autonomously computable* with respect to an extractor $\llbracket E \rrbracket$ if for every document in the database \mathbb{D} , applying $\llbracket E \rrbracket$ to $Repl(g, U)(D)$ can be computed from the update expression, the update relation, the extraction formula that defines the extractor, and the extracted relation.⁵

There is an important distinction between the problems of updating traditional relational views and updating materialized extractions. Extracted relations contain pairs of offsets from input documents, not document content. Thus, as happens for spans within regions S_3 and S_5 in Figure 3, replacing a string of one length by a string of another length somewhere in the document will cause a span further down in the document to shift, even if the content of that span is unaffected.

More generally, given a document D and the corresponding updated document D' , if span S in D is disjoint from all spans produced by the well-defined update spanner $\llbracket g \rrbracket$, $ReplSpan(g, U)(S)$ is shifted from S by an amount that is dependent on the lengths of all spans in the update relation that precede S in D and the lengths of the strings specified by U , as captured by Algorithm 2.

Algorithm 2: Shift a span.

Input: span $S = [i, j)$, update relation R , update variable X , replacement specification U

Output: span $S' = [i', j') = ReplSpan(g, U)(S)$

Precondition: R contains no span that overlaps S or another span in R

shift, $l \leftarrow 0$;

for $t \in R$ **do**

$[m, n) \leftarrow t[X]$; /* span for the update variable */

if $m < i$ **then**

$l \leftarrow computeLength(U, t)$; /* using each $|S(u_i)|$ from t */

shift $\leftarrow shift + (n - m) - l$

end

end

return $[i - shift, j - shift)$

Definition 4.3. Update expression $Repl(g, U)$ is *pseudo-irrelevant* with respect to an extractor $\llbracket E \rrbracket$ if for every document D and updated document $D' = Repl(g, U)(D)$, $\llbracket E \rrbracket(D') = \{S' \mid \exists S \in \llbracket E \rrbracket(D), S' = ReplSpan(g, U)(S)\}$ (i.e., extracted spans merely shift).

Thus, a pseudo-irrelevant update is a special case of an autonomously computable update. By definition, if an update expression is irrelevant with respect to an extractor, then it is also pseudo-irrelevant with respect to that spanner (all shifts are of length 0).

5 CATEGORIZING DOCUMENT UPDATES

We wish to identify whether an update is pseudo-irrelevant with respect to a given extractor, independently of input documents. The

⁵We note in passing that updating the span relation when adding documents to or deleting documents from \mathbb{D} can be performed autonomously. For new documents, we run the extractors on those documents only (the text of which forms part of the update) and then add the extracted tuples to the materialized view without any need to reference documents in \mathbb{D} . When deleting documents, we simply delete all tuples with corresponding document identifiers from the extracted relation.

essence of our approach is to inspect various kinds of *overlap* between an update expression and an extractor. The proposed process verifies some *sufficient* conditions for pseudo-irrelevant updates.

By examining Definition 4.3, we can deduce three possibilities that could cause an update (changing document D to become D') to fail to be pseudo-irrelevant with respect to an extractor: (1) a span s extracted from D fails to be extracted from D' ; (2) a span s' extracted from D' does not correspond to any extracted span from D prior to the update; and (3) a span s extracted from D changes to become s' extracted from D' , but it is not a simple shift.

Example 5.1. Consider the wh-word extractor specified by

$$\Sigma^* \left(\left[_ \right] Q \{ \text{wh} \} \gamma_{lc}^+ \right) \left(\left[_ \right] (\Sigma - \left[_ \right] - \left[_ \right] - \left[_ \right])^* \left[_ \right] \right) \Sigma^*$$

If an update replaces semicolons by periods, a wh-question that matched before the update might no longer match the specification, and therefore a tuple might be deleted from the extracted relation. If it instead replaces periods by semicolons, a newly formed wh-question might be created, and thus a tuple might be inserted into the extracted relation. If it replaces “ere” by “y”, “Where” could be replaced by “Why” and thus the extracted relation might change by more than a simple shift.

We leave it to future work to determine under what conditions an update that overlaps extracted spans (e.g., changing “Who” to “Why” in the example) or their contexts (e.g., changing periods to exclamation points for this extractor) happens to be pseudo-irrelevant. Instead, we determine when there can be no overlap, and then under which further conditions an update is pseudo-irrelevant. Section 5.1 describes how to identify the contexts of extracted spans. Section 5.2 describes how to identify spans in an updated document D' that have had their contents changed by an update to D . Finally Section 5.3 describes a sound, but not complete, mechanism to determine whether an update expression, specified by the update formula and replacement specifier, is pseudo-irrelevant with respect to a document spanner specified by an extraction formula.

5.1 Contextualization of Extraction Formulas

For an extraction formula to match a document, certain strings must appear either inside or outside the regions marked by capture variables. If a document is updated, we wish to know whether the extracted content or any of the strings that specify required contextual information is disrupted in any way. The tests we have implemented compare updated regions to regions marked by a capture variable, so, given the specification of a document spanner E with exposed variable v , we define a corresponding regular formula $C_v(E)$ in which all contextual expressions are also marked with capture variables. For any document D , $\llbracket C_v(E) \rrbracket$ produces the same spans for the capture variable v as does $\llbracket E \rrbracket$, but it also captures spans of characters that must remain unchanged. Each uncovered subexpression in $C_v(E)$ is of the form X_i^* where $L(X_i)$ includes *unigrams* (i.e., strings of length 1) only, essentially restricting the uncovered subexpression to strings over a restricted alphabet $\hat{\Sigma}_i \subseteq \Sigma$. For example, in Example 5.1, $\left[_ \right] \text{wh} \gamma_{lc}^*$ is already covered by the capture variable Q , but the blanks before and after that pattern are a part of the context than cannot change, as is the question mark at the end; they should be covered. On the other hand, the expressions

Algorithm 3: Capture contexts for $v \in SVars(E)$.

