

3rd MUMIA Training School
FORTH, Heraklion, Crete – Greece
21-25 July 2014



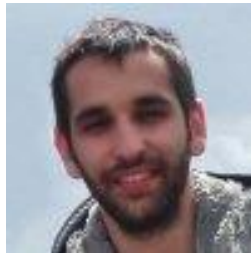
BRIDGING THE WEB OF DOCUMENTS WITH THE WEB OF DATA AT SEARCH TIME

Yannis Tzitzikas
University of Crete and
FORTH-ICS

ACKNOWLEDGEMENTS

- Research done mainly with my students (more in the acknowledgements)

- Pavlos Fafalios



- Panagiotis Papadakos

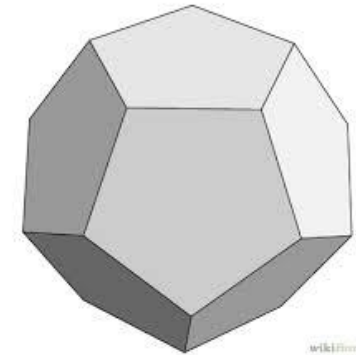
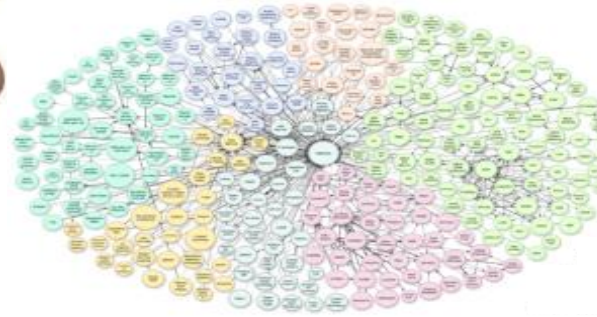


OUTLINE

- Introduction and Background (60')
 - Exploratory Search
 - Semantic Web and Linked Data
 - Faceted Exploration
- Bridging the Web of Documents with the Web of Data (50')
 - Possible Ways
 - Focus on doing this at search time
 - Case studies: presented as a series of 9 milestones
- Synopsis, Challenges, Discussion (10')
- References
- Acknowledgements



“Study without desire spoils the memory, and it retains nothing that it takes in.” Leonardo da Vinci



1. INTRODUCTION AND BACKGROUND

Exploratory Search
Semantic Web and Linked Data
Faceted Exploration

4

WHAT USERS USUALLY WANT/DO WHEN SEARCHING?

Kinds of information needs

- *Precision-oriented*
 - Locate one resource or/and its attributes
e.g. Find the telephone of a store
- *Recall-oriented*
 - Locate a **set** of resources
e.g. Medical information seeking, travel planning

Over 60% of web search queries are ***recall-oriented***
[Broder 02, Rose and Levinson 04]

RECALL-ORIENTED INFORMATION NEEDS

- In Recall-Oriented Information Needs:
 - the users require >1 hit
 - essentially such needs correspond to **decision tasks**
- Examples of Recall-oriented information needs
 - Booking
 - Product-buying
 - Bibliography search
 - Patent Search
 - Medical Search
 -

QUESTIONS



How many of you (raise your hand):

- have bought your mobile phone by looking at the first hit(s) that google returned to your query?
- have ever booked the 1st hotel returned by google?
- have entirely read the first paper that Scholar Google returned to your query?
- have made at least once a search about a medical issue?
 - keep your hand raised if you had submitted to Google only one query
 - keep your hand raised if you have read only the first 5 hits of that single query that you submitted to Google, and eventually that was enough for fully satisfying your information need

EXPLORATORY SEARCH



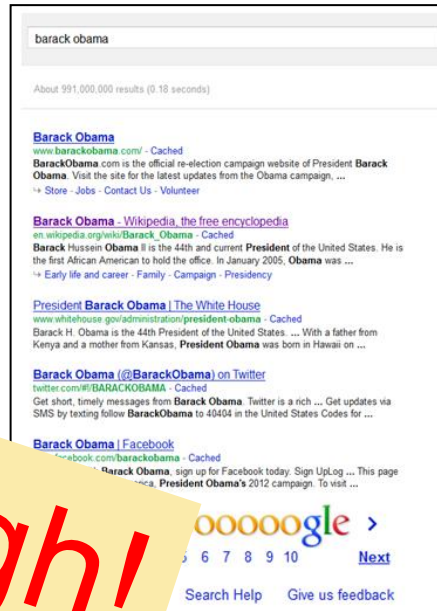
Wikipedia:

“**Exploratory search** is a specialization of information exploration which represents the **activities carried out by searchers** who are either:

- a) **unfamiliar with the domain** of their goal (i.e. need to learn about the topic in order to understand how to achieve their goal)
- b) **unsure about the ways to achieve their goals** (either the technology or the process)
- c) or even **unsure about their goals in the first place**.

Consequently, exploratory search covers **a broader class of activities** than typical **information retrieval**, such as **investigating, evaluating, comparing, and synthesizing**, where new information is sought in a defined conceptual area; **exploratory data analysis** is another example of an information exploration activity. Typically, therefore, such users generally **combine querying and browsing** strategies to foster learning and investigation.”

THEREFORE...



○ Ranking is not enough for exploratory search

By GARY MARCHIONINI

EXPLORATORY SEARCH: FROM FINDING TO UNDERSTANDING

Research tools critical for exploratory search success involve the creation of new interfaces that move the process beyond predictable fact retrieval.

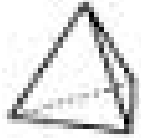
From the earliest days of computers, search has been a fundamental application that has driven research and development. For example, a paper published in the inaugural year of the *IBM Journal* 36 years ago outlined challenges of text retrieval that continue to the present [4]. Today's data storage and retrieval applications range from database systems that manage the bulk of the world's structured data to Web search engines that provide access to petabytes of text and multimedia data. As computers have become consumer products and the Internet has become a mass medium, searching the Web has become a daily activity for everyone from children to research scientists.

SOME COMMON REQUIREMENTS FOR EFFECTIVE EXPLORATORY SEARCH

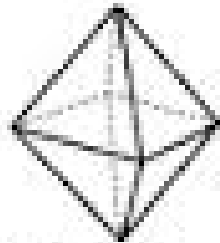


- Allow easy and fast access even to **low ranked** hits
- Allow browsing and inspecting the found hits in **groups** (according to various criteria)
- Offer **overviews** of the search results
 - Compute and show descriptions and **count** information for the various groups, or other **aggregated** values
- Allow **gradual** restriction/ranking of the search results

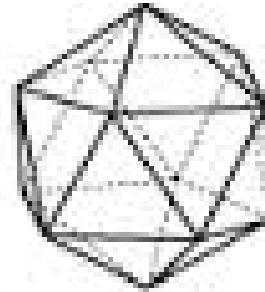
FACETED SEARCH



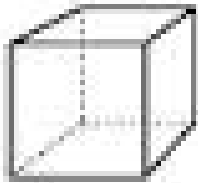
Tetrahedron



Octahedron



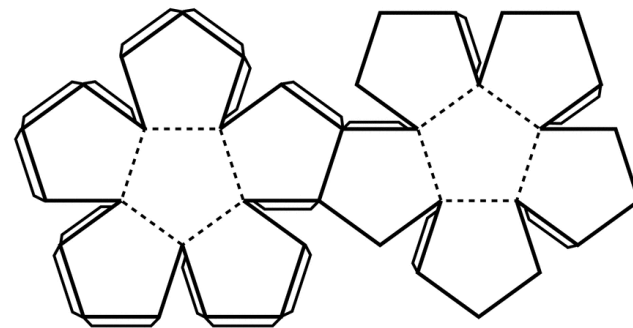
Icosahedron



Cube



Dodecahedron



FACETED SEARCH/EXPLORATION

Faceted Exploration is a **widely** used interaction scheme **for Exploratory Search**

A short (and rather informal) definition:

FE is a **session-based** interactive method for **query formulation** (commonly over a multidimensional information space) through simple clicks offering

- ✓ an overview of the result set (groups and count information)
- ✓ never leading to empty results sets

- The access paradigm supported is a conceptual exploration,
 - far more frequent in "search" tasks than the retrieval by exact specification supported by search engines and database queries.
- Simple and easily understood by users.



FACETED SEARCH

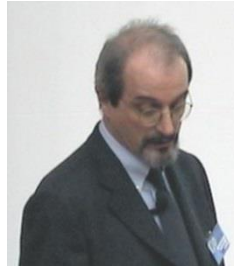
Wikipedia:

- **Faceted search**, also called **faceted navigation** or **faceted browsing**, is a technique for accessing information organized according to a faceted classification system, allowing users to explore a collection of information by applying multiple filters. A faceted classification system classifies each information element along multiple explicit dimensions, called facets, enabling the classifications to be accessed and ordered in multiple ways rather than in a single, pre-determined, taxonomic order.
- Facets correspond to properties of the information elements. They are often derived by analysis of the text of an item using entity extraction techniques or from pre-existing fields in a database such as author, descriptor, language, and format. Thus, existing web-pages, product descriptions or online collections of articles can be augmented with navigational facets.
- Within the academic community, faceted search has attracted interest primarily among library and information science researchers, and to some extent among computer science researchers specializing in information retrieval

A PART OF MY INCOMPLETE FACEBOOK OF FACETED CLASSIFICATION/BROWSING/SEARCH



S. R. Ranganathan
(1892-1972
Faceted
Classification



Giovanni Maria
Sacco
Dynamic Taxnomies



Sébastien Ferré
University Rennes 1



Marti A. Hearst
University of
California, Berkeley

- Let's now see some examples from some widely used systems

EXAMPLE: EBAY

Related : olympus om lens olympus camera olympus digital voice recorder olympus dm olympus e500 olympus camera charger ... ☐ Include description

Categories

- Cameras & Photography** (113,553)
 - Camera & Photo Accessories (73,325)
 - Lenses & Filters (18,488)
 - Digital Cameras (5,129)
 - Film Photography (3,164)
 - Flashes & Accessories (3,069)
 - More ▾
- Business, Office & Industrial** (6,187)
 - Medical/ Lab Equipment (4,441)
 - Office Equipment & Supplies (1,227)
 - Electrical & Test Equipment (289)
 - Industrial Supply/ MRO (45)
 - Other Business & Industrial (10)
 - More ▾

[See all categories](#)

Condition [see all](#)




- ☐ New (114,023)
- ☐ Used (21,216)
- ☐ Not specified (2,341)

All listings Auction Buy it now

Sort: Best Match View:

137,645 results for olympus [★ Save search](#)

Worldwide

	Olympus Trip AF MD 35mm Point and Shoot Film Camera	59s left Today 12:20	£1.25 2 bids Postage not specified
	Olympus SZ-14 Black Digital Camera Top-rated seller		£79.99 Buy it Now Postage not specified
	Olympus PEN E-PM1 Black + 14-42mm Lens + FL-LM1 Flash Top-rated seller		£159.99 Buy it Now + £7.99 postage

EXAMPLE OF FDT: BOOKING.COM

Your Search
Sými
1 Night (Sept 12 - Sept 13)
2 adults
[Change Search](#)

[Show map](#)

Map data ©2013 Google

Filter by:

Price (per night)
☒ € 0 - € 49 active
☐ € 50 - € 99
☐ € 100 - € 149
☐ € 150 - € 199
☐ € 200 +
Star Rating
☐ 2 stars (1)
☐ 3 stars (3)
☐ 4 stars (2)
☒ Unrated active
Hotel Type
☐ Apartments (5)
☐ Vacation Homes (2)
☐ Villas (1)
☐ Guesthouses (3)
Review Score
☐ Wonderful: 9+ (2)
☐ Very good: 8+ (5)
☐ Good: 7+ (7)
☐ Pleasant: 6+ (9)
☐ No rating (2)

7 out of 25 properties are available in and around Sými
Showing 1 - 7

Sort by: Recommended Stars Location Price Review Score
[List](#) [Grid](#)

Kirilos Studios
Sými • [Show map](#)
Reservation possible without a credit card
Latest booking: 10 hours ago
Studio - Split Level
Only 2 left
€ 40
[Book now](#)

Wonderful 9.1
Score from 11 reviews

Sými Garden Studios
Sými • [Show map](#)
Latest Booking: September 8
Studio
Only 2 left
€ 50
Studio (4 Adults)
Only 4 left
€ 70
Studio with Sea View and Harbour View Last one!
Last chance!
Only 1 left
€ 70
[Book now](#)

Good 7.3
Score from 27 reviews

Villa Pinots!
Sými • [Show map](#)
1 person is looking at this guesthouse.
Reservation possible without a credit card
Latest booking: 1 hour ago
Double or Twin Room
Only 2 left
€ 55
Double or Twin Room Breakfast included
Only 2 left
€ 60
[Book now](#)

Excellent 8.9
Score from 11 reviews

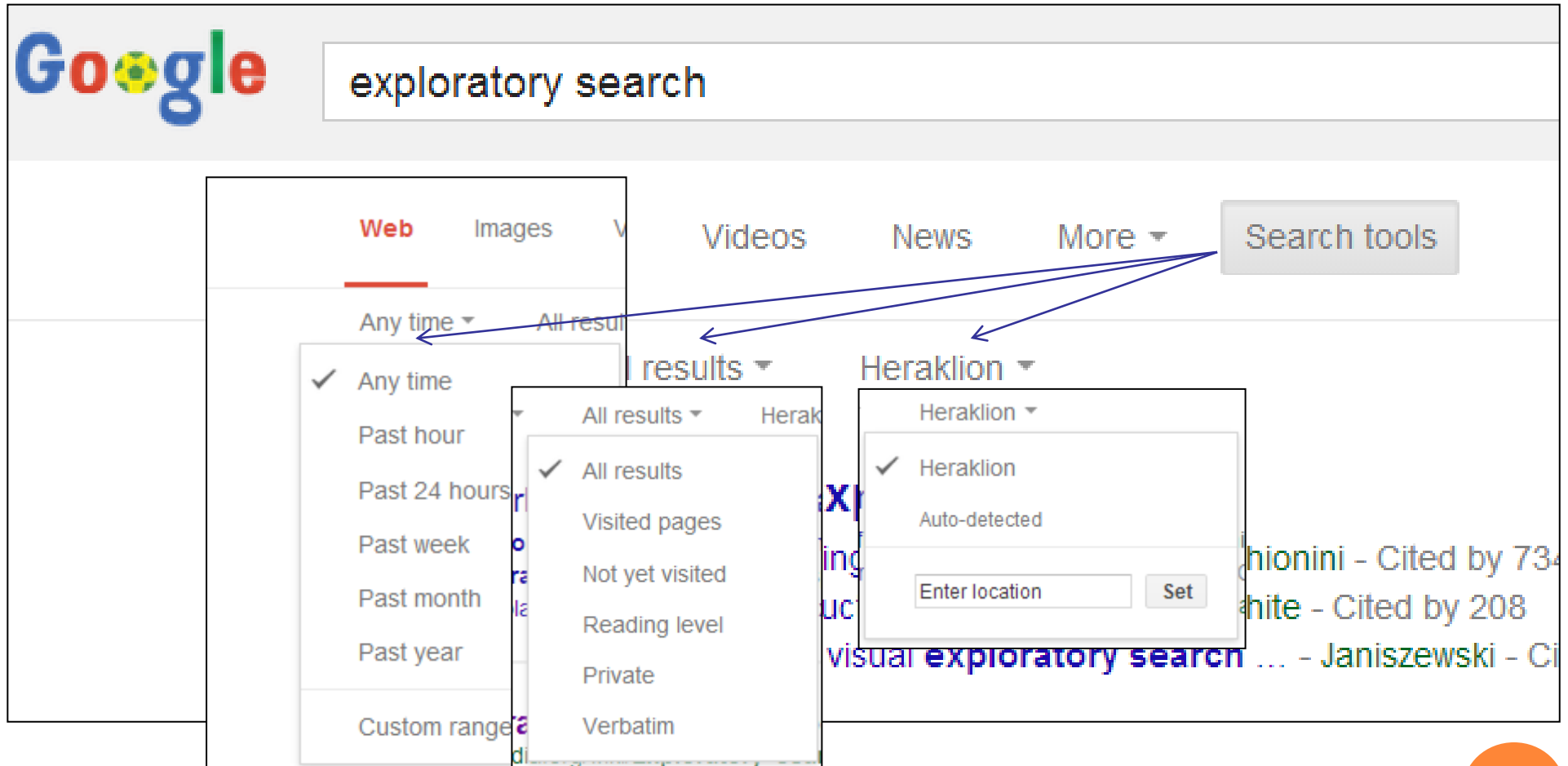
Kokona
Sými • [Show map](#)
Reservation possible without a credit card
Latest Booking: yesterday
Double or Twin Room
Only 4 left
€ 58
Double or Twin Room Breakfast included
Only 4 left
€ 68
[Book now](#)

Very good 8.5
Score from 155 reviews



Grace Hotel & Studios
Sými • [Show map](#)
Latest booking: 22 hours ago
Studio with Sea View
Only 2 left
€ 60
[Book now](#)

EXAMPLE: GOOGLE SEARCH

(LIMITED FUNCTIONALITY: NO COUNT INFORMATION)



EXAMPLE: SCHOLAR GOOGLE



Scholar About 1,430,000 results (0.06 sec)

Articles

Case law

My library

Any time

Since 2014

Since 2013

Since 2010

Custom range...

Sort by relevance

Sort by date

☒ include patents

☒ include citations

[Exploratory search: from finding to understanding](#)
[G Marchionini](#) - Communications of the ACM, 2006 - dl.acm.org
From the earliest days of computers, **search** has been a fundamental application that has driven research and development. For example, a paper published in the inaugural year of the IBM Journal 36 years ago outlined challenges of text retrieval that continue to the present [4]. ...
Cited by 734 Related articles All 18 versions Web of Science: 150 Import into BibTeX Save More

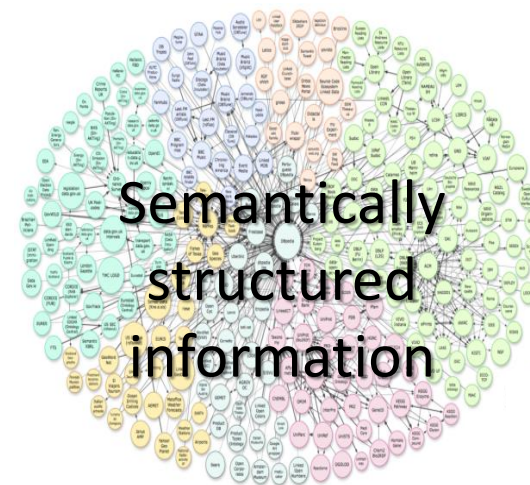
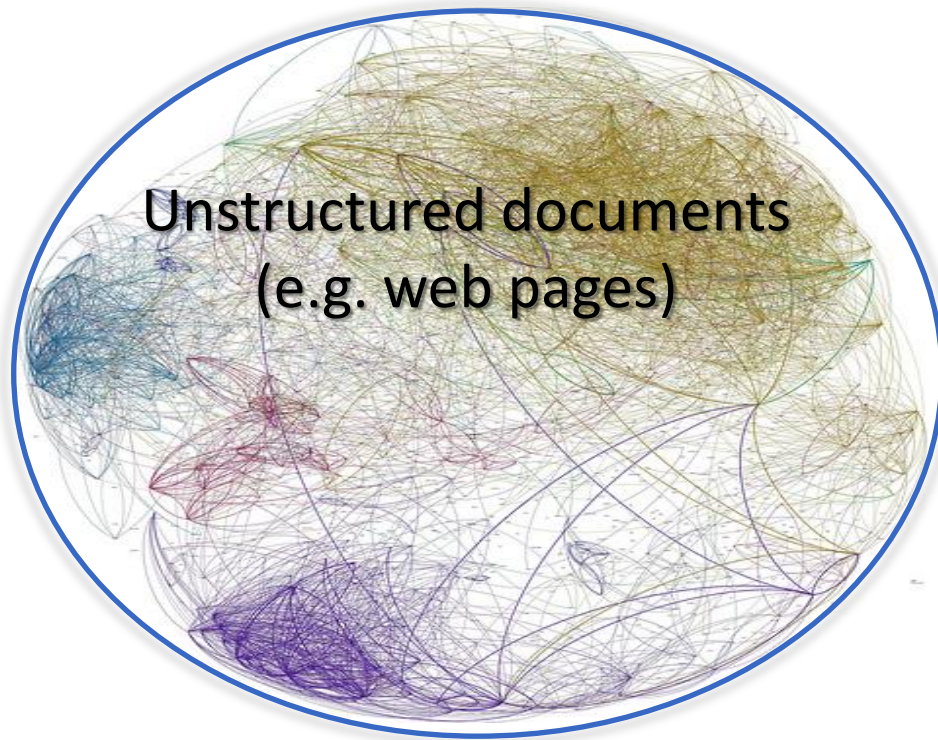
[\[HTML\] Supporting exploratory search. introduction. special issue. communications of the ACM](#)
[RW White](#), [B Kules](#), [SM Drucker](#) - Communications of the ACM, 2006 - eprints.soton.ac.uk
Online **search** has become an increasingly important part of the everyday lives of most computer users. **Search** engines, bibliographic databases, and digital libraries provide adequate support for users whose information needs are well defined. However, there are ...
Cited by 208 Related articles All 7 versions Web of Science: 10 Import into BibTeX Save More

[\[PDF\] The influence of display characteristics on visual exploratory search behavior](#)
[C Janiszewski](#) - Journal of Consumer Research, 1998 - JSTOR
Visual information **search** is a combination of two distinct types of behavior. Goal-directed **search** behavior occurs when consumers use stored **search** routines to collect information in a deliberate manner. In contrast, **exploratory search** behavior occurs when consumers are ...
Cited by 298 Related articles All 9 versions Web of Science: 99 Import into BibTeX Save More

MORE ON FE

- We will return to define more formally the interaction of faceted exploration, after first making a short introduction of the **Semantic Web** and **Linked Data**.
 - After that we will show how FE can be applied on such data

RELATIVELY RECENTLY A ... SECOND “WORLD” STARTS GROWING





SEMANTIC WEB AND LINKED DATA

23

THE SEMANTIC WEB VISION

- The **Semantic Web** is an evolving extension of the WWW where the **content** can be expressed **not only in natural language** but also in **formal languages** (e.g. RDF/S, OWL) that can be read and used by software agents, permitting them to **find**, **share** and **integrate** information more easily
- *Imagine that the objective is the collaborative creation and evolution of a **world wide distributed graph** (about everything 😊)*
 - *We could say that this graph resembles the structure of an EntityRelationship Diagram (recall Databases)*

SEMANTIC WEB TECHNOLOGIES

- For achieving the Semantic Web vision, the recent years several technologies have been emerged, many of them are international (W3C) standards
- These technologies include:
 - **knowledge representation** languages (e.g. RDF/S, OWL) and formats for exchanging knowledge
 - **query** languages (π.χ. SPARQL),
 - **rule languages** and **inference engines**
 - Techniques for constructing **mappings** for integrating/harmonizing schemas and data
 - Technologies for **mining** structure knowledge from texts
 - Various APIs.
- **Linked Open Data** Are based on these technologies.



The Semantic Web Technology Stack (not a piece of cake...)

Most apps use only a subset of the stack

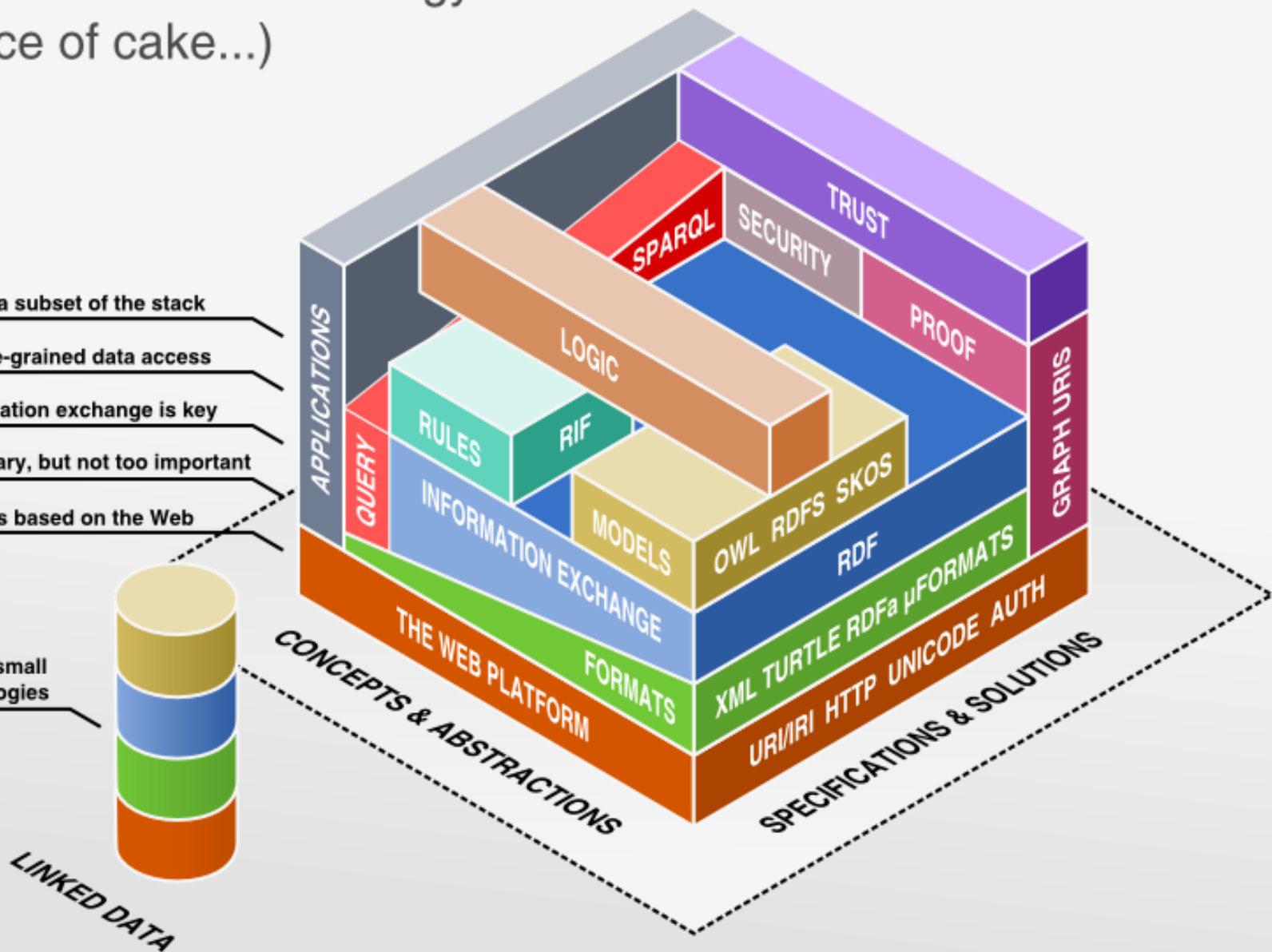
Querying allows fine-grained data access

Standardized information exchange is key

Formats are necessary, but not too important

The Semantic Web is based on the Web

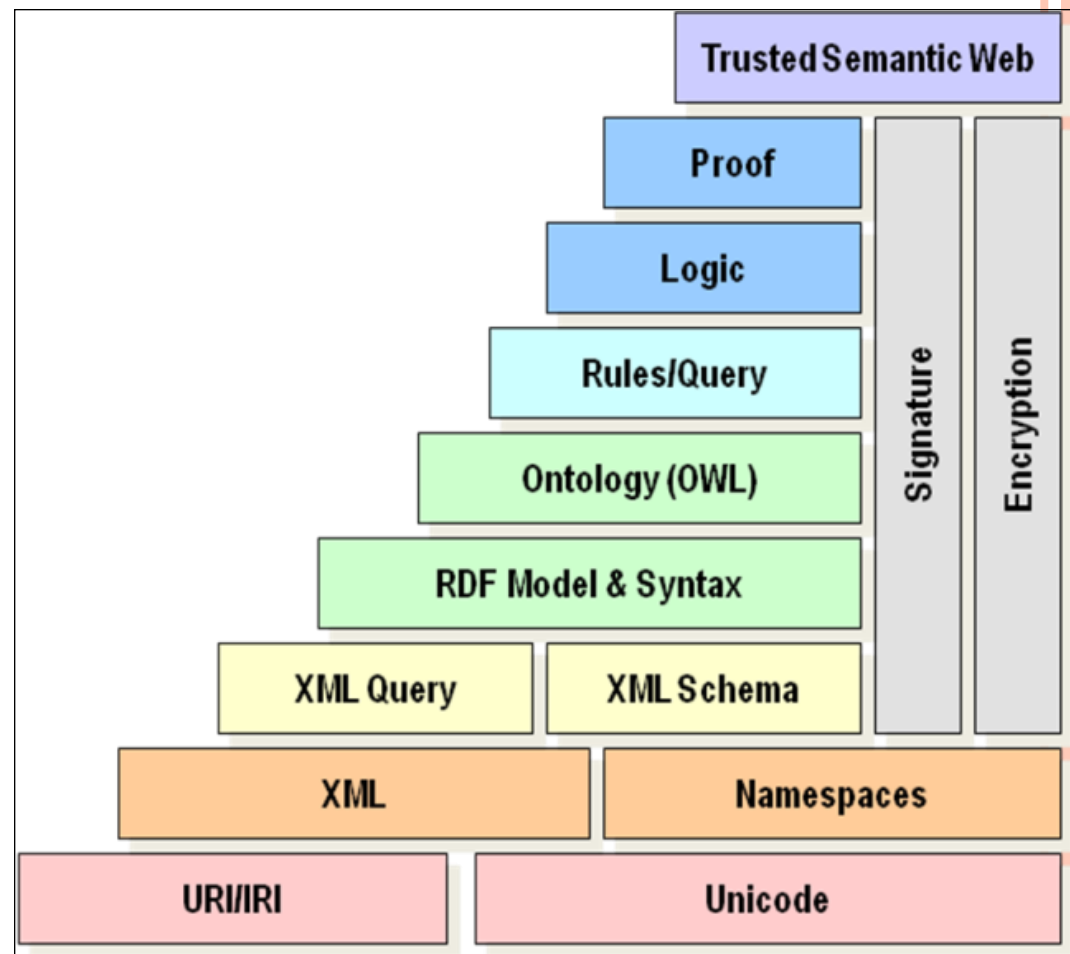
Linked Data uses a small
selection of technologies



THE **TECHNOLOGY STACK** OF SEMANTIC WEB

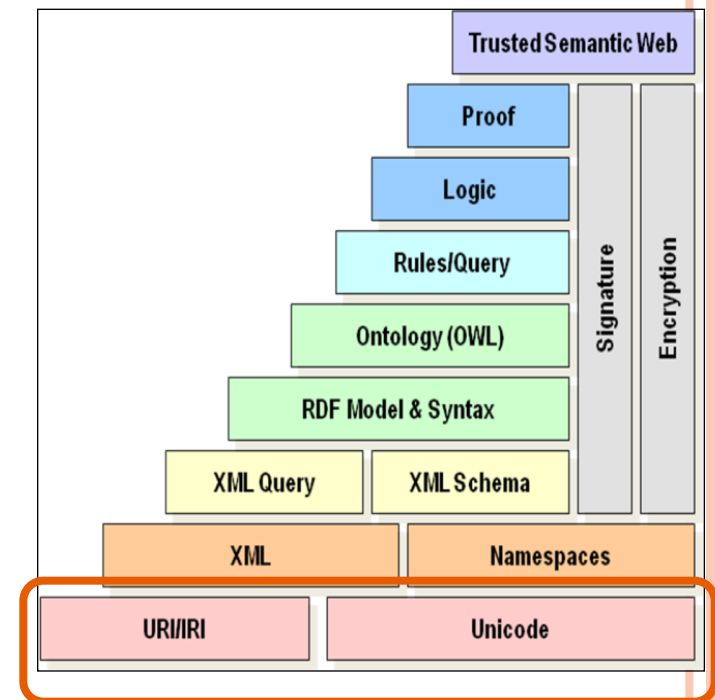
By starting from the World Wide Web where the content is represented as **hypertext** (i.e. pages containing texts and links to other pages), let's now give an overview of the basic technologies on which the Semantic Web is based.

