

Keyword Search over Data Service Integration for Accurate Results

Vidmantas Zemleris

CMS Computing dept., CERN / LSIR, IC, EPFL
Supervised by: prof. K.Aberer, dr. R.Gwadera (EPFL);
dr. V.Kuztesov (Cornell), dr. P.Kreuzer (CERN)

15th April 2013



Outline

- 1 Introduction
 - Problem statement
 - State-of-the-Art
- 2 Implementation
- 3 Conclusions & Future work
- 4 Back-up/extra

Outline

1 Introduction

- Problem statement
- State-of-the-Art

2 Implementation

3 Conclusions & Future work

4 Back-up/extra

Preliminaries (1/2)

Virtual data service integration (EII)

- *lightweight **virtual** integration* (vs. data-warehousing, publish-subscribe)
- usually queried with structured languages, e.g. SQL, YQL, etc
- growing # of sources and applications: corporate, governmental, mashups...
 - ▶ e.g. Yahoo's YQL, Google Fusion Tables, ...

How it works?

- process the query & send requests to services
- consolidate the results:
 - ▶ namings & dataformats (XML, JSON, ..)
 - ▶ apply filters, aggregations, service composition

Preliminaries (1/2)

Virtual data service integration (EII)

- *lightweight **virtual** integration* (vs. *data-warehousing, publish-subscribe*)
- usually queried with structured languages, e.g. SQL, YQL, etc
- growing # of sources and applications: corporate, governmental, mashups...
 - ▶ e.g. *Yahoo's YQL, Google Fusion Tables, ...*

How it works?

- process the query & send requests to services
- consolidate the results:
 - ▶ namings & dataformats (XML, JSON, ..)
 - ▶ apply filters, aggregations, service composition

Preliminaries (2/2)

Example query (in DASQL used at CMS, CERN)



Remark: this is close to boolean retrieval + (aggregation XOR projections).

The problem: it is **tiring/overwhelming** to:

- learn a query language
- remember how exactly data is structured and named

Could this be extended by Keyword Queries? e.g.

- *list sizes of RelVal datasets where number of events > 1000*
- *avg(dataset size) Zmmg 'number of events' > 1000*

Preliminaries (2/2)

Example query (in DASQL used at CMS, CERN)



Remark: this is close to boolean retrieval + (aggregation XOR projections).

The problem: it is **tiring/overwhelming** to:

- learn a **query language**
- remember how exactly **data is structured and named**

Could this be extended by Keyword Queries? e.g.

- *list sizes of RelVal datasets where number of events > 1000*
- *avg(dataset size) Zmmg 'number of events' > 1000*

Preliminaries (2/2)

Example query (in DASQL used at CMS, CERN)



Remark: this is close to boolean retrieval + (aggregation XOR projections).

The problem: it is **tiring/overwhelming** to:

- learn a **query language**
- remember how exactly **data is structured and named**

Could this be extended by Keyword Queries? e.g.

- *list sizes of RelVal datasets where number of events > 1000*
- *avg(dataset size) Zmmg 'number of events' > 1000*

Problem statement

Given:

- schema terms (entity and field names)
- value terms
 - ▶ values listing (for some fields)
 - ▶ constrains, e.g. regexps, mandatory service inputs
- query: $KWQ = (kw_1, kw_2, \dots, kw_n)$

Task: interpret each $kw_i \in KWQ$ as:

- schema term (result type; projections; or field name in a predicate)
- values term (a value in a predicate) [$\text{pred} := \text{field_name op. value}$]
- (aggregation) operator, or *unknown*.



Challenges

- no direct access to the data, use:
 - ▶ bootstrapping values lists (available only for some fields)
 - ▶ rely on regexps otherwise
- no fully predefined schema
 - ▶ some field names are unclean (directly from XML, JSON responses)
 - ▶ bootstrap list of fields through queries and maintain it...

State of the art

- Nature of Keyword queries:
 - ▶ ambiguous: *structured queries as results*
 - ▶ nearby keywords are related
- “Keyword Search over EII” received not much attention:
 - 1 KEYMANTIC - rule-based + heuristics
 - 2 KEYRY - HMM-based
 - 3 Natural Language to compose services (farther...)

Closely related works

① (I do not need that much detail:)

② KEYMANTIC

- ① based on metadata, keywords are scored as potential matches to *schema terms* (entity names and their attributes, using some entity matching techniques) or as potential *value* matches (by checking any available constraints, such as the regular expressions imposed by the database or data-services).
- ② Then, these scores are combined, and refined by heuristics that increase the scores of query interpretations with the nearby keywords having related labels assigned.
- ③ Finally, these labels are interpreted as SQL queries.