Input: extraction formula E , $v \in SVars(E)$
Output: modified extraction formula $C_v(E)$
Precondition: v is exposed, $v_i \notin SVars(E)$ for all i

```
1  $\mathcal{D} \leftarrow \Delta(E, v)$ ; /* get the disjunctive form using Algorithm 1 */
2  $D_c \leftarrow \text{emptyList}()$ ; /* prepare for list of covered expressions */
3 forall  $d \in \mathcal{D}$  do
4    $\mathcal{T}_c \leftarrow \text{cover}(\text{ExpressionTree}(d), 1)$ ; /* start cover with  $v_1$  */
5    $\mathcal{T}_m \leftarrow \text{mergeConsecVars}(\mathcal{T}_c)$ ; /*  $v_i\{e_i\}v_j\{e_j\} \rightarrow v_k\{e_i e_j\}$  */
6    $D_c.add(\mathcal{T}_m)$  /* keep covered expr for disjunct */
7 end
/* ensure each disjunct includes all  $v_i$  so  $C_v(E)$  is functional */
8 result  $\leftarrow \epsilon$ ;
9 forall  $d \in D_c$  do
10   $\mathcal{T}_v \leftarrow \text{findSubtree}(d, v)$ ; /* find node with  $v$  */
11  forall  $y \in \text{getAllVars}(D_c) \setminus \text{getAllVars}(d)$  do
12    /* change  $v\{\dots\}$  to  $v\{y\{\dots\}\}$  in this disjunct */
13     $\text{setParent}(\text{setParent}(\mathcal{T}_v.child, \text{newNode}(y)), \mathcal{T}_v)$ 
14  end
15  result  $\leftarrow \text{result} \vee \text{toRegExp}(d)$  /* add disjunct to  $C_v(E)$  */
16 end
17 return result
```

Function $\text{cover}(\mathcal{T} : \text{expression tree}, i : \text{int}) : \text{expression tree}$
Output: cover contexts in a subtree, starting with v_i

```
18  $\mathcal{T}_c \leftarrow \text{emptyTree}()$ ;
19 if  $\mathcal{T}.root == *$  then
20   if  $\text{matchesUnigrams}(\mathcal{T}.root.left)$  then
21      $\mathcal{T}_c = \mathcal{T}$  /*  $\hat{\Sigma}^*$  does not get covered */
22   else
23      $\mathcal{T}_c \leftarrow \text{setParent}(\mathcal{T}.root, \text{newNode}(v_i))$ 
24   end
25 end
26 if  $\mathcal{T}.root == \vee$  then
27    $\mathcal{T}_c \leftarrow \text{setParent}(\mathcal{T}.root, \text{newNode}(v_i))$ 
28 end
29 if  $\text{varNode}(\mathcal{T}.root)$  then
30    $\mathcal{T}_c = \mathcal{T}$  /* already covered */
31 end
32 if  $\mathcal{T}.root == \bullet$  /* traverse both sides */
33 then
34    $\mathcal{T}.left \leftarrow \text{cover}(\mathcal{T}.left, i)$ ;
35    $i' = \text{maxIndex}(\mathcal{T}.left)$ ;
36   if  $i \leq i'$  then
37      $i = i' + 1$  /* set  $i$  to the largest index so far */
38   end
39    $\mathcal{T}.right \leftarrow \text{cover}(\mathcal{T}.right, i)$ ;
40    $\mathcal{T}_c \leftarrow \mathcal{T}$ 
41 else
42    $\mathcal{T}_c \leftarrow \text{setParent}(\mathcal{T}.root, \text{newNode}(v_i))$  /* unigram */
43 end
44 return  $\mathcal{T}_c$ 
45 end
```

$(\Sigma - \textcircled{\cdot} - \textcircled{!} - \textcircled{?})^*$ and Σ^* are examples of $\hat{\Sigma}^*$ that can be freely updated by any string in their languages.

LEMMA 5.2. *Given an extraction formula E and exposed variable $v \in SVars(E)$, Algorithm 3 returns $C_v(E)$ in quadratic time.*

PROOF. Using induction on the height of the expression tree for E , it is straightforward to show that Algorithm 3 covers all subexpressions except for the closure of those that satisfy the predicate *matchesUnigrams*. Testing this predicate requires merely checking that no node in the subtree is labelled \bullet or $*$.