The current technology stack of **Semantic Web** can be illustrated as follows:



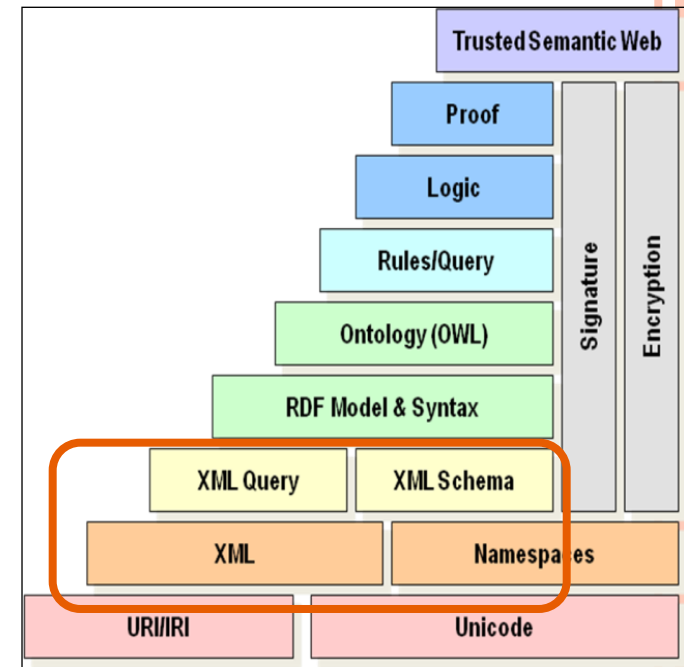
TECHNOLOGICAL STACK: URIs & UNICODE

- At the lowest level we have the **URI** (Uniform Resource Identifiers) for identifying and locating resources, and **UNICODE** for representing characters from various natural languages



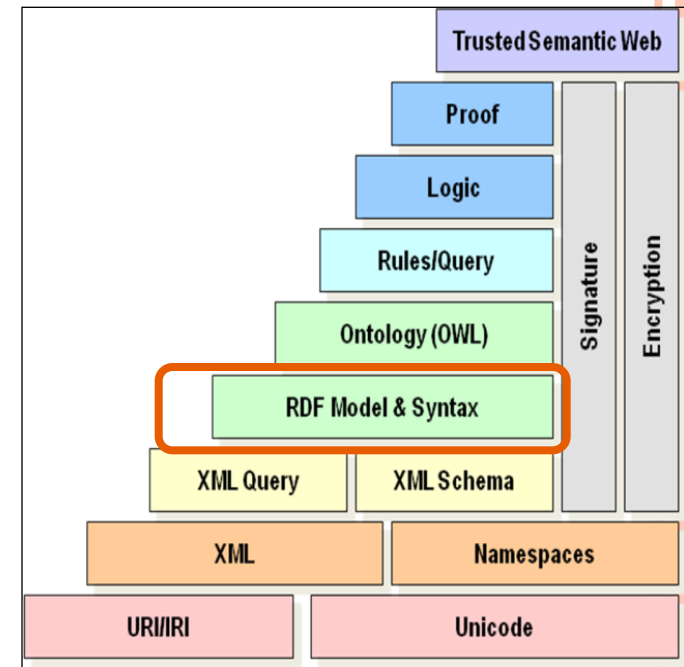
TECHNOLOGICAL STACK: XML

- XML provides syntax for having structured documents, but does not impose any constraint regarding the semantics of these documents.
- XML Schema offers a method to restrict the structure that XML documents can have.
- For querying XML documents, we have XPath and XQuery.



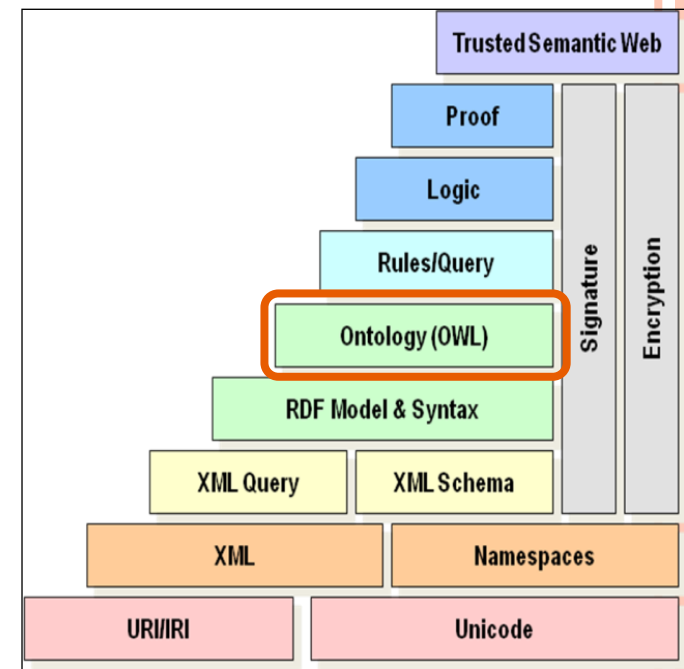
TECHNOLOGICAL STACK: **RDF** & **RDFS**

- **RDF (Resource Description Framework)** is a structurally object-oriented model for representing objects (resources) and associations between them.
- It allows expressing content in form of **triples** (*subject, predicate, object*), and sets of triples actually form a semantic network/graph.
- These triples can be expressed in various formats (**TriG, N3 RDF/XML**), some of them are based on XML (RDF/XML)
- **RDFS** allows defining the vocabulary to be used in RDF and semantic relationships between the elements of this vocabulary



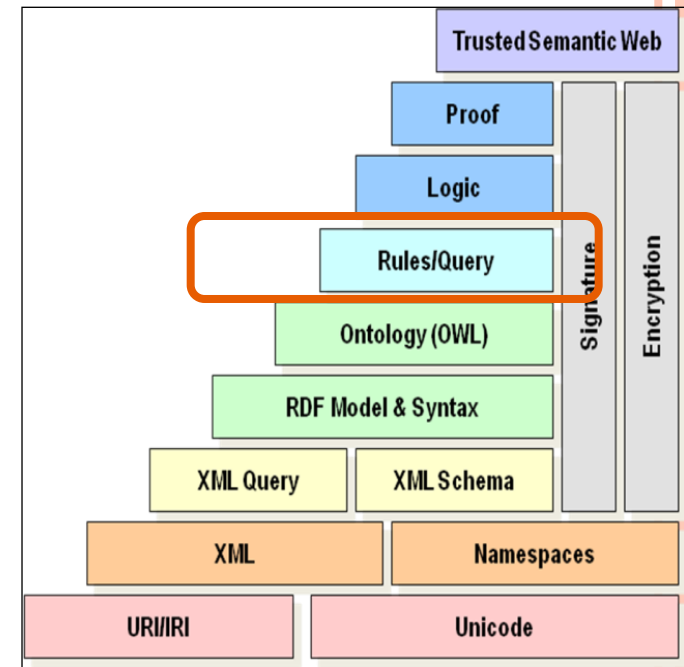
TECHNOLOGICAL STACK: OWL

- The Web Ontology Language (**OWL**) is a family of knowledge representation languages or ontology languages for authoring ontologies or knowledge bases. The languages are characterized by formal semantics and RDF/XML-based serializations for the Semantic Web.



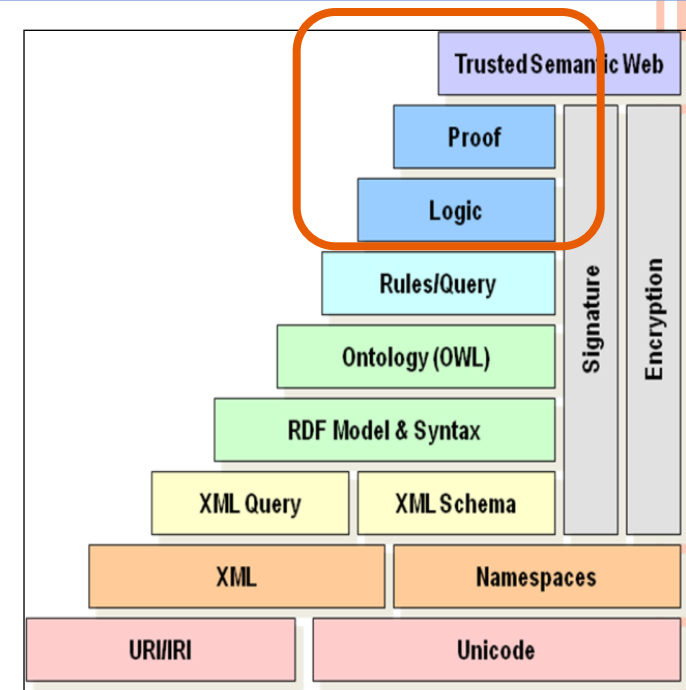
TECHNOLOGICAL STACK: SPARQL & SWRL

- For exploiting the structured content that has been represented using RDF/S, there are **query languages** and **rule languages**
- Specifically, **SPARQL** (SPARQL Protocol and RDF Query Language) is a query language for knowledge expressed in RDF/OWL.
- **SWRL** (Semantic Web Rule Language) allows expressing inference rules (essentially Horn rules).



TECHNOLOGICAL STACK: LOGIC, PROOF, TRUST

- The layers *Logic* and *Proof* concern the enrichment of the expressiveness of the representation languages
- Finally the *Trust* layer concerns trust issues, e.g. digital signatures (for proving that one particular person has written or agrees with a particular document or sentence, as well as trust networks allowing users to define who they trust (and so on), eventually yielding trust networks (*Web of Trust*)



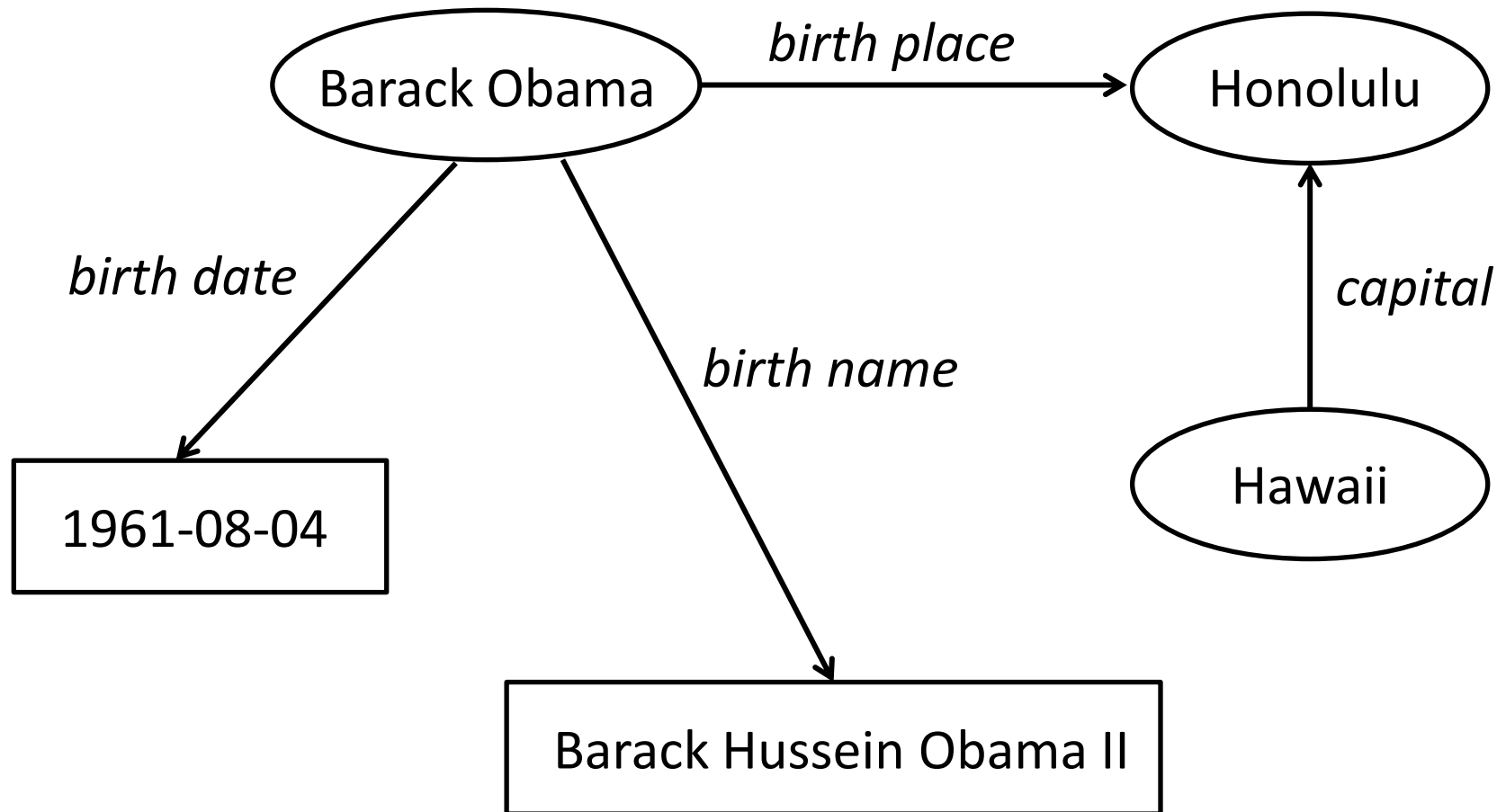
STANDARDIZATION, PERSISTENT STORAGE

- The part of the stack that has been implemented and **standardized** allows **describing resources** (of digital objects mainly) using **classes** and **properties** that have been defined in **ontologies** (expressed using **RDF/S** or **OWL**) which are accessible from the network which in turn can be connected (one ontology can extend classes and properties of other ontologies).
- For the **persistent storage** of such contents (independently of whether it corresponds to data or metadata) in the form of connected semantic network (semantic graph), there are tools, usually referred to as **triplestores**. We could call them **semantic data bases** since they receive as input RDF/S data, support integrity constraints, enable querying (SPARQL support), and concurrent access.

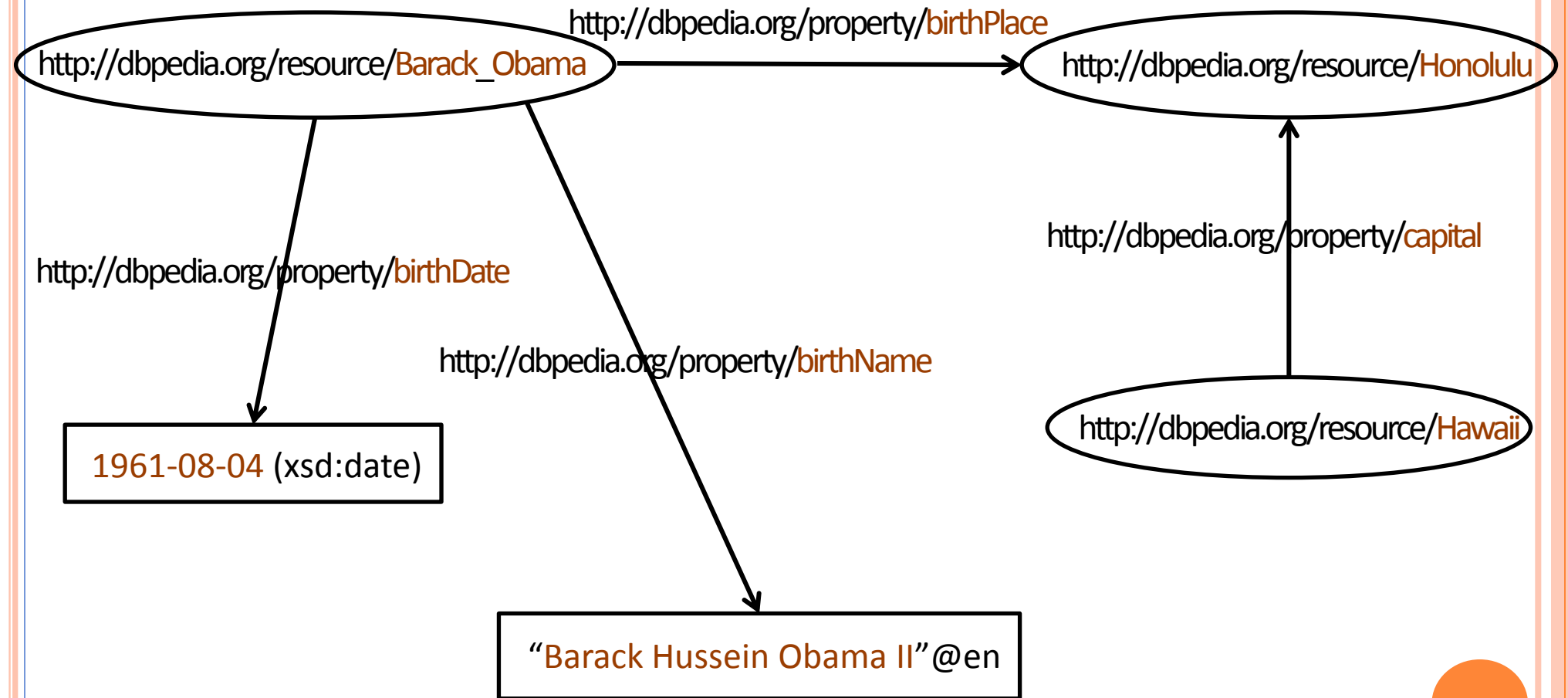


EXAMPLE OF AN RDF GRAPH

Directed Labeled (Multi)Graph



RDF GRAPH



RDF GRAPH REPRESENTATION

In XML (RDF/XML):

```
<?xml version="1.0" encoding="utf-8" ?>
<rdf:RDF
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:dbpprop="http://dbpedia.org/property/"
  <rdf:Description rdf:about="http://dbpedia.org/resource/Barack_Obama">
    <dbpprop:birthDate rdf:datatype="http://www.w3.org/2001/XMLSchema#date">
      1961-08-04
    </dbpprop:birthDate>
    <dbpprop:birthPlace rdf:resource="http://dbpedia.org/resource/Honolulu" />
    <dbpprop:birthName xml:lang="en">
      Barack Hussein Obama II
    </dbpprop:birthName>
  </rdf:Description>
  <rdf:Description rdf:about="http://dbpedia.org/resource/Hawaii">
    <dbpprop:capital rdf:resource="http://dbpedia.org/resource/Honolulu" />
  </rdf:Description>
</rdf:RDF>
```

In N-Triples:

```
<http://dbpedia.org/resource/Barack_Obama> <http://dbpedia.org/property/birthPlace> <http://dbpedia.org/resource/Honolulu> .
<http://dbpedia.org/resource/Barack_Obama> <http://dbpedia.org/property/birthDate> "1961-08-04"^^<http://www.w3.org/2001/XMLSchema#date> .
<http://dbpedia.org/resource/Barack_Obama> <http://dbpedia.org/property/birthName> "Barack Hussein Obama II"@en .
<http://dbpedia.org/resource/Hawaii> <http://dbpedia.org/property/capital> <http://dbpedia.org/resource/Honolulu> .
```




As of September 2011

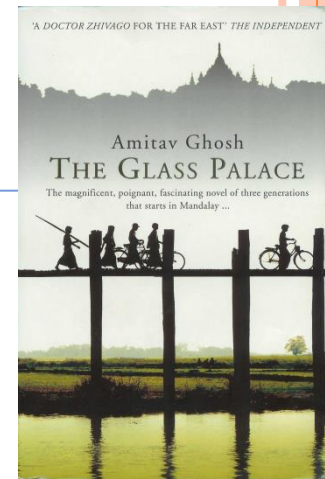
EXAMPLE: LINKED DATA SCENARIO



This example will show how data from two *bookstores* can be **exposed** and published in a structured way, how they can be **integrated**, how they can be **queried**, how they can be **enriched** with other knowledge (e.g. context or domain knowledge) from other sources (e.g. wikipedia).

LINKED DATA SCENARIO

A SIMPLIFIED BOOKSTORE DATA (DATASET “A”)



Suppose the data of the bookstore A are stored in a relational database. The relational tuples related to the book of our running example could be

Books

ISBN	Author	Title	Publisher	Year
0006511409X	id_xyz	The Glass Palace	id_qpr	2000

Authors

ID	Name	Homepage
id_xyz	Ghosh, Amitav	http://www.amitavghosh.com

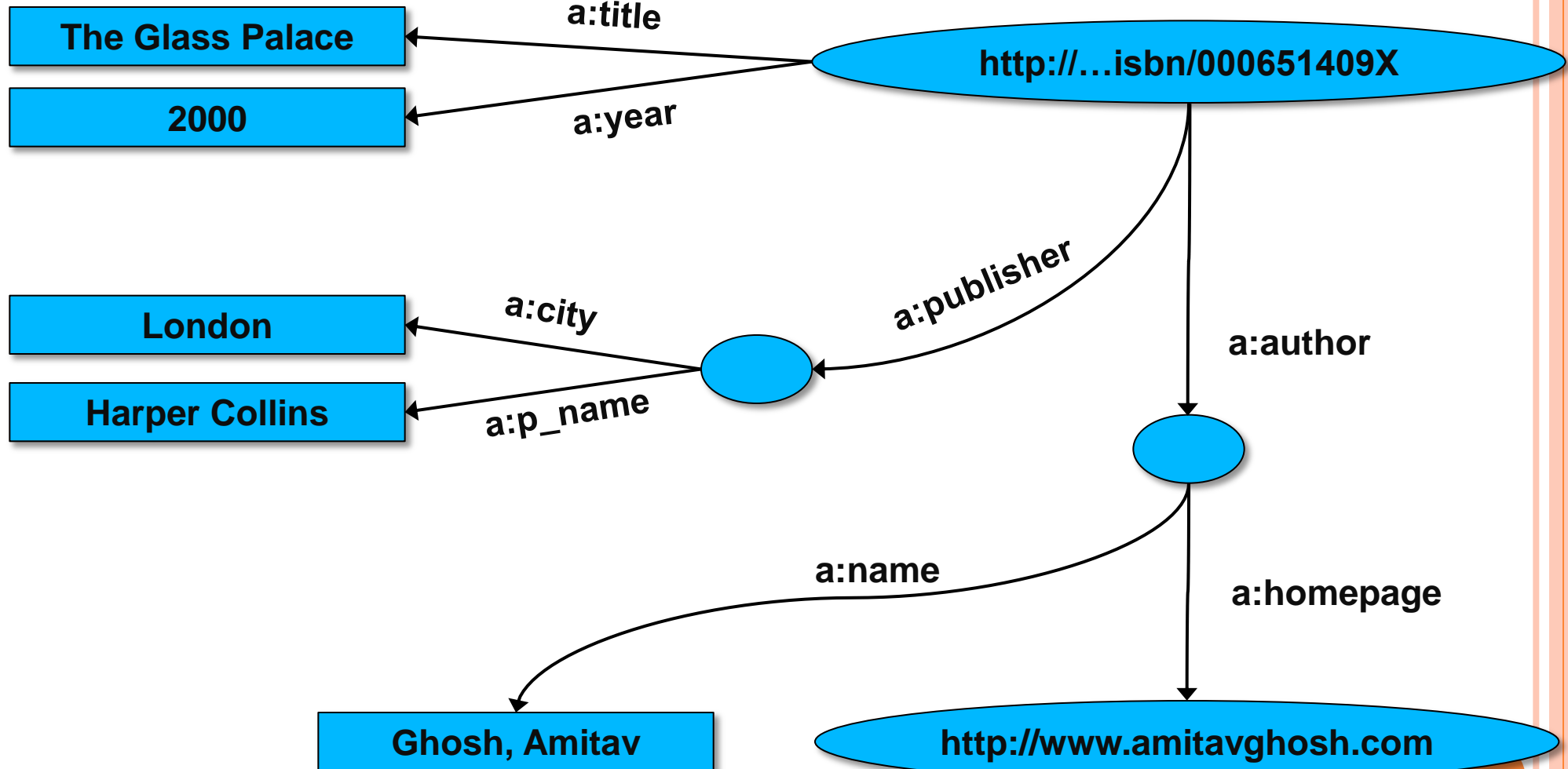
Publishers

ID	Publisher's name	City
id_qpr	Harper Collins	London



LINKED DATA SCENARIO

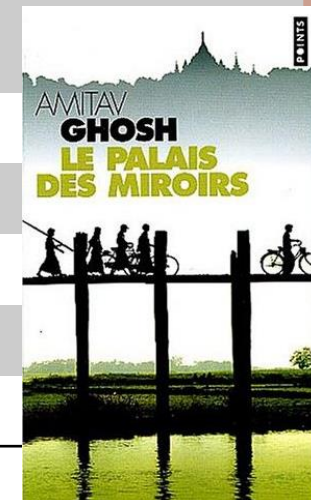
1ST: **EXPORT** YOUR DATA AS A **GRAPH**



LINKED DATA SCENARIO

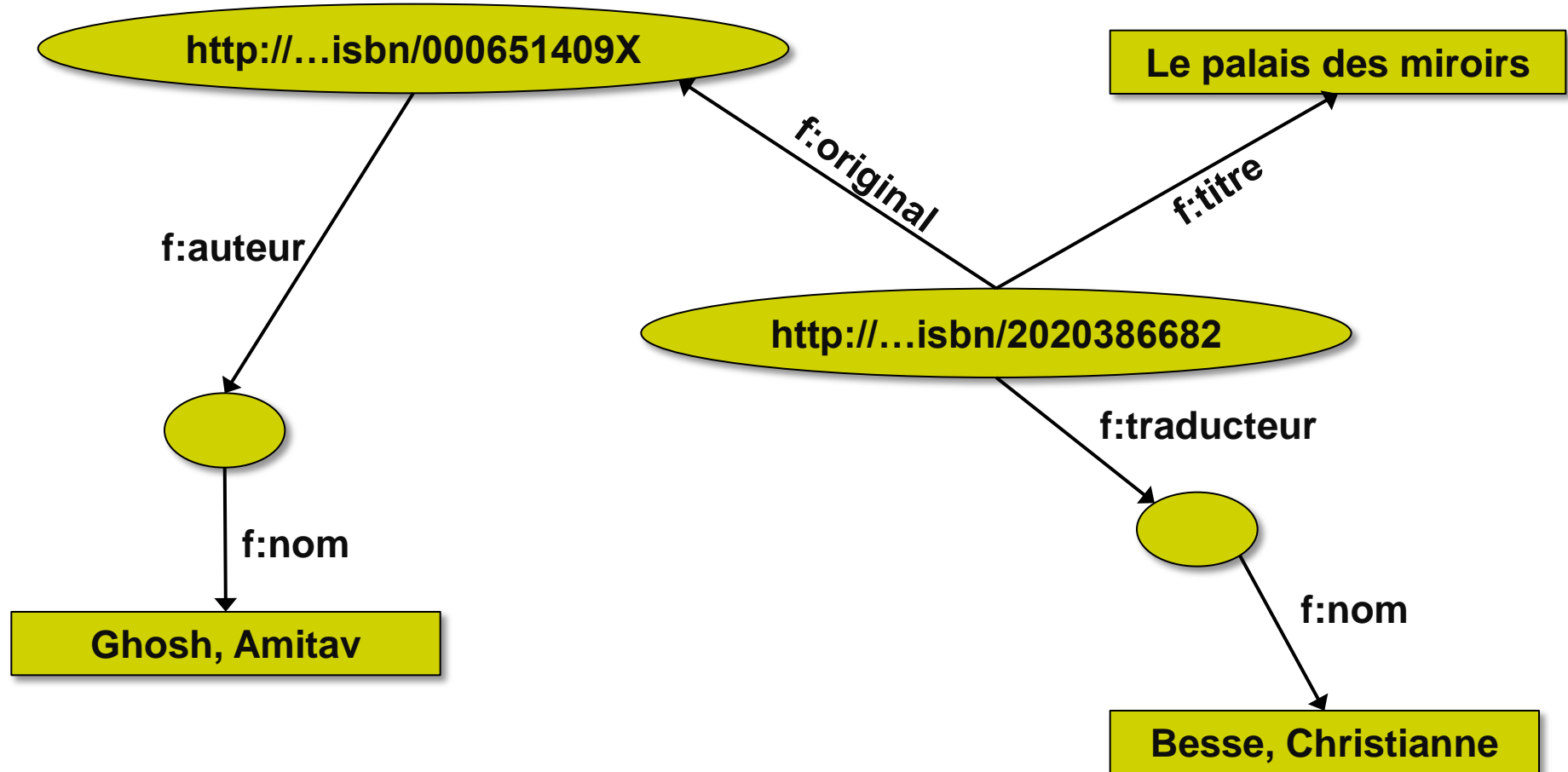
ANOTHER BOOKSTORE DATA (DATASET “F” FROM FRANCE)

