



# EXPLORATORY SEARCH

Research Prototypes for Information Systems Laboratory  
FORTH-ICS

Contact person: Yannis Tzitzikas

Last modification: August 1, 2016

copyright: FORTH-ICS

1



# OUTLINE

- Context
- A few words about Exploratory Search
- An overview of research prototypes from ISL
  - m1: faceted search using static and dynamic metadata
  - m2: instant overview search
  - m3: entity mining over documents
  - m4: configurability of entity mining
  - m5: application in professional search
  - m6: application in an research infrastructure
  - m7: preference-enriched faceted search
  - m8: top-k semantic graphs

# CONTEXT

- Users should enjoy their data

# EXPLORATORY SEARCH



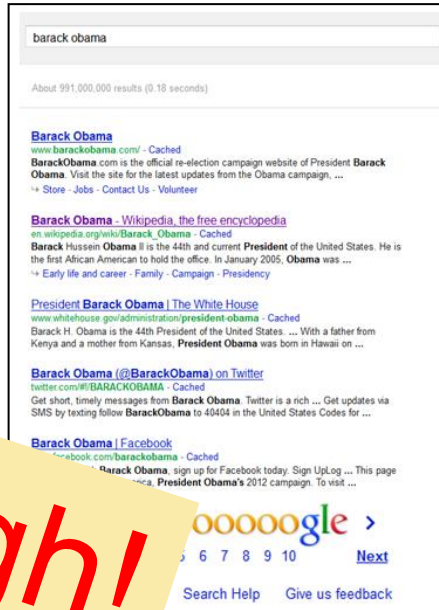
## *Wikipedia:*

“**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.*

**F**

rom 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.

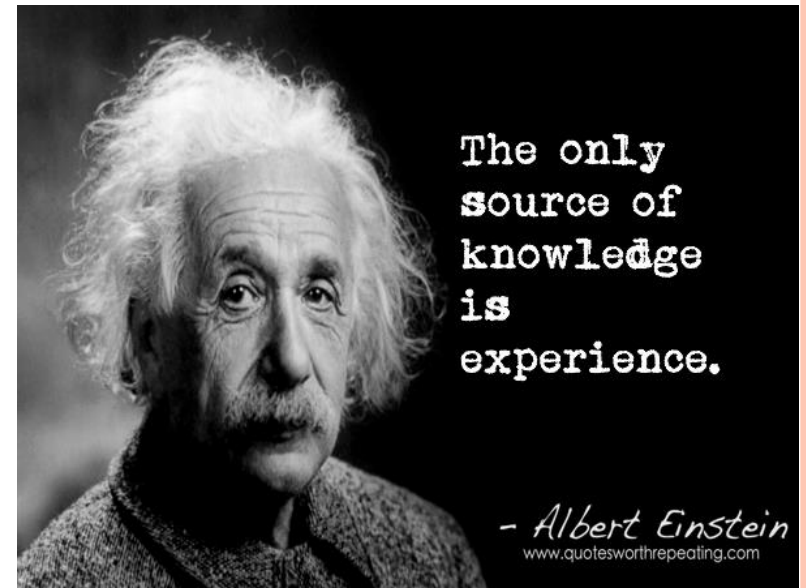
# SOME COMMON REQUIREMENTS FOR EFFECTIVE EXPLORATORY SEARCH



- 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

# RESEARCH PROTOTYPES FROM ISL RELATED TO EXPLORATORY SEARCH

- They are presented like a “story” organized in milestones that correspond to activities of ISL (Information Systems Laboratory) of FORTH-ICS (2009-now)





# MILESTONE 1. THE MITOS WSE (2009)

- MITOS is WSE built from scratch. 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)

The screenshot shows the MITOS search interface with a search bar and 'Advanced Search' button. Below the search bar, it displays 'Faceted taxonomies results' and a list of search results. The results are categorized into faceted taxonomies: 'By clustering', 'By domain', and 'By filetype'. The 'By clustering' faceted taxonomy is highlighted with a red dot and a blue arrow pointing to a larger view of it.

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

**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)

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

# 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 (ms)

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

8558 initial results

## 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](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](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 -

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)

By domain

- ⊕ gr (606)

By date

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

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.

[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  
[http://www.ics.forth.gr/isl/services.html](#) - 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 in RDF/XML

Results per page 10

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.html](http://www.ics.forth.gr/news/ian_cernockv_lecture.html) - 1241420696000 - 20KB

