## COnnected KEywords

#### Grégory Smits, Olivier Pivert and Virginie Thion

IRISA, Campus de Beaulieu, 35042 Rennes, France {gregory.smits, olivier.pivert, virginie.thion}@irisa.fr

RCIS, May 13-15 2015, Athens, Greece Institut de recherche en informatique et systemes aléatoires



.



Cooperative Querying Approach

## Giving Access to Structured Data



#### SQL QUERY LANGUAGE





INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

G. Smits et al.

COnnected KEywords

## Giving Access to Structured Data





G. Smits et al

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

2/20

## Motivations



More expressive query interface:

keywords and connectives,

 constrained but rich query language, (SPJ queries, aggregates, some nested queries, ...) 'movies Clint E. Gene H.' 'genre of movies entitled Scarface' 'movies genre is western released in 2014'

'actors appearing in Dogma'

- tool to help domain experts:
  - query their corporate databases,
  - manage the COKE language.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

## Outline

### Introduction

- 1.1 Problem
- 1.2 Motivating example

### 2 Data Model and Language

- 2.1 Schema Graph Model
- 2.2 Query Graph Model and Searchable Vocabulary
- 2.3 a Constrained Keyword-based Query Language
- 3 Cooperative Querying Approach
- 3.1 KW Query Definition
- 3.2 KW Query Interpretation
- 3.3 Translation into SQL



G. Smits et al

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

Translation of the Relational Schema into a Graph Schema

- capture the data semantics,
- have a more manageable model,
- make it possible to apply graph-based algorithms.



- tables, attributes and definition domains represented by nodes,
- primary to foreign key constraints form links between table nodes,
- domain nodes are linked to their attribute nodes themselves linked to their table nodes.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

G. Smits et al.

COnnected KEywords

## Query Graph Model: mapping between the model and SQL

Directed graph representing the elements that may appear in a user query:

- table, attribute, domain and function nodes,
- projection, predicate, selection, join and transformation edges,
- subquery edges.

#### Definition

A query is covered by the COKE system if it forms a connected subgraph of the query graph.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

## Query Graph Model: mapping between the model and SQL

'name of producers, title and production year of movies with Michael Keaton in the role of Batman'



SQL> SELECT title, production\_year, P.name FROM producer P join movies M on M.pid = P.id join mov\_cast MC on MC.mid = M.id join characters C on C.id = MC.cid join actors A on A.id = MC.aid WHERE a.name = 'Keaton' and C.name = 'batman'; 'number of movies produced by the producer of the movie entitled Dogma'



SQL> SELECT count(\*)
FROM movies M join mov\_prod MP on M.id = MP.mid
 join producers P on P.id = MP.pid
WHERE P.name = ( SELECT P.name
 FROM producers P2 join mov\_prod MP2 on
 P2.id = MP2.id
 join movies M2 on M2.id = MP2.mid
 WHERE M2.title = 'Dogma');



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

G. Smits et al.

## Searchable Vocabulary: mapping keywords to elements of the query graph



#### Definition

A keyword is a word or a group of words attached to a searchable element of the query graph:

- table, attribute, domain and function nodes,
- projection, predicate, join, selection and transformation edges,
- subqueries edges,
- and some paths (predicate + selection, joins, joins + selection).

## **∮**≸IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

G. Smits et al

COnnected KEywords

# Searchable Vocabulary: mapping keywords to elements of the query graph



Extract of the labelled query graph

## **∮**≸IRISA

G. Smits et al.

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

## COKE: a Constrained Keyword-based Query Language

A COKE query is a keyword-based expression of a selection-projection-join SQL query:

- with a restricted vocabulary,
- and a constrained syntax (projection part + selection part).

Backus Naur form grammar: COKE\_query ::= {Projection\_Statement}\* {Selection\_Statement}\* Projection\_Statement ::=  $F^{f} [e^{f}] A^{k} e^{\pi} R^{k} | F^{f} [e^{f}] R^{k} | [Attributes [e^{\pi}]] R^{k} | Attributes$ Attributes :=  $A^{k}_{l} \{, A^{k}_{l}\}* [and A^{k}_{l}]$ Selection\_Statement ::=  $\{e^{\bowtie} | p^{\bowtie^{*}}\} R^{k} e^{\sigma} A^{k}_{l} e^{\theta} D^{k}_{l} | A^{k}_{l} e^{\pi} R^{k} e^{\theta} D^{k}_{l} | \{e^{\bowtie} | p^{\bowtie^{*}}\} R^{k} p^{\sigma\theta} D^{k}_{l} | A^{k}_{l} e^{\theta}$  $D^{k}_{l} | R^{k} p^{\sigma\theta} D^{k}_{l} | p^{\sigma\theta} D^{k}_{l} | p^{\bowtie^{*}\sigma\theta} D^{k}_{l} | R^{k} D^{k}_{l} | D^{k}_{l} | e^{\eta^{op}} D^{k}_{l}.$ 



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

G. Smits et al.

## The COKE Syntax Illustrated

#### COKE query = a projection part + a selection part

#### Projection statement:

- basic form: "movies", "title year movies", "number movies", etc.,
- more expressive syntax: "title, length and genre of movies", "number of movies", etc.

Selection statement (joins, restrictions, some nested queries):

- basic form: "Clint E. Gene H.", "actor Clint E.", "producer Clint E. actor Gene H.", etc.,
- more expressive syntax:. "by a producer whose name is Clint E.", "produced by Clint E. with Gene H.", "with an actor of Dogma", etc.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

## COKE: a Cooperative Querying Approach

A constrained keyword-based query language with:

- 1. a restricted vocabulary,
- 2. a limited syntax,
- 3. and semantic expectations (connected subgraph of the query graph).