A	B	C	D
1	ID	Titre	Traducteur
2	ISBN 2020286682	Le Palais des Miroirs	\$A12\$
3			
4			
5			
6	ID	Auteur	
7	ISBN 0-00-6511409-X	\$A11\$	
8			
9			
10	Nom		
11	Ghosh, Amitav		
12	Besse, Christianne		



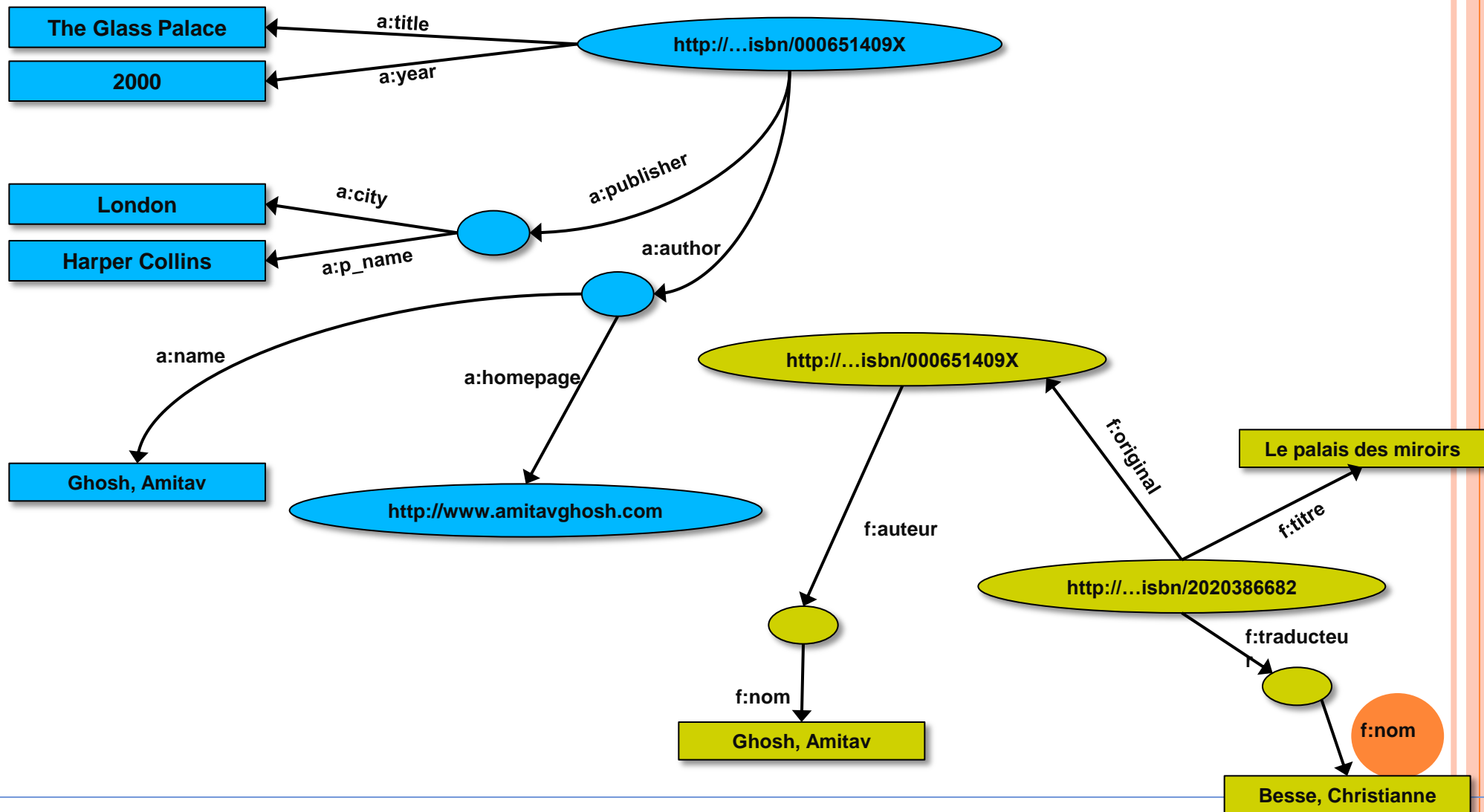
LINKED DATA SCENARIO

2ND: **EXPORT** YOUR SECOND SET OF DATA



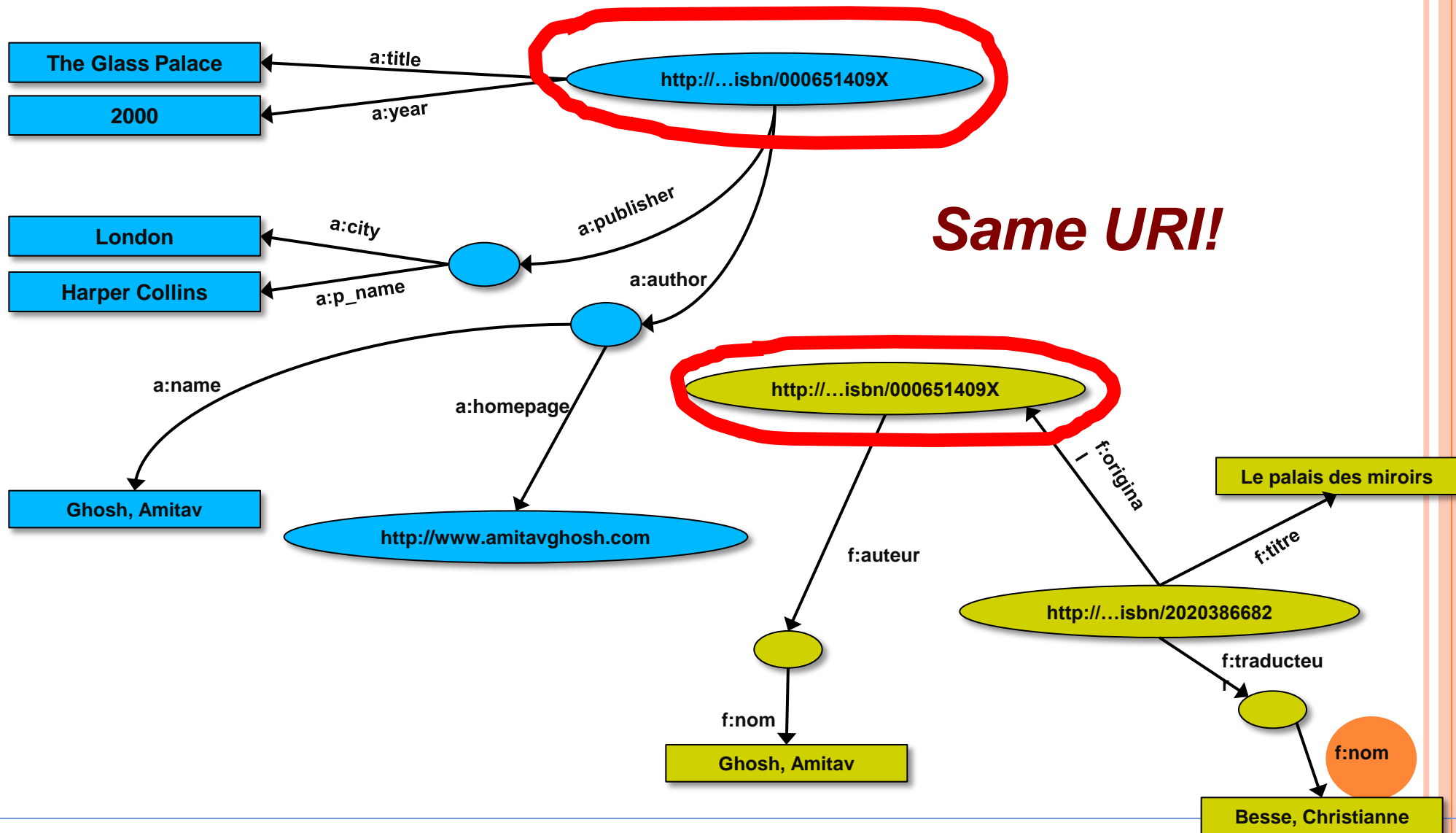
LINKED DATA SCENARIO

3RD: START **MERGING** YOUR DATA



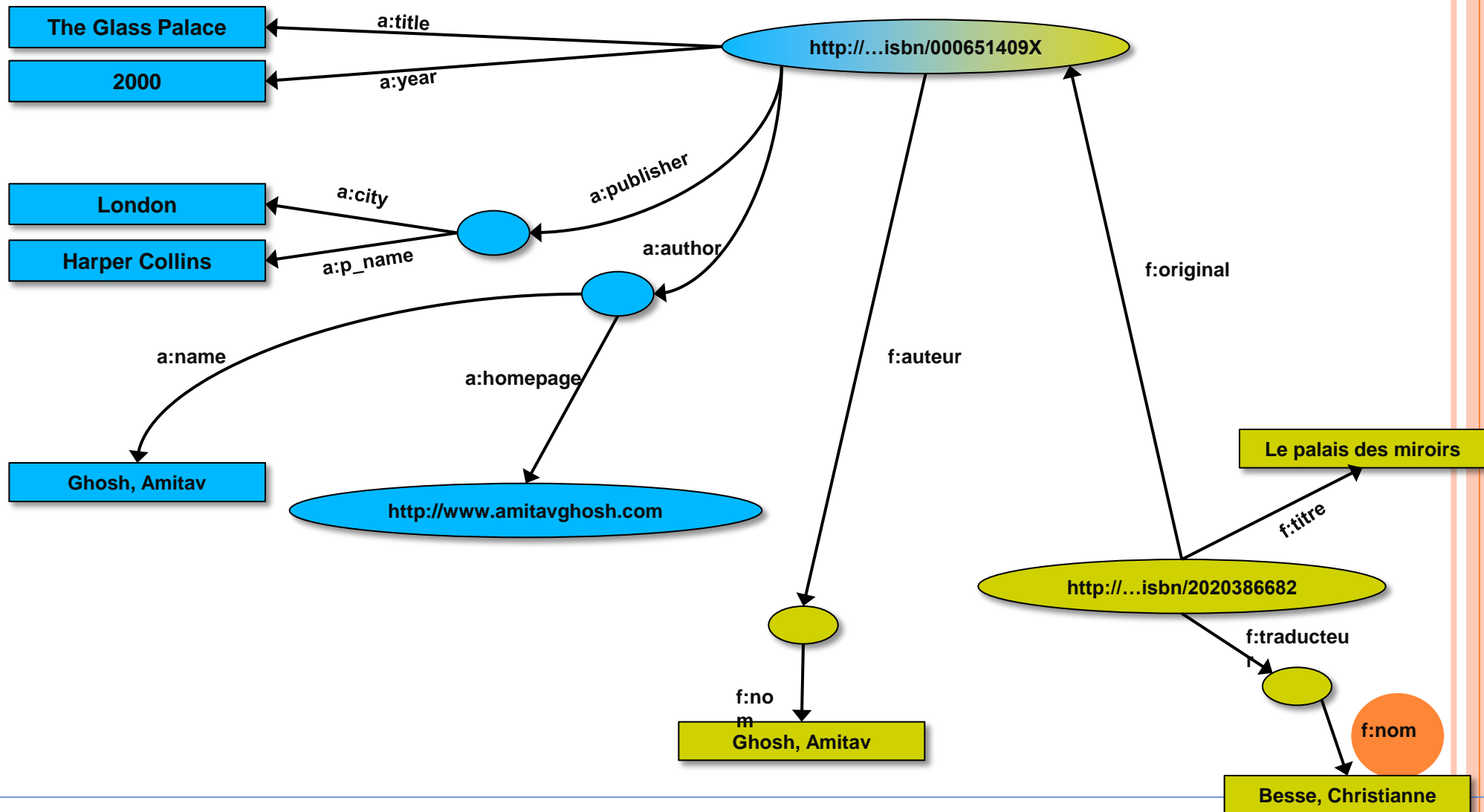
LINKED DATA SCENARIO

3RD: START MERGING YOUR DATA (CONT)



LINKED DATA SCENARIO

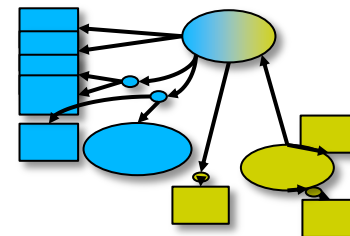
3RD: START MERGING YOUR DATA



LINKED DATA SCENARIO

START MAKING **QUERIES**...

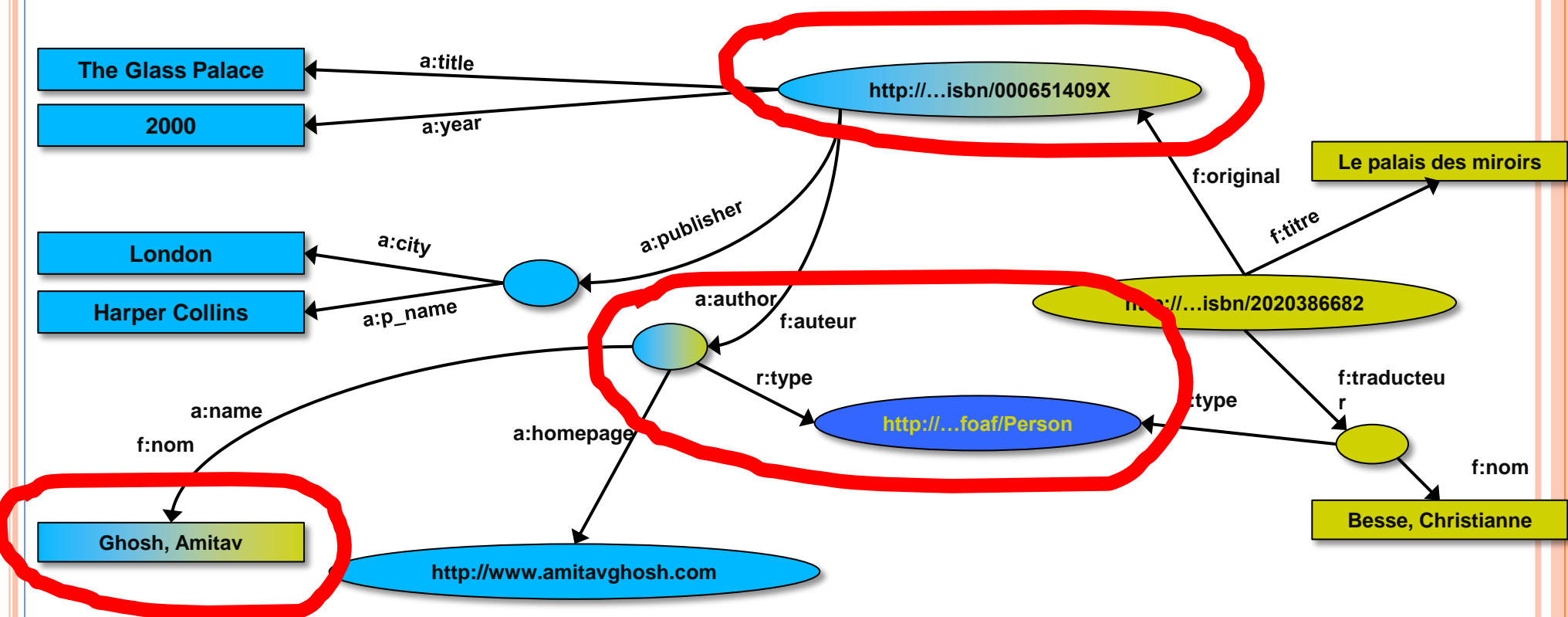
- User of data “F” can now ask queries like:
 - “*give me the title of the original*”
- This information is not in the dataset “F”, but it can be retrieved by merging with dataset “A”



- We “feel” that `a:author` and `f:auteur` should be the **same**. But an automatic merge does not know that! We can add some extra information to the merged data:
 - `a:author same as f:auteur`
 - both identify a “**Person**” a term that a community may have already defined (a “Person” is uniquely identified by his/her name and, say, homepage (it can be used as a “category” for certain type of resources))

LINKED DATA SCENARIO

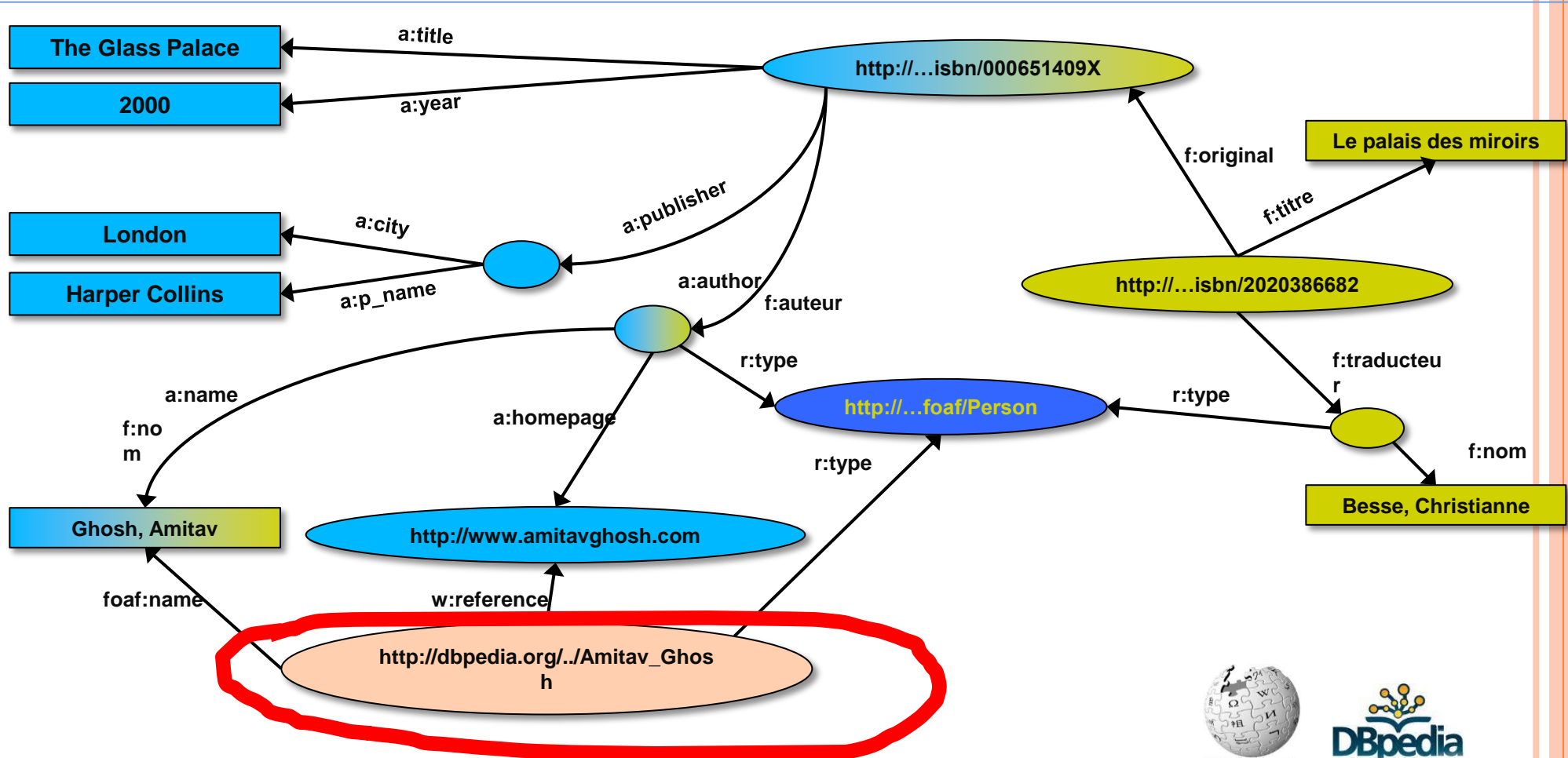
3RD REVISITED: USE THE EXTRA KNOWLEDGE



- User of dataset “F” can now query: “*give me the home page of the original’s ‘auteur’*”. The information is not in datasets “F” or “A”, but was made available by **merging** datasets “A” and datasets “F” and adding three simple extra statements as an extra “glue”

LINKED DATA SCENARIO

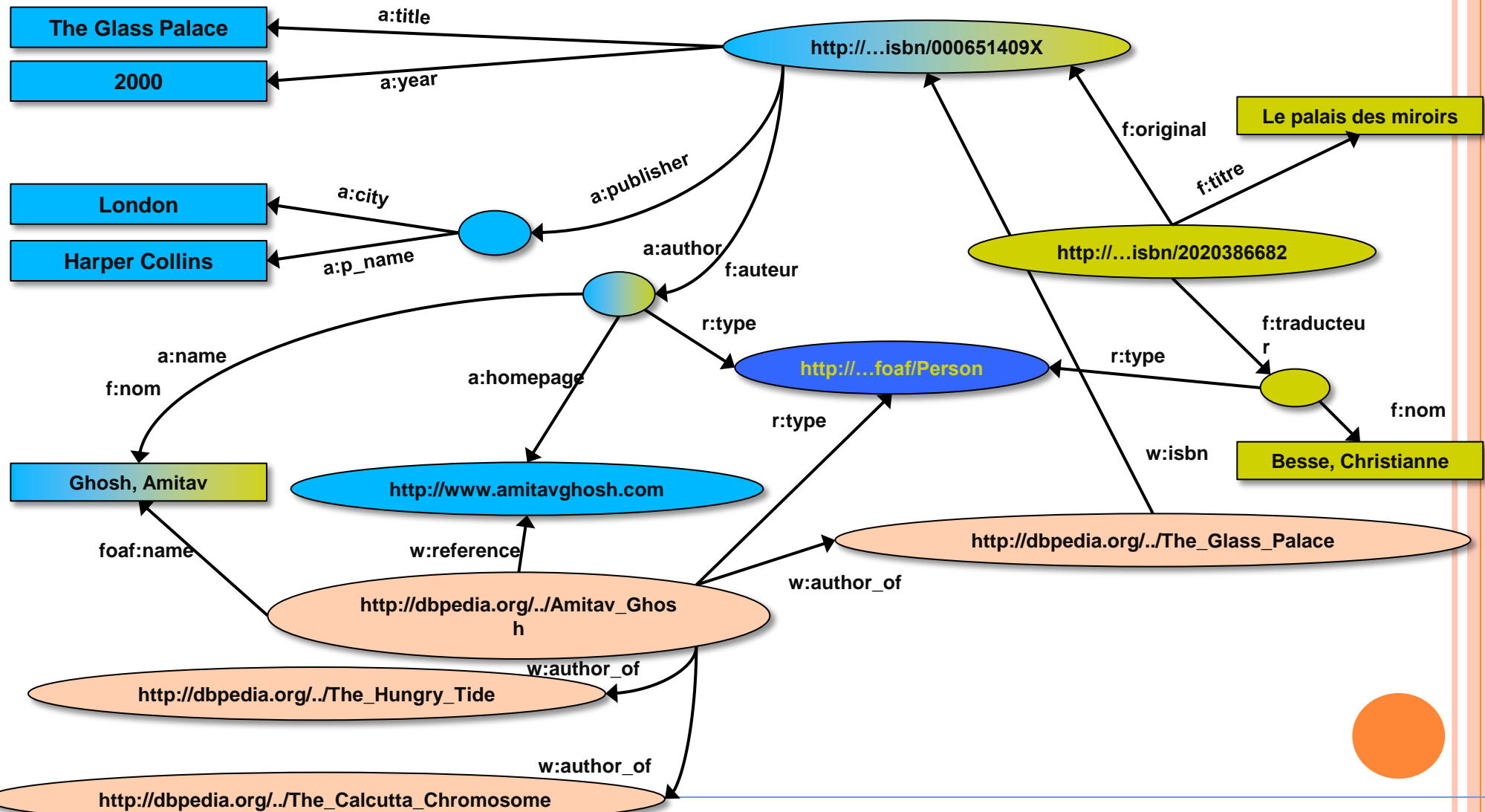
MERGE WITH WIKIPEDIA DATA



Using, e.g., the “**Person**”, the dataset can be combined with other sources. For example, data in Wikipedia can be extracted using dedicated tools

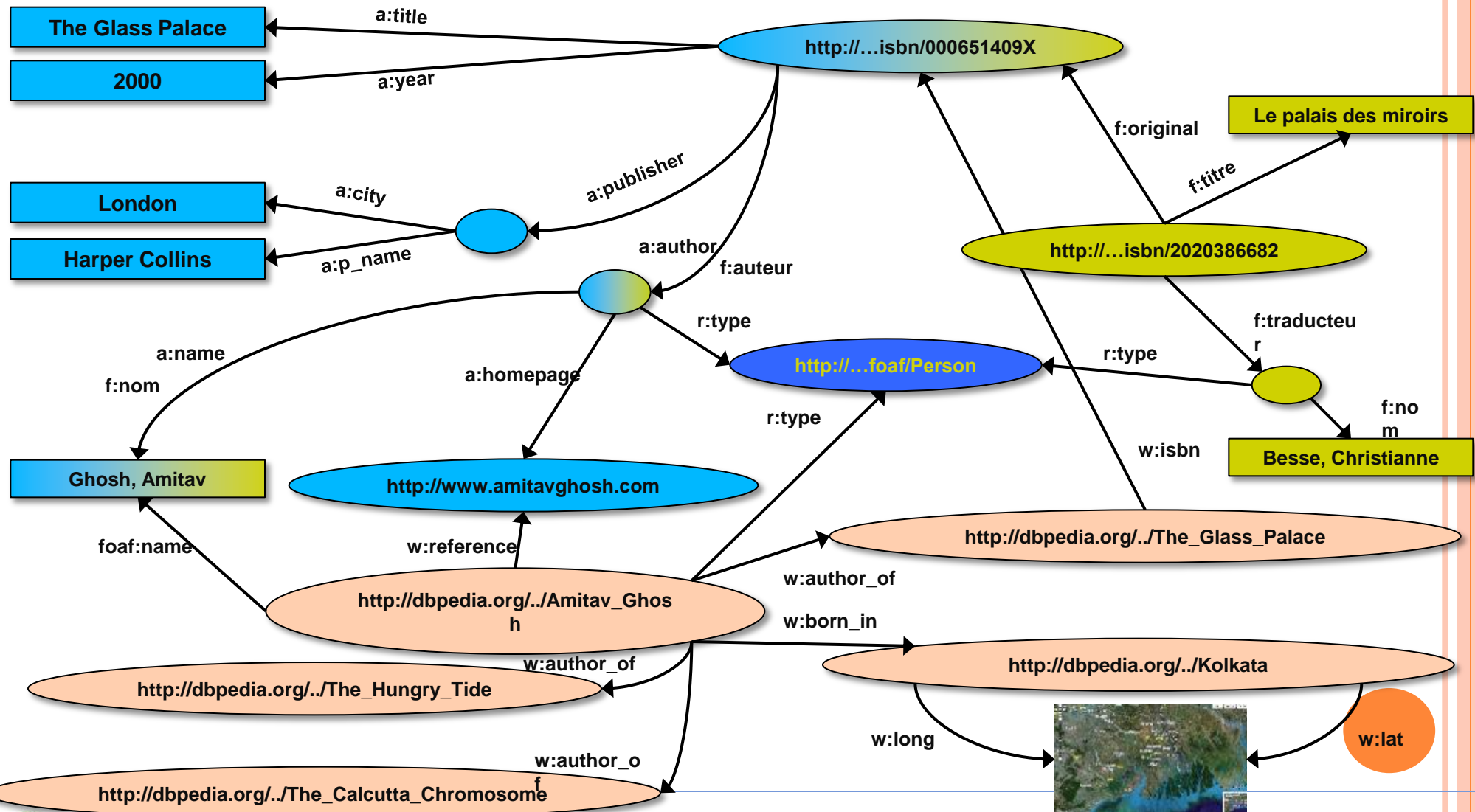
LINKED DATA SCENARIO

MERGE WITH WIKIPEDIA DATA



LINKED DATA SCENARIO

MERGE WITH WIKIPEDIA DATA



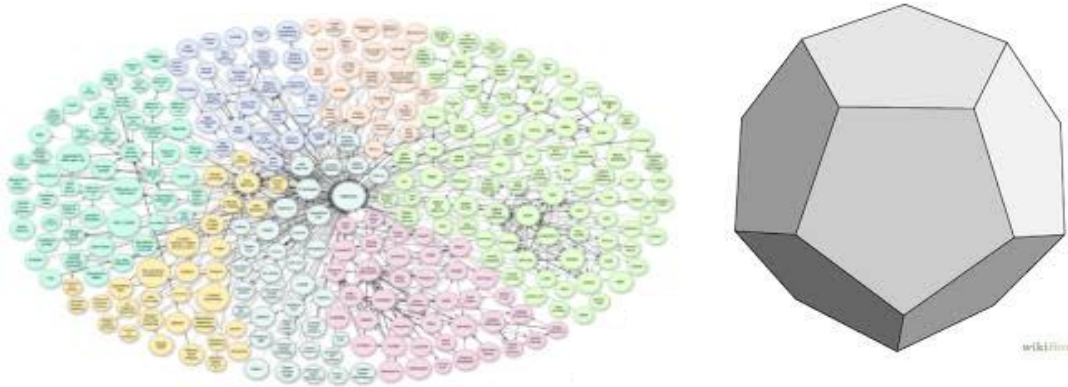
QUERIES AND THE DIFFICULTY OF QUERY FORMULATION

- The availability of Semantic Data allows formulating complex queries. E.g.:
 - *scientists who have worked in an institute in Germany and are known for their work in fuzzy set theory*
- however to formulate such queries one has to know not only SPARQL, but also the vocabulary used. E.g. (sketch):

```
SELECT ?x
WHERE {
    ?x type scientist .
    ?x ns:knownFor "fuzzySetTheory" .
    ?x ns:hasWorkedIn ?y .
    ?y type institution .
    ?y ns:locatedAt "Germany" .
}
```



QUERYING THE SEMANTIC WEB / LOD





- Faceted exploration can allow users to satisfy their information needs without having to be aware of the **employed terminology, contents, or query language** of the sources
- An example from an application offering faceted exploration to DBpedia follows

**German scientists
known for their
work in the
field of evolution**



GUIDED EXPLORATION - AN EXAMPLE



search powered by  neofonie

[About Neofonie](#) [About DBpedia](#) [Imprint](#) [Help](#)

enter search terms...