Complexity Analysis: Given an input expression E , the corresponding expression tree can be constructed in linear time in the input size, and projection via a top-down traversal of the expression tree can also be performed in linear time. Normalization has time complexity that is linear in the input length (Lemma 3.5), and *cover()* and *mergeConsecVars()* require a top-down traversal over the expression tree, again in linear time. For an input expression E with $m < |E|$ occurrences of v , the time complexity of the algorithm is thus $O(m * |E|)$. \square

To be most effective in identifying contexts, we wish to identify as many instances of X_i^* as we can, since they indicate portions of the document that can be (almost) freely updated. To this end, standard rewrite rules for regular expressions can be applied, even in the presence of capture variables. For example, if R_i are any regular expressions with capture variables and $\hat{\Sigma}$ is any expression matching a unigram, distributive laws can be used to pull X_i^* outside disjunctions:

$$\begin{aligned} (\hat{\Sigma}^* R_1) \vee (\hat{\Sigma}^* R_2) &\rightarrow \hat{\Sigma}^* (R_1 \vee R_2) \\ (R_1 \hat{\Sigma}^*) \vee (R_2 \hat{\Sigma}^*) &\rightarrow (R_1 \vee R_2) \hat{\Sigma}^* \end{aligned}$$

Example 5.3. Consider the following expression:

$$\gamma_{jb} = (\Sigma^* \langle \text{journal} \rangle \vee \Sigma^* \langle \text{booktitle} \rangle) T \{ \gamma_t \}$$
$$(\langle \text{journal} \rangle \Sigma^* \vee \langle \text{booktitle} \rangle \Sigma^*)$$

Applying Algorithm 3 to γ_{jb} with variable T produces

$$C_T(\gamma_{jb}) = v_1 \{ (\Sigma^* \langle \text{journal} \rangle \vee \Sigma^* \langle \text{booktitle} \rangle) \} T \{ \gamma_t \}$$
$$v_2 \{ (\langle \text{journal} \rangle \Sigma^* \vee \langle \text{booktitle} \rangle \Sigma^*) \}$$

After rewriting with the distributive laws, the algorithm produces

$$C_T(\gamma_{jb'}) = \Sigma^* v_1 \{ (\langle \text{journal} \rangle \vee \langle \text{booktitle} \rangle) \} T \{ \gamma_t \}$$
$$v_2 \{ (\langle \text{journal} \rangle \vee \langle \text{booktitle} \rangle) \} \Sigma^*$$

marking the portions that are required to be in any matching document, while leaving matches to arbitrary text (i.e., Σ^*) unmarked.

Other simple rewrite rules, such as removing superfluous closures $((R^*)^* \rightarrow R^*$, where R is any regular expression) and removing superfluous disjunctions $(L(R) \subseteq L(\hat{\Sigma}^*) \implies \hat{\Sigma}^* \vee R \rightarrow \hat{\Sigma}^*)$, also provide means to expose more instances of $\hat{\Sigma}^*$.

5.2 Post-Update Spanner

We use information provided by the update expression to construct a *post-update spanner* that identifies regions in updated documents corresponding to updated spans. First we need to ensure that an update does not modify any of the context used to identify which spans to update.

Definition 5.4. An update spanner $\llbracket g \rrbracket$ is *durable* if it is well-defined and spans marked by the update variable are disjoint from spans marked by context variables after applying $\llbracket C_{UVar(g)}(g) \rrbracket$ to any document D .

PROPOSITION 5.5. *Testing whether an update spanner is durable can be performed in polynomial time.*

PROOF. Again the proof has two parts, one for each property. Testing for being well-defined is straightforward. Then we show that detecting whether an updated span might overlap with a span marked by a context variable can be done by constructing a specific spanner and testing for emptiness in polynomial time.

Let v represent $UVar(g)$. Testing whether $\llbracket g \rrbracket$ is well-defined requires polynomial time (Proposition 3.7), and $C_v(g)$ can be constructed in polynomial time (Lemma 5.2).

Let v_i represent the i^{th} variable marking the context in $C_v(g)$ and construct a spanner $p_i(g)$:

$$p_i(g) = \pi_X(\rho_{v \rightarrow X}(\llbracket C_v(g) \rrbracket)) \bowtie \llbracket \Gamma_{(X \cap Y)} \rrbracket \bowtie \pi_Y(\rho_{v_i \rightarrow Y}(\llbracket C_v(g) \rrbracket))$$

where $\Gamma_{(X \cap Y)}$ is again the disjunction of the fifth through thirteenth basic relationships in Table 1. Therefore, the spanner $p_i(g)$ matches all documents that can be updated while spans marked by the update variable have at least one subspan in common with the set of spans marked as context. Therefore, if for all i , $p_i(g) = \emptyset$, then $\llbracket g \rrbracket$ is durable. Constructing $p_i(g)$ and testing for emptiness requires polynomial time, and the number of cover variables is less than the size of the input. \square

Given a specification of an update spanner, we derive a new spanner to match documents that result from an update. Some substrings in an updated document come from outside the regions matched by the update variable, and others come from the replacement specifier U , either from substrings explicitly contained in U or from the use of back-references to substrings in the original document. Thus, U implicitly describes a language $L(U)$ comprising the set of possible replacement values, but it uses back-references to variables in g . We define a regular language $\diamond(U, g)$ that describes a slightly broader

space of replacement values: $\diamond(U, g) \supseteq L(U)$ and $\diamond(U, g) \setminus L(U)$ is fairly small.