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



# CONT.

## ○ Evaluation with Users (main results) :

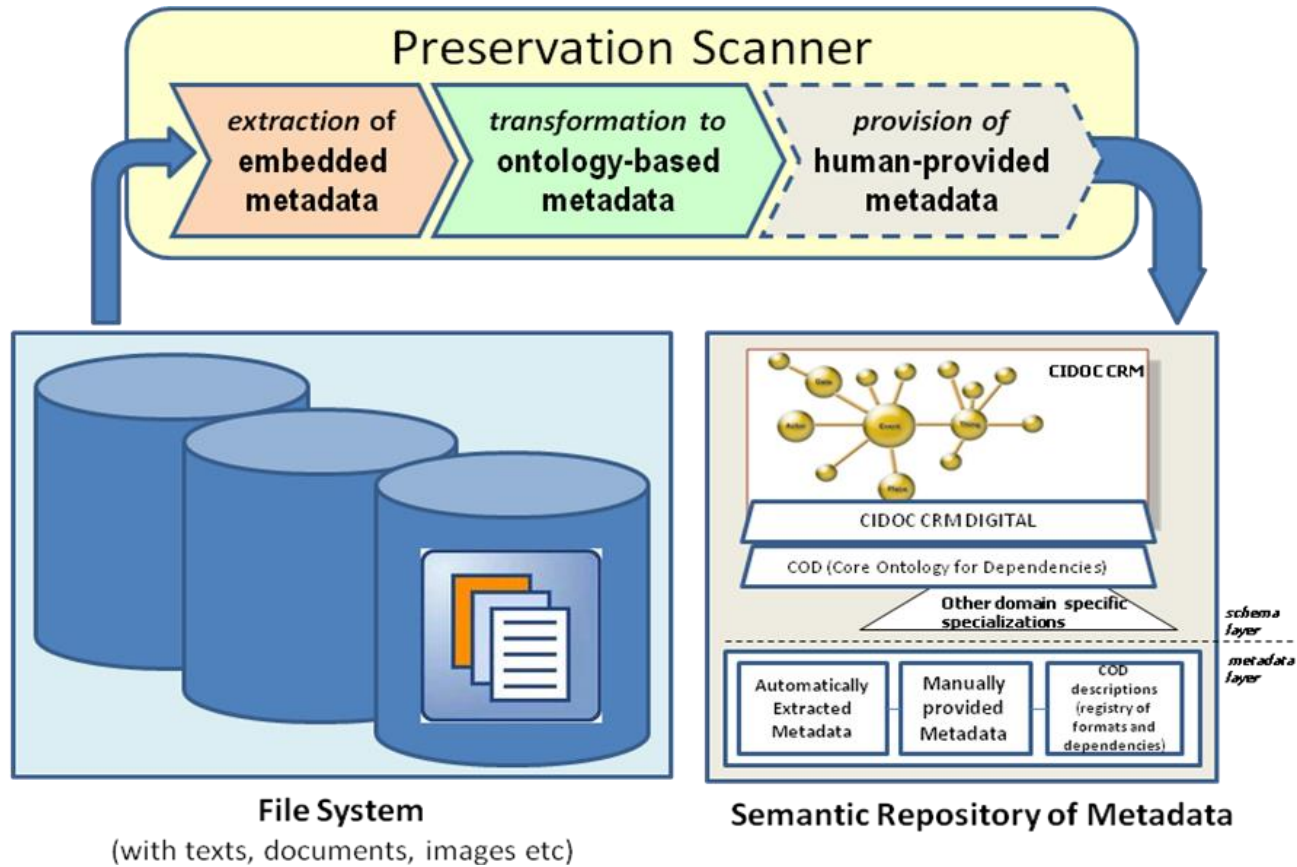
- Faceted search, combining dynamically and statically mined metadata
  - lead to much improved task completeness with much less user interactions
  - was more preferred by the users (advanced and plain ones) and lead to greater satisfaction, than plain clustering or faceted interfaces

## ○ Most Important Related Publications

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

# PRESCAN

- ◆ **PreScan**: Automated extraction of file-embedded metadata from file systems.



# MORE ABOUT PRESCAN

## ○ Features

- Automatic Scanning of file systems
  - Automatic Format Identification and Extraction of Embedded Metadata
  - Support for Human-entered/edited Metadata
  - Periodic Re-Scannings without losing the human-provided metadata
  - Referential Integrity services
- An early version of this tool is described in the paper [http://users.ics.forth.gr/~tzitzik/publications/Tzitzikas\\_2009\\_MEDES.pdf](http://users.ics.forth.gr/~tzitzik/publications/Tzitzikas_2009_MEDES.pdf)



# 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 autocompletion 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)

Cluster Label Tree of the top suggestion "marilyn"



» 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](http://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))

[Marilyn Agency - Wikipedia, the free encyclopedia](#) - 3

A **Marilyn** is a model agency in Paris. **Marilyn** Agency est la plus importante agence de mannequin F

Search interface showing the input "mari" and a dropdown list of suggestions: "marilyn", "marilyn monroe", "mario games". A yellow box highlights the suggestions list with the text "List of query's suggestions." Below the suggestions are radio buttons for "SET (default)", "PET", "STIE", and "PTIE".

First page of results of the top suggestion "marilyn"

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

(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)



» 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.ht  
(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  
e.