First | Previous | Next | Last

▼ item type

start typing...

Person (2)

Scientist (2)

▼ nationality

start typing...

Germany (2)

▼ is known for

start typing...

Evolution (2)

Homology (biology) (1)

▼ died in

start typing...

Tübingen (1)

▼ born in year year

start typing...

from... to...

1896 (1)

1826 (1)

▼ born in

start typing...

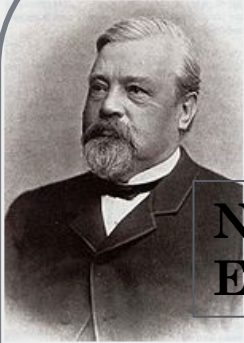
Würzburg (1)

Hanover (1)

Your Filters [Reset Filters](#)

Results 1 to 6 of 485

item type [Scientist](#) nationality [Germany](#) is known for [Evolution](#)



Karl Gegenbaur
Carl Gegenbaur (21 August 1826 - 14 June 1903), "Karl Gegenbaur - Encyclopaedia Britannica" (biography), Encyclopædia Britannica, 2006, Britannica.com webpage: Britannica-KarlG. also Karl Gegenbaur, was a German anatomist and professor who demonstrated that the field of comparative anatomy offers important evidence supporting of the theory of evolution.

New State's Extension

"German scientists known for their work in the field of evolution"

evolution of corals and cephalopods. Schindewolf was on the faculty at the University of Marburg from 1919 until 1927. He then he became director of the Geological Survey of Berlin. In 1948 he became a professor at the University of Tübingen, where he retired as professor emeritus in 1964.

Transition
markers

SURVEYING FACETED EXPLORATION APPROACHES

We use the term *BA* to refer to a *browsing approach*. We can survey the works (*BAs*) that have been applied or proposed according to various aspects:

1. ***Characteristics of the underlying information space.*** The structuring of the underlying information base is an important aspect since each case requires tackling different difficulties.
2. ***Configuration.*** Some approaches can be applied without requiring any form of configuration or application design (regarding the browseable information space), while some others require configuration steps, e.g. specify the contents and structuring of the browsable part through the view-based approach over a DB or an RDF repository. Since the browsable part of the information source is defined by a query, its structure may be different from that of the original source.
3. ***State Space.*** In general we can view the interaction as a state space consisting of *states* and *transitions*, therefore we can characterize, or comparatively evaluate, two *BAs* by comparing their state spaces, e.g. by identifying properties which are satisfied by their state space.

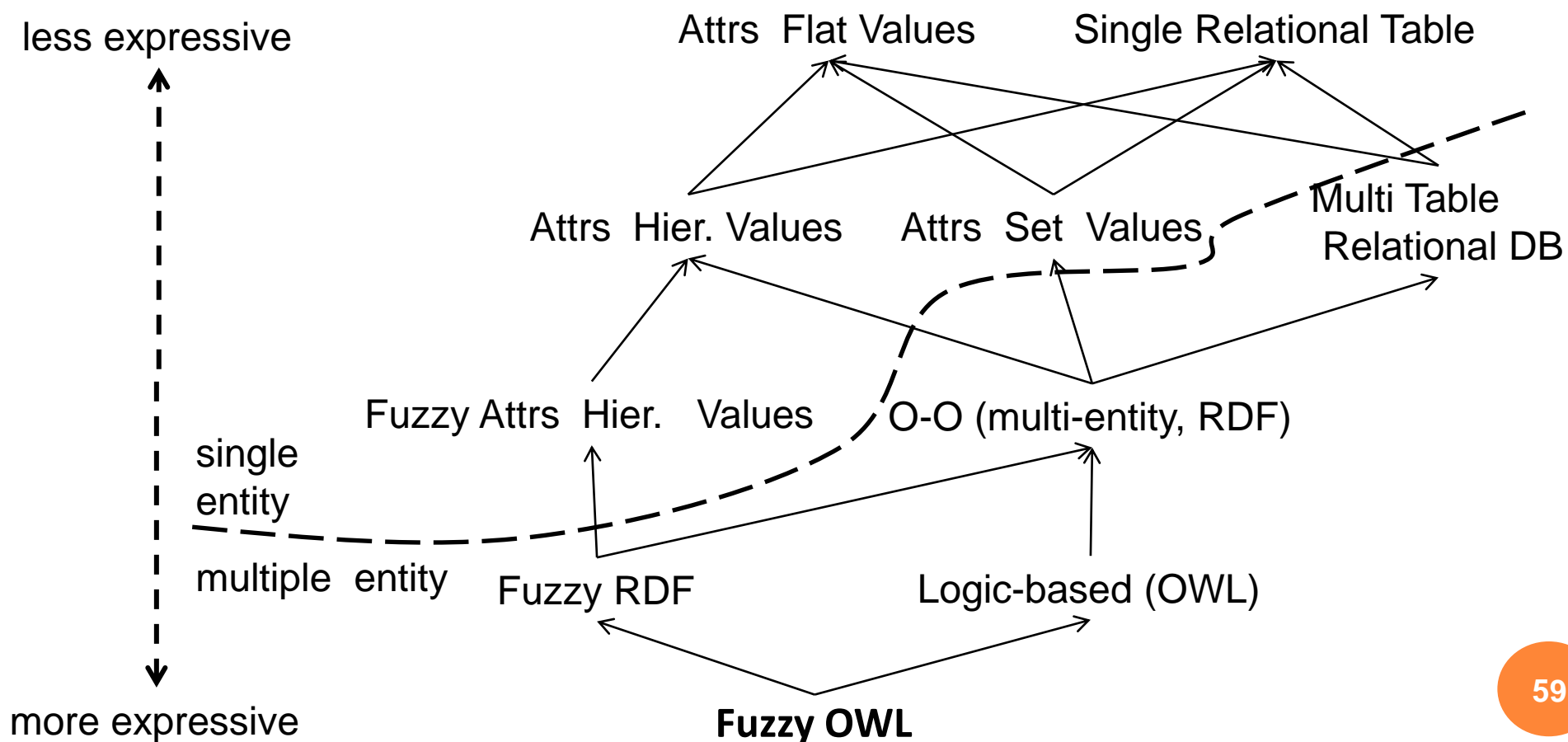
SURVEYING FACETED EXPLORATION APPROACHES

1. *THE UNDERLYING INFORMATION SPACE*

- Some browsing approaches are applicable to simple structures, while others to complex information structures (e.g. OWL-based KBs). There are several options, some of them follow:
 - **attribute-value pairs** with **flat** values (e.g. name=Yannis),
 - **attribute-value pairs** with **hierarchically organized values** (e.g. location=Crete),
 - **set-valued attributes** (either flat or hierarchical) (e.g. accessories={ABS, ESP}),
 - **multi-entity** (or object-oriented) (e.g. RDF, linked open data), and relational databases.
 - Furthermore, we could have **fuzziness** and we can consider this as an independent aspect (e.g. there are fuzzy extensions of the RDF model).
- Therefore one important aspect is how the underlying information is structured.

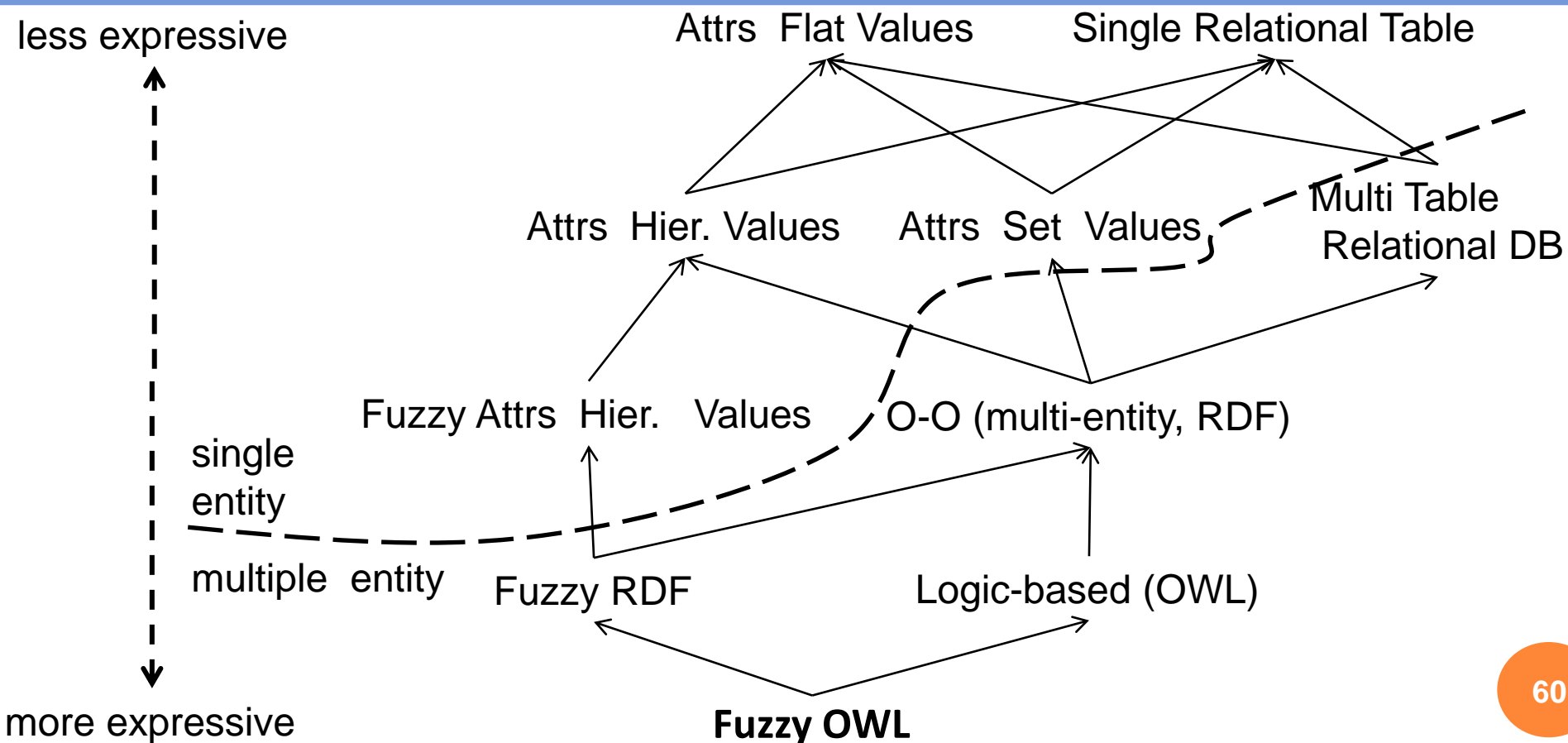
CONT.

Here we can see the above categories **organized hierarchically** where an option X is a (direct or indirect) child of an option Y if whatever information can be expressed in Y can also be expressed in X.



CONT.

The **value of this diagram** is that it depicts the fact that if a browsing approach is appropriate for an option X, then certainly it is appropriate for all options which are parents of X.



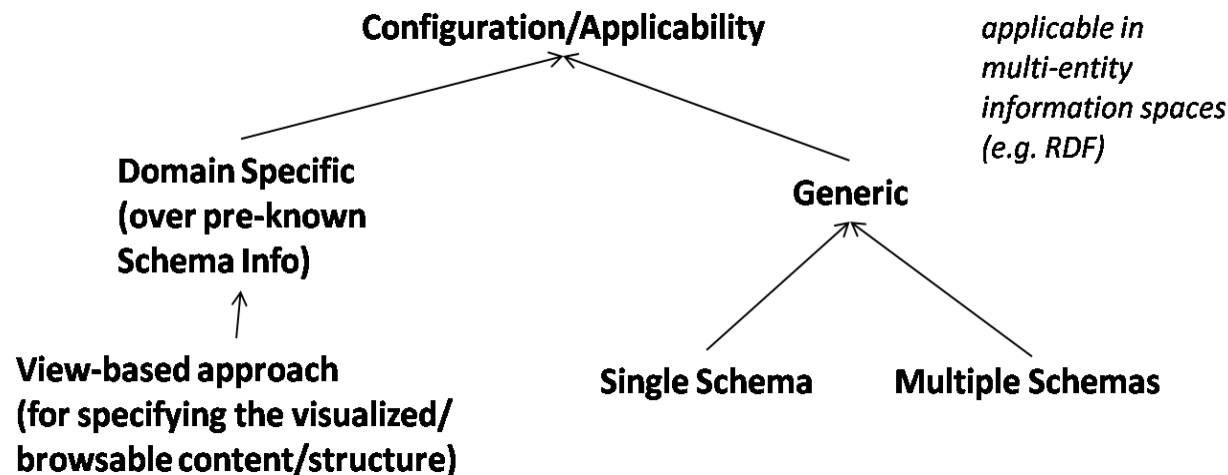
2. CONFIGURATION

○ *View-based* approach

- The structure of the browsable part has to be explicitly specified through configuration steps (e.g. by using a query language or logic rules)

○ *Generic* approach

- No configuration requirements w.r.t the underlying information space



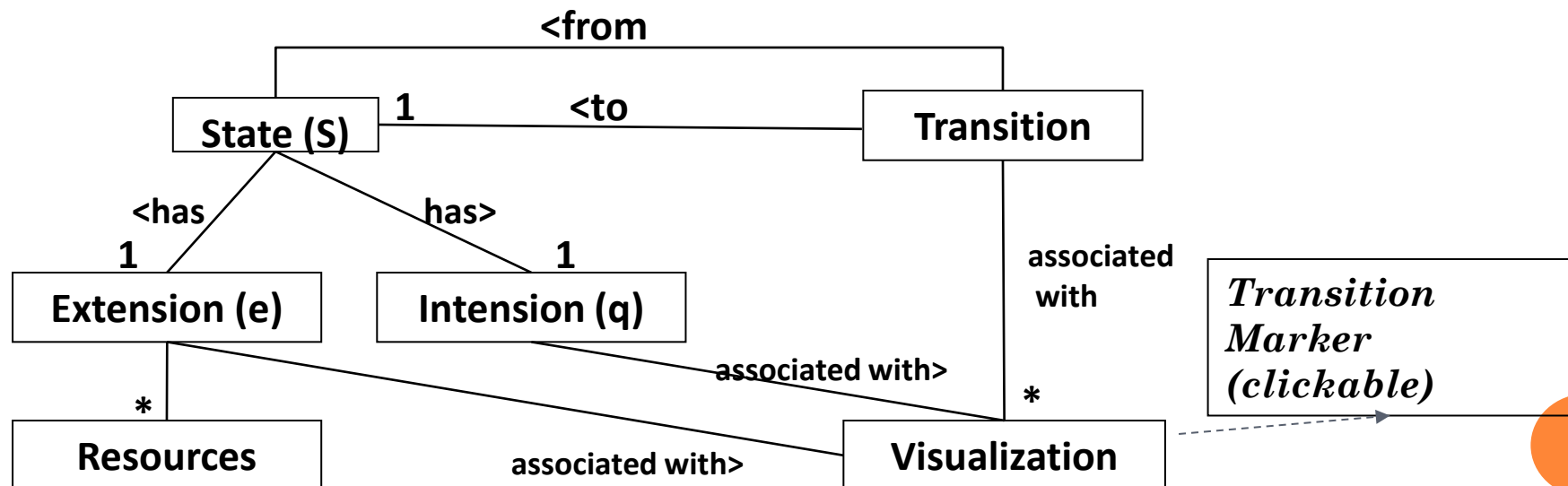
SURVEYING FACETED EXPLORATION APPROACHES

3. STATE SPACE (STATES, TRANSITIONS AND TRANSITION MARKERS)

- **State Space.** In general we can view the interaction as a state space consisting of **states** and **transitions**, therefore we can characterize, or comparatively evaluate, two BAs by comparing their state spaces, e.g. by identifying properties which are satisfied by their state space.

CONT. STATES

- A **state** has an **extension** (a set of items displayed, e.g. WS results), an **intension** (i.e. condition or query; satisfied by the extension), probably a name, and a number of **transitions** each leading to a different state.
- In addition each state has a visualization format for its (a) extension, (b) intension, as well as (c) its transitions (e.g. a tree-control, a list, a table). In any case, each transition has a clickable **transition marker** signifying the existence of the transition. Usually these markers are enriched with information regarding the target state.



CONT. STATES

Tennis Player (8 entities)

- Roger Federer (3)
- Rafael Nadal (2)
- Steffi Graf (1)
- MANSOUR BAHRAMI (2)
- Maria Sharapova (1)
- Elena Baltacha (1)
- Guillermo Vilas (1)
- David Ferrer (1)

Country (4 entities)

- Emirates (1)
- America (1)
- Argentina (1)
- United States (1)

» Tennis Player:

- MANSOUR BAHRAMI ■
- Roger Federer ■
- Steffi Graf ■

5 results: [reset](#)

[Best Tennis Players - Alistaday - A List a Day](#)
all time? We look at the top ten players in the sport's history, including Steffi Graf and Rafael Nadal.
[http://www.alistaday.com/sports/tennis/best-tennis-players/](#) - find its entities

[Top 10 Tennis Players of All Time - International Business ...](#)
Roger Federer is the greatest male tennis player of all time.
[http://www.ibtimes.com/top-10-tennis-players-all-time-704359](#) - find its entities

[Best Male Tennis Player - Top Ten List - TheTopTens.com](#)
Based on over 2,000 votes, Roger Federer is ranked number 1 out of 48 choices. Agree? Disagree? Place your vote on the top 10 list of Best Male Tennis Player.

Intension (as visualized)

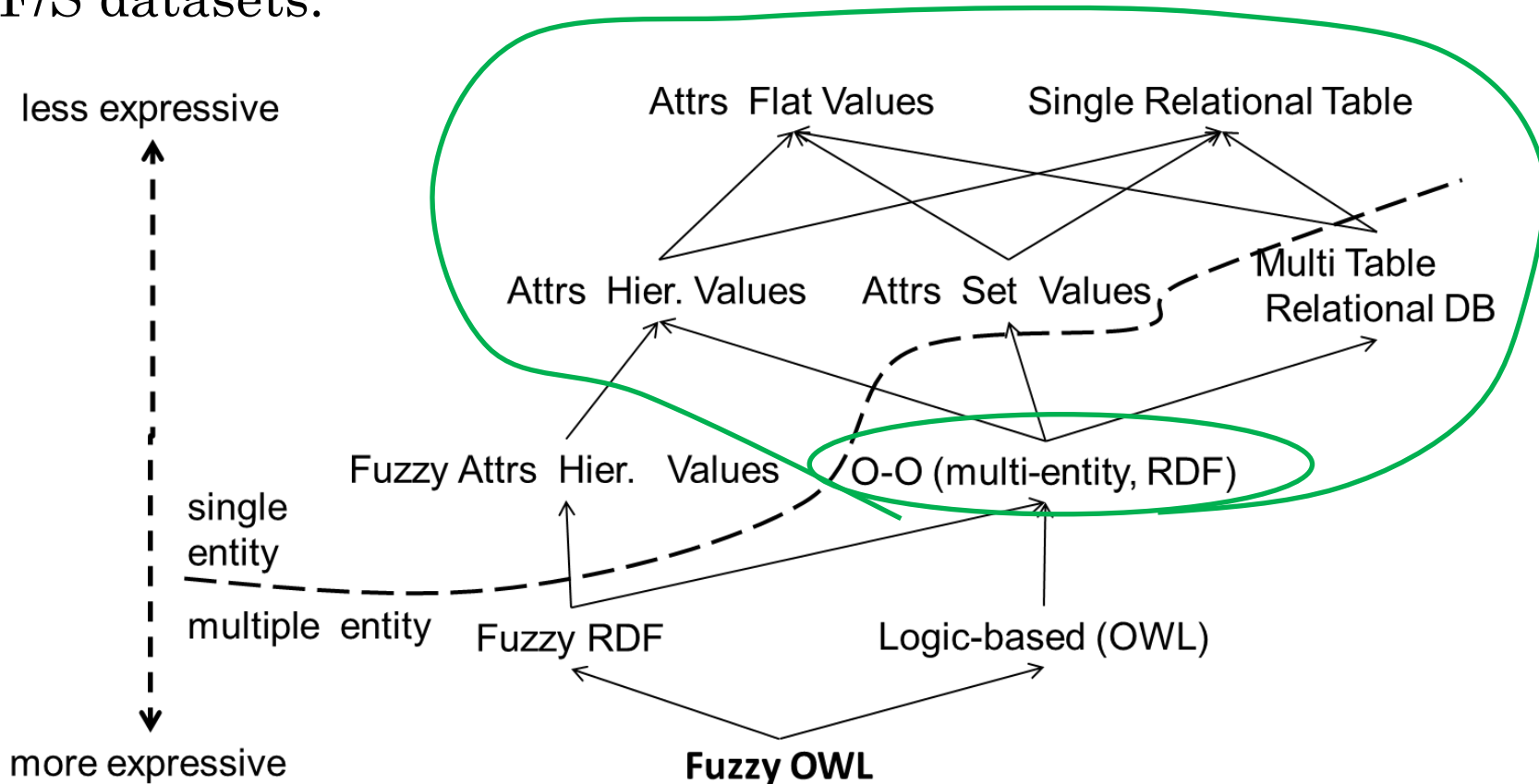
state

extension

Transition markers


A “GENERIC” INTERACTION MODEL FOR EXPLORING RDF/S DATASETS

- Now we will see an interaction model for faceted exploration over RDF/S datasets.



A GENERIC INTERACTION MODEL FOR EXPLORING RDF/S DATASETS

- Objective: Through simple clicks the user can reach states whose extension corresponds to the answer of complex queries.
- Example:



*Japanese cars for sale which
are driven by persons who
work at FORTH and know a
person who knows Bob*

- Source: It is a simplified version of the model for browsing Fuzzy RDF that is described in
 - Nikos Manolis, Yannis Tzitzikas: Interactive Exploration of Fuzzy RDF Knowledge Bases. ESWC (1) 2011: 1-16

SUPPORTED KINDS OF TRANSITIONS

- Supported transitions
 - Class-based
 - Property-based
 - Property-path based
 - Entity Type Switch

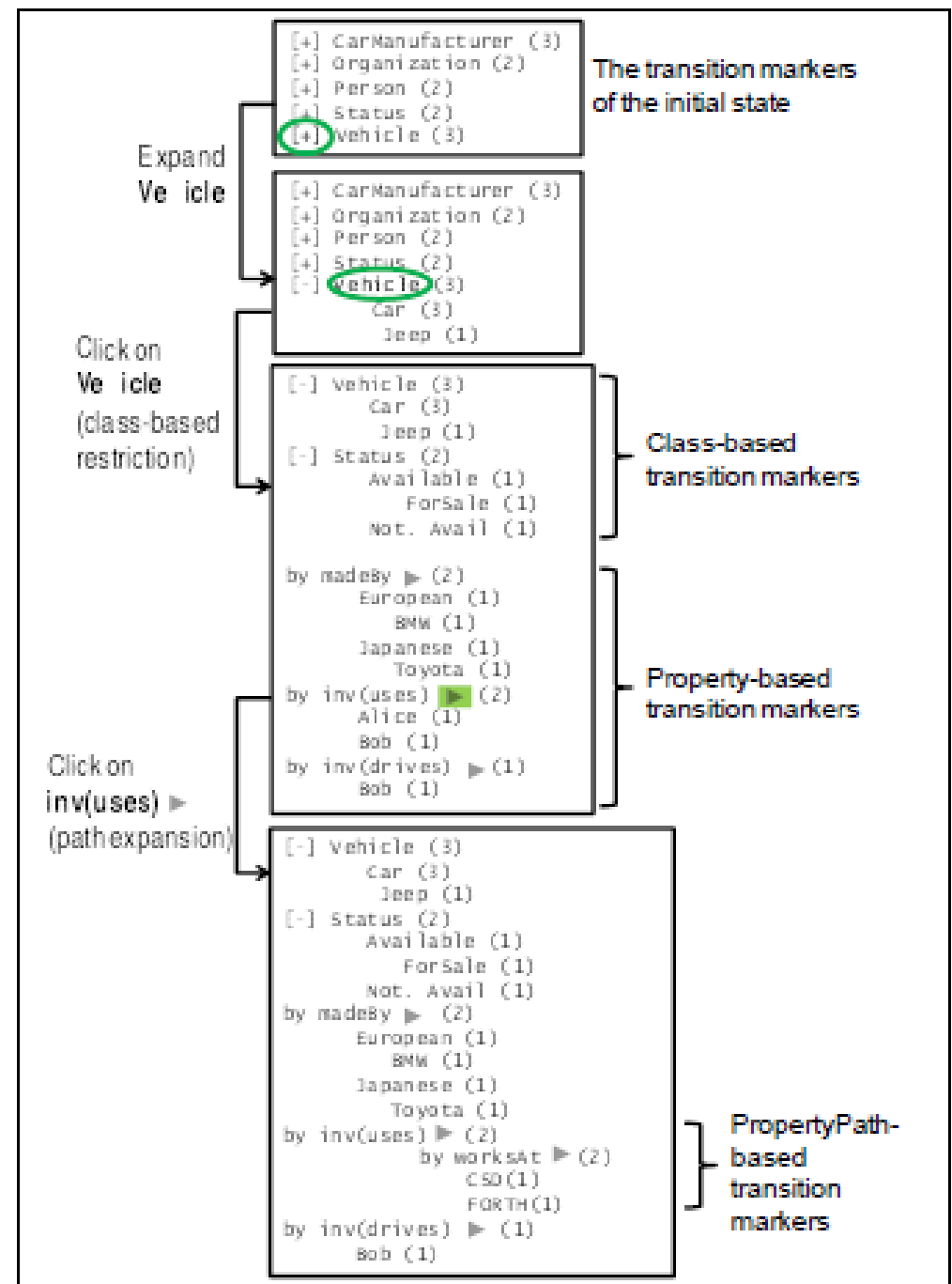


Fig. 10. Sketch of the GUI part for transition markers

NOTATIONS (FOR RDF/S)

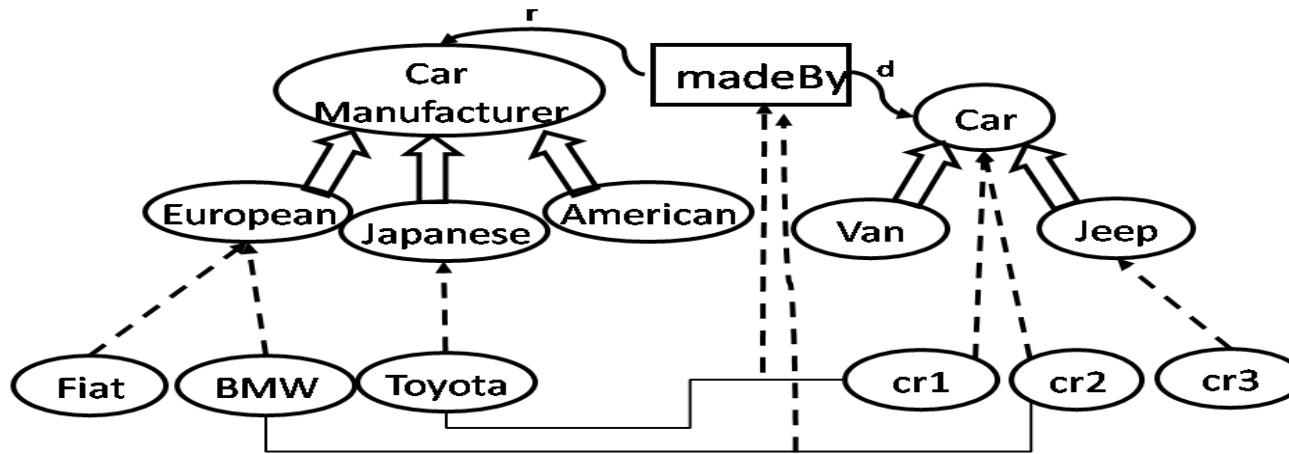
- A RDF/S KB is defined by a set of **RDF triples** of the form:

(subject, predicate, object)



- The **closure** $C(K)$ of a KB K contains all the triples explicitly asserted or inferred (based on the RDF/S semantics) from a KB
- The **Schema** of a RDF/S KB K is a 6-tuple $\Gamma = \langle C, Pr, domain, range, \leq_{cl}^*, \leq_{pr}^* \rangle$
 - C : Classes, Pr : Properties, $domain$ and $range$ of properties, $subclassOf$ (among C), $subPropertyOf$ (among Pr)
 - **Instance notations** for a KB:
 - Instances of a class $c \in C$: $inst(c) = \{ o \mid (o, type, c) \in C(K) \}$
 - Class instance triples: **(o type c)**
 - Instance of a property $p \in Pr$: $inst(p) = \{ (o, p, o') \mid (o, p, o') \in C(K) \}$
 - Property instance triples: **(o p o')**

CONT.