An interactive query construction and interpretation process to:

- help users express their information needs with covered queries,
- interactively determine the exact meaning of the queries.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE



Cooperative Querying Approach 000000

13/20

## Query Definition at the Lexical Level



## Autocompletion and locked query field to ensure lexical coverage:

index of the whole vocabulary.



#### Definition

The lexical analyzer (made easy by autocompletion) produces mappings between keywords and elements of the query graph.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

G. Smits et al.

## Query Interpretation at the Syntactic Level

Parsing to identify valid projection and selection statements

#### Definition

A syntactic analysis is a mapping from adjacent keywords to projection/selection graph patterns.

 $\begin{aligned} \mathsf{COKE\_query} &:= \{\mathsf{Projection\_Statement}\}^* \{\mathsf{Selection\_Statement}\}^* \\ \mathsf{Projection\_Statement} &:= F^f \left[ e^f \right] A_l^k \ e^\pi \ R^k \ | \ F^f \left[ e^f \right] R^k \ | \ [Attributes \left[ e^\pi \right] \right] R^k \ | \ Attributes \\ \mathsf{Attributes} &:= A_l^k \ \{, \ A_l^k\}^* \ [and \ A_l^k] \\ \mathsf{Selection\_Statement} &::= \{ e^{\bowtie} | p^{\bowtie^*} \} \ R^k \ e^\sigma \ A_l^k \ e^\theta \ D_l^k \ | \ A_l^k \ e^\pi \ R^k \ e^\theta \ D_l^k \ | \ \{ e^{\bowtie} | p^{\bowtie^*} \} \ R^k \ p^{\sigma\theta} \ D_l^k \ | \ A_l^k \ e^\theta \\ D_l^k \ | \ R^k \ p^{\sigma\theta} \ D_l^k \ | \ p^{\bowtie^*\sigma\theta} \ D_l^k \ | \ R^k \ D_l^k \ | \ P^k \ e^{\eta^{op}} \ D_l^k. \end{aligned}$ 

#### Example:

 $\langle$  ('title production year of movie',  $\{A_i^k\} + e^{\pi} R^k$ ) ('with an actor named Mickael Keaton',  $e^{\triangleright \triangleleft} R p^{\sigma \theta} D$ ) ('in the role of Batman',  $p^{\triangleright \dashv} \sigma^{\theta} D_k^k$ ) $\rangle$ .

#### € IRISA G. Smits et al.

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

## Interactive Syntactic Disambiguation

Some syntactic patterns may be ambiguous or uncovered:

- ask for a confirmation of the right interpretation,
  - expansion of ambiguous patterns.
- show syntactically uncovered parts of the query.



## **∮** §IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

COnnected KEywords

## Discovering the Meaning of the Query

#### Definition

A query is semantically consistent if its interpretation forms a connected subgraph of the query graph.

In case of isolated statements (e.g. 'title of movies name Clint E.'), re-establish the semantic consistency:

- identify the subject (relation node) of each statement,
  - $\blacktriangleright$  projection statement  $\rightarrow$  the concerned relation,
  - ► selection statement → the relation on which a restriction is applied (selection edge).
- determine meaningful joins,
  - link an isolated subject node to a previously activated subject,
  - favor the shortest join path between two subjects.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

G. Smits et al

COnnected KEywords

## Discovering the Meaning of the Query (example)

"title genre and production year of movies actor named Gene H. producer Clint E."

Subject relation nodes: {movies, actors, producers}.

- projection statement,
   P<sup>1</sup> "title, genre and
   production year of movies"
- selection statements, S<sup>1</sup> "actor named Gene H." S<sup>2</sup> "producer Clint E."





INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

## Translation to SQL

Depth-first traversal of the connected subgraph:

- starting with projection statements,
- translation rules for each graph pattern,
- ▶ incremental filling of the SELECT, FROM and WHERE clauses.



SQL> SELECT movie.title, movies.genre, movie.production year FROM movies, mov cast, actors, mov prod, producers WHERE movies.idM = mov cast.idM and mov cast.idA = actors.idA and actors.name = 'Gene Hackman' and movies.idM = mov prod.idM and mov prod.idP = producers.idP and producers.name = 'Clint Eastwood'; Translation rules:

- $A \pi R \rightarrow \text{append } A \text{ to}$ SELECT and R to FROM,
- *R*<sup>1</sup> ⋈ *R*<sup>2</sup> append a natural join between *R*<sup>1</sup> and *R*<sup>2</sup> to *FROM*,

18/20

•  $A \theta D$  append  $A\theta D$  to WHERE.

#### **∮**§IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

G. Smits et al

## Experimentations

Implementation on top of the MusicBrainz DB:

- manual definition of the query graph and vocabulary.
- Expressivity assessment:
  - translation of the QALD corpus (100 NL queries):
    - ▶ 81 queries successfully translated into COKE queries,
    - ▶ 8 concern data not available in the DB version of MusicBrainz,
    - ▶ 11 uncovered (6 Boolean queries and 5 relying on indirect knowledge).
- A first attempt of precision assessment (precision 100%):
  - translation by a "COKE" queries expert,
  - using unambiguous statements, no join deduction,
  - basic queries (80/81 queries contain 1 projection and 1 selection).

Efficiency assessment:

► NL approach: avg. 31.6s, COKE: avg. 0.65s.

G. Smits et al.

## Conclusion and Perspectives

A keyword-based query approach taking connectives into account:

- constrained keyword-based language,
- interactive query definition and disambiguation strategy.

A tool to help domain experts query their corporate DBs

Perspectives:

- make COKE queries closer to natural language,
- automatic acquisition of the vocabulary,
- handle more complex queries (aggregates, groups),
- real users feedbacks in a CRM applicative context.



INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRE

G. Smits et al.

COnnected KEywords