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

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

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



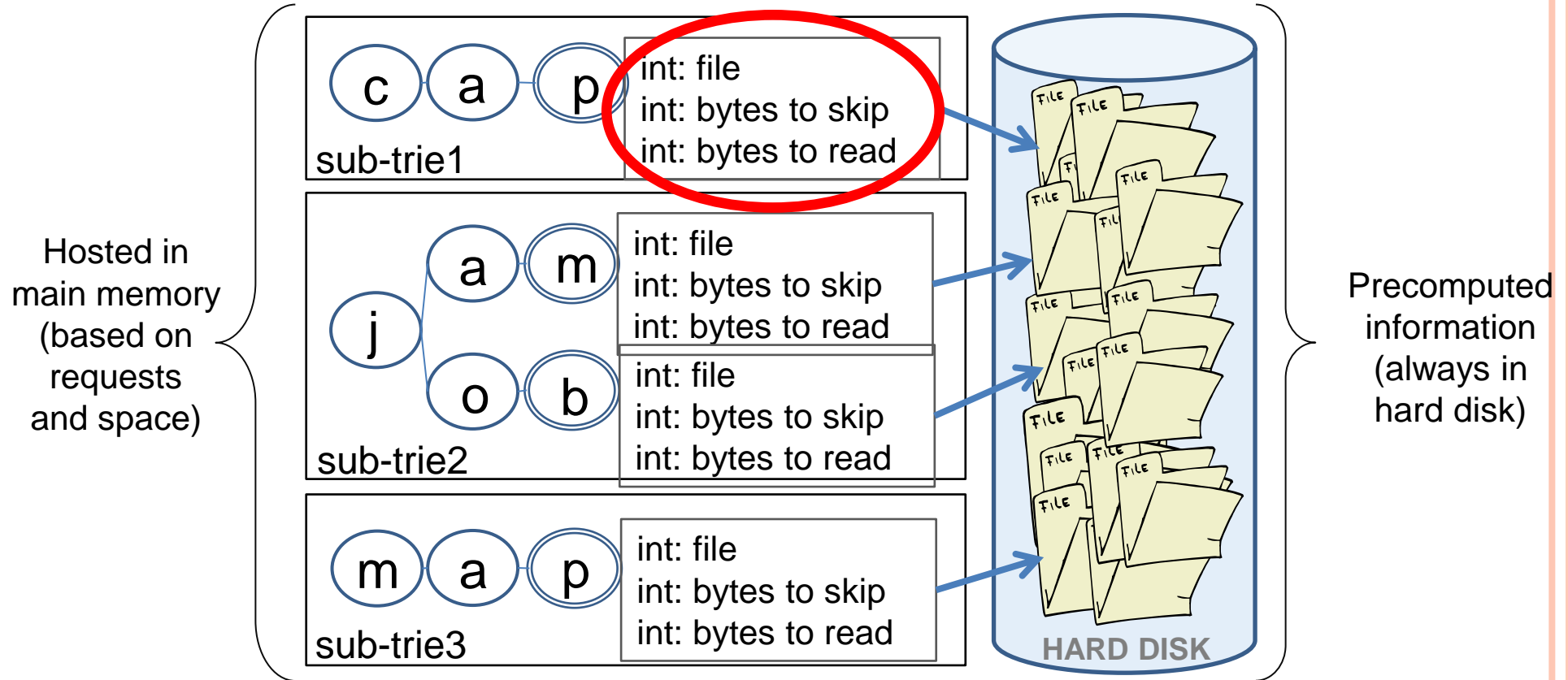
# IOS (INSTANCE OVERVIEW SEARCH)

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 the search term 'tim'. The search results are organized into several sections:

- Search Input:** A search bar containing 'tim' with a dropdown menu showing suggestions: 'tim curry', 'tim montgomery', 'tim powers', and 'tim berners-lee' (highlighted).
- Entity Mining (a):** A list of search results for 'tim' with red boxes around the text 'find its entities' for each result:
  - Tim Berners-Lee - Wikipedia, the free encyclopedia
  - Tim Berners-Lee Bio. A graduate of Oxford University, Tim Berners-Lee invented the World Wide Web, an internet-based hypermedia initiative for global information ...
  - Tim Berners-Lee | Facebook
  - Tim Berners-Lee - Mahalo.com
- Faceted Search (b):** A sidebar with three facets:
  - Person (521 entities):** Tim Berners (30), Robert Cailliau (9), John Poole (7), Sir Timothy John (7), Mark Fischetti (6), Queen Elizabeth II (6), Albert Hofmann (5). A 'show all' button is present.
  - Date (566 entities):** 1980 (12), today (12), 2010 (11), 8 June 1955 (11), March 1989 (7), 25 December 1990 (6), December 2004 (6), 1993 (5). A 'show all' button is present.
  - Organization (416 entities):** MIT (20).
- Metadata-based Faceted Search (b):** A sidebar with three facets:
  - 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).
- Results Clustering (c):** A sidebar showing a hierarchical view of results for 'Tim Berners-Lee(100)':
  - web(48)
  - inventor(22)
  - server(7)
  - sir(19)
  - internet(11)
  - data(9)
  - scientist(6)
  - world(37)
  - news(5)
  - access(7)
  - net(5)
  - linked data(4)
  - information(6)

# 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***

# 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 two red-bordered boxes: one containing "100 results to mine" and another with a checked checkbox for "mine only snippets". The search results are categorized into "Person" (1427 entities) and "Organization" (842 entities). Under "Person", a detailed entry for Nicolas Sarkozy is highlighted with a red border, showing his photo, name, title, birth date, birth place, profession, and links to his website and Wikipedia page. Under "Organization", a list of organizations is shown, including Harvard, White House, Congress, University of Hawaii, and Columbia University. On the left side, there are two text snippets with red-bordered boxes around the "find its entities" links: one for Barack Obama from Wikipedia and another for Barack Obama from Facebook.

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

Y. Tzitzikas, Athens DB, Athens 2012

# CONT

The screenshot shows an entity mining interface. At the top left is the logo for 'entity mining ios' with a magnifying glass over the 'i'. A search bar contains 'barack obama' and a 'Search' button. Below the search bar, it shows '100 results to mine' and a checkbox for 'mine only snippets'. A red box highlights the text 'Results of selected entities: reset'. The main results area shows several snippets, with a red box highlighting 'Joe Biden (13)' and 'John McCain (8)' under the 'Person' category. Another red box highlights 'White House (22)' under the 'Organization' category. A yellow box on the left contains a bullet point: 'Exploitation for restricting the focus'.