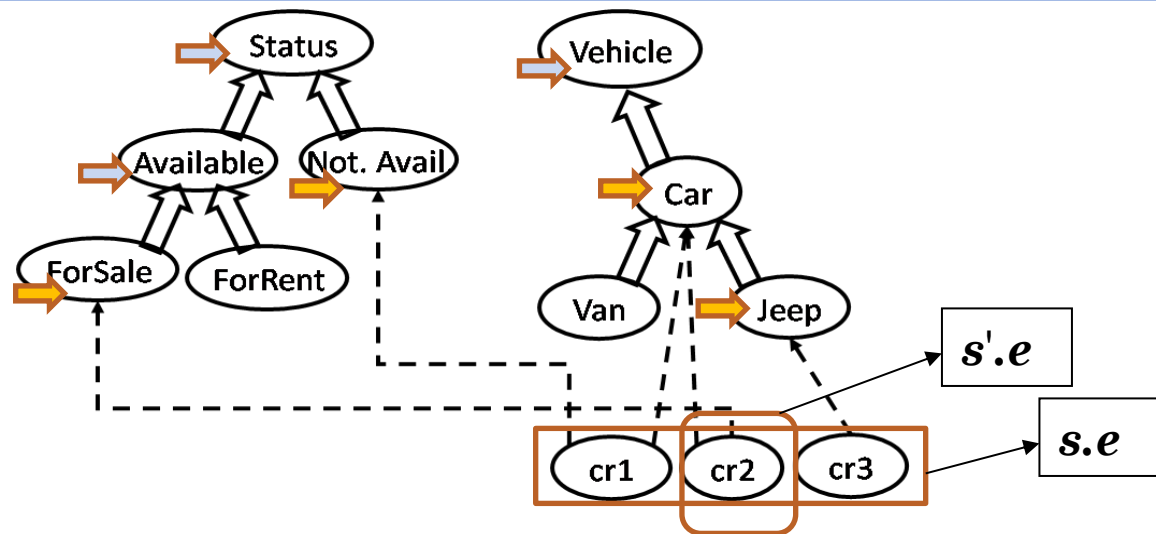
- $\text{inst}(\text{Jeep}) = \{\text{cr3}\}$
- $\text{inst}(\text{Car}) = \{\text{cr1}, \text{cr2}, \text{cr3}\}$
- $\text{Inst}(\text{Car}, \text{madeBy}, \text{CarManufacturer}) = \{$
 $(\text{cr1}, \text{madeBy}, \text{Toyota})$
 $(\text{cr2}, \text{madeBy}, \text{BMW})$
 $\}$

AUXILIARY DEFINITIONS: RESTRICTIONS AND JOINS

Assuming an initial state we need to tackle the following:

- ❖ How the **transitions markers** (*tms*) available to that state are computed
- ❖ How the **new state's extension** is computed after selecting a *tm*
- To define the above we need to define the notion of **restriction** and **join**
 $p \in PR \cup PR^{-1}$, E : set of resources, $vset$: a set of resources or literals
- **Restrictions** : Given a set E
 1. $Restrict(E, p : v) = \{ e \in E \mid (e, p, v) \in inst(p) \}$
 2. $Restrict(E, p : vset) = \{ e \in E \mid \exists v' \in vset \text{ and } (e, p, v') \in inst(p) \}$
 3. $Restrict(E, c) = \{ e \in E \mid e \in inst(c) \}$
- **Joining values**: Compute values which are linked with the elements of E
 - $Joins(E, p) = \{ v \mid \exists e \in E \text{ and } (e, p, v) \in inst(p) \}$

CLASS-BASED TRANSITIONS



Candidate tms: $TM_{cl}(s.e) = \{ c \in C \mid Restrict(s.e, c) \neq \emptyset \}$

Clicking on a $c \in TM_{cl}(s.e)$ then $s'.e = Restrict(s.e, c)$

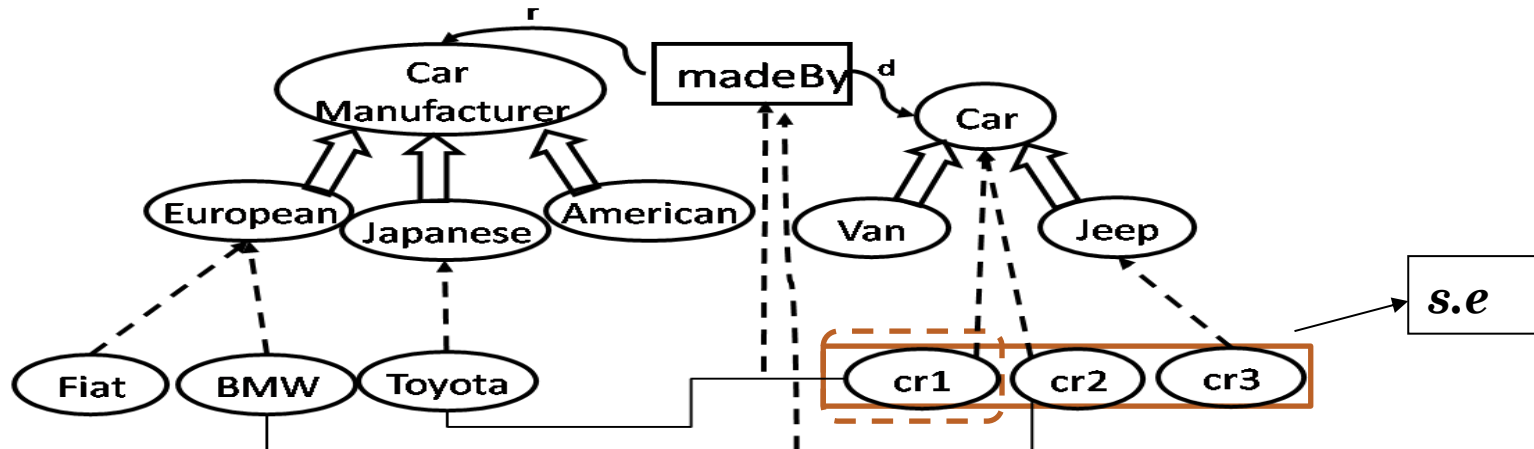
[-] Vehicle(3)
 Car(3)
 Jeep(1)
 [-] Status(2)
 Available(1)
 ForSale(1)
 Not Available(1)



[-] Vehicle(1)
 Car(1)
 [-] Status(1)
 Available(1)
 ForSale(1)

PROPERTY-BASED TRANSITIONS

- The goal: Restrict the current state's extension according to a property value v (resource or literal). E.g. All cars made by Toyota



- Candidate properties: $\text{Props}(s) = \{ p \in \text{Pr} \cup \text{PR}^{-1} \mid \text{Joins}(s.e, p) \neq \emptyset \}$

Candidate tms for a $p \in \text{Props}(s)$: $\text{Joins}(s.e, p) = \{ v \mid e \in s.e \text{ and } (e, p, v) \in \text{inst}(p) \}$

by madeBy(2)
BMW(1)
Toyota(1)

OR

by madeBy(2)

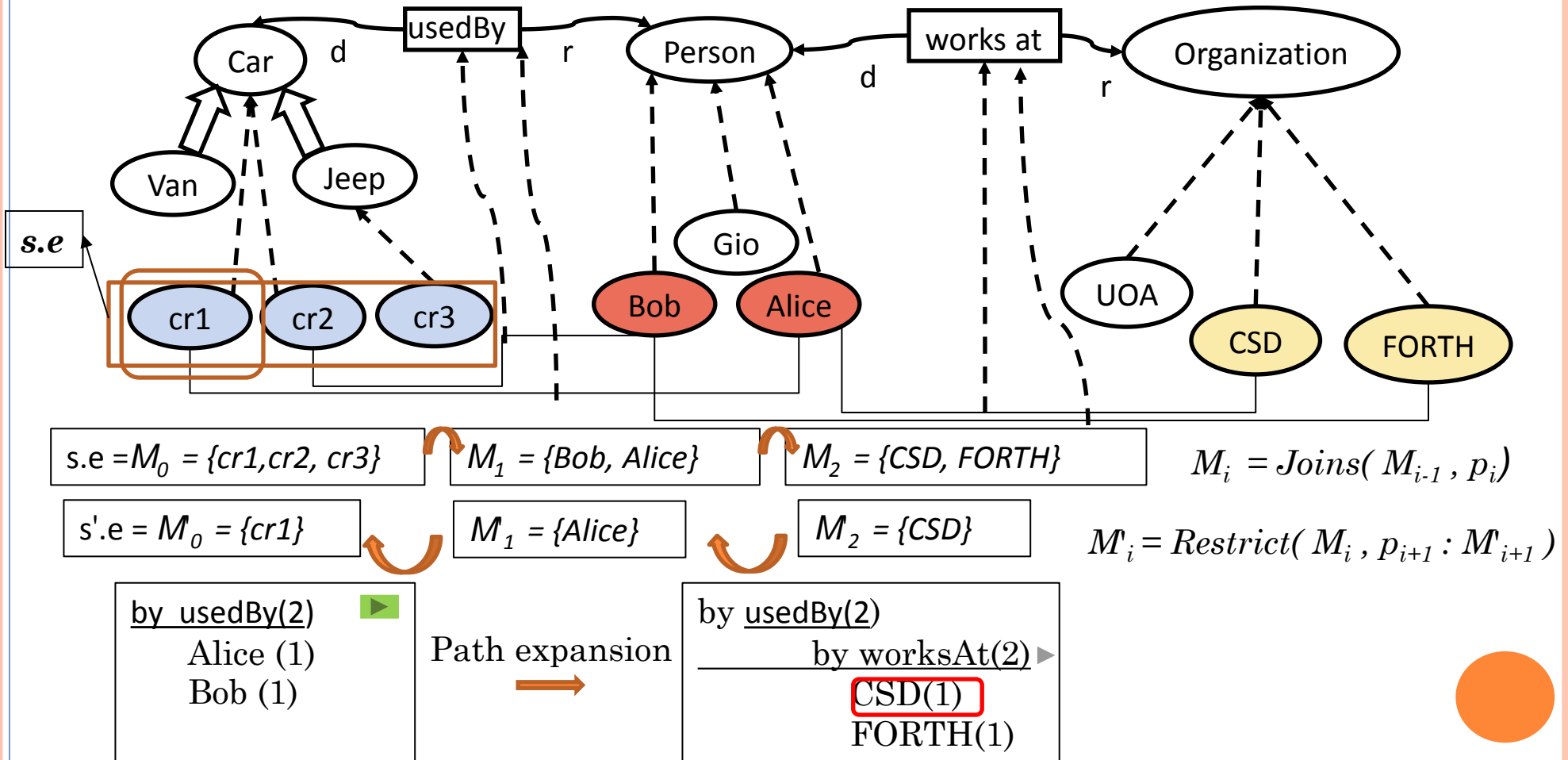
European
BMW(1)
Japanese (1)
Toyota(1)

$TM_{cl}(\text{Joins}(s.e, p))$

PROPERTY PATH-BASED TRANSITIONS

- Let p_1, \dots, p_k be a sequence of properties

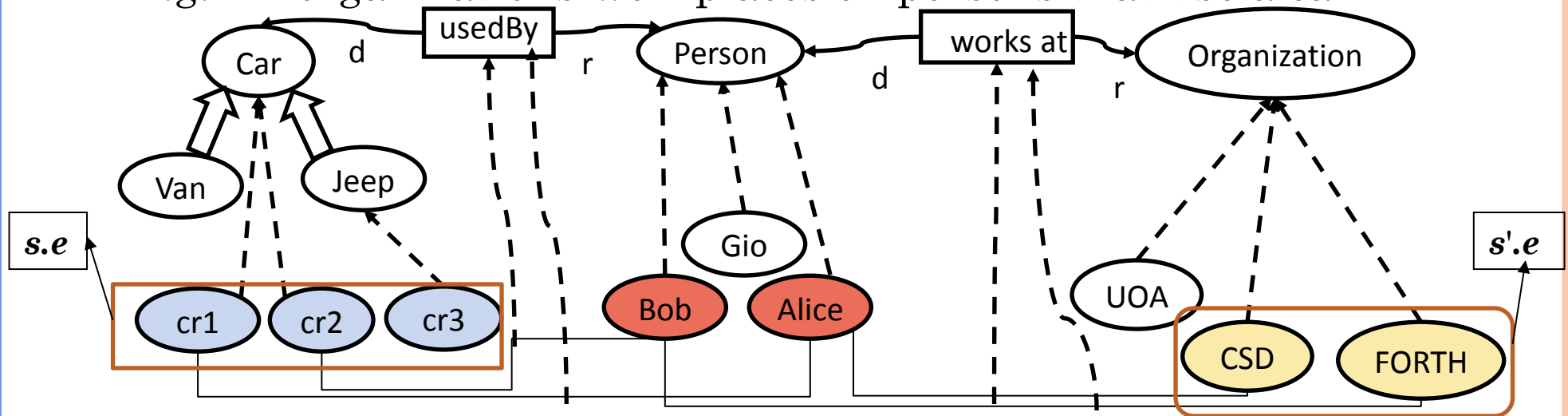
E.g. All cars used by persons working at CSD



ENTITY TYPE SWITCH TRANSITIONS

Allow users to move to a state whose extension is the current set of tms

E.g. All organizations-workplaces of persons that use a car



$s.e = M_0 = \{cr1, cr2, cr3\}$

$M_1 = \{Bob, Alice\}$

$M_2 = \{CSD, FORTH\}$

$s'.e = M_2 = \{CSD, FORTH\}$

by inv(uses)(2)

all
Alice (1)
Bob (1)

Path expansion

by inv(uses)(2)

by worksAt(2)

all
CSD(1)
FORTH(1)

APPLICABILITY

- To apply it, one has to implement Restrict and Join using the technology used for the underlying dataset
 - (e.g. over SPARQL, over SQL, over WhatEverQL, ...)
- For more see the related publication
 - It also captures the fuzzy aspect

SYNOPSIS OF THE FIRST PART



- A significant percentage of information needs are recall-oriented and this justifies the need for exploratory search services
- Faceted Search/Exploration is an effective model for exploratory search (now de facto standard)
- There are many variation of FE according to: structural complexity of the information space, configurability, supported state space
- We have seen a model for FE that captures the majority of datasets that exist



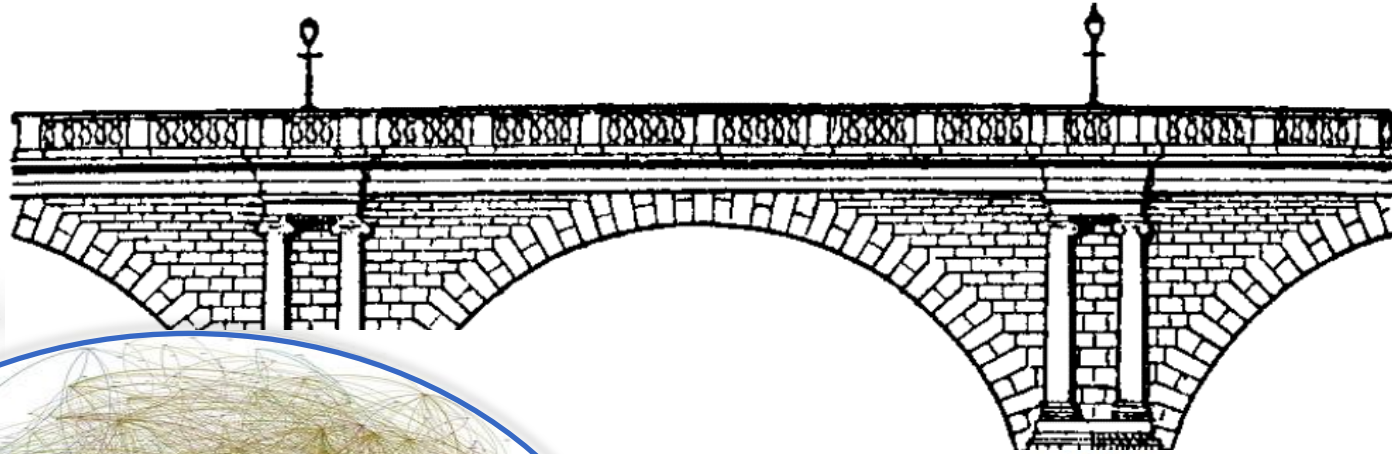
3. BRIDGING THE WEB OF DOCUMENTS WITH THE WEB OF DATA AT SEARCH TIME

79

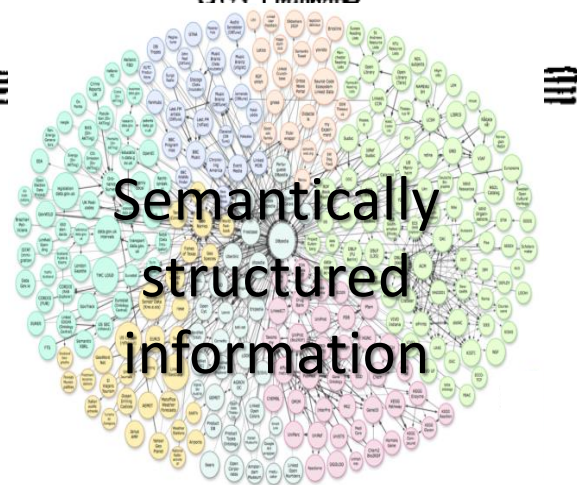
Possible Ways

Focus on doing this at *Search Time*

AN INTERESTING QUESTION



Unstructured documents
(e.g. web pages)

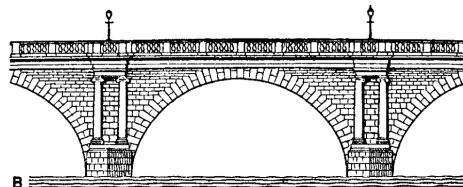


THIS QUESTION IS INSTANCE OF A BIGGER QUESTION

- How to integrate the results and/or foster collaboration between different communities



IR
Community



NLP
Community



Database
Community

AI
Community



BRIDGING THE WEB OF DOCUMENTS WITH THE WEB OF DATA

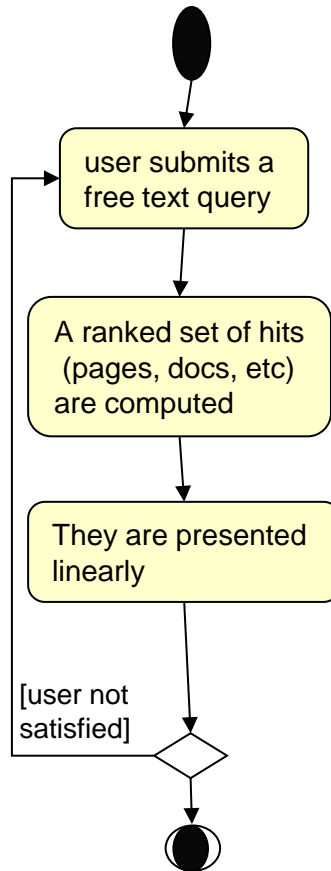
- Why
 - Both exist (why not exploit both?)
- When and Where
 - There is plethora of options. Let's consider the direction: How the searching over the Web of Documents can benefit from the existence of the Web of Data. A few options follow:
 - At indexing time
 - E.g. for the disambiguation of words
 - At query formulation time
 - E.g. autocompletion/query expansion/term_recommendation services that exploit the web of data, the various vocabularies, ontologies, etc.
 - At query evaluation time
 - Exploitation of how words/concepts are connected in the ranking formula
 - After query evaluation
 - For semantically post-processing of search results
 - ...

OUR FOCUS IN THIS TALK

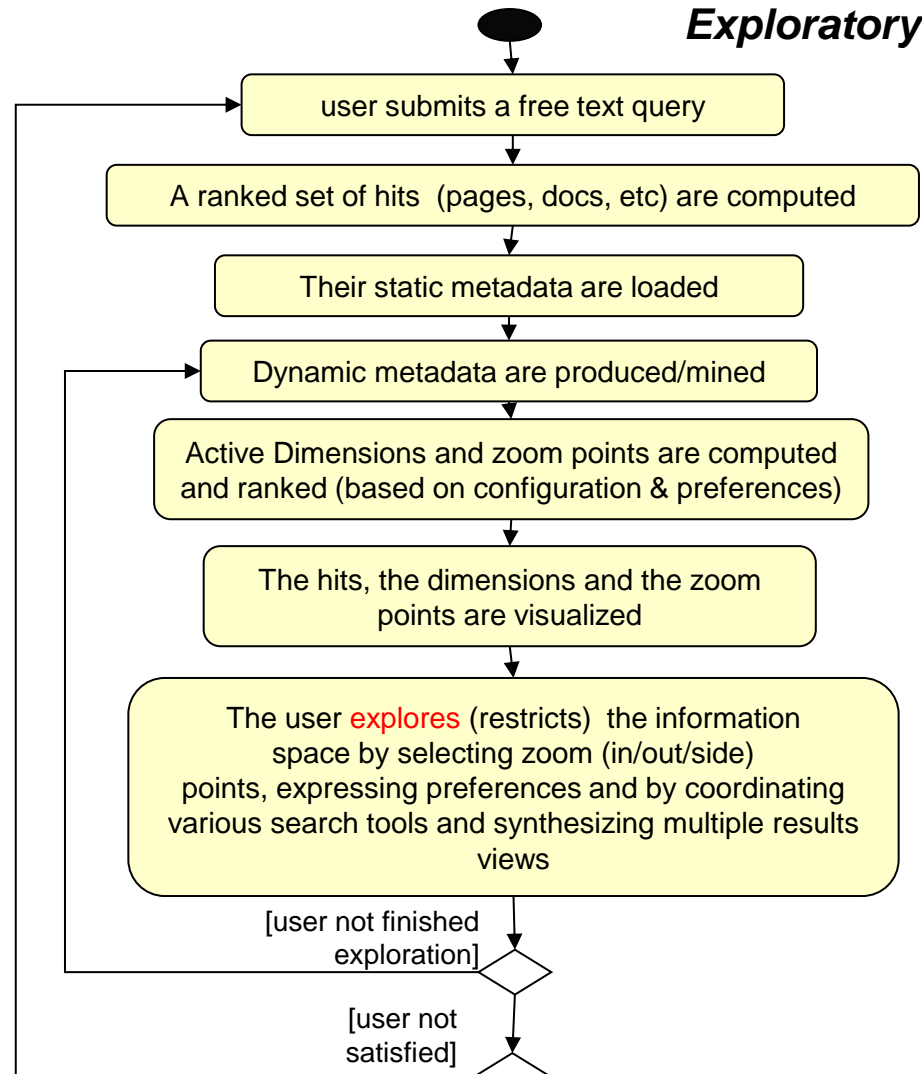
- *Don't change the way users search for information (keyword queries in Web Search Engines).*
- *Try to exploit LOD to offer value-adding services*

A PROCESS FOR **EXPLORATORY** SEARCH

Web searching today

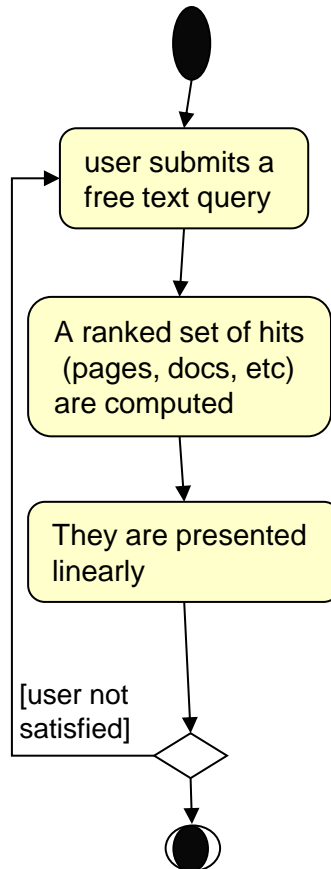


Exploratory Searching

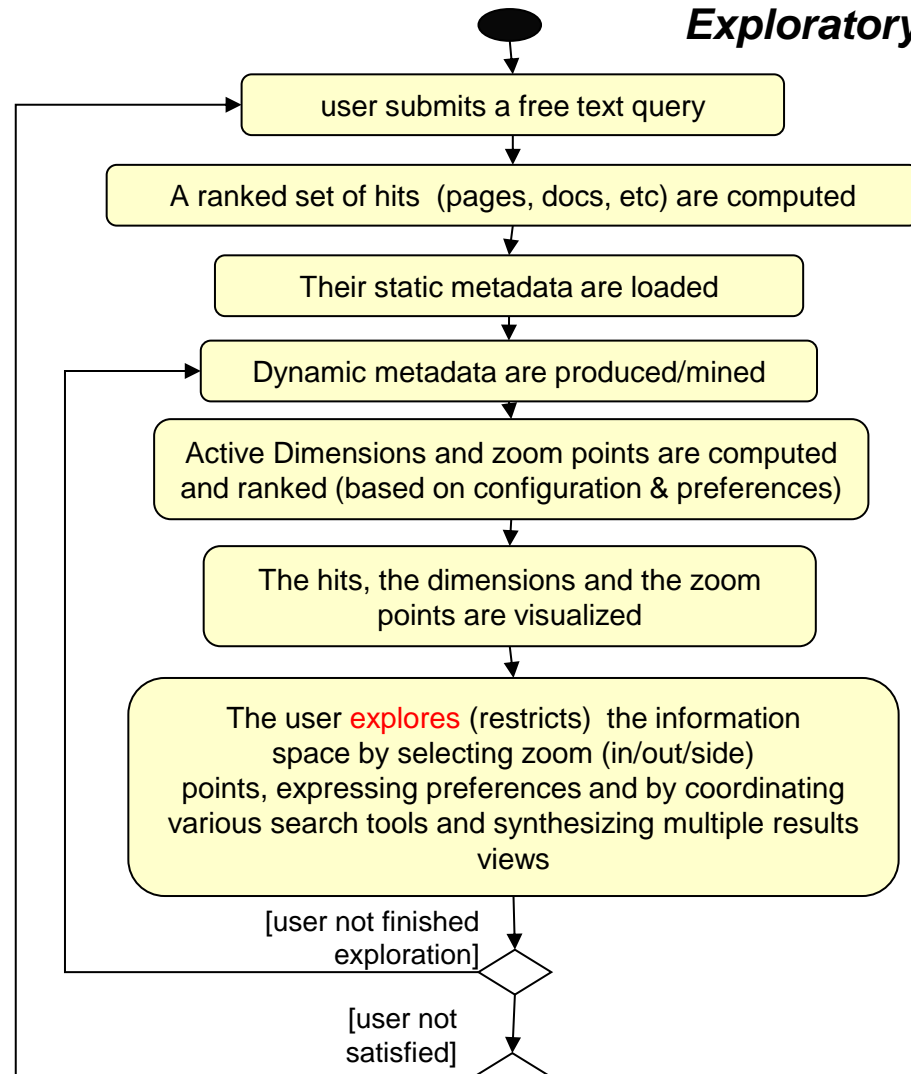


A PROCESS FOR **EXPLORATORY** SEARCH

Web searching today



Exploratory Searching



BENEFITS (OF THIS PROCESS)

- Does not change the way users search for information
- Can be applied over existing search systems/interfaces
- Provides overviews of the results (not only of the top hits)
- By clicking on a facet term (metadata value, mined entity, etc) the user can see the related hits even if they are low ranked
- Allows restricting the answer gradually

*Next we will see various applications of this process
(these applications concern Web Search, Exploratory Search and
LOD)*

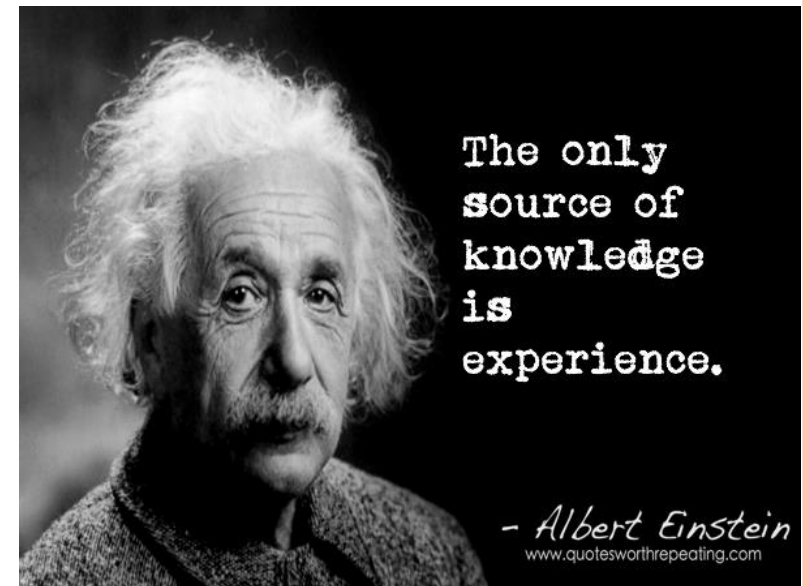


4. CASE STUDIES

87

CASE STUDIES OUTLINE

- As Case Studies I will present you a kind of “story” organized in 9 milestones
- These 9 milestones correspond to activities of ISL (Information Systems Laboratory) of FORTH-ICS.
 - They correspond to the period 2009-2014
 - They are presented in (almost) chronological order.