LEMMA 5.6. *Given a replacement specifier U and a well-defined spanner specified by g , Algorithm 4 outputs $\diamond(U, g)$ with time complexity $O(|g| + |U|)$.*

PROOF. The proof is by construction. Given g and

$$U = u_1 \cdots u_i \$ (v_i) u_{i+1} \cdots u_j \$ (v_j) u_{j+1} \cdots u_n$$

Algorithm 4 scans U and g and substitutes each back-reference $\$(v_i)$ with the disjunction of all expressions specified with that variable in g : ($exp_1^{v_i} \vee exp_2^{v_i} \vee \cdots$), omitting other capture variables enclosed by v_i .⁶ The update spanner specified by g is functional, which implies that in every run exactly one of the expressions marked as variable v_m in g is encountered, and in $\diamond(U, g)$ the corresponding run will encounter its associated back-reference in U . The resulting regular expression is:

$$\diamond(U, g) = u_1 \cdots u_i \bullet (exp_1^{v_i} \vee exp_2^{v_i} \vee \cdots) \bullet u_{i+1} \cdots u_j \bullet (exp_1^{v_j} \vee exp_2^{v_j} \vee \cdots) \bullet u_{j+1} \cdots u_n$$

where $exp_k^{v_m}$ is a regular expression marked by $v_m \in SVars(g)$.

Complexity Analysis: To build $\diamond(U, g)$, Algorithm 4 scans U to retrieve all back-referenced variables and then the expression tree of g is scanned, substituting back-references in U with associated regular expressions taken from g . The complexity is $O(|U| + |g|)$. \square

Definition 5.7. Given an update spanner specified by g and a replacement specifier U , a spanner that matches the updated documents and marks all updated spans is called a *post-update spanner* and is specified by $\nabla(g, U)$.

Example 5.8. For the update expression in Example 3.2, the corresponding post-update spanner is

$$\nabla(\mathcal{Y}_{DOI}, \text{"DOI: (C)"}) = \Sigma^* \langle \text{ee} \rangle F \{ \text{"DOI: } (\Sigma - \langle \rangle)^* \langle \rangle \} \langle \text{ee} \rangle \Sigma^* \vee \Sigma^* \langle \text{ee} \rangle F \{ \text{"DOI: } (\Sigma - \langle \rangle)^* \langle \rangle \} \langle \text{ee} \rangle \Sigma^*$$

⁶Eliminating capture variables inside v_i does not affect the set of strings marked by v_i .

Algorithm 4: Define space of replacement values $\diamond(U, g)$.

Input: update expression $Repl(g, U)$
Output: $\diamond(U, g)$
 $\mathcal{T} \leftarrow ExpressionTree(g);$
 $bkrefs[] \leftarrow get_bkrefs(U);$ /* get all back-references from U */
 $result_strs \leftarrow \{\};$ /* empty dictionary */
forall $bkref$ **in** $bkrefs$ **do**
 /* retrieve all substrings associated with $bkref$ */
 $subexps[] \leftarrow get_subexp(\mathcal{T}, bkref);$
 /* concatenate all substrings separated by disjunction symbols */
 $result_str \leftarrow "\vee".join(toRegExp(subexp[]));$
 /* enclose with parentheses and add to the dictionary */
 $result_strs[bkref] \leftarrow "(" + result_str + ")"$
end
/* substitute each $bkref$ with corresponding string in dictionary */
 $\diamond(U, g) \leftarrow replace(U, result_strs);$
return $\diamond(U, g)$

Algorithm 5: Build post-update spanner $\nabla(g, U)$.

Input: update expression $Repl(g, U)$, update variable x
Output: $\nabla(g, U)$
Precondition: $\llbracket g \rrbracket$ is durable
1 $\mathcal{D} \leftarrow \Delta(g, x);$ /* get the disjunctive form using Algorithm 1 */
2 $\nabla(g, U) \leftarrow \emptyset;$
 /* process each disjunct g_i in $\Delta(g, x)$ */
3 **forall** $g_i \in \mathcal{D}$ **do**
4 $\mathcal{T} \leftarrow ExpressionTree(g_i);$
 /* get subexpressions on left/right of update variable */
5 $\mathcal{T}_l \leftarrow get_left(\mathcal{T}, x);$
6 $\mathcal{T}_r \leftarrow get_right(\mathcal{T}, x);$
7 $\nabla(g, U) \leftarrow \nabla(g, U) \vee (toRegExp(\mathcal{T}_l) \bullet x \{ \diamond(U, g_i) \}$
8 $\bullet toRegExp(\mathcal{T}_r))$
9 **end**
10 **return** $\nabla(g, U)$

LEMMA 5.9. For a durable update spanner $\llbracket g \rrbracket$, Algorithm 5 outputs $\nabla(g, U)$ in quadratic time in the input length.

PROOF. The proof is by contradiction, using Figure 3 for intuition.

Let D be a document that has m regions marked by $x = UVar(g)$ when processed by $\llbracket g \rrbracket$ and D' denote the result of updating D by $Repl(g, U)$. Let us assume that there are some *problematic regions* of D that are marked by x but do not have correctly marked corresponding regions when running $\llbracket \nabla(g, U) \rrbracket$ on D' . Using Figure 3 as a guide, let s_4 be the *leftmost* problematic region of D . Because the update spanner is durable, all subexpressions in g identified as contexts, as well as any instances of $\hat{\Sigma}^*$ that match to the left of s_4 , must correctly match the regions preceding s_4 (e.g., s'_1 , s'_2 , and s'_3 are correctly identified). Based on Lemma 5.6, the region s_4 must match $\diamond(U, g)$. Thus the substring from the start of D to the end of s_4 ($D_{[d_0, d_4]}$) corresponds correctly to the substring from the start of D' to the end of s'_4 ($D'_{[d'_0, d'_4]}$). Therefore, s_4 cannot be problematic, and so there cannot be any problematic region.