**entity mining ios**

barack obama

100 results to mine  
 mine only snippets

**Results of selected entities: reset**

**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)

Barack Obama  
BarackObama.com is the official re-election campaign website of President Barack 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

Record - Barack Obama  
my was losing more than  
d quickly to pass the

American Recovery ...  
<http://www.barackobama.com/record> - find its entities

News for barack+obama Barack Obama - Wikipedia, the free encyclopedia  
Barack Hussein Obama II is the 44th and current President of the United

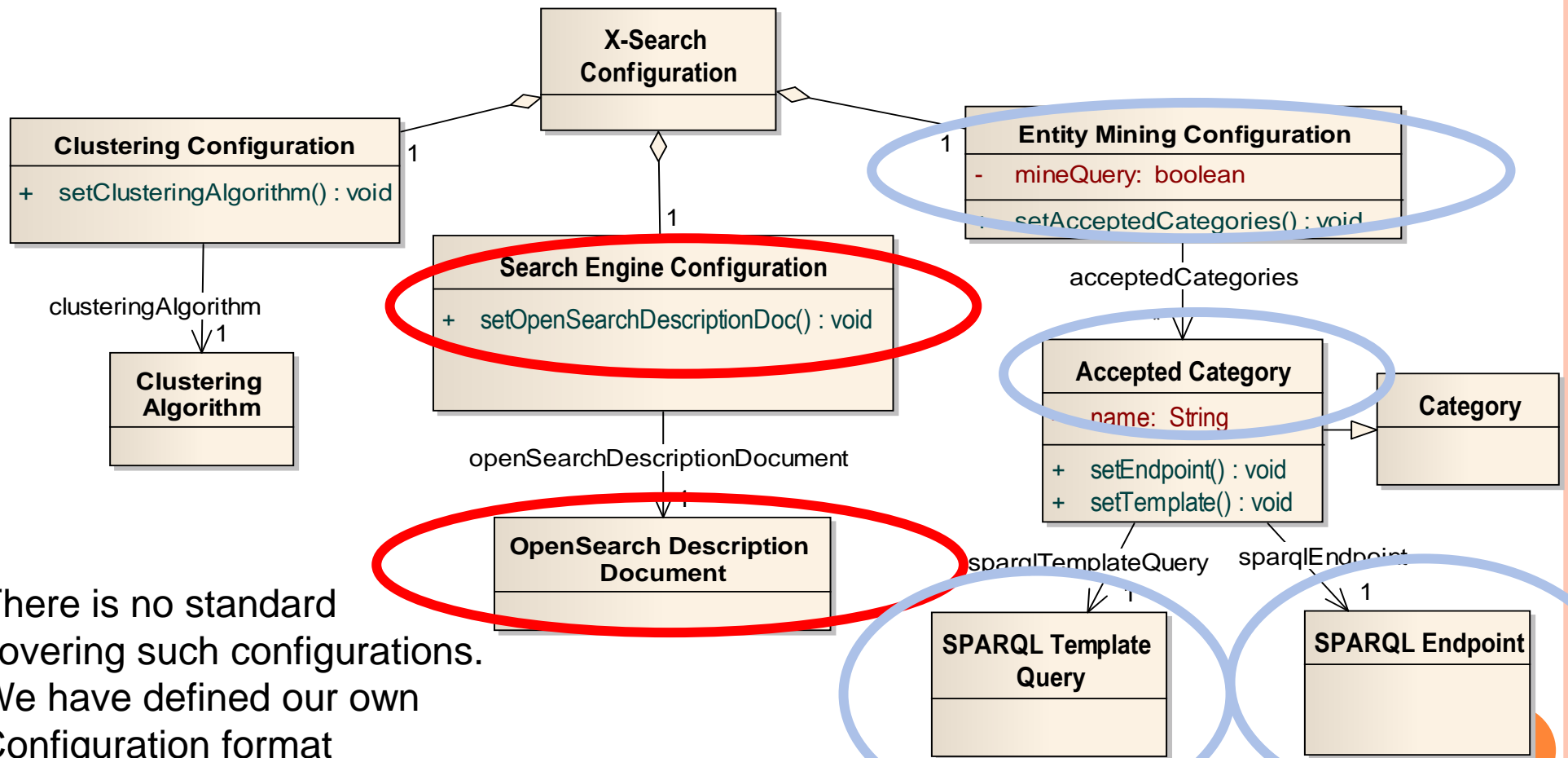
- Exploitation for restricting the focus

# 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 This is the first .. “gap”





top tennis players

Search

about X-ENS | admin configuration

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  
 e tennis player and tennis tournament statis  
 opens with ...

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

**top hits**

### Semantic Entity Enrichment (close)

#### Properties of: Andy Roddick

##### Description

Andrew Stephen "Andy" Roddick (born Aug 30, 1982) is an American professional tennis player and a former World No. 1. He is..

##### Depiction



##### BirthPlace

Omaha, Nebraska

##### BirthDate

1982-08-30

## Country (11 entities)

- India (8)
- Canada (2)

**Entity**

# 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 Salampassis

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 'Powered by X-Search' and lists various entities such as 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 'International Patent Classification (IPC) (111 entities)' with a list of codes and counts.
- Cluster Explorer:** Shows 'Powered by X-Search' and a hierarchical cluster of terms: 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 'Results: 200' and a list of results. The first result is 'Use of nadolol for inhibiting the onset of migraine' by SQUIBB & SON & INC, with IPCs: [A61K31/22, A61K31/21, A61K31/135] EP-0350080-A2, 1. The second result is 'Method for treating migraine symptoms with ibup' by MCNEIL PPC INC.
- Details:** Provides detailed information for the selected result: 'USE OF NADOLOL FOR INHIBITING THE ONSET OF MIGRAINE HEADACHES', Year: 1990, Publication number: EP0350080, Kind: A3, and Application number: EP89115176.