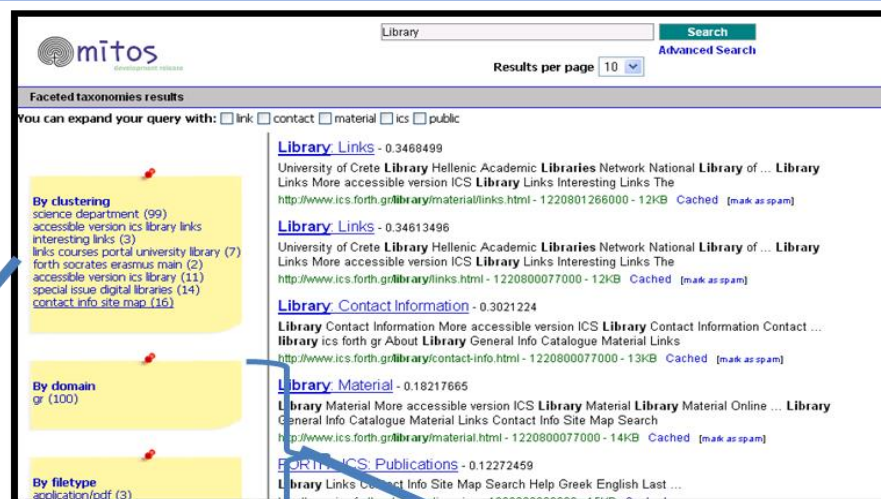


MILESTONE 1. THE MITOS WSE (2009)

- MITOS is WSE build from scratch and collaboratively with the students of the Computer Science Department of the University of Crete.
- Apart from the classical WSE functionality, Mitos offers faceted search over the results of the submitted queries.
 - It supports facets corresponding to metadata attributes of the web pages (static metadata), as well as facets corresponding to the outcome of snippet-based clustering algorithms (a kind of dynamic metadata).
 - The user can then restrict his/her focus gradually, by interacting with the resulting multidimensional structure through simple clicks.



THE MITOS WSE (2009)



Dimension based on dynamic
(query-dependent) metadata (of
the top ranked hits)

By clustering

- ▶ architecture (6)
- ⊕ contact (10)
- ⊕ content (10)
- ▶ copyright notice csd (13)
- ▶ course content english (7)
- ▶ csd (14)
- ▶ department (8)
- ⊕ forth (38)
- ▶ health telematics network (6)
- ▶ ics (39)
- ▶ information (11)
- ▶ network (8)
- ▶ physical (5)
- ▶ science (22)
- ▶ ΙΤΕ ΤΕΧΝΙΚΕΣ αναφορές (22)

REST (3981)

Dimensions based on static
metadata

N. Manolis and Y. Tzitzikas (ESWC'11)

By domain

- ⊕ gr (4067)

By date

- ⊕ 2008 (479)
- ⊕ 2007 (694)
- ⊕ 2006 (1340)
- ⊕ 2005 (184)
- ⊕ 2004 (106)
- ⊕ 2003 (82)
- ⊕ 2002 (88)
- ⊕ 2001 (28)
- ⊕ 2000 (13)
- ⊕ 1999 (4)
- ⊕ 1998 (6)
- ⊕ 1997 (1)
- ▶ Unknown (1042)

By filetype

- ▶ application/msword (16)
- ▶ application/pdf (1476)
- ▶ application/vnd.ms-powerpoint (29)
- ▶ text/html (2546)

By language

- ▶ Any (UTF-8) (18)
- ▶ Greek (1209)
- ▶ Latin-1 (Europe, Latin America, Caribbean, Canada, Africa) (944)
- ▶ Latin-2 (Central and Eastern Europe) (4)
- ▶ Unknown (1892)

Y. Tzitzikas, Panel@ExploreDB, Athens
2014

90

THE MITOS WSE (2009)



information systems laboratory

Search

Advanced Search

Results per page 10

Faceted taxonomies with on-demand clustering results

8558 Results 1 - 10 from 8558 for information systems laboratory (ms)

By clustering

- ▶ activities (39)
- ▶ biomedical informatics laboratory (10)
- ▶ decision support systems (2)
- ▶ dimitris (7)
- ▶ events (43)
- ▶ forth (25)
- ▶ history (35)
- ▶ ics (85)
- ▶ informatics (64)
- ▶ laboratories (2)
- ▶ projects (53)
- ▶ publications (47)
- ▶ seminars (26)
- ▶ support (5)
- ▶ yannis (3)

REST (8458)

By domain

- ▶ gr (8558)

By date

- ▶ 2009 (606)
- ▶ 2008 (668)
- ▶ 2007 (3113)

Information Systems Laboratory - 0.916549

information systems laboratory ics isl img isl src images buttons isl ...
http://www.ics.forth.gr/isl/publications/by_name.jsp?Person_ID=7 - 0 - 8KB
Cached - Similar pages [mark as spam]

FORTH - ICS: Information Systems Laboratory - 0.8703994

information systems information retrieval systems database workflow
management systems semantically rich ... forth ics information systems
laboratory information systems laboratory head laboratory prof
<http://www.ics.forth.gr/isl/index.html> - 1173087253000 - 17KB Cached - Similar
pages [mark as spam]

Information Systems Laboratory - 0.8702642

information systems laboratory ics isl panos constantopoulos muse
multimedia accessible version ... ics isl isl centre cultural informatics history
events activities projects publications
http://www.ics.forth.gr/isl/publications/by_year.jsp?Year_of_publication=1987 - 0 -
17KB Cached - Similar pages [mark as spam]

Information Systems Laboratory - 0.8692224

information systems laboratory null ics isl accessible version ics isl isl ...
centre cultural informatics history events activities projects publications
seminars people links

publication=null - 0 -

information systems laboratory ics isl panos constantopoulos office
document retrieval multos ... accessible version ics isl isl centre cultural
informatics history events activities
http://www.ics.forth.gr/isl/publications/by_year.jsp?Year_of_publication=1986 - 0 -

A user wants to get
information about
**Information
Systems
Laboratory**

8558 initial results

We can focus on “By date”
facet, clicking the “2009” label.

(CONT.)



information systems laboratory

Search

Advanced Search



Results in RDF/XML

Results per page 10

Faceted taxonomies with on-demand clustering results

606 Results 1 - 10 from 606 for information systems laboratory. (15729 ms)

By clustering

- ▶ athanasios mouchtaris (3)
- ⊕ communication (15)
- ▶ distributed (13)
- ▶ dynamic (13)
- ⊕ forth (11)
- ▶ home page (3)
- ⊕ networks (13)
- ▶ **news (5)**
- ▶ oikonomou (2)
- ▶ page (5)
- ▶ presentation (3)
- ▶ publications (6)
- ▶ spring (3)
- ▶ tziritas (2)
- ▶ ire (3)

REST (558)

We can further limit the results, by selecting one of the clusters (they were recomputed for the new focus)

The results of the selected group are loaded in the results' panel and all facets are updated.

By domain

- ⊕ gr (606)

By date

- ⊖ 2009 (606)
- ⊕ June (71)
- ⊕ May (80)
- ⊕ April (212)

[Information Systems Laboratory: Seminars](#) - 0.28782406

challenge succeed transition traditional **information systems information** retrieval **systems** database workflow ... management **systems** semantically rich large scale adaptive **information systems systems** characterized
<http://www.ics.forth.gr/isl/services.html> - 1244639664000 - 21KB [Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[Information Systems Laboratory: Seminars](#) - 0.276

seminars seminars ics isl in
matics subjects developed
including greek
ml - 1244123830000 - 16KB

[Information Systems](#) - 0.19771457

ems CS 463 **Information Retrieval**
Teaching Material Lectures and Program
s Links

<http://www.csd.uoc.gr/~hy463/2007/en/grades.html> - 1241012788000 - 2KB
[Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[CS-463 Information Retrieval Systems](#) - 0.18768528

CS 463 **Information Retrieval Systems** CS 463 **Information Retrieval**
Systems ... Course **Information** Teaching Material Lectures and Program
Exercises and Assignments Grades

<http://www.csd.uoc.gr/~hy463/2007/en/announcements.html> - 1241012778000 - 2KB [Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[CS-463 Information Retrieval Systems](#) - 0.18636355

CS 463 **Information Retrieval Systems** CS 463 **Information Retrieval**
Systems Spring ... **Information** Teaching Material Lectures and Program

(CONT.)



information systems laboratory

Search

Advanced Search

Results per page 10



Results in RDF/XML

Faceted taxonomies with on-demand clustering results

5 Results 1 - 5 from 5 for information systems laboratory. (86 ms)

By clustering
news (5)

By domain
gr (5)

By date
2009 (5)
June (2)
May (2)
January (1)

By filetype
text/html (5)

By language
Latin-1 (Europe, Latin America,
Caribbean, Canada, Africa) (5)

[FORTH - ICS: News](#) - 0.053278793

information greek information greek information greek information greek
information greek information ... greek information greek information
greek information greek information greek information greek
<http://www.ics.forth.gr/news/news-prev.html> - 1241420849000 - 54KB [Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[FORTH - ICS: News](#) - 0.040669773

laboratories publications services library links contact info site map search
help ...
<http://www.ics.forth.gr/news.html> - 1244102246000 - 23KB [Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[FORTH - ICS: Welcome Note by the Director of ICS-FORTH -](#)
0.028257346

zoomin ics announcements news press releases **laboratories** publications
services library links ...
<http://www.ics.forth.gr> - 1232022089000 - 19KB [Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[FORTH - ICS: News](#) - 0.02504004

technical aspects multimodal **systems** tams department informatics
university hamburg germany university ... **information** science university
pennsylvania professor head cognitive informatics **laboratory** laval university
<http://www.ics.forth.gr/news/lectures-prev.html> - 1244102254000 - 100KB
[Cached](#) - [Similar pages](#) [\[mark as spam\]](#)

[FORTH - ICS: Lectures](#) - 0.01583067

announcements news **laboratories** publications services library links contact
info site map ...
http://www.ics.forth.gr/news/ian_cernockv_lecture.htm - 1241420696000 - 20KB

With only 2 clicks, we
have limited the
results to 5 hits.

CONT.

○ Evaluation with Users (main results) :

- Faceted search, combining dynamically and statically mined metadata
 - lead to much improved task completeness with much less user interactions
 - was more preferred by the users (advanced and plain ones) and lead to greater satisfaction, than plain clustering or faceted interfaces

○ Most Important Related Publications

- [ECDL'09] P. Papadakos, S. Kopidaki, N. Armenatzoglou and Y. Tzitzikas. Exploratory Web Searching with Dynamic Taxonomies and Results Clustering. In ECDL 2009
- [WISE'09] S. Kopidaki, P. Papadakos, and Y. Tzitzikas. STC+ and NM-STC: Two novel online results clustering methods for web searching. In WISE 2009
- [J. KAIS 2012] P.Papadakos, S.Kopidaki, Nikos Armenatzoglou and Y. Tzitzikas On exploiting Static and Dynamically mined Metadata for Exploratory Web Searching , KAIS Journal, 2012

MILESTONE 2. DURING TYPING?



Q instant overview search

- Then we questioned ourselves:
 - ***why not offering this functionality **during query typing**, i.e. a kind of **richer autocomplete service**?***
- This resulted to what we called **Instant Overview Search (IOS)**.
- The idea:
 - For the **frequent queries**, **pre-compute and store** not only the first page of results, but **also the analysis of these hits**
- Technical challenge
 - Since the amount of information that has to be stored for each query is higher (and obviously does not fit in main memory) we devised a **partitioned trie-like index** for efficiency (plus a dedicated **cache**)

IOS (INSTANCE OVERVIEW SEARCH), 2011-2012



This is Pavlos Fafalios
(I supervise his PhD)

*We want to find information about the life of Marilyn
Monroe
(and probably its connection to Pavlos?)
However, we are not sure for the spelling of her name.
So, we start typing "mari".*

(CONT.)

[home page](#) • [visit uoc-csd](#) • [visit ics-forth](#)

» clusters overview:

- marilyn(99)
 - monroe(33)
 - quotes(4)
 - biography(3)
 - manson(8)
 - home(8)
 - free(7)
 - encyclopedia(3)
 - quotes(5)
 - photo(5)
 - collection(5)
 - news(5)
 - encyclopedia(4)
 - biography(4)
 - new(10)
 - york(4)
 - gavin rossdale(3)
 - images(5)
 - andy(3)
 - online(3)



Search interface showing the input field with "mari", a search button, and a dropdown list of suggestions: "marilyn", "marilyn monroe", and "mario games". Below the suggestions are radio buttons for "SET (default)", "PET", "STIE", and "PTIE". A yellow arrow points to the suggestions list.

List of query's suggestions.

» first page results overview:

Marilyn (singer) - Wikipedia, the free encyclopedia - 0

Peter Robinson (born 3 November 1962), better known as **Marilyn**, is a British pop singer who achieved international fame in the 1980s with his hit song ...Early life - Blitz years - Career - Recent activity
[en.wikipedia.org/wiki/Marilyn_\(singer\)](http://en.wikipedia.org/wiki/Marilyn_(singer))

Marilyn Monroe - Wikipedia, the free encyclopedia - 1

Marilyn Monroe born Norma Jeane Mortenson, but baptized Norma Jeane Miller - Somethings Got to Give - Some Like It Hot
en.wikipedia.org/wiki/Marilyn_Monroe

Marilyn (hill) - Wikipedia, the free encyclopedia - 2

A **Marilyn** is a mountain or hill in the United Kingdom, Ireland or Isle of ...
[en.wikipedia.org/wiki/Marilyn_\(hill\)](http://en.wikipedia.org/wiki/Marilyn_(hill))

First page of results of the top suggestion "marilyn"

We can continue typing the query. Instantly new suggestions are shown

Cluster Label Tree of the top suggestion "marilyn"

(CONT.)

[home page](#) • [visit uoc-csd](#) • [visit ics-forth](#)

» clusters overview:

- marilyn monroe(100)
 - quotes(8)
 - news(10)
 - online(5)
 - photos(7)
 - collection(6)
 - images(7)
 - death(6)
 - gallery(5)
 - photo(6)
 - video(4)
 - free(3)
 - pictures(5)
 - encyclopedia(3)
 - biography(5)**
 - links(3)

By clicking a label, the results of the specific cluster are loaded in the results panel.



We selected the suggestion “marilyn monroe”. The results’ first page and cluster label tree for this suggestion were loaded immediately.

» first page results overview:

[Marilyn Monroe Biography from Who2.com](#)

Marilyn Monroe's sex appeal talent and untimely death combined to make her an enduring star and one of Hollywood's most recognizable icons. Early in.

<http://www.who2.com/marilynmonroe.html>

[Marilyn Monroe: Biography from Answers.com](#)

Marilyn Monroe Actor Born: 1 June 1926 Birthplace: Los Angeles California Died: 4 August 1962 (drug overdose) Best Known As: Hollywood's most.

<http://www.answers.com/topic/marilyn-monroe>

[The Marilyn Pages-Marilyn Monroe biography and images](#)

life of Marilyn Monroe Norma Jean. ... Marilyn Monroe. The Marilyn Pages have moved to ellensplace.net/marilyn.html (If you are not taken to the new ...

<http://www.ionet.net/~jellenc/marilyn.html>

[The Marilyn Pages-Marilyn Monroe biography and images](#)

life of Marilyn Monroe Norma Jean. ... NOTE FOR AOL USERS ♦ Site Awards for The Marilyn Pages ♦ to ellen's

<http://www.ellensplace.net/marilyn.html>

IOS (INSTANCE OVERVIEW SEARCH), 2011-2012

We can exploit this technique for any kind of pre-processing of search results (e.g. metadata-based faceted search, snippet-based clustering, entity mining, etc)

The screenshot displays the IOS interface with search results for 'tim' and 'tim berners-lee'. The interface is divided into several sections:

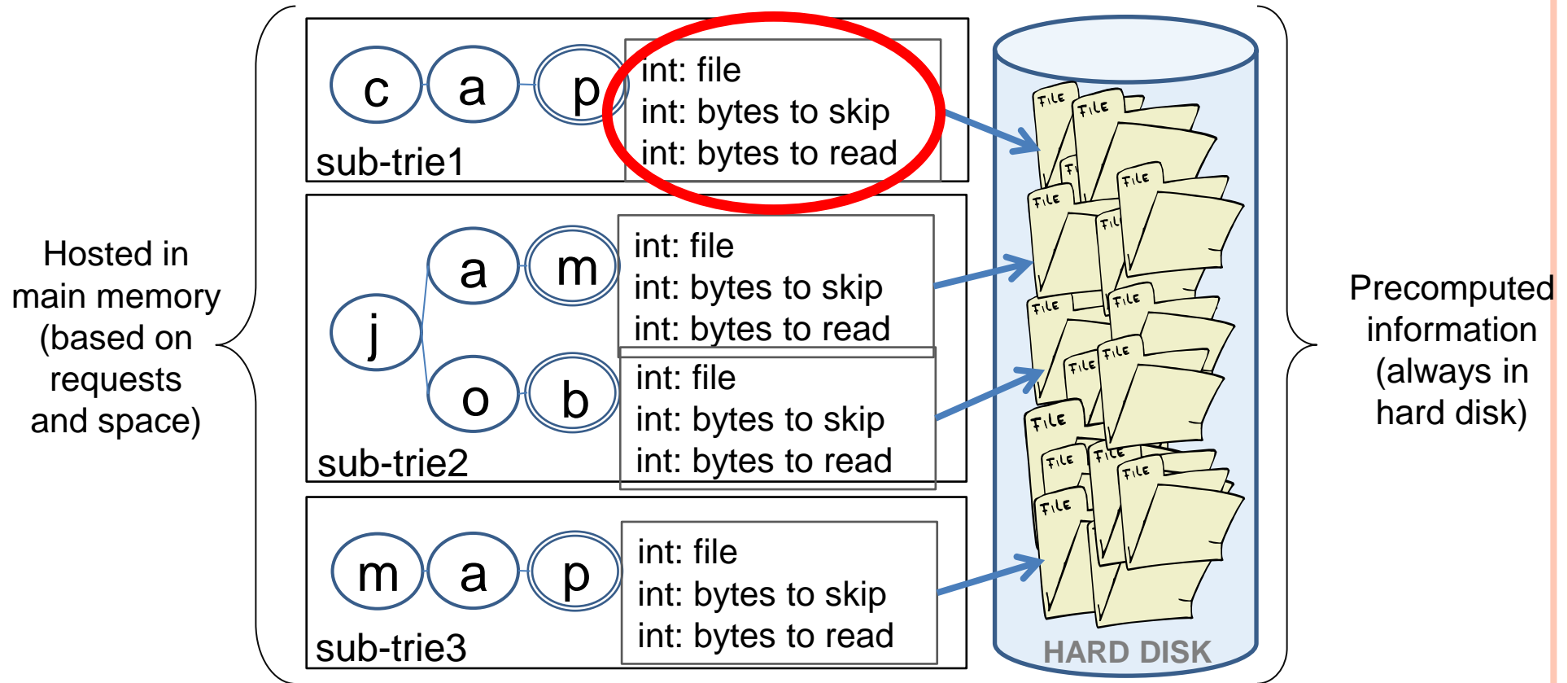
- Search Bar:** Contains the text 'tim' and a dropdown menu with suggestions: 'tim curry', 'tim montgomery', 'tim powers', and 'tim berners-lee' (highlighted).
- Search Results:** Displays results for 'tim berners-lee' with links to Wikipedia, Tim Berners-Lee Bio, Facebook, and Mahalo.com. Each link has a red box labeled 'find its entities'.
- Faceted Search:** A sidebar on the right shows faceted search results for 'tim berners-lee' (100 entities). The facets are:
 - By date:** 2011 (32), 2010 (29), 2009 (17), 2008 (8), Unknown (14).
 - By filetype:** application/msword (11), application/pdf (16), application/ms-powerpoint (3), text/html (70).
 - By language:** Any (UTF-8) (27), English (4).
- Entity Mining:** A section on the left shows entity mining results for 'tim' and 'tim berners-lee'. It lists entities like 'Person' (521 entities), 'Date' (566 entities), and 'Organization' (416 entities). Red boxes labeled 'show all' are present next to each entity list.

a) Entity mining

b) Metadata-based faceted search

c) Results clustering

IOS INDEXES



Average Retrieval Time \approx **135ms**

*Experiments over a server running on a **modest personal computer**, with a synthetic query log of **1 million** distinct queries and synthetic precomputed information of **1 Terabyte***

100

CONT

○ Key results

- A partitioned trie-based index structure that can efficiently support recommendations for millions of distinct queries even with modest hardware
 - One can provide instant access to large amount of data, utilizing the existing resources, without requiring more hardware
- A hybrid caching policy (70% static and 30% dynamic) seems to be the more appropriate choice yielding a throughput increment of around 80% and a 25% speedup

○ Demo

- <http://www.ics.forth.gr/isl/ios>
 - Select the system “Instant Entity Mining + Clustering (over Bing)”

○ Related Publications

- [WISE'11] P. Fafalios and Y. Tzitzikas, Exploiting Available Memory and Disk for Scalable Instant Overview Search, 12th International Conference on Web Information System Engineering (WISE 2011), Sydney, Australia, October 2011
- [WWW'12] P. Fafalios, I. Kitsos and Y. Tzitzikas, Scalable, Flexible and Generic Instant Overview Search, 21st International Conference on World Wide Web, (WWW 2012), Demo Paper, Lyon, France, April 2012

MILESTONE 3. ENTITY MINING AND LOD?

- Then we questioned ourselves:
 - ***why not exploiting LOD in the context of entity mining of the search results?***
- Motivation
 - LOD contains plenty of information about **Named Entities** (their names, attributes, relationships with other entities, etc)
- Output
 - IOS Entity Mining
 - LOD is used as source for Named Entity Recognition
 - LOD is used for providing more information about the identified entities



The screenshot shows the IOS Entity Mining interface. At the top left is the logo with the text "entity mining" and "ios". A search bar contains the text "barack obama" and a "Search" button. Below the search bar, there are options: "100 results to mine" and a checkbox for "mine only snippets". The results are categorized into "Person" (1427 entities) and "Organization" (842 entities). Under "Person", a detailed entry for Nicolas Sarkozy is highlighted, showing his photo, name, title, birth date, birth place, profession, and web site. Under "Organization", a list of organizations is shown with their entity counts. A yellow box on the left contains a list of features: "Automatically connects knowledge with documents at query time", "No preprocessing", and "No indexing". A red box highlights the "find its entities" link for the Wikipedia entry for Barack Obama. Another red box highlights the "find its entities" link for the Facebook entry for Barack Obama.

entity mining
ios

barack obama

Search

100 results to mine
☐ mine only snippets

Person (1427 entities)

Barack Obama - Wikipedia, the free encyclopedia
Barack Hussein Obama II (born August 4, 1961) is the 44th and current President of the United States. He is the first African American to hold the office. Obama ...
http://en.wikipedia.org/wiki/Barack_Obama - find its entities

Barack Obama
Barackobama.com is the official re-election campaign website of President Barack Obama. Visit the site for the latest updates from ...

Organization (842 entities)

Harvard (14)
White House (18)
Congress (14)
University of Hawaii (10)
Columbia University (8)

Nicolas Sarkozy
Current President of France
Birth date: 1955-01-28
Birth place: Paris, France
Profession: Lawyer
Web site: <http://www.sarkozy.fr>
Page: http://en.wikipedia.org/wiki/Nicolas_Sarkozy

Quote of the day: "This is a good first step, but it is only a step. Congress needs to pass the rest of my American Jobs Act so that can create jobs and put ..."
<https://www.facebook.com/barackobama> - find its entities

- Automatically connects knowledge with documents at query time
- No preprocessing
- No indexing

CONT

The screenshot shows a web interface for 'entity mining'. At the top left is a logo with the text 'entity mining' and 'ios' (where 'i' is a magnifying glass over a globe). A search bar contains 'barack obama' and a 'Search' button. Below the search bar, it says '100 results to mine' and a checkbox for 'mine only snippets'. On the left, a red-bordered box contains the text 'Results of selected entities: reset'. Below this, there are three search results, each with a title, a snippet, and a link to 'find its entities'. The first result is 'Barack Obama' from 'BarackObama.com'. The second is 'About Barack Obama — Barack Obama' from 'PresidentObama speaking'. The third is 'News for barack+obamaBarack Obama - Wikipedia'. On the right, there are two panels. The top panel is titled 'Person (1427 entities)' and lists names with counts and a small icon. 'Joe Biden (13)' and 'John McCain (8)' are highlighted with red boxes. The bottom panel is titled 'Organization (842 entities)' and lists 'Harvard (14)' and 'White House (22)', with the latter highlighted by a red box. A 'show all' link is at the bottom right of the 'Person' panel.

entity mining
ios

barack obama

Search

100 results to mine
☐ mine only snippets

Results of selected entities: reset

Barack Obama
BarackObama.com is the official re-election campaign website of PresidentBarack Obama. Visit the site for the latest updates from the Obama campaign, ...
<http://www.barackobama.com/> - find its entities

About Barack Obama — Barack Obama
Barack Obama is the 44th President of the United States of America. PresidentObama speaking. President Obama was born in Hawaii on August 4th 1961 to a ...
<http://www.barackobama.com/record> - find its entities

News for barack+obamaBarack Obama - Wikipedia, the free encyclopedia
Barack Hussein Obama II is the 44th and current President of the United States of America. He was born on August 1, 1961, in Honolulu, Hawaii, to a Kenyan mother and a ...
<http://www.barackobama.com/record> - find its entities

Person (1427 entities)

- Barack Obama (16)
- Michelle Obama (19)
- George W. Bush (16)
- Ann Dunham (15)
- Craig Robinson (15)
- Joe Biden (13)**
- John McCain (8)**
- Kennedy (9)
- Sarkozy (8)
- Clinton (6)

show all

Organization (842 entities)

- Harvard (14)
- White House (22)**
- Congress (14)

- Exploitation for restricting the focus

CONT.

○ Some results

- Real-time NEM over snippets is feasible and yields about 1.2 entities per snippet
- NEM over contents is more time consuming, but mines much more entities
- String similarity between the query and the entity name does not improve entity ranking (in our setting)
- The top-10 entities derived from snippet mining are quite different from those derived from contents mining ($< 30\%$ Jaccard similarity)

○ Related Publications

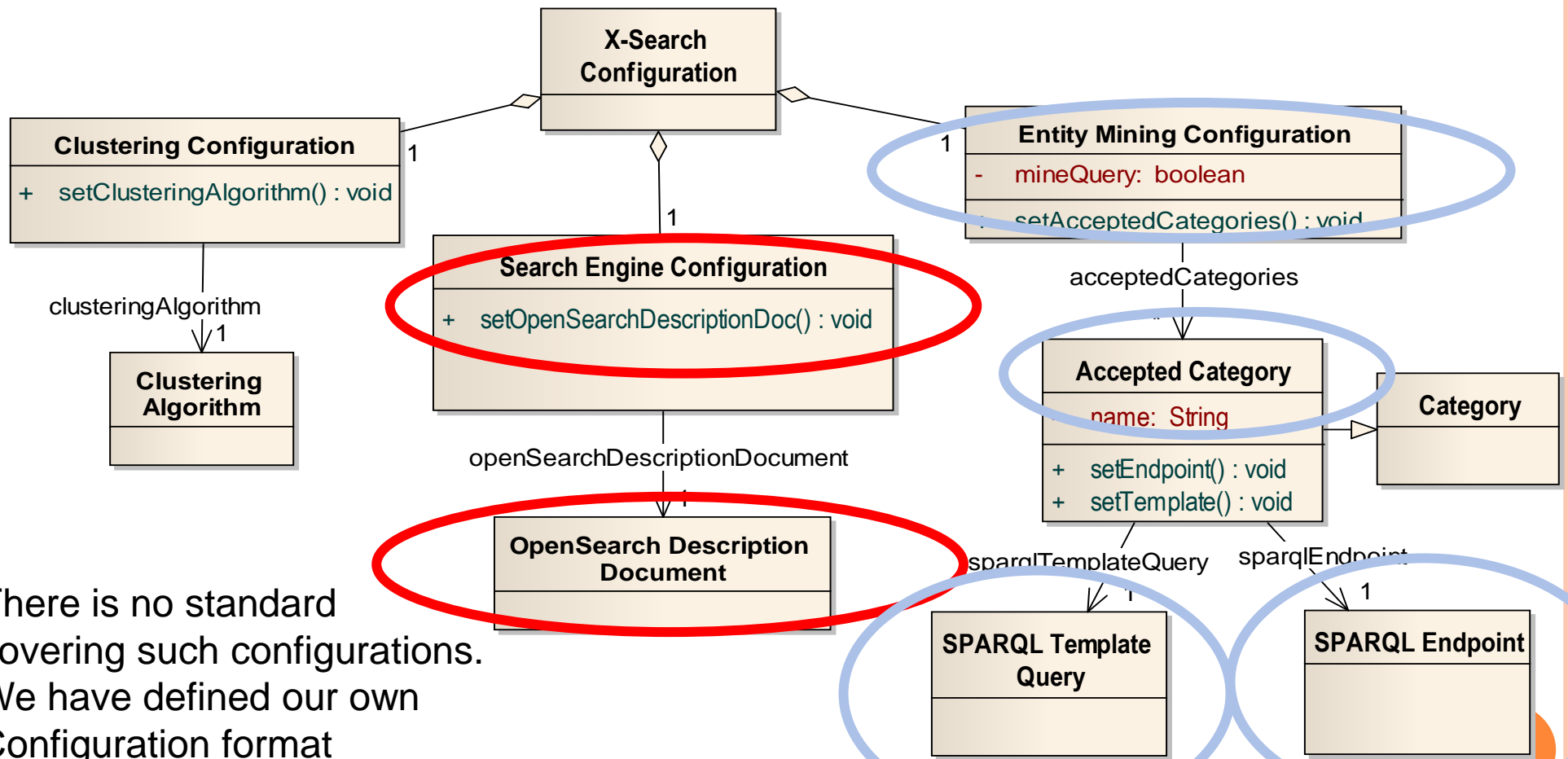
- [IRFC'12] P. Fafalios, I. Kitsos, Y. Marketakis, C. Baldassarre, M. Salampasis and Y. Tzitzikas, *Web Searching with Entity Mining at Query Time*, 5th Information Retrieval Facility Conference (IRFC 2012), Vienna, Austria, July 2012

MILESTONE 4. CONFIGURABILITY (AND LOD)


- Then we questioned ourselves:
 - ***why not allowing the user to configure himself the entities of interest by exploiting LOD (again in the context of entity mining of the search results)?***
- Outcome
 - X-ENS (e**X**plore **EN**tities in **S**earch)
- Related Publications
 - [SIGIR'13] P. Fafalios and Y. Tzitzikas, X-ENS: Semantic Enrichment of Web Search Results at Real-Time, 36th International ACM SIGIR Conference, Demo Paper, Dublin, Ireland, 28 July - 1 August 2013



XSEARCH-CONFIGURABILITY: THE CONCEPTUAL MODEL



There is no standard covering such configurations. We have defined our own Configuration format



200
results to mine

Tennis Player (39 entities)

- Roger Federer (14)
- Rafael Nadal (7)
- Novak Djokovic (5)
- Andy Roddick (4)
- serena williams (4)
- Maria Sharapova (3)
- Andy Murray (2)
- Tsvetana Pironkova (2)
- Urszula Radwanska (2)
- Vania King (2)


Tennis - ATP World Tour - Home

... photos, video, behind-the-scenes footage of the tennis player and tennis tournament statistics. It opens with ...