③ KEYRY incorporate users feedback through training an Hidden Markov Model's (HMM) tagger taking keywords as its input.

- ① It uses the List-Viterbi [?] algorithm to produce the top-k most probable tagging sequences (where tags represent the “meaning” of each keyword).
- ② This is interpreted as SQL queries and presented to the users.
- ③ The HMM

Uniqueness of this implementation

- ① no assumptions on input query
 - ▶ plain keywords vs. full-sentence
 - ▶ still can use patterns if present (phrases, predicates/conditions)
- ② implements a specific real-world use-case
 - ① different selection of entry points / scoring heuristics and entity matching
 - ② scoring
 - ③ specific query language
- ③ first open-source implementation
 - ▶ the code will be further maintained

Outline

- 1 Introduction
 - Problem statement
 - State-of-the-Art
- 2 Implementation
- 3 Conclusions & Future work
- 4 Back-up/extra

Implementation Overview

- ❶ *tokenizer*: clean up; identify patterns
- ❷ identify and score “*entry points*”
 - ❶ string matching [for entity names]
 - ❷ IR (IDF-based)[unclean fieldnames]
 - ❸ known values
 - ❹ regular expressions on allowed values
- ❸ combine *entry points*
 - ❶ consider various *entry points* permutations (called *configurations*)
 - ❷ promote ones respecting keyword dependencies or other heuristics
 - ❸ interpret as structured queries

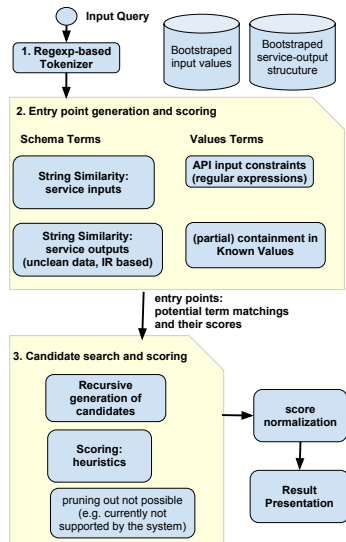


Figure 1: Query processing

Project priorities and constraints

Due to the constraints on the project duration, a number of items had to be excluded from the implementation: question answering approaches with deep language processing; complex service orchestration (feeding of outputs into inputs of other services, which is anyway not directly supported by the EII system and the service performance is not adequate for this¹); and lastly the performance is of lower priority, as the end user's perceived performance is still dominated by services taking minutes to respond, and the performance was already covered by the earlier works.

¹this due to issues with data service performance and unavailability of basic capabilities such as pagination or sorting of their results; we do not control the data services, so a number of suggestions for the providers have been proposed (see appendix ??); second, these improvements would take a considerable effort to be implemented, pushing this far beyond the scope of this project

Example query processing

dataset sizes RelVal 'number of events > 1000'

RelVal -> 1.0, value: group=RelVal

RelVal -> 0.7, value: dataset=*RelVal*

datasets -> 0.9, schema: dataset

Chunks:

'number of events>1000'

--> 0.93, filter: dataset.nevents>1000

--> 0.93, filter: file.nevents>1000

'datasets; sizes'

--> 0.99, projection: dataset.size

'sizes'

--> 0.41, projection: dataset.size

... and some more with lower scores...

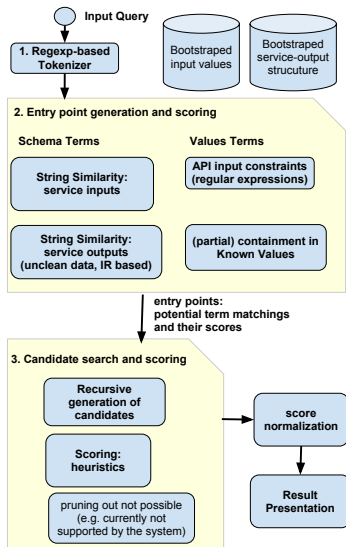


Figure 2: Query processing

Step 1: Tokenizer

1 Clean-up

- ▶ remove extra spaces, normalize formatting
- ▶ recognize simple unambiguous expressions

2 Split into tokens on these regular expressions:

- 1 [terms] operator value (e.g. "number of events">10, dataset=Zmm)
- 2 terms in quotes (e.g. "magnetic field")
- 3 individual terms

Future - aggregates, e.g. avg(number of events)

– more exact than ambiguous “avg number events”

– easier than currently strict syntax: | avg(dataset.nevents)

Step 2: Entry point Generation and Scoring (1/2)