Two large red circles are drawn over the Entities Explorer and Cluster Explorer panels, highlighting the X-Search powered features.

# 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” ( $\approx 1.5$  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

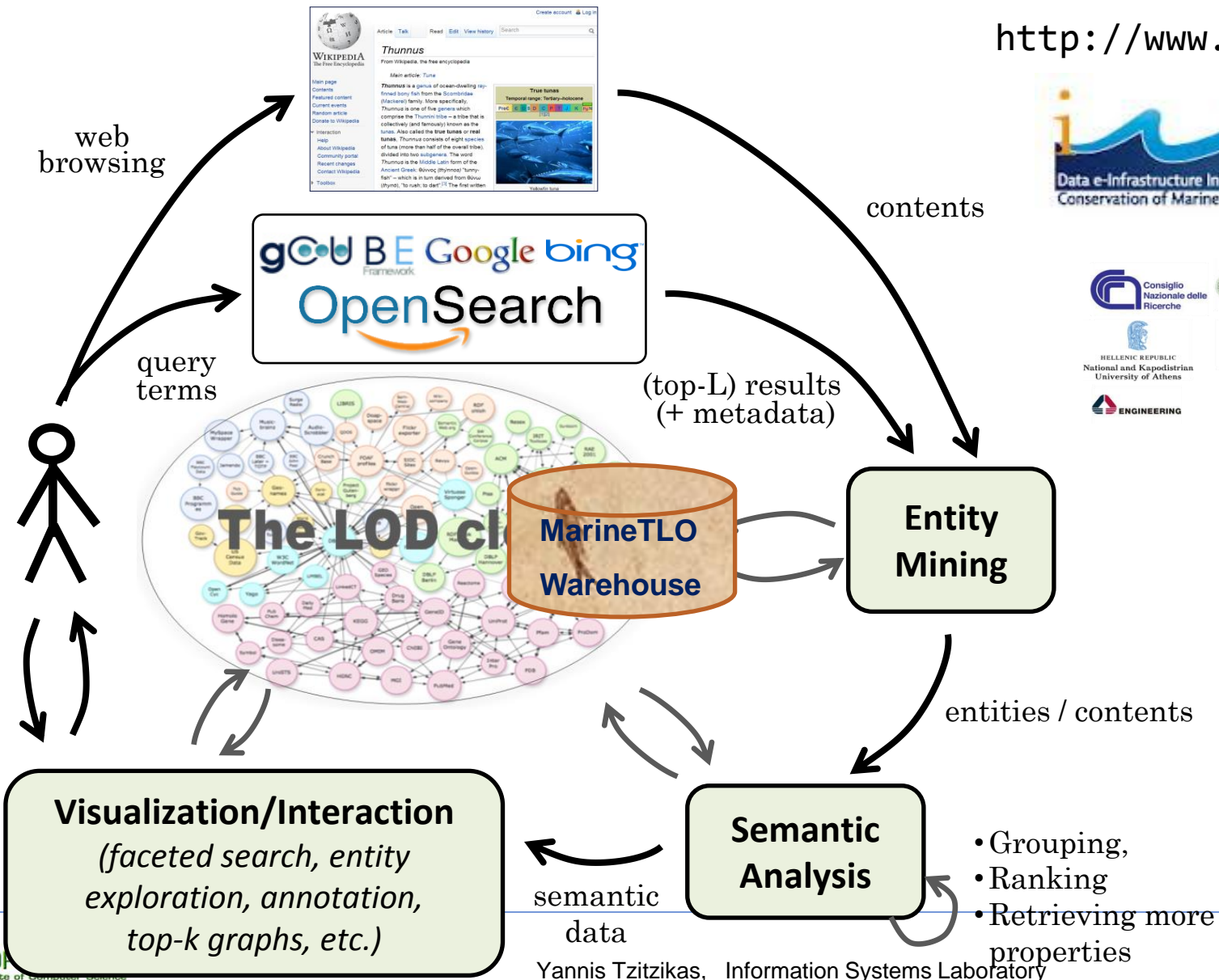
# APPLYING IN THE CONTEXT OF AN INFRASTRUCTURE

- 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

# XSEARCH:

## SEMANTIC POST-PROCESSING OF SEARCH RESULTS

<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

---

- URI: [http://dbpedia.org/resource/Yellowtail\\_amberjack](http://dbpedia.org/resource/Yellowtail_amberjack) (open)
- Value: Yellowtail amberjack

**Semantic Entity Exploration**

Properties of: [Yellowtail\\_amberjack](#)