<http://www.atpworldtour.com/> - find its entities

Semantic Entity Enrichment (close)

Properties of: Andy Roddick

Description	Depiction
Andrew Stephen "Andy" Roddick (born Aug 30, 1982) is an American professional tennis player and a former World No. 1. He is..	
BirthPlace	BirthDate
Omaha, Nebraska	1982-08-30

Country (11 entities)

- India (8)
- Canada (3)

top hits

Entities

MILESTONE 5. PROFESSIONAL SEARCH SYSTEMS?

- Then we questioned ourselves:
 - ***why not applying and testing this in the context of a professional search system?***
- Outcome
 - Application in **patent search**. Missing relevant documents is unacceptable in patent search (*recall oriented search procedure*). Retrieval of all relevant documents is usually necessary
 - Patents contain plenty of named entities of various kinds
 - *Companies, Countries, Persons, Product types, Laws, etc*
 - Inclusion of **PerFedPat** System
 - In collaboration with Mike Salampasis

The screenshot displays the ezDL web application interface. The top navigation bar includes 'File', 'Tools', 'Perspectives', and 'Help'. The main interface is divided into several panels:

- Advanced Query:** Contains search filters for Full Text/Abstract, Title, Publication number, Application number, Priority number, Year, Applicant(s), Inventor(s), European Classification (ECLA), International Patent Classification (IPC), and U.S. Classification. The search term 'migraine' is entered in the Full Text/Abstract field.
- Entities Explorer:** Displays a list of entities related to the search, including International Patent Classification (IPC), Inventor, Applicant, European Classification (ECLA), Disease, Publication Year, Publication Country, Application Year, Application Country, Publication Number, Drug, and Chemical Substance. It also shows the International Patent Classification (IPC) with 111 entities.
- Cluster Explorer:** Displays a list of clusters related to the search, including text:migraine(50), migraine(46), behandlung(29), traitement(28), migräne(24), treatment(23), verwendung(13), treating(14), method(12), prevention(8), verfahren(8), vorbeugung(8), pain(7), douleur(6), and traiter(7).
- Results, Details:** Shows the search results, including the title 'Use of nadolol for inhibiting the onset of migraine', the applicant 'SQUIBB & SON & INC', and the IPCs [A61K31/22, A61K31/21, A61K31/135]. It also shows the publication number 'EP0350080-A2, 1'.

Two red circles highlight the 'Entities Explorer' and 'Cluster Explorer' panels. The 'Results, Details' panel shows the search results, including the title 'Use of nadolol for inhibiting the onset of migraine', the applicant 'SQUIBB & SON & INC', and the IPCs [A61K31/22, A61K31/21, A61K31/135]. It also shows the publication number 'EP0350080-A2, 1'.

PERFEDPAT (CONT)

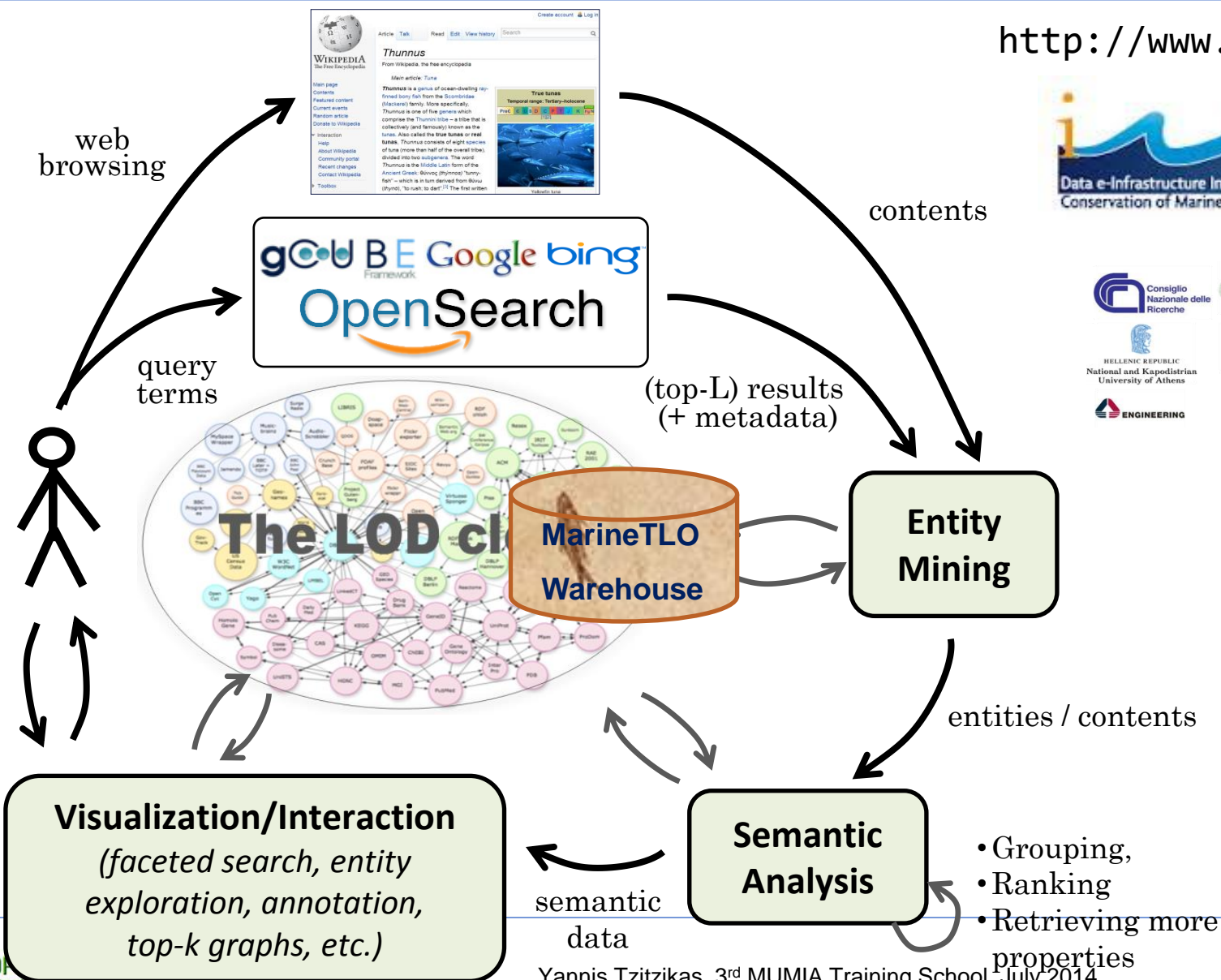
- The proposed functionality:
 - offers a tight integration of different search tools with the main retrieval engine,
 - connects the search results (i.e. patents) with data and knowledge,
 - can be exploited by any patent search system (i.e. it acts as a service over a ranked list of results)
 - The time that we have to pay is proportional to the number of the top results that we want to “explore” (≈ 1.5 sec / 100 results)
- Related Publications
 - P. Fafalios, M. Salampasis and Y. Tzitzikas, Exploratory Patent Search with Faceted Search and Configurable Entity Mining, 1st International Workshop on Integrating IR technologies for Professional Search, in conjunction with ECIR'13, Moscow, Russia, March 2013

MILESTONE 6

APPLYING IN THE CONTEXT OF AN INFRASTRUCTURE

- Then we questioned ourselves:
 - ***why not applying this in another domain of professional search in the context of a **real and operating EU research infrastructure**?***
- Outcome
 - X-Search in the context of the ongoing **iMarine Research Infrastructure** project

<http://www.i-marine.eu/>



EXAMPLE: X-SEARCH DEPLOYED IN AN OPERATIONAL RESEARCH INFRASTRUCTURE (2012-NOW)

Semantically Enriched Results

Query: tuna
In Collections: FIGIS

Mined Entities

- FAOCountry(24)
 - Republic of ... (1)
 - Viet nam(1)
 - Venezuela(2)
 - Yugoslavia(2)
 - Senegal(1)
- Species(8)
 - eastern Paci... (1)
 - yellowtail a... (1)
 - Ara(1)
 - pantropical... (1)
 - Indo-Pacific... (1)
- WaterAreas(3)
 - Mediterranea... (1)

Object Metadata

Thunnus albacares (Bonnaterre, 1788) - Fact sheet

Yellowfin **tuna**... (Venezuela), Ca bo Vang (Viet nam), **Tuna** zutoperka (Yugoslavia)... There are important yellowfin **tuna** fisheries throughout tropical and subtropical seas. The most... major surface fishing techniques for yellowfin **tuna** in the Pacific, even though this method

Textual Clustering

- Root(15)
 - fact sheet(27)
 - thunnus(8)
 - stenella(4)
 - linnaeus fac...(3)
 - axis(2)
 - tengraulis...(1)
 - dax fact s...(1)
 - purus lin...(1)

Semantic Entity Exploration

- URI: <http://www.fao.org/figis/flod/entities/codeentity/3e6d22db-1f06-437d-ac4a-9d3c8b895bf5> (open)
- Value: yellowtail amberjack

Semantic Entity Exploration

Properties of: Yellowtail_amberjack

Type	SameAs
Animal (open)	Seriola lalandi (open)
Thing (open)	
Species (open)	
FLODSpecies (open)	
Fish (open)	
CodeEntity (open)	
Eukaryote (open)	
Animal (open)	
Fish (open)	

Subject
Category:Fish of the Red Sea (open)
Category:Fish of the Indian Ocean (open)
Category:Seriola (open)

BinomialAuthority	Class
Georges Cuvier (open)	Actinopterygii (open)
Achille Valenciennes (open)	

Family
Carangidae (open)

Genus	Kingdom	Order	Phylum
Seriola (open)	Animal (open)	Perciformes (open)	Chordate (open)

Depiction	Thumbnail
Seriola lalandi.jpg (open)	200px-Seriola lalandi.jpg (open)

Object Metadata

Thunnus thynnus (L)

Atlantic bluefin **tuna**... fisheries. Off Sicily, no oceanic but seasonal tolerate a... for almost catches of northern b

Result of Entity Mining

Result of Textual Clustering

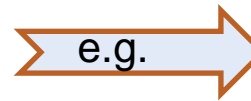
Yannis Tzitzikas, Panel@ExploreDB

MILESTONE 7 PARALLELIZATION

- Then we questioned ourselves:
 - ***Can we do the same over the full contents of the search hits? Downloading takes time. The processing also is expensive***
- Outcome
 - Investigation of how the task can be partitioned to several machines
 - A MapReduce-based parallelization of downloading and entity mining task for exploiting the resources of cloud
- Key points
 - How to tackle the uncertainty (of hits size) for achieving load balancing and reaching the ideal speedup
 - Two processes for carrying out this task

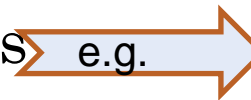
LEVELS OF FUNCTIONALITY

- **L0**: Categories + Entities



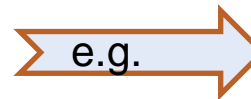
Location
Athens
Greece

- **L1**: L0 + Count Info for Entities



Location
Athens (3)
Greece (5)

- **L2**: L1 + Rank Entities



Location
Greece (5)
Athens (3)



Ranked
Entities

- **L3**: L2 + Doc List for each entity

- **L4**: L3 + Semantic Enrichment

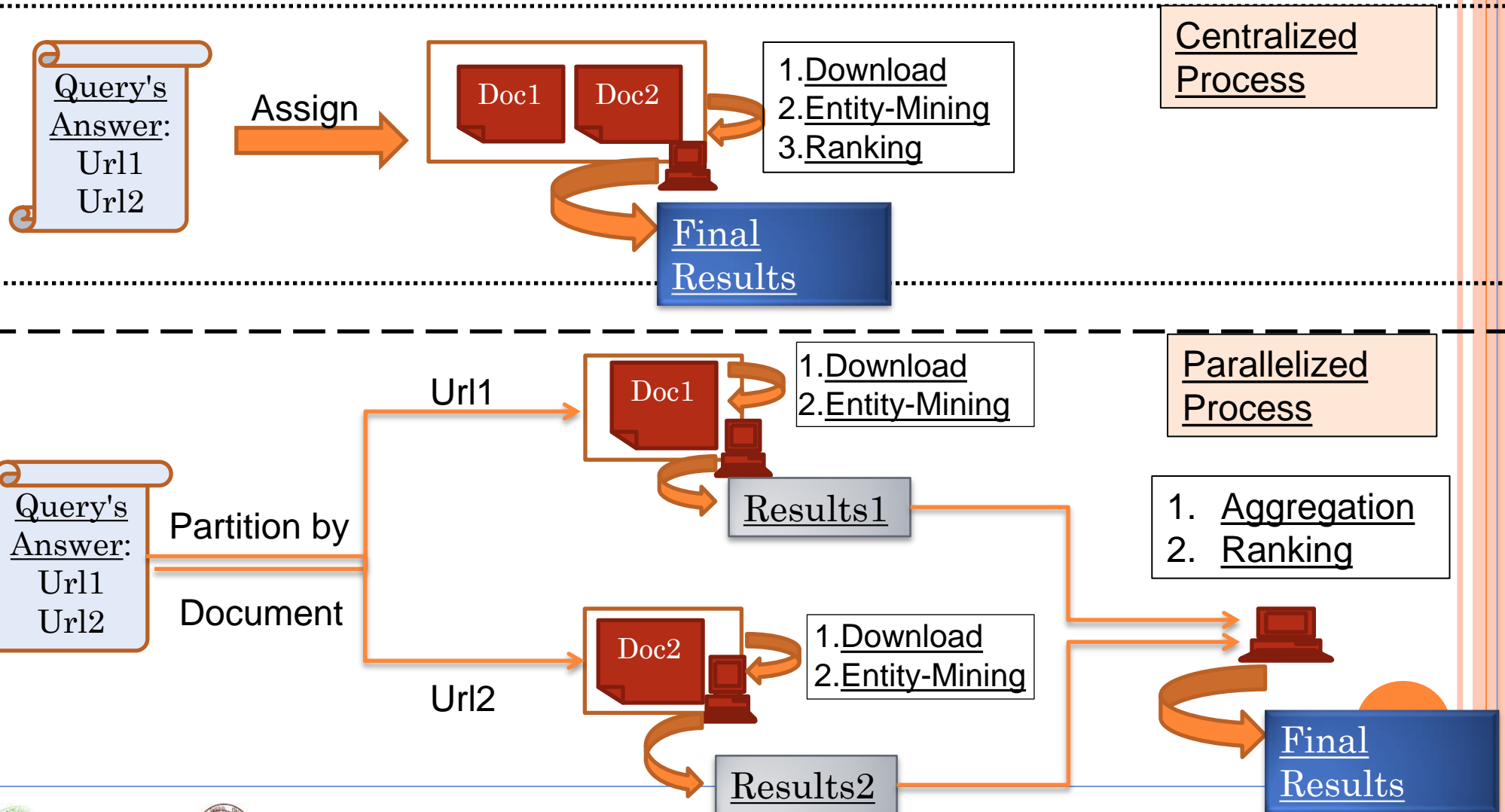


Location
Greece (5) LOD_URL
Athens (3) LOD_URL

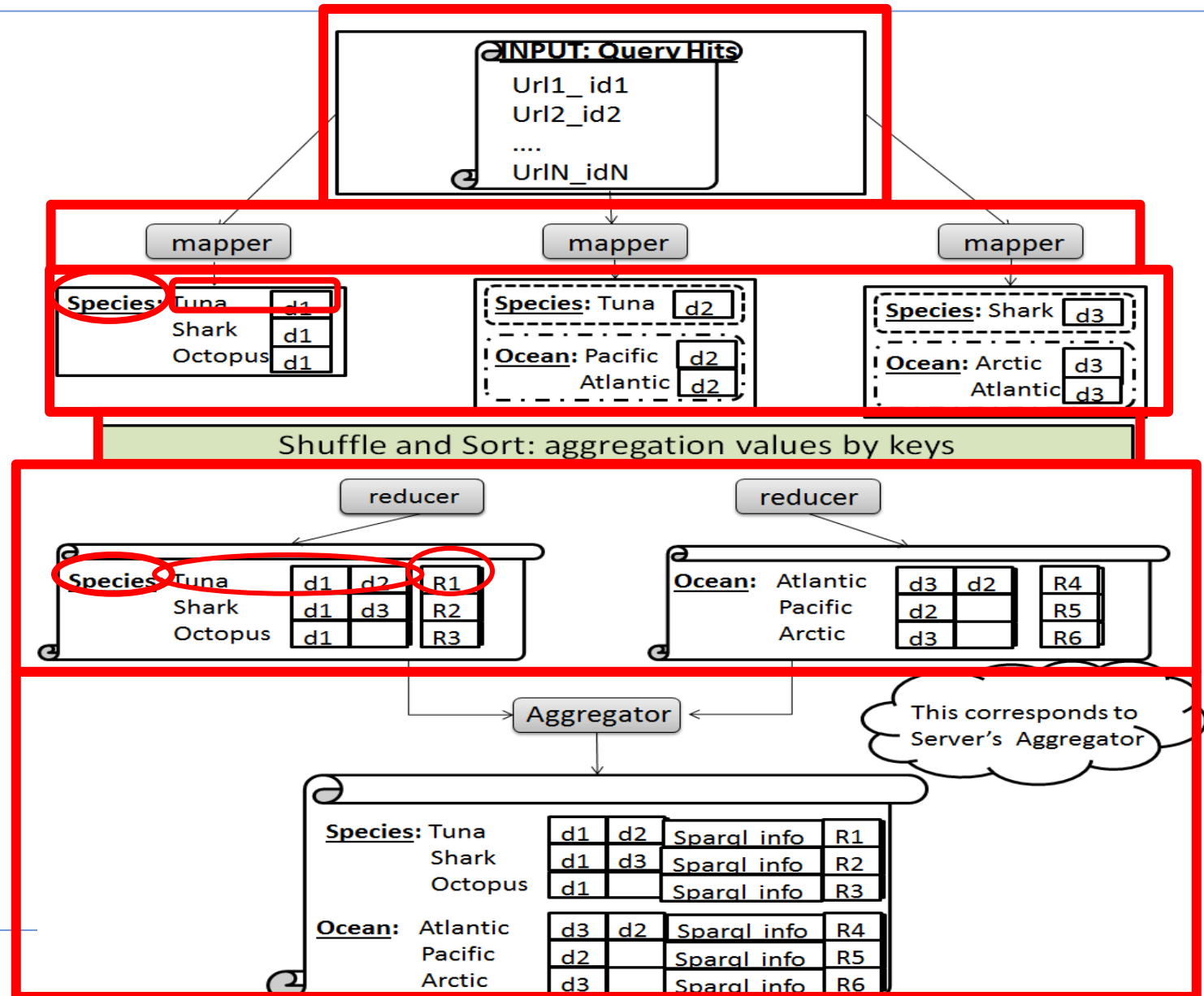


THE PARALLELIZATION PROCEDURE

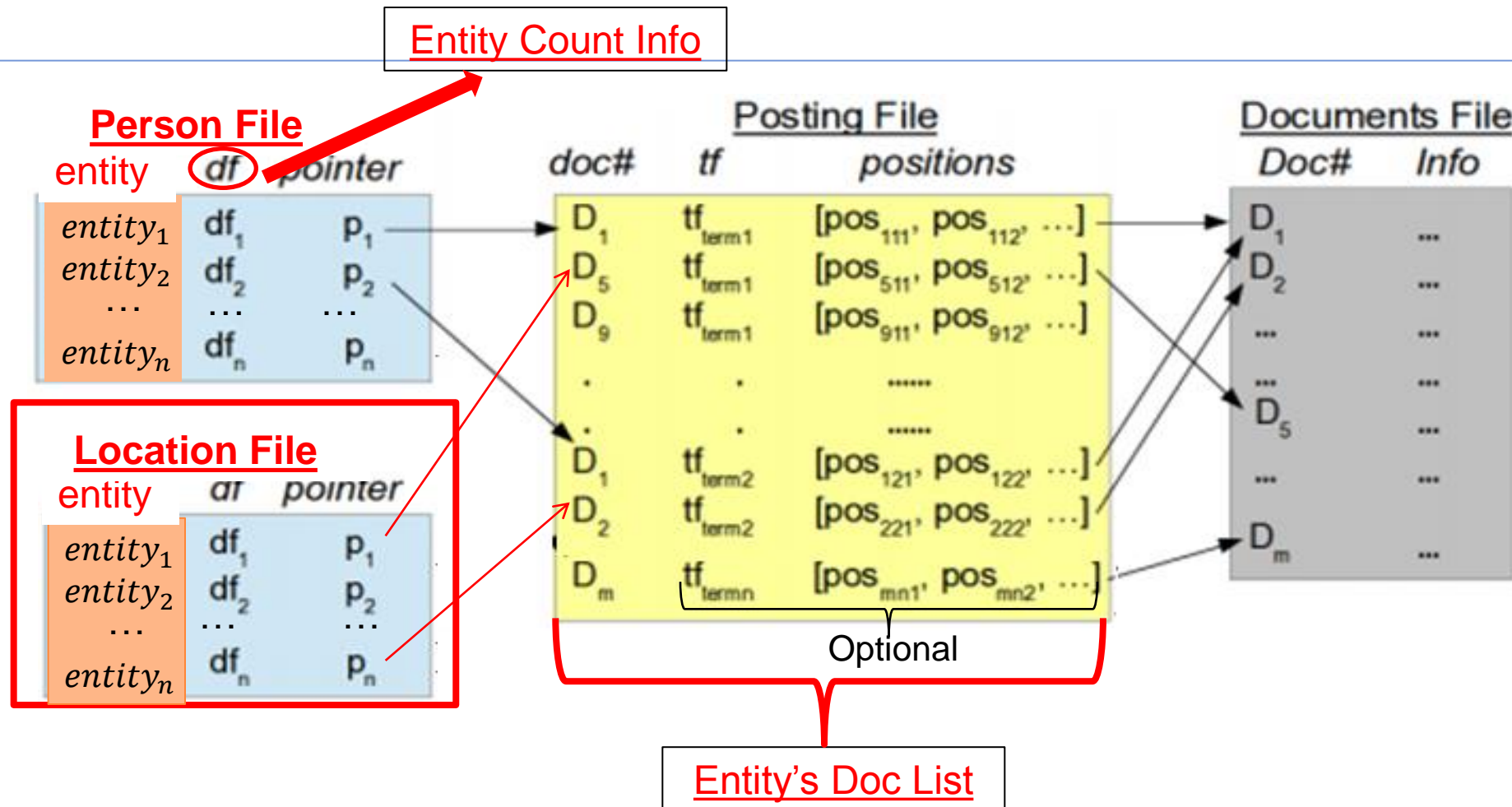
- Main idea: *partition by documents (hit)*



EXAMPLE OF DISTRIBUTED NEM PROCESSING USING MAPREDUCE



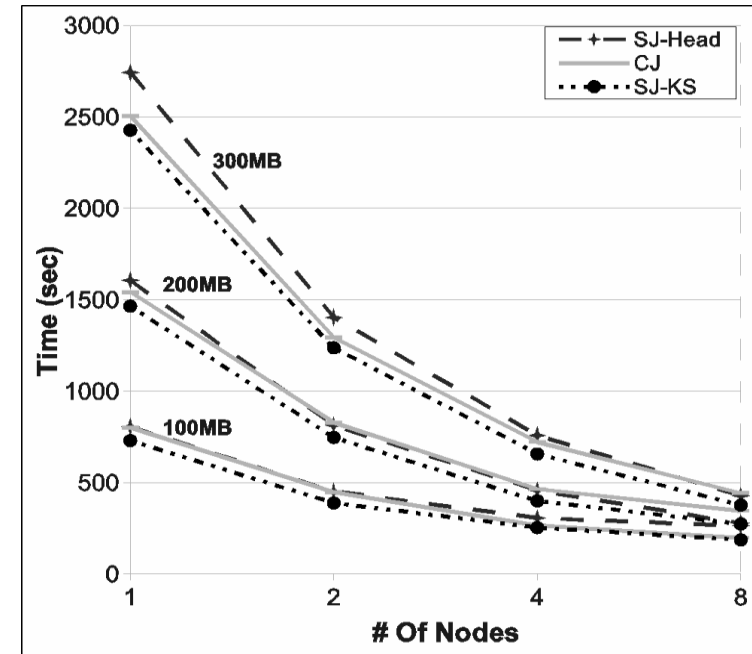
AN ANALOGY TO INVERTED FILES



- Note that our task is much more CPU and memory intensive
- Requires document's downloading

EXPLOITING MORE THAN ONE MACHINES [J. DAPD 2014]

- Two models of evaluation were investigated
- A thorough evaluation of the parameters that affect performance were conducted
- We reached a speedup close to the ideal (according to Amdahl's law)!



○ Related Publications

- I. Kitsos, K. Magoutis and Y. Tzitzikas, Scalable Entity-based Summarization of Web Search Results using MapReduce, Journal on Distributed and Parallel Databases (DAPD), 32(3), 2014



WHAT NEW CAR SHOULD I BUY?