Matching schema terms

- did not work well: string edit-distance, semantic similarity

$$\text{similarity}(A, B) = \begin{cases} 1, & \text{if } A = B \\ 0.9, & \text{if } \text{lemma}(A) = \text{lemma}(B) \\ 0.7, & \text{if } \text{stem}(A) = \text{stem}(B) \\ 0.6 \cdot \text{dist}(\text{stem}(A), \text{stem}(B)), & \text{otherwise} \end{cases}$$

Matching multi-word unclean schema terms

- some terms are non-informative → **IDF needed**
- use Information Retrieval library (Whoosh) with BM25F scoring
- create *virtual documents* each representing “a field of an entity”
 - ▶ fully-qualified machine readable (e.g. block.replica.creation_time)
 - ★ tokenized+stemmed (e.g. creation_time)
 - ★ context - tokenized+stemmed parent (e.g. block; replica)
 - ▶ human readable title, if any (e.g. “Creation time”)

Step 2: Entry point Generation and Scoring (2/2)

Matching Value terms:

- Regular expression (regexp) can result in false positives:
 - ▶ it do not guarantee that a certain value exists
 - ▶ regexp could be loosely defined
 - ▶ thus, regexp matches are scored lower than other methods
- Known values
 - ▶ automatically bootstrapped
 - ▶ decreasing score: full match, partial match, and matches of keywords containing wildcards
 - ▶ If keyword's value matches a regular expression, but is not contained in the known values list and the accepted values of the given field are considered to be static (not changing often), we exclude this very likely false match that reduces the false positives.

Step 3: Answer candidate scoring: formulas

$$score_avg = \frac{\sum_{i=1}^{|KWQ|} \left(score(tag_i|kw_i) + \sum_{h_j \in H} h_j(tag_i|kw_i; tag_{i-1,...,1}) \right)}{N_non_stopword} \quad (1)$$

$$score_prob = \sum_{i=1}^{|KWQ|} \left(\ln(score(tag_i|kw_i)) + \sum_{h_j \in H} h_j(tag_i|kw_i; tag_{i-1,...,1}) \right) \quad (2)$$

- $score(tag_i|kw_i) \sim$ likelihood of kw_i to be tag_i
- $h_j(tag_i|kw_i; tag_{i-1,...,1})$ - the score boost returned by heuristic h_j given a tagging so far (all tags are not needed).

Step 3: Answer candidate scoring: heuristics

- Relationships between keywords:
 - ▶ promoting such combinations where nearby keywords refer to related schema terms (e.g. entity name and it's value)
 - ▶ balance between taking the keyword or leaving it out (the one that we are unsure about)
 - ▶ boost important keywords (different parts of speech are of different importance, e.g. stop-words are less useful than nouns)
- Qualities of Data Integration System:
 - ▶ promote data service inputs over filters on their results: 1) it is more efficient, especially when this is possible; 2) there are much more of possible entities to filter, so more false matches are expected there, while the service inputs shall cover large part of cases
 - ▶ if some keyword can be matched as the requested entity, and mapping of other keywords fits the service constraints
 - ▶ if requested entity and a filter condition is the same (a small increase, a common use-case is retrieving an entity given it's "primary key" identifier or a wild-card)
 - ▶ for being able to execute the query, the service constraints must be satisfied; still it could useful to the interpretations that achieve high

Tuning the scoring parameters

- ① tuned individual components to “sufficient” level
 - ▶ unit tests and manual testing
- ② fine-tune the whole system (by hand)
 - ▶ use keyword queries by written users or developer for evaluation
 - ▶ important variables to be tuned:
 - ★ weights for regexps, **etc**
 - ★ likelihood of not taking a keyword
 - ★ BM25F “field” and “query” weights (for IR matching multi-word terms)

Evaluation: Accuracy

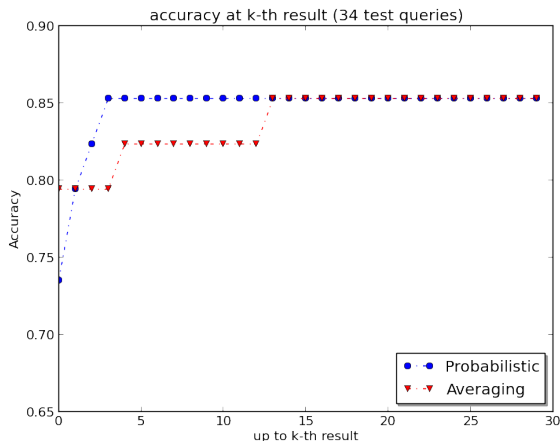


Figure 3: Accuracy comparison of the two scoring methods at kth result

- accuracy of 85% @ 4th suggestion
- testing set is limited - need more live feedback