Complexity Analysis: The time complexity of creating $\Delta(g, x)$ and any expression tree and to traverse the expression tree to extract right and left contexts is $O(|g|)$, and the complexity of constructing $\diamond(U, g_i)$ is $O(|g_i| + |U|)$. Finally, if g has m disjuncts in $\Delta(g, x)$, m is $O(|g|)$ and thus the time complexity of the algorithm is $O(|g|^2 + |g| * |U|)$. \square

5.3 Detecting Pseudo-Irrelevance for Spanners

An update is pseudo-irrelevant with respect to a core spanner S if it is pseudo-irrelevant with respect to each primitive extractor defining S :

THEOREM 5.10. Let $\mathbb{E}_{Repl(g, U)}$ represent the set of extractors for which $Repl(g, U)$ is pseudo-irrelevant. $\mathbb{E}_{Repl(g, U)}$ is closed under union, projection, natural join, and string selection.

PROOF. We need to show that if $\llbracket E_1 \rrbracket \in \mathbb{E}_{Repl(g, U)}$, $\llbracket E_2 \rrbracket \in \mathbb{E}_{Repl(g, U)}$, and $Y \subseteq SVars(E_1)$, then

- (1) $\llbracket E_1 \cup E_2 \rrbracket \in \mathbb{E}_{Repl(g, U)}$ when $\llbracket E_1 \rrbracket$, $\llbracket E_2 \rrbracket$ are union compatible,
- (2) $\llbracket \Pi_Y E_1 \rrbracket \in \mathbb{E}_{Repl(g, U)}$,
- (3) $\llbracket E_1 \bowtie E_2 \rrbracket \in \mathbb{E}_{Repl(g, U)}$, and
- (4) $\llbracket \zeta_Y^R E_1 \rrbracket \in \mathbb{E}_{Repl(g, U)}$.

Each of these can be proven by contradiction. \square

Definition 5.11. Let $\llbracket S \rrbracket$ be a spanner defined by update formula g or a post-update spanner specified by $\nabla(g, U)$. Given an extraction formula E with exposed variable v , $\llbracket E \rrbracket$ is *independent* of $\llbracket S \rrbracket$ if for every document D , $\llbracket C_v(E) \rrbracket(D)$ includes no span that overlaps a span in $\llbracket S \rrbracket(D)$. Otherwise, we say that $\llbracket E \rrbracket$ depends on $\llbracket S \rrbracket$.

PROPOSITION 5.12. Independence can be verified in polynomial time.

PROOF. As before, we show that detecting whether an updated span might overlap a span marked by a context variable can be done by constructing a specific spanner and testing for emptiness in polynomial time.

First rename variables such that $UVar(g) \notin SVars(C_v(E))$ and $X, Y \notin (SVars(C_v(E)) \cup \{UVar(g)\})$, and let $v_i \in SVars(C_v(E))$.

Then construct the following spanner:

$$\begin{aligned} depends(E, g, v_i, UVar(g)) &= \pi_X(\rho_{v_i \rightarrow X}(\llbracket C_v(E) \rrbracket)) \bowtie \llbracket \Gamma_{X \cap Y} \rrbracket \\ &\quad \bowtie \pi_Y(\rho_{UVar(g) \rightarrow Y}(\llbracket g \rrbracket)) \end{aligned}$$

Again, $\Gamma_{X \cap Y}$ is the disjunction of the fifth through thirteenth basic relationships in Table 1, capturing overlapping regions. The spanner $depends(E, g, v_i, UVar(g)) \neq \emptyset$ if there exists at least one span in $\pi_{v_i}(\llbracket C_v(E) \rrbracket)$ that touches at least one span in $\pi_{UVar(g)}(\llbracket g \rrbracket)$. Therefore, if for all such v_i , $depends(E, g, v_i, UVar(g)) = \emptyset$, then $\llbracket E \rrbracket$ is independent of $\llbracket g \rrbracket$. There are at most $|E|$ such tests, each of which takes polynomial time, proving the claim. \square

Definition 5.13. Given an extraction formula E with exposed variable v , an update expression $Repl(g, U)$ respects the alphabets in $C_v(E)$ if, for every expression of the form $\hat{\Sigma}^*$ not covered in $C_v(E)$, the update neither deletes nor inserts a symbol in $(\Sigma - \hat{\Sigma})$.

If $Repl(g, U)$ does not respect the alphabets in $C_v(E)$, tuples might be inserted into or deleted from an extracted relation $\llbracket E \rrbracket(D)$, as illustrated by Example 5.1.

THEOREM 5.14. A durable update spanner $\llbracket g \rrbracket$ with the replacement specifier U is pseudo-irrelevant with respect to an extractor $\llbracket E \rrbracket$ if $\llbracket E \rrbracket$ is independent of both $\llbracket g \rrbracket$ and $\llbracket \nabla(g, U) \rrbracket$ and $Repl(g, U)$ respects the alphabets in $C_v(E)$.