<p><b>Type</b></p> <ul style="list-style-type: none"> <li>Animal (open)</li> <li>Thing (open)</li> <li>Species (open)</li> <li>FLODSpecies (open)</li> <li>Fish (open)</li> <li>CodeEntity (open)</li> <li>Eukarvote (open)</li> <li>Animal (open)</li> <li>Fish (open)</li> </ul>	<p><b>SameAs</b></p> <ul style="list-style-type: none"> <li>Seriola lalandi (open)</li> </ul>		
<p><b>Subject</b></p> <ul style="list-style-type: none"> <li>Category:Fish_of_the_Red_Sea (open)</li> <li>Category:Fish_of_the_Indian_Ocean (open)</li> <li>Category:Seriola (open)</li> </ul>			
<p><b>BinomialAuthority</b></p> <ul style="list-style-type: none"> <li>Georges Cuvier (open)</li> <li>Achille Valenciennes (open)</li> </ul>	<p><b>Class</b></p> <ul style="list-style-type: none"> <li>Actinopterygii (open)</li> </ul>		
<p><b>Family</b></p> <ul style="list-style-type: none"> <li>Carangidae (open)</li> </ul>			
<p><b>Genus</b></p> <ul style="list-style-type: none"> <li>Seriola (open)</li> </ul>	<p><b>Kingdom</b></p> <ul style="list-style-type: none"> <li>Animal (open)</li> </ul>	<p><b>Order</b></p> <ul style="list-style-type: none"> <li>Perciformes (open)</li> </ul>	<p><b>Phylum</b></p> <ul style="list-style-type: none"> <li>Chordate (open)</li> </ul>
<p><b>Depiction</b></p> <ul style="list-style-type: none"> <li>Seriola lalandi.jpg (open)</li> </ul>	<p><b>Thumbnail</b></p> <ul style="list-style-type: none"> <li>200px-Seriola lalandi.jpg (open)</li> </ul>		

Result of Entity Mining

Result of Textual Clustering

Y. Tzitzikas, Panel@ExploreDB 2014







# WHAT NEW CAR SHOULD I BUY?



**JALOPNIK**

CHART BY JASON TORCHINSKY

© 2013

# MILESTONE 7 PREFERENCES

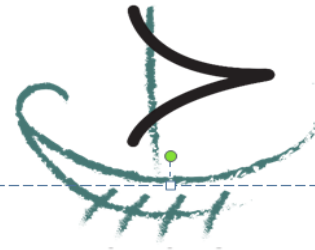
- Then we questioned ourselves:
  - ***What about the ordering of facets, terms and objects? Should the user only restrict the focus? Why not allowing the user to change the order based on his/her preferences?***
- Outcome
  - A framework for preferences over multi-dimensional and hierarchical information 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<sup>1,2</sup> and Y. Tzitzikas<sup>1,2</sup>