Presenting the results to the user

Query: Zmmg number of events>10

Are searching for: [dataset](#), [file](#), [block](#), [run](#), [status](#), [see all](#)

color coding:

input predicates - cheap

filters on outputs - expensive

entity to return

```
0.52 file dataset=*Zmmg* | grep file.name, file.nevents>10 debug
0.52 dataset dataset=*Zmmg* | grep dataset.name, dataset.nevents>10 debug
0.52 block dataset=*Zmmg* | grep block.name, block.nevents>10 debug
0.28
0.28 Explanation:
find Block name (i.e. block.name) for each block where dataset=*Zmmg* AND Number of events (i.e. block.nevents) >
0.12 dataset dataset=*Zmmg* | grep dataset.nevents, dataset.name, dataset.nfiles>10 debug
0.12 dataset dataset=*Zmmg* | grep dataset.nevents, dataset.name, dataset.nblocks>10 debug
0.08 dataset dataset=*event* | grep dataset.name, dataset.nblocks>10 debug
0.08 file dataset=*event* | grep file.name, file.nevents>10 debug
0.08 run dataset=*event* | grep run.run_number, run.nlumis>10 debug
```

Showing only top 10 suggestions. [see all](#)

Autocompletion prototype

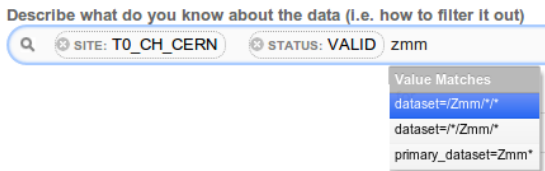


Figure 5: prototype of auto-completion based interface

an interesting feature of this is, it could seamlessly provide feedback for:

- entity matching techniques
- (or even keyword search itself)

Using the feedback for self-improvement

First, the implicit feedback from auto-completion could be useful for evaluating and improving the quality of the entity matching (learned edit-distance metrics, updating the weights of different matching metrics). Second, the user's selections in auto-completion fields could serve as training data for machine learning-based algorithms (and it is of better quality because user is selecting auto-completed values for separate terms), however it is important to gather sufficiently large sets of high quality feedback, to avoid over-fitting the machine learning models.

Also users implicit feedback (clicking on the link), could be useful, however it is of limited quality (the user may click on it just for figuring out what the query returns even if it was not directly related with his information needs expressed as the the input keywords). An alternative to this could be asking users if the result was what they were asking for, when they see the results of actual structure query (this could be a little bit overwhelming, but looks quite optimal).

Additional problem is that what was so far modelled by the sequential machine learning algorithms, such as HMM, was not directly the structured query, however the query configurations

Outline

- 1 Introduction
 - Problem statement
 - State-of-the-Art
- 2 Implementation
- 3 Conclusions & Future work
- 4 Back-up/extra

Conclusions

- growing popularity of data-service integration systems, increase need accessing data easily
 - ▶ one successful approach can be keyword search proposing structured queries
 - ▶ keyword search in EII (i.e. then no complete access is available) received fairly little attention...
- discussed a real-world case and **first open-source** implementation
 - ▶ able to use patterns
 - ▶ **users quite like the idea..**
 - ▶ **most of performance problems may be deep in underlying services...**
 - ▶ the system is going to be further supported
 - ▶ **this will contribute to improve the overall efficiency of the physics analysis program by the CMS Collaboration**

Future work

- better interpretation of patterns in queries
- tuning the accuracy based on users' feedback
- exploring the auto-completion further
- exploring the Machine Learning approaches once more data is gathered?
- technical
 - ▶ large parts of keyword search can be moved to client-side
 - ▶ performance improvements

Outline

- 1 Introduction
 - Problem statement
 - State-of-the-Art
- 2 Implementation
- 3 Conclusions & Future work
- 4 Back-up/extra

Complete list of project deliverables

① keyword search engine and related components

- ▶ implementation of entity matching techniques & heuristics
- ▶ code for bootstrapping of: 1) allowed values, 2) fields in service results
- ▶ tuning the system's parameters
- ▶ prototype of advanced auto-completion input widget
- ▶ slight relaxation of DASQL: prototype of "simple service orchestration" even then the existing fields are not known in advance - gives more power and simplifies the keyword search

② log analysis and data service performance benchmarking at CMS

- ▶ proposed solutions for data service providers

③ user surveys, presentations and tutorials at the CMS Collaboration

- ▶ constant cooperation with a selected group of ~5+ users for feedback

Data integration war-stories: Dataservice Performance

References

- prototype online: <https://docs-bulk-tool.cern.ch/das/>