PROOF. Consider any document D with m regions marked by $\llbracket g \rrbracket$ to be updated, like s_2 and s_4 in Figure 3. Let D' denote D 's updated peer, which has regions marked by $\llbracket \nabla(g, U) \rrbracket$, including those like s'_2 and s'_4 in Figure 3. Assume that some spans of D are marked by $\llbracket E \rrbracket$ to be extracted. To prove that an update is pseudo-irrelevant with respect to an extractor, we must prove three properties (as explained earlier):

- (1) For every extracted region in D , there should be a corresponding region in D' that will be extracted, i.e., no span disappears from the extracted relation after the update. Since $\llbracket E \rrbracket$ is independent of $\llbracket g \rrbracket$, the only regions that can be updated are those matching $\hat{\Sigma}^*$, where $\hat{\Sigma} \subseteq \Sigma$. However, since $Repl(g, U)$ respects the alphabets in $C_v(E)$, any update expression matching $\hat{\Sigma}^*$ includes a replacement specification U that also matches $\hat{\Sigma}^*$. Thus every span found before the update must have a corresponding span found after the update.
- (2) For every region extracted from D' , there should be a corresponding region extracted from D , i.e., no new span appears in the extracted relation because of the update. Since $\llbracket E \rrbracket$ is independent of $\llbracket \nabla(g, U) \rrbracket$, the only regions that could have been updated are those matching $\hat{\Sigma}^*$, where $\hat{\Sigma} \subseteq \Sigma$. As before, since $Repl(g, U)$ respects the alphabets in $C_v(E)$, any update expression matching $\hat{\Sigma}^*$ includes a replacement specification U that also matches $\hat{\Sigma}^*$. Thus every span found after the update must have a corresponding span found before the update.
- (3) For extracted spans, $ReplSpan$ should be realized as a *shift* function. Because of independence, the lengths of the extracted regions remain the same after the update, since no update can occur inside an extracted region. But other regions preceding an extracted region and matching $\hat{\Sigma}^*$ can be updated, which makes the starting offset of an extracted region move back or forth. Since the length of the extracted region is fixed, the end offset moves by the same amount, which is indeed a shift of the extracted region. \square

random samples), ranging from 2.3 hours to 43.4 hours and with the final extractor requiring more than 125 hours. Undoubtedly, various optimization techniques adopted by mature products such as SystemT can accelerate the extraction process, and these were not implemented on our research-based system. We note, however, that Chen et al. [6] report that the fastest extractor they tested on a Wikipedia corpus of 35 MB requires 100 seconds (using different hardware and software), which implies an extraction time of more than 6 hours when scaled to the size of DBLP.⁷ Elsewhere, Shen et al. [32] report a complex, but optimized extractor taking 61 minutes for a corpus smaller than 2% of the size of DBLP! Thus, we are convinced that our extraction times are indicative of the times needed for extractors more generally.

Table 5 shows how many seconds are required for verifying pseudo-irrelevancy for each update, displaying the minimum, average, and maximum times over all five primitive DBLP extractors. The table also shows the minimum, average, and maximum times required for shifting the tables extracted for the first seven extractors. It is far faster to verify pseudo-irrelevancy than to re-execute the corresponding extractor: verification times are at most 3.5 minutes per primitive extractor *independently of the database size*, and simply shifting extracted relations is far faster than re-extracting them, especially when the extractions are highly selective.

Table 5: Verification and shift times (sec) for all updates.

Update	Verification			Shift		
	min'm	avg	max'm	min'm	avg	max'm
Hashtags	20	64.0	127	0.2±0.4	10.2	39.2±76.8
URLs	36	108.4	212	1.1±0.1	67.1	239.3±1.8
Dates	1	33.6	64	1.5±0.2	78.7	273.5±6.5
DOIs	18	58.2	110	1.2±0.2	72.4	252.4±7.6

Chen et al. [6] have proposed re-using extracted data when it has clearly not been affected by an update. The approach requires three steps: (1) determining which pieces of documents have been updated, (2) re-executing the extractor on portions of the updated documents where changes to the extracted data *or to its context* (specified as a window around each extracted region) might change the results of extraction, and (3) copying over previous extractions guaranteed to be unaffected by the update, but shifting their offsets so as to be able to determine overlaps with the next update.

We note first that incremental re-extraction is orthogonal to determining pseudo-irrelevancy. An update might be provably pseudo-irrelevant to $\llbracket E \rrbracket$ even if it changes data very near each extracted region, such as is true for γ_{date} with respect to S_{Mod} . On the other hand, if an update is not provably pseudo-irrelevant, such as γ_{URL} with respect to Q_{RN} , incremental extraction can be applied rather than naively re-executing the extractor on the complete document.

To determine the cost of using the approach proposed by Chen et al., we measure the time to execute a simple matcher for the DBLP corpus on the same hardware used in our other experiments. In particular, by running Unix's `diff` against 20% of the DBLP corpus we find that the program requires between 27 ms and 68 ms (mean

⁷In practice, document-at-a-time extraction can easily be parallelized, so a system with 32 processors could reduce overall elapsed time to 12 minutes or so.

= 49 ms) per file. This implies that the first step alone requires over 42 hours to execute. Even if the computation is distributed across 32 processors running in parallel, the execution time is excessive for updates that can be determined to be pseudo-irrelevant.

7 ADDITIONAL RELATED WORK

Document Spanners. Researchers have addressed many problems using the document spanner model, including how to deal with documents with missing information [25] and how to eliminate inconsistencies from extracted relations [11]. Others have studied the complexity of evaluating spanners and computing the results of various algebraic operations over span relations [2, 13, 28, 29].

