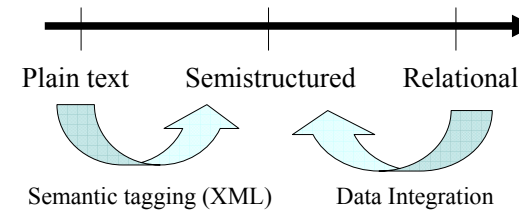


Formal languages problems in semistructured data management

S. Afonin
Moscow State University

Semistructured data

- **Semistructured data** (SSD) are the data with irregular, rapidly changing or even unknown structure.



FDPW-2005

SSD Formal models

- Data structure: tree, ordered tree, directed graph
- Integrity constraints: data schemes (e.g. XML DTD), additional constraints (path equivalencies)
- Query language: functional, regular path queries (CRPQ, XPath)

FDPW-2005

Graph representation

The directed edge-labelled graph $D = \langle V, \Sigma, E \rangle$ where

V is the set of vertices, and

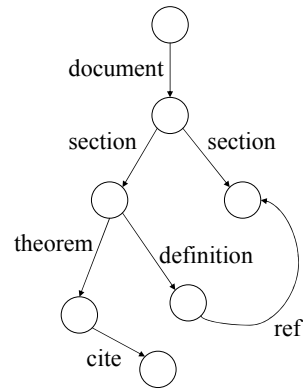
$E \subseteq \langle V \times \Sigma \times V \rangle$ is the set of Σ -labelled edges is called a *semistructured database*.

Vertices correspond to domain-specific objects, edges represent relations between these objects.

Edge labels reflect relation type.

FDPW-2005

Graph representation



```

\begin{document}
\section{A}\label{secA}
\section{B}
\begin{definition}
see also \ref{secA}
\end{definition}
\begin{theorem}\cite{Ivanov}}
\end{theorem}
\end{document}
    
```

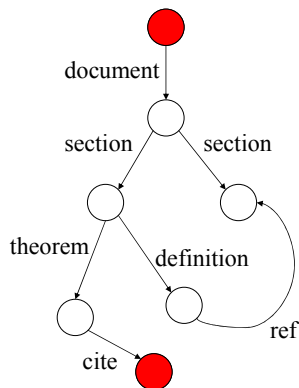
FDPW-2005

Regular path queries

- A *query* is a regular expression (language)
- Query *result* is the set of all pairs of the graph vertices such that there exists at least one path between them labeled by a word from a given language
- Example: $X^*.subsection.(theorem+definition).cite Y$
- Conjunctive regular path queries:
 - $X^*.theorem Y, Y cite Z$

FDPW-2005

SSD Examples



Query:

```

*.section.(theorem+definition).cite
*(sub)*section.(theorem+definition).cite
    
```

The whole graph may need to be searched during query evaluation, which is inefficient

FDPW-2005

View-based query processing

- Let us assume that we know the results for the queries E_1, \dots, E_k .
- Can this data be used during evaluation of an arbitrary query?
- Is it possible to compute the result of a query using the result of view queries only?

FDPW-2005

Definitions

- An **alphabet** is a finite non-empty set of symbols
- A **word** is a finite sequence of symbols; ϵ - is the **empty word**
- A **language** is an arbitrary set of words

FDPW-2005

Language operations

- Union
 - $L_1 + L_2 = \{ u \mid u \text{ in } L_1 \text{ or } u \text{ in } L_2 \}$
- Concatenation
 - $L_1.L_2 = \{ uv \mid u \text{ in } L_1 \text{ and } v \text{ in } L_2 \}$
- Iteration
 - $L^* = \{ \epsilon \} + L + L.L + L.L.L + \dots$
- Regular languages is the minimal class of languages which contains *singletons* and closed under union, concatenation and iteration.

FDPW-2005

Language representation problem

- For finite set $\{E_1, \dots, E_k\}$ of regular languages, and a subset T of language operations $\{+, \cdot, *\}$ is it decidable whether or not a given language R may be constructed from $\{E_i\}$ using a finite number of operations from T ?
 - $T = \{\text{concatenation}\}$ - language factorization
 $R = a^* + a^*ba^*b(a+b)^*$ $E_1 = a^* + a^*b(a+b)^*b$ $E_2 = (bab)^*$
 $R = E_1.E_2.E_2.E_1.E_1$
 - $T = \{+, \cdot, *\}$ - maximal rewriting
- Language representation problem is decidable (K.Hashiguchi, 1982)

FDPW-2005

Representation & query evaluation

- $R = a^* + a^*ba^*b(a+b)^*$
 $E_1 = a^* + a^*b(a+b)^*b$ $E_2 = (bab)^*$
 $R = E_1.E_2.E_2.E_1.E_1$
- E_1 ■ E_2

(1, 2)	(1, 3)	E_1	E_2	E_2	E_1	E_1
(1, 4)	(2, 1)	1	--> 2	--> 1	--> 3	--> 1
(2, 3)	(3, 2)			...		
(3, 1)	(2, 2)			...		

FDPW-2005

Semigroups of regular languages

- Regular languages are closed under concatenation
- $(S, \cdot) = \langle E_1, \dots, E_k \rangle$ - finitely generated semigroup
- A query R may be represented in terms of $\{E_i\}$ if and only if the language R belongs to (S, \cdot)
- The membership problem is decidable (K.Hashiguchi)

FDPW-2005

Query rewriting under constraints

- Theorem** The set of all possible representations of an element R is a regular language.

$$\square E = (e + b^*a)(b + ab^*ab^*a)^* \quad (S, \cdot) = \langle E \rangle \quad R = E.E.E$$

$$R = (E.E.E).E^*$$

- Open Problem** For given set Q_1, \dots, Q_n of “the most popular queries” find the minimal number of semigroup generators E_1, \dots, E_k , such that all the languages Q_i belongs to the semigroup

FDPW-2005

Language equations

- For given regular languages L (query) and E (view) find a regular language X such that $L = EX$
- If the equation has a solution then it has the unique maximal solution, which is a regular language (L.Kari)
- Maximal solution contains all solutions
- Minimal solution contains no solutions
- If the equation has a solution then it has at least one minimal solution (L.Kari)

FDPW-2005

Maximal and minimal solutions

- Maximal solution drawback: redundancy
 - $a^*b = a^*.X$ $X = a^*b$ is the maximal solution
 $X = b$ is a solution as well
- Infinitely many minimal solutions
 - $a^* = (e+a).X$
 - $X = (aa)^*$ and $X = e + a(aa)^*$ are minimal solutions
 - $X = \{0, 2, 4, \dots, 2n, 2n+1, \quad 2n+3, 2n+4, 2n+6, \dots\}$
 $2n+2$

FDPW-2005

Open Problems

- Is there exists an algorithm for minimal solution finding?
- Is there exist regular languages L and E such that the equation $L=EX$ has no regular minimal solutions?

FDPW-2005

Thank you

FDPW-2005