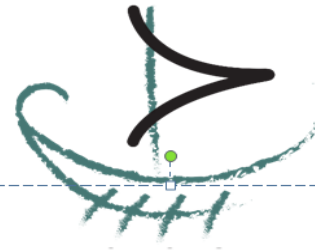
MILESTONE 8 PREFERENCES

- Then we questioned ourselves:
 - ***What about the ordering of facets, terms and objects? Should the user only restrict the focus? Why not allowing the user to change the order based on his/her preferences?***
- Outcome
 - A framework for preferences over multi-dimensional and hierarchical information spaces
 - An extension of the interaction model of faceted search with preferences
 - The Hippalus system that realizes it

SYSTEM: HIPPALUS (2013)




Hippalus: Preference-enriched Faceted Exploration



P. Papadakos^{1,2} and Y. Tzitzikas^{1,2}

- Allows faceted browsing and also supports **Preferences**
 - User actions specify the ranking of the information space
 - Gradual preference specification
 - Automatic resolution of conflicts
 - Different preference composition modes
 - E.g. if the user defines the desired ordering wrt each dimension, then the first block of the ranked objects is the skyline

HIPPALUS: INTERACTION OVER A KB OF 50 CARS

 Hippalus

Preference-Enriched Faceted Exploratory System

Facets

+ Acceleration (43)

+ Body_Type (50)

+ Doors (50)

+ Drive_System (50)

+ Engine_Power (50)

+ Engine_Torque (48)

+ Engine_Volume (50)

+ Fuel_Cons_city (43)

+ Fuel_Cons_highway (43)

+ Fuel_Tank (46)

+ Fuel_Type (50)

+ Gears (50)

+ ID (50)

+ Manufacturer (50)

+ Model (50)

+ Price (50)

+ Speed (47)

+ Transmission (50)

+ Trunk (40)

+ Vehicle_Type (50)

+ Weight_Empty (39)

+ Year (50)

In focus: 50 objects Number of buckets: 1

1

- Alfa-Romeo-8C-ID3
- Alfa-Romeo-Brera-ID1
- Alfa-Romeo-MiTo-ID2
- Audi-A3-ID4
- Audi-S8-ID5
- Audi-TT-ID6
- BMW-1-ID7
- BMW-3-ID8
- BMW-7-ID9
- Citroen-C1-ID10
- Citroen-C3-ID11
- Fiat-Bravo-ID12
- Fiat-Punto-ID13
- Ford-Fiesta-ID14
- Ford-Ka-ID15
- Hyundai-i10-ID16
- Hyundai-i30-ID17
- Kia-Ceed-ID18
- Lancia-Delta-ID19
- Mazda-3-ID20
- Mazda-RX-8-ID21
- Mercedes-Benz-A-ID22
- Mercedes-Benz-C-ID23
- Mercedes-Benz-C-ID25
- Mercedes-Benz-SL-ID24
- Mitsubishi-Colt-ID26
- Mitsubishi-X-Trail-ID27
- Nissan-Micra-ID28
- Nissan-Navara-ID29
- Opel-Astra-ID30

Preference Actions

Clear

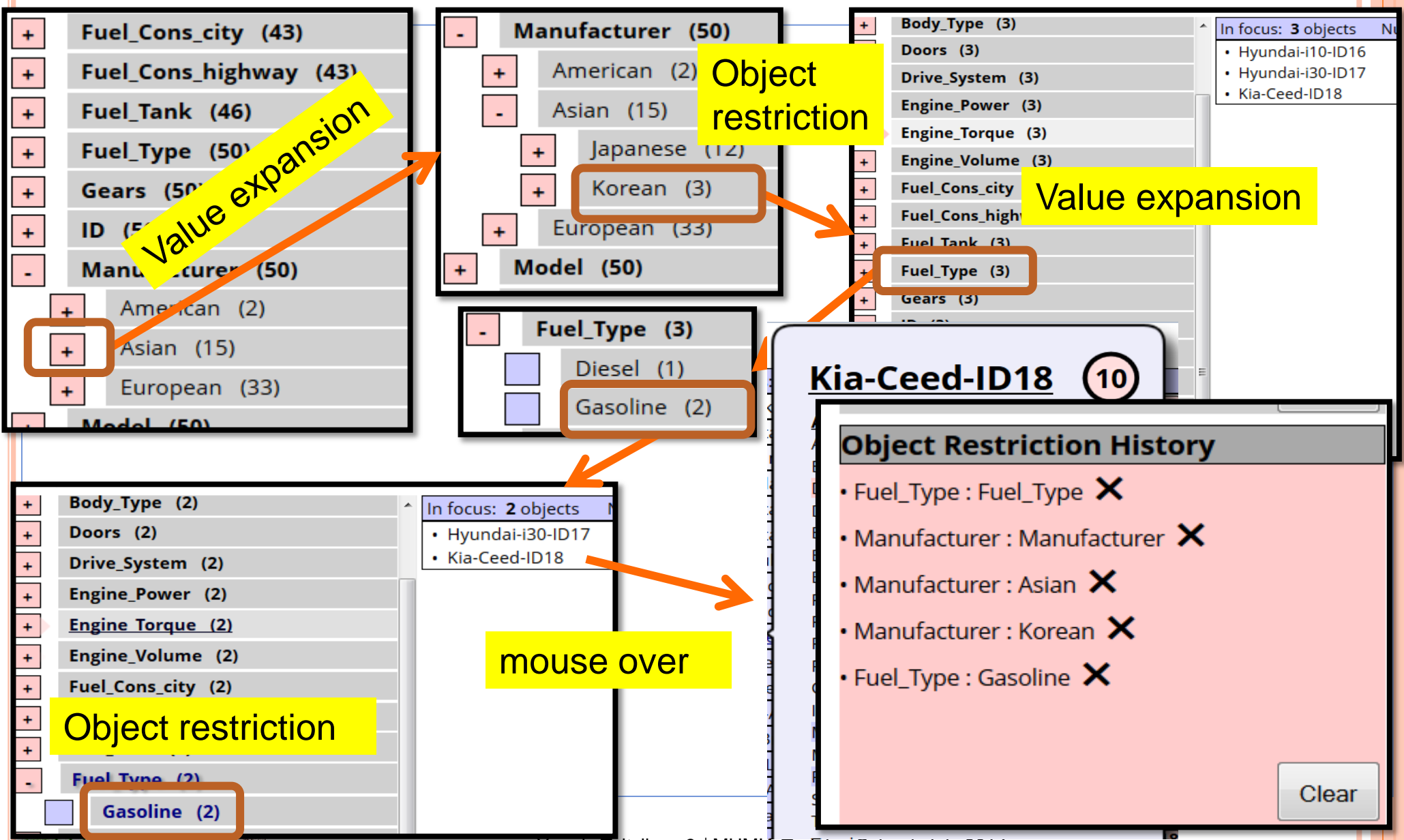
Composition: Combination

Interesting Objects

Clear

Object Restriction History

HIPPALUS: FDT INTERACTIONS



HIPPALUS: PREFERENCE ACTIONS

Cars ordered with priority on manufacturer

The screenshot displays the HIPPALUS interface with several components:

- Model (50):** A dropdown menu showing car models like Fiat-Punto, Hyundai-i30, etc.
- Preference Actions:** A section with a text input field containing "1) objects order: prefer term Manufacturer...Korean to Manufacturer...European" and a close button (X).
- Facets:** A list of facets for filtering results:
 - Acceleration (8):** +
 - Body_Type (8):** +
 - Doors (8):** - (highlighted with a yellow box and labeled "Object restriction")
 - Drive_System (8):** +
 - Engine_Power (8):** +
 - Engine_Torque (8):** +
- Level 2: Price_Euros:** A dropdown menu showing price ranges like "15538 (1)".
- Results:** A list of car models with their IDs and bucket numbers, ordered by manufacturer priority. The list is divided into two sections: "In focus: 50 objects" (top) and "In focus: 8 objects" (bottom).

Section	Object	Bucket
In focus: 50 objects	Hyundai-i30-ID17	1
	Hyundai-i10-ID16	2
	Kia-Cee-ID18	3
In focus: 8 objects	Peugeot-207-ID33	1
	BMW-1-ID7	2
	BMW-3-ID8	3
	Alfa-Romeo-Brera-ID1	4
	Audi-TT-ID6	5
	Saab-9-3-ID36	6
	Alfa-Romeo-8C-ID3	7
	Mercedes-Benz-SL-ID24	8
Continuation of results	Peugeot-207-ID33	17
	Lancia-Delta-ID19	18
	Mercedes-Benz-A-ID22	19
	Volkswagen-Scirocco-ID48	20
	Audi-A3-ID4	21
	Skoda-Octavia-ID39	22
	Volvo-C30-ID50	23

CONT.

Evaluation

- Over a collection of 50 cars
- With the preference-enriched, all users completed successfully all tasks leading to
 - ideal scores for Precision and Recall!
 - on average in 1/3 of the time!
 - on average with 1/3 of the actions!
- None of the users completed successfully all tasks with the plain interface
- All users (either plain or experts) preferred the preference-enriched interface

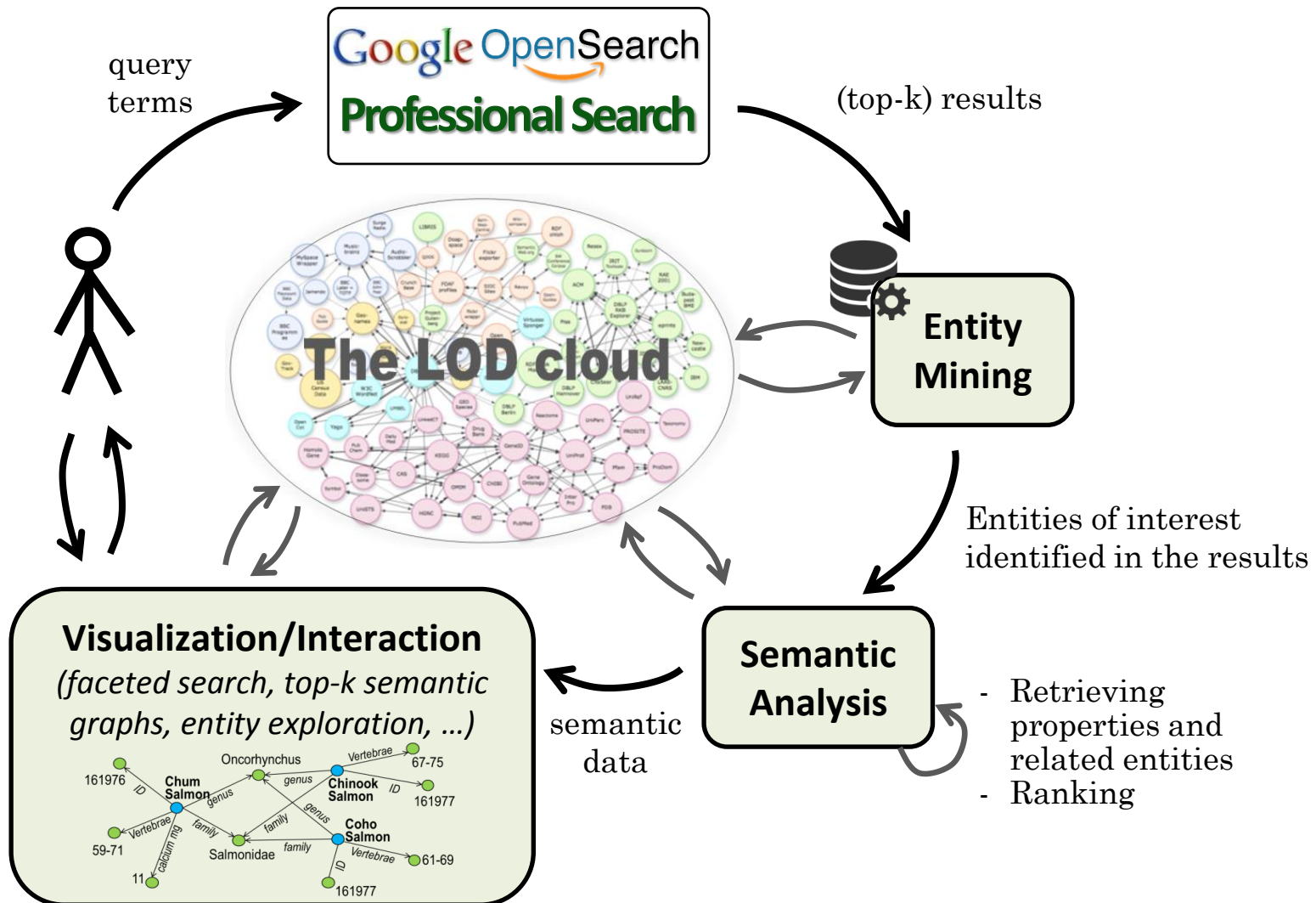
More information in the publications

- [J. FI 13] Yannis Tzitzikas and Panagiotis Papadakos. Interactive Exploration of Multidimensional and Hierarchical Information Spaces with Real-Time Preference Elicitation. In *Journal FUNDAMENTA INFORMATICA*, 2013
- [ExploreDB'14] Panagiotis Papadakos, Yannis Tzitzikas: Hippalus: Preference-enriched Faceted Exploration. In EDBT/ICDT Workshops 2014

MILESTONE 9 FROM DIMENSIONS TO GRAPHS

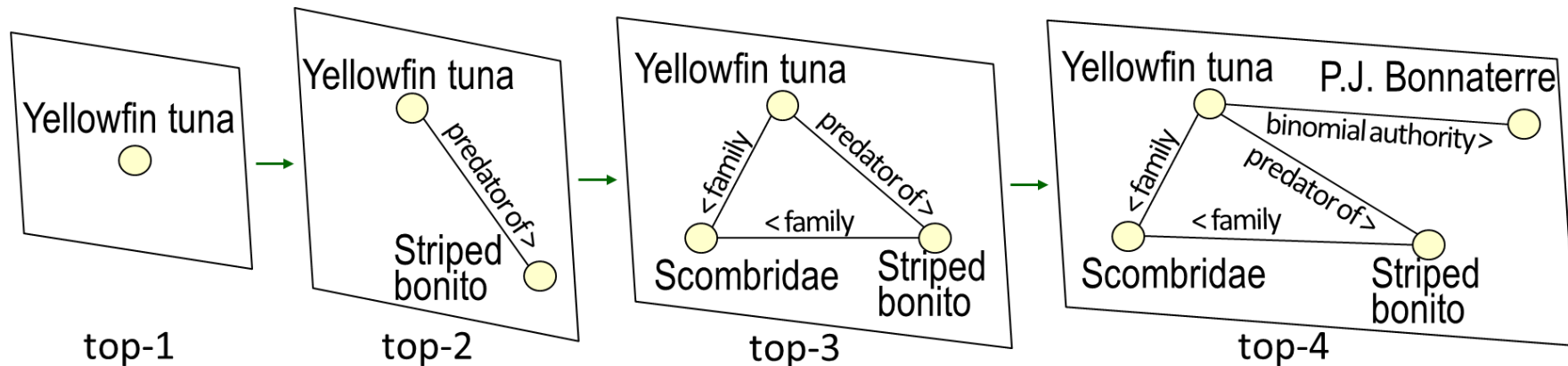
- Then we questioned ourselves:
 - ***So far we have seen services for getting and exploiting multidimensional spaces over the search results. But what if the notion of dimension cannot be defined, or in case there are too many? What can be done without having to configure entity types?***
- Outcome
 - A semantic post-processing of results that does not yield a multidimensional space but a **graph**.
- Challenges
 - Graph construction and exploitation for identifying the important (useful for the user) nodes and relationships

TOP-K SEMANTIC GRAPHS



TOP-K SEMANTIC GRAPHS (CONT.)

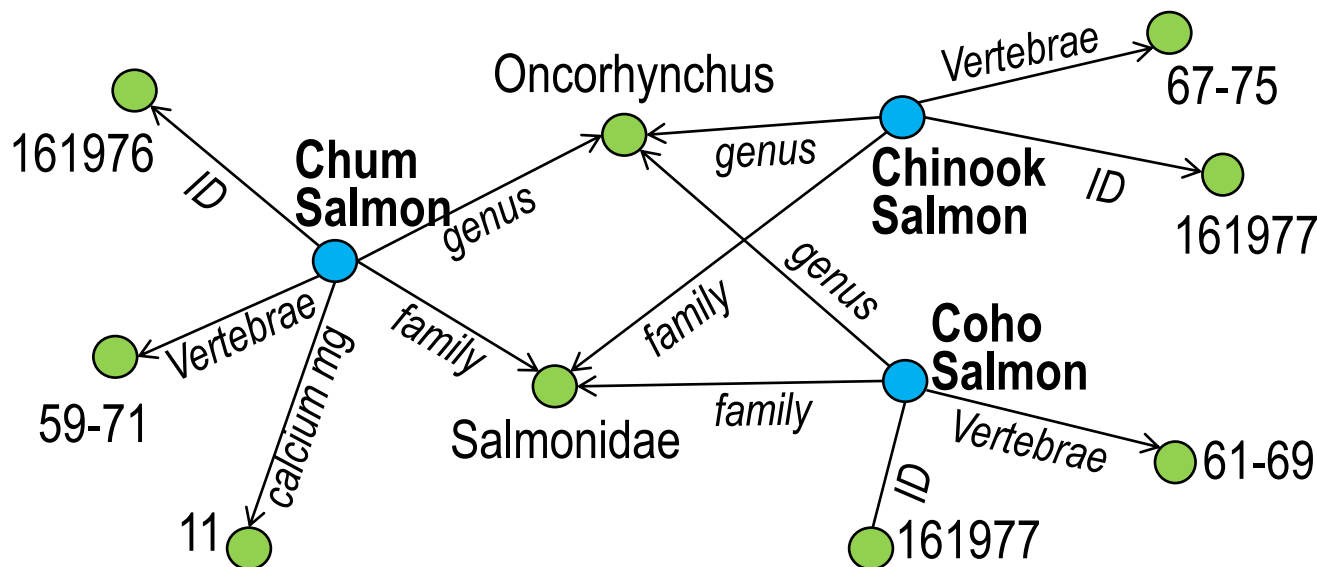
- The system can return the top-K graph for any K from 1 to number of nodes produced
 - Vertices: the K most highly ranked nodes
 - Edges: the edges that connect the K most highly ranked nodes
- The user is free to increase or reduce the value of K
- Example (from a real domain):



TOP-K SEMANTIC GRAPH

This graph

- can **complement** the query answer with useful information regarding the connectivity of the identified entities
- allows users to **instantly inspect** information that may lie in different places and that may be laborious and time consuming to locate
- provides useful information about the **context** of the identified entities
- allows the users to get a **more sophisticated overview** and to make better sense of the results



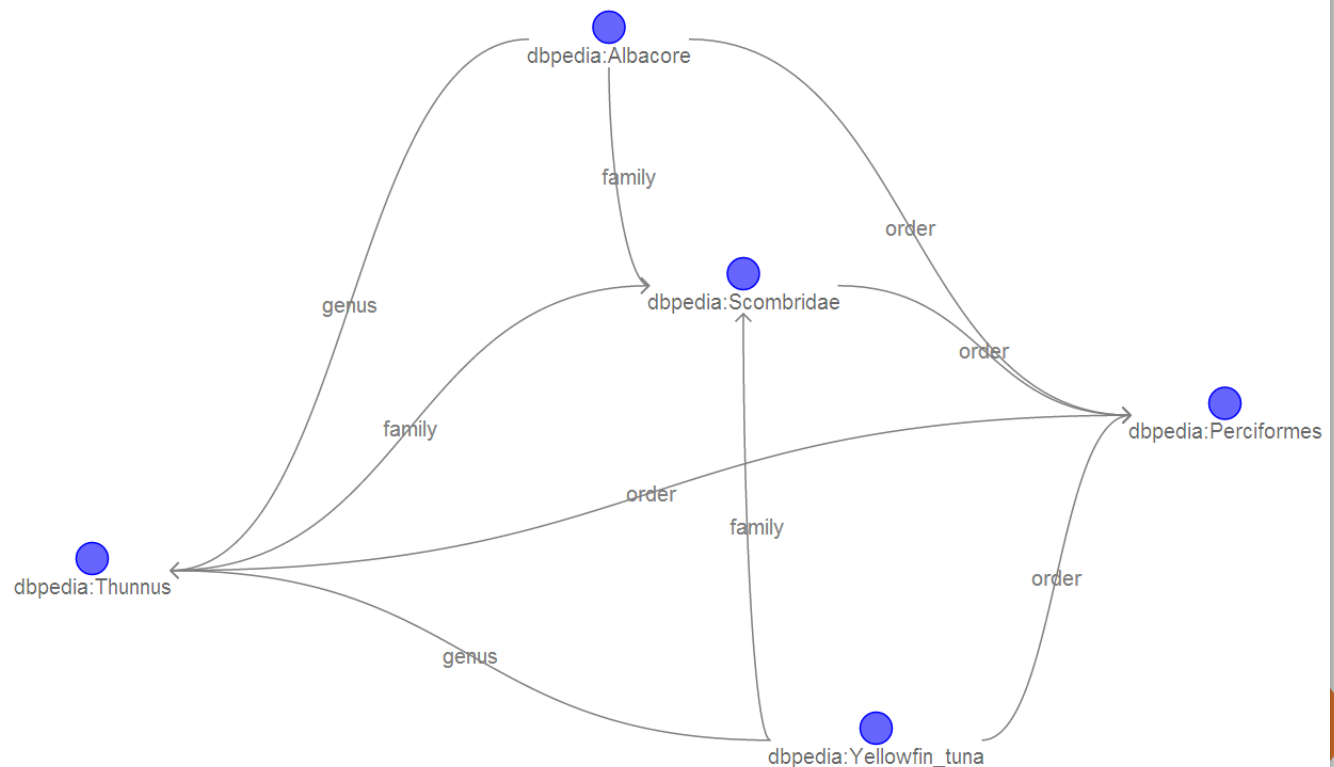
PROTOTYPE (2014)

<http://139.91.183.72/x-ens-2/>

TOP-5 LIST

1. dbpedia:Yellowfin_tuna
2. dbpedia:Perciformes
3. dbpedia:Scombridae
4. dbpedia:Albacore
5. dbpedia:Thunnus

TOP-5 GRAPH



CONT.

○ Evaluation (main results)

- Usefulness – Survey for the marine domain
 - *The majority of participants believe that the appearance of a graph of semantic information related to the search results can help them during an exploratory search process*
- Effectiveness – Comparative evaluation of ranking schemes:
 - The proposed PageRank-based ranking scheme produces more preferred ranking compared to other link analysis-based algorithms
- Efficiency – Case study over online DBpedia
 - The exploitation of LOD can be supported at query-time
 - For up to 100 detected entities we can offer the proposed functionality at real-time, even if we query an online KB (like DBpedia)
- The major bottleneck is the reliability and performance of online SPARQL endpoints
 - We expect this limitation to get overcome in the near future
 - In the meanwhile, we can use caching / indexing / dedicated warehouses / distributed infrastructure

○ Related Publications:

- P. Fafalios and Y. Tzitzikas, Post-Analysis of Keyword-based Search Results using Entity Mining, Linked Data and Link Analysis at Query Time, IEEE 8th International Conference on Semantic Computing (ICSC'14), Newport Beach, California, USA, June 2014

5. SYNOPSIS AND DISCUSSION (30')

135

Synopsis
Challenges

SYNOPSIS

- We have discussed information needs of exploratory nature
- We have seen the basics of faceted exploration
- We have seen an overview of semantic technologies
- We have seen ways to exploit semantic datasets during searching (exploratory search) with emphasis on doing this at search time

○ So what's next ?

DIRECTIONS & CHALLENGES

○ Ubiquity.

- Faceted browsing of search results and gradual restriction should be possible for any kind of query, for any domain and with no predetermined facets.
- In other words, methods that bypass the need for explicit configuration (regarding facets, entities types, LOD sources) are required.
 - This is why we currently study M9
 - Then we also have to define the interaction model over such graphs

DIRECTIONS & CHALLENGES (2)

○ **Fusion of Structured and Unstructured Content.**

- The exploitation of LOD in the exploratory search process is promising, e.g. for Named Entity Recognition and disambiguation. However, the fusion of structured and unstructured content requires more work.
 - E.g. should (and if yes how) the results of the semantic post processing should affect the search hits (a kind of semantic feedback)

DIRECTIONS & CHALLENGES (3)

○ User Control.

- Explicit, user-provided and controllable preference management is beneficial for supporting a transparent decision making process.
- We believe that the framework supported by the Hippalus system is a first step towards this direction.

DIRECTIONS & CHALLENGES (4)

○ Evaluation.

- We need easy to follow methods for evaluating the effectiveness of exploratory search methods, and easily reproducible evaluation results.
- Although the classical IR has well established methodologies for evaluation, things are not so clear and straightforward in interactive IR (IIR).



QUESTIONS AND EXERCISES

142

144

REFERENCES AND LINKS

Organized on the basis of the case studies

REFERENCES AND LINKS

- Faceted Search and Dynamic Taxonomies
 - Sacco, Giovanni Maria; Tzitzikas, Yannis (Eds.), *Dynamic Taxonomies and Faceted Search: Theory, Practice, and Experience*, Series: The Information Retrieval Series , Vol. 25, 2009
 - Giovanni Maria Sacco: **Dynamic Taxonomies: A Model for Large Information Bases.** IEEE Trans. Knowl. Data Eng. **12**(3): 468-479 (2000)
- Browsing Approach for (plain and fuzzy) RDF
 - Sébastien Ferré's Publications
 - Nikos Manolis, Yannis Tzitzikas: **Interactive Exploration of Fuzzy RDF Knowledge Bases.** ESWC (1) 2011: 1-16



REFERENCES (M1)

○ Mitos-related (**Web Searching and Faceted Search**)

- [PCI'08] Panagiotis Papadakos, Yannis Theoharis, Yannis Marketakis, Nikos Armenatzoglou and Yannis Tzitzikas. Mitos: Design and Evaluation of a DBMS-based Web Search Engine. In PCI 2008.
- [FIND'08] Yannis Tzitzikas, Nikos Armenatzoglou, and Panagiotis Papadakos. FleXplorer: A Framework for Providing Faceted and Dynamic Taxonomy-based Information Exploration. In DEXA/FIND 2008.
- [ECDL'09] Panagiotis Papadakos, Stella Kopidaki, Nikos Armenatzoglou and Yannis Tzitzikas. Exploratory Web Searching with Dynamic Taxonomies and Results Clustering. In ECDL 2009.
- [WISE'09] Stella Kopidaki, Panagiotis Papadakos, and Yannis Tzitzikas. STC+ and NM-STC: Two novel online results clustering methods for web searching. In WISE 2009.
- [J. KAIS 2012] Panagiotis Papadakos, Stella Kopidaki, Nikos Armenatzoglou and Yannis Tzitzikas. On exploiting Static and Dynamically mined Metadata for Exploratory Web Searching. In KAIS Journal 2012.

REFERENCES (M2)

○ Instance Overview Search-related

- Pavlos Fafalios, Yannis Tzitzikas: Exploiting Available Memory and Disk for Scalable Instant Overview Search. WISE 2011: 101-115
- P. Fafalios, I. Kitsos and Y. Tzitzikas, Scalable, Flexible and Generic Instant Overview Search, Proceedings of the 21st International Conference on World Wide Web (demo paper), WWW 2012, Lyon, France, April 2012

○ Links to Online Prototypes

- <http://www.ics.forth.gr/ios>

REFERENCES (M3+M4)

- Semantic-Post Processing of Search Results with Entity Mining and LOD
 - P. Fafalios, I. Kitsos, Y. Marketakis, C. Baldassarre, M. Salampasis and Y. Tzitzikas, Web Searching with Entity Mining at Query Time, Proceedings of the 5th Information Retrieval Facility Conference, IRF 2012, Vienna, July 2012
 - Pavlos Fafalios, Yannis Tzitzikas: X-ENS: semantic enrichment of web search results at real-time. SIGIR 2013, Dublin, Ireland, 28 July - 1 August 2013
- Links to Online Prototypes
 - <http://62.217.127.118/x-ens/>

REFERENCES (M5)

- Semantic-Post Processing of Search Results with Entity Mining in **Patent Search**
 - P. Fafalios and Y. Tzitzikas, “Exploratory Professional Search through Semantic Post-Analysis of Search Results”, In “Professional Search in the Modern World”, Lecture Notes in Computer Science (LNCS), Springer, 2014 (accepted for publication as a State-of-the-Art volume in LNCS)
 - P. Fafalios, M. Salampasis and Y. Tzitzikas, *Exploratory Patent Search with Faceted Search and Configurable Entity Mining*, 1st International Workshop on Integrating IR technologies for Professional Search, in conjunction with ECIR'13, Moscow, Russia, March 2013
- Links to Online Prototypes
 - <http://www.perfedpat.eu/index.php/download-perfedpat>
 - <http://139.91.183.72/x-search-metadata-groupings/>

REFERENCES (M6)

- Semantic-Post Processing of Search Results with Entity Mining in **Marine Search**
 - P. Fafalios and Y. Tzitzikas, “Exploratory Professional Search through Semantic Post-Analysis of Search Results”, In “Professional Search in the Modern World”, Lecture Notes in Computer Science (LNCS), Springer, 2014 (*accepted for publication as a State-of-the-Art volume in LNCS*).
- Links to Online Prototypes
 - <http://62.217.127.118/x-search/>
 - <http://62.217.127.118/x-search-fao/>

REFERENCES (M7)

- Semantic-Post Processing of Search Results with Entity Mining in **over the Cloud**
 - I. Kitsos, K. Magoutis and Y. Tzitzikas, Scalable Entity-based Summarization of Web Search Results using MapReduce, Journal on Distributed and Parallel Databases (DAPD), 2014

REFERENCES (M8)

- Extending Faceted Search with **Preferences**
 - [WISE'12] Panagiotis Papadakos, Yannis Tzitzikas and Dafni Zafeiri: An Interactive Exploratory System with Real-Time Preference Elicitation. Demo paper, In *Web Information Systems Engineering - WISE 2012*, Pafos, Cyprus, Volume 7651, p. 817-820, November 2012.
 - [J. FI 13] Yannis Tzitzikas and Panagiotis Papadakos. Interactive Exploration of Multidimensional and Hierarchical Information Spaces with Real-Time Preference Elicitation. In *Journal FUNDAMENTA INFORMATICA*, Volume 122, Issue 4, pp 357-399, 2013.
 - [ExploreDB'14] Panagiotis Papadakos, Yannis Tzitzikas: Hippalus: Preference-enriched Faceted Exploration. EDBT/ICDT Workshops 2014: 167-172
- Video Demonstration available at <http://www.youtube.com/watch?v=Cah-z7KmlXc>
- Links to Online Prototypes <http://www.ics.forth.gr/isl/Hippalus>

REFERENCES (M9)

- **Graph-resulting** Semantic Post-Processing of Search
 - P. Fafalios and Y. Tzitzikas, Post-Analysis of Keyword-based Search Results using Entity Mining, Linked Data and Link Analysis at Query Time, IEEE 8th International Conference on Semantic Computing (ICSC'14), Newport Beach, California, USA, June 2014
- Prototype (configured for the marine domain)
 - <http://139.91.183.72/x-ens-2/>

ACKNOWLEDGEMENTS

- Apart from Panagiotis Papadakos and Pavlos Fafalios, other students have also contributed to the “story” that was presented:
 - Nikos Armenatzoglou
 - Stella Kopidani
 - Nikos Manolis
- These slides include material from
 - The LOD integration scenario was taken from
 - http://www.csee.umbc.edu/courses/graduate/691/spring14/01/notes/01_introduction/01motivating_example.pptx

Thanks for you attention