Static Analysis of Programs Using Regular Languages. We use an extended form of finite-state automata to determine whether an update expression is pseudo-irrelevant with respect to a document spanner. Similar static analyses of regular expressions have been used in diverse areas, including access control mechanisms for XML database systems [27] and conflicting gestures in multi-touch environments [22]. We have also used finite automata to statically analyze extractors specified by JAPE [8] in the context of updating extracted views [19].

Materialized View Maintenance. In the relational setting, a view definition, expressed in relational algebra, along with tuples to be inserted or removed from base tables provides information that is used to maintain the content of a view without recomputing it from scratch. Blakeley et al. address this problem where the update involves insertion or deletion of a set of tuples into/from one base relation at a time [3, 4]. They give sufficient and necessary conditions on the irrelevancy of an update, which is stronger than our work proposing only sufficient conditions. For a differential view update, Blakeley et al. additionally require the content of associated base relations before the updates.

8 CONCLUSIONS

Given a program defined as a document spanner and an update specification, we determine sufficient conditions for autonomously re-computing which spans of an updated document are extracted. In particular, we propose three sufficient conditions for pseudo-irrelevancy of updates with respect to an extraction program: *durability*, *independence*, and *respect for alphabets*. We prove that we require time and space that are polynomial in the size of the extraction program and the update specification to perform the five required tests to determine that the revised extracted relation can be computed autonomously. We also designed some practical extractors and conducted experiments on two real-world datasets to conclude that the runtime overhead imposed by our verification is small in practice when compared to re-evaluating extractors, even if the re-evaluation is performed incrementally. Furthermore, because it uses static analysis, verification is independent of database size.

ACKNOWLEDGMENTS

We gratefully acknowledge financial assistance received from the University of Waterloo and NSERC, the Natural Sciences and Engineering Research Council of Canada.