- 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

- 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

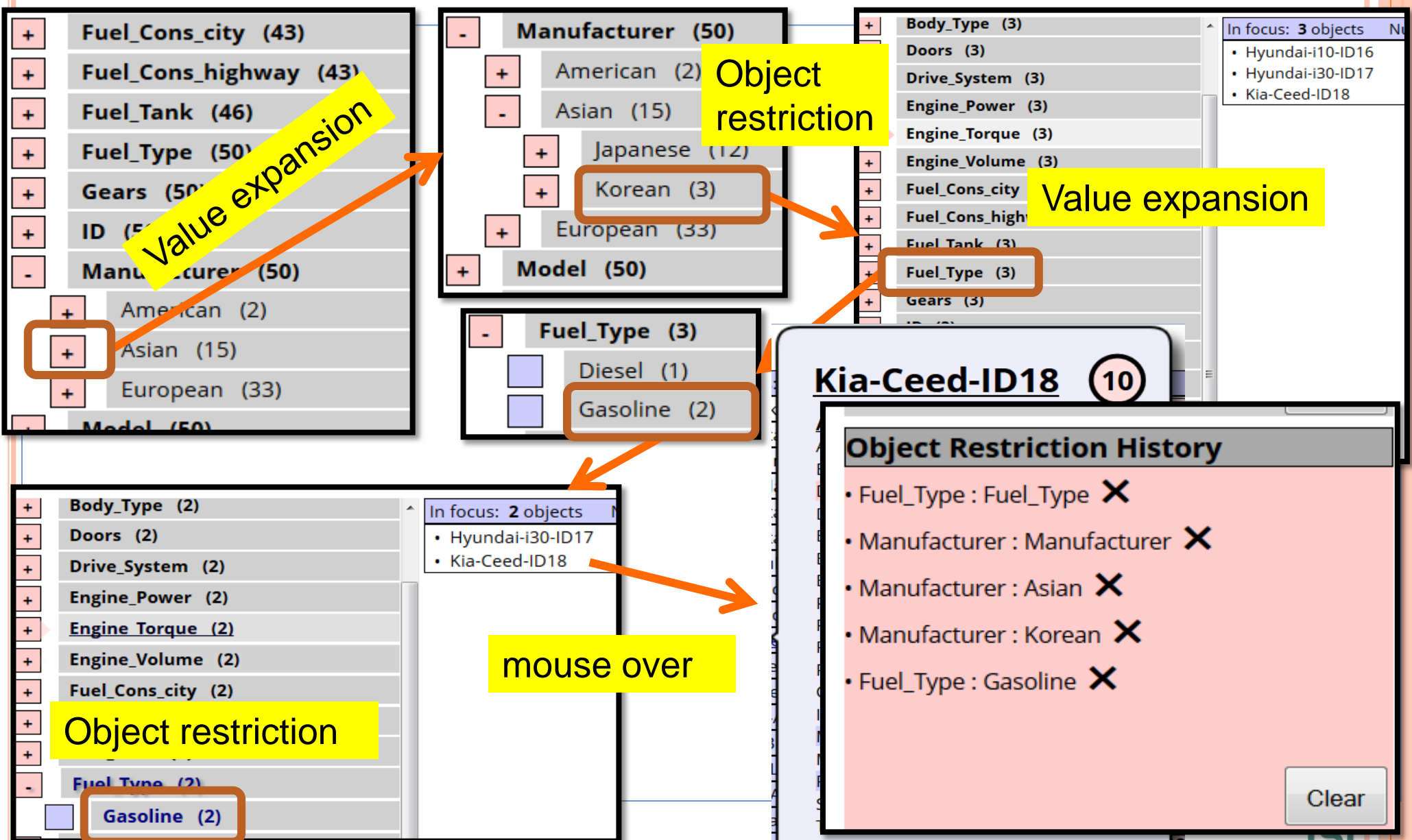
**Interesting Objects**

Clear

**Object Restriction History**



# HIPPALUS: FDT INTERACTIONS



# HIPPALUS: PREFERENCE ACTIONS

Cars ordered with priority on manufacturer

+ Model (50)

In focus: 50 objects    Number of buckets: 33

Preference Actions

In focus: 50 objects    Number of buckets: 50

Facets

In focus: 8 objects    Number of buckets: 8

+ Acceleration (8)

• Peugeot-207-ID33 (1)

+ Body\_Type (8)

• BMW-1-ID7 (2)

- Doors (8)

• BMW-3-ID8 (3)

2 (8)

• Alfa-Romeo-Brera-ID1 (4)

+ Drive\_System (8)

• Audi-TT-ID6 (5)

+ Engine\_Power (8)

• Saab-9-3-ID36 (6)

+ Engine\_Torque (8)

• Alfa-Romeo-8C-ID3 (7)

Level 2: Price\_Euros

• Mercedes-Benz-SL-ID24 (8)

1538 (1)

• 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)

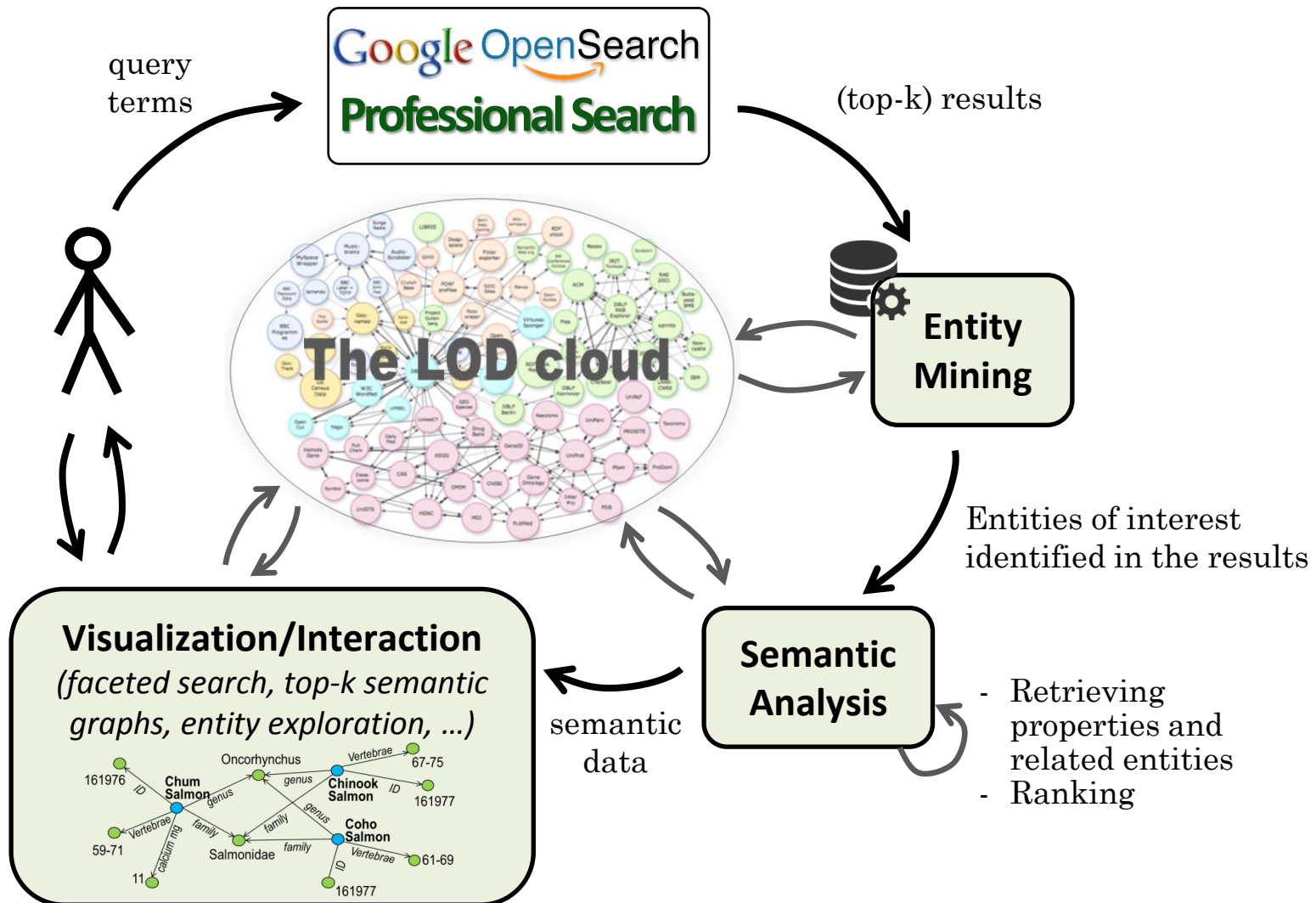
Object restriction

2 (8)