REFERENCES

- [1] James F. Allen. 1983. Maintaining Knowledge about Temporal Intervals. *Commun. ACM* 26, 11 (1983), 832–843.
- [2] Antoine Amarilli, Pierre Bourhis, Stefan Mengel, and Matthias Niewerth. 2019. Constant-Delay Enumeration for Nondeterministic Document Spanners. In *Proc. 22nd International Conference on Database Theory, ICDT (LIPIcs)*, Pablo Barceló and Marco Calautti (Eds.), Vol. 127. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, Lisbon, 22:1–22:19.
- [3] José A. Blakeley, Neil Coburn, and Per-Åke Larson. 1989. Updating Derived Relations: Detecting Irrelevant and Autonomously Computable Updates. *ACM Trans. Database Syst.* 14, 3 (1989), 369–400.
- [4] José A. Blakeley, Per-Åke Larson, and Frank Wm. Tompa. 1986. Efficiently Updating Materialized Views. In *Proc. 1986 ACM SIGMOD International Conference on Management of Data*, Carlo Zaniolo (Ed.). ACM Press, Washington, DC, 61–71.
- [5] Amit Chandel, P. C. Nagesh, and Sunita Sarawagi. 2006. Efficient Batch Top-k Search for Dictionary-based Entity Recognition. In *Proc. 22nd International Conference on Data Engineering, ICDE*, Ling Liu, Andreas Reuter, Kyu-Young Whang, and Jianjun Zhang (Eds.). IEEE Computer Society, Atlanta, 28:1–28:10.
- [6] Fei Chen, AnHai Doan, Jun Yang, and Raghu Ramakrishnan. 2008. Efficient Information Extraction over Evolving Text Data. In *Proc. 24th International Conference on Data Engineering, ICDE*. IEEE Computer Society, Cancun, 943–952.
- [7] Rada Chirkova and Jun Yang. 2012. Materialized Views. *Found. Trends Databases* 4, 4 (2012), 295–405.
- [8] Hamish Cunningham, Diana Maynard, and Valentin Tablan. 2000. *JAPE: a Java annotation patterns engine*. Technical Report CS-00-10. Dept. Comp. Sci., Univ. Sheffield. 28 pages.
- [9] The dblp team. 2022. Statistics. Retrieved November 7, 2022 from <https://dblp.org/statistics/index.html>
- [10] Johannes Doleschal, Benny Kimelfeld, Wim Martens, Yoav Nahshon, and Frank Neven. 2019. Split-Correctness in Information Extraction. In *Proc. 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS*, Dan Suciu, Sebastian Skritek, and Christoph Koch (Eds.). ACM, Amsterdam, 149–163.
- [11] Ronald Fagin, Benny Kimelfeld, Frederick Reiss, and Stijn Vansummeren. 2014. Cleaning inconsistencies in information extraction via prioritized repairs. In *Proc. 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. ACM, Snowbird UT USA, 164–175.
- [12] Ronald Fagin, Benny Kimelfeld, Frederick Reiss, and Stijn Vansummeren. 2015. Document Spanners: A Formal Approach to Information Extraction. *J. ACM* 62, 2 (2015), 12:1–12:51.
- [13] Fernando Florenzano, Cristian Riveros, Martín Ugarte, Stijn Vansummeren, and Domagoj Vrgoc. 2020. Efficient Enumeration Algorithms for Regular Document Spanners. *ACM Trans. Database Syst.* 45, 1 (2020), 3:1–3:42.
- [14] Dominik D. Freydenberger and Sam M. Thompson. 2020. Dynamic Complexity of Document Spanners. In *Proc. 23rd International Conference on Database Theory, ICDT (LIPIcs)*, Carsten Lutz and Jean Christoph Jung (Eds.), Vol. 155. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, Copenhagen, 11:1–11:21.
- [15] Robert J. Gaizauskas, Hamish Cunningham, Yorick Wilks, Peter J. Rodgers, and Kevin Humphreys. 1996. GATE: An Environment to Support Research and Development in Natural Language Engineering. In *Proc. Eighth International Conference on Tools with Artificial Intelligence, ICTAI*. IEEE Computer Society, Toulouse, 58–66.
- [16] Ashish Gupta, H. V. Jagadish, and Inderpal Singh Mumick. 1996. Data Integration using Self-Maintainable Views. In *Proc. Advances in Database Technology - EDBT'96, 5th International Conference on Extending Database Technology (Lecture Notes in Computer Science)*, Peter M. G. Apers, Mokrane Bouzeghoub, and Georges Gardarin (Eds.), Vol. 1057. Springer, Avignon, 140–144.
- [17] Panagiotis G. Ipeirotis, Eugene Agichtein, Pranay Jain, and Luis Gravano. 2007. Towards a query optimizer for text-centric tasks. *ACM Trans. Database Syst.* 32, 4 (2007), 21:1–21:46.
- [18] Alpa Jain, Panagiotis G. Ipeirotis, and Luis Gravano. 2008. Building query optimizers for information extraction: the SQOUT project. *SIGMOD Rec.* 37, 4 (2008), 28–34.
- [19] Besat Kassaie and Frank Wm. Tompa. 2019. Predictable and Consistent Information Extraction. In *Proc. DocEng '19: ACM Symposium on Document Engineering*. ACM, Berlin, 14:1–14:10.
- [20] Besat Kassaie and Frank Wm. Tompa. 2020. Detecting Opportunities for Differential Maintenance of Extracted Views. (2020), 19 pages. arXiv:2007.01973
- [21] Besat Kassaie and Frank Wm. Tompa. 2020. A Framework for Extracted View Maintenance. In *Proc. DocEng '20: ACM Symposium on Document Engineering*. ACM, Virtual Event, 16:1–16:4.
- [22] Kenrick Kin, Björn Hartmann, Tony DeRose, and Maneesh Agrawala. 2012. Proton: multitouch gestures as regular expressions. In *ACM Conf. on Human Factors in Computing Systems, CHI '12*. ACM, Austin, 2885–2894.
- [23] Michael Ley. 2009. DBLP - Some Lessons Learned. *Proc. VLDB Endow.* 2, 2 (2009), 1493–1500.
- [24] Yunyao Li, Frederick Reiss, and Laura Chiticariu. 2011. SystemT: A Declarative Information Extraction System. In *Proc. 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - System Demonstrations*. ACL, Portland, OR, 109–114.
- [25] Francisco Maturana, Cristian Riveros, and Domagoj Vrgoc. 2018. Document Spanners for Extracting Incomplete Information: Expressiveness and Complexity. In *Proc. 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. ACM, Houston TX USA, 125–136.
- [26] Andrea Morciano. 2016. *Engineering a runtime system for AQL*. Master's thesis. École Polytechnique de Bruxelles, Université Libre de Bruxelles.
- [27] Makoto Murata, Akihiko Tozawa, Michiharu Kudo, and Satoshi Hada. 2006. XML access control using static analysis. *ACM Trans. Inf. Syst. Secur.* 9, 3 (2006), 292–324.
- [28] Liat Peterfreund, Dominik D. Freydenberger, Benny Kimelfeld, and Markus Kröll. 2019. Complexity Bounds for Relational Algebra over Document Spanners. In *Proc. 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS*, Dan Suciu, Sebastian Skritek, and Christoph Koch (Eds.). ACM, Amsterdam, 320–334.
- [29] Liat Peterfreund, Balder ten Cate, Ronald Fagin, and Benny Kimelfeld. 2019. Recursive Programs for Document Spanners. In *Proc. 22nd International Conference on Database Theory, ICDT (LIPIcs)*, Pablo Barceló and Marco Calautti (Eds.), Vol. 127. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, Lisbon, 13:1–13:18.
- [30] Frederick Reiss, Sriram Raghavan, Rajasekar Krishnamurthy, Huaiyu Zhu, and Shivakumar Vaithyanathan. 2008. An Algebraic Approach to Rule-Based Information Extraction. In *Proc. 24th International Conference on Data Engineering, ICDE*. IEEE Computer Society, Cancun, 933–942.
- [31] Jonathan Schler, Moshe Koppel, Shlomo Argamon, and James W. Pennebaker. 2006. Effects of Age and Gender on Blogging. In *Proc. 2006 AAAI Spring Symp. on Computational Approaches to Analyzing Weblogs*. AAAI, Stanford, 199–205.
- [32] Warren Shen, AnHai Doan, Jeffrey F. Naughton, and Raghu Ramakrishnan. 2007. Declarative Information Extraction Using Datalog with Embedded Extraction Predicates. In *Proc. 33rd International Conference on Very Large Data Bases, VLDB*. ACM, Vienna, 1033–1044.
- [33] Danny Sullivan. 2020. *A reintroduction to our Knowledge Graph and knowledge panels*. Google. Retrieved December 7, 2022 from <https://blog.google/products/search/about-knowledge-graph-and-knowledge-panels/>