# MILESTONE 8 FROM DIMENSIONS TO GRAPHS

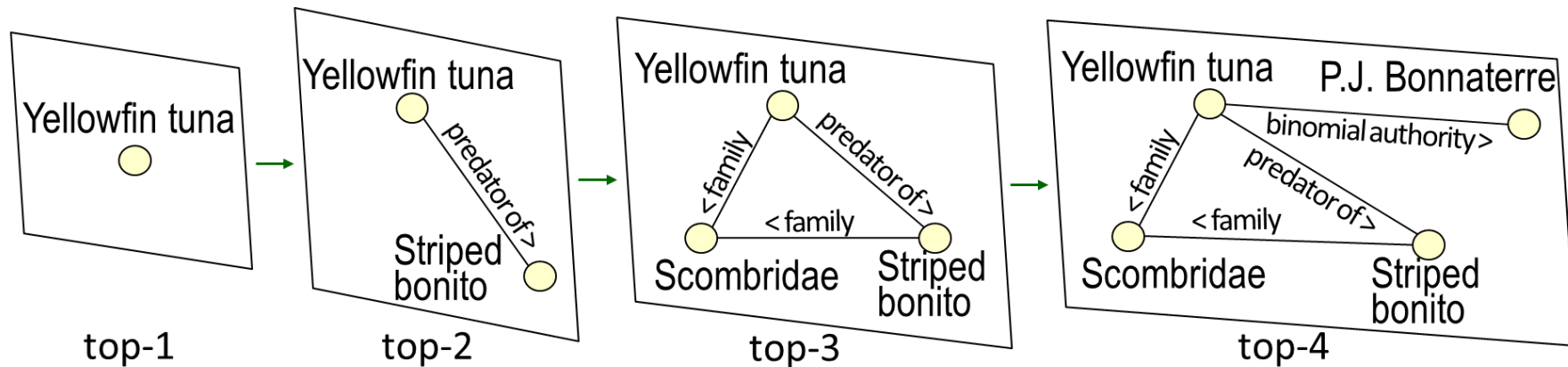
- Then we questioned ourselves:
  - ***So far we have seen services for getting and exploiting multidimensional spaces over the search results. But what if the notion of dimension cannot be defined, or in case there are too many? What can be done without having to configure entity types?***
- Outcome
  - 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.)

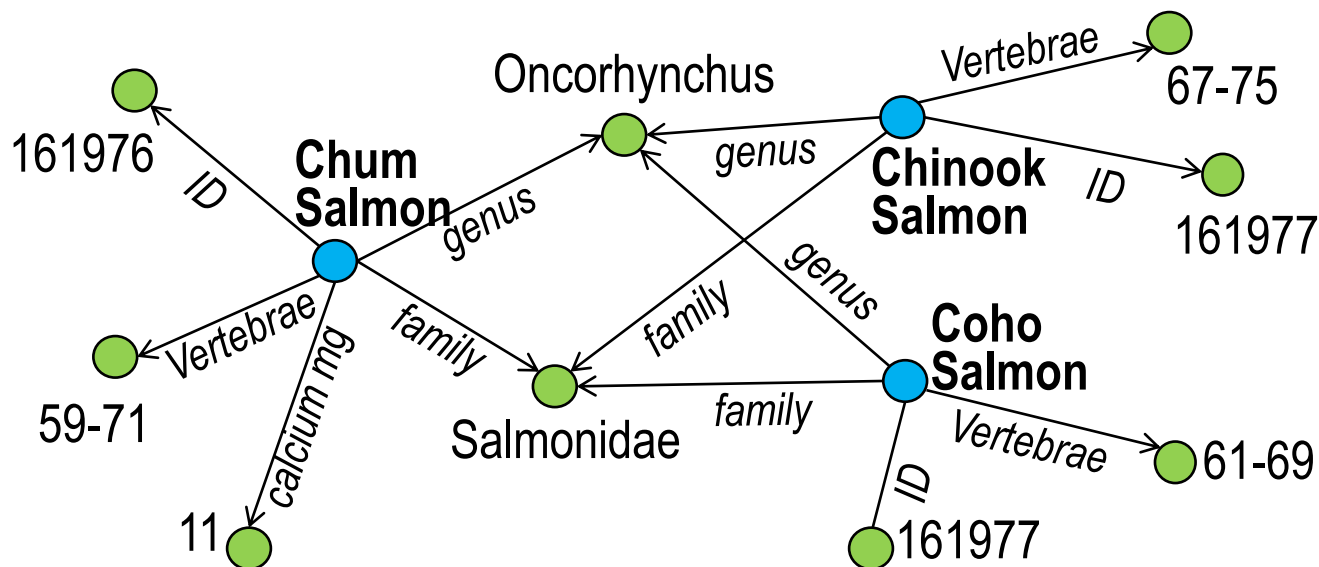
- The system can return the top-K graph for any K from 1 to number of nodes produced
  - Vertices: the K most highly ranked nodes
  - Edges: the edges that connect the K most highly ranked nodes
- The user is free to increase or reduce the value of K
- Example (from a real domain):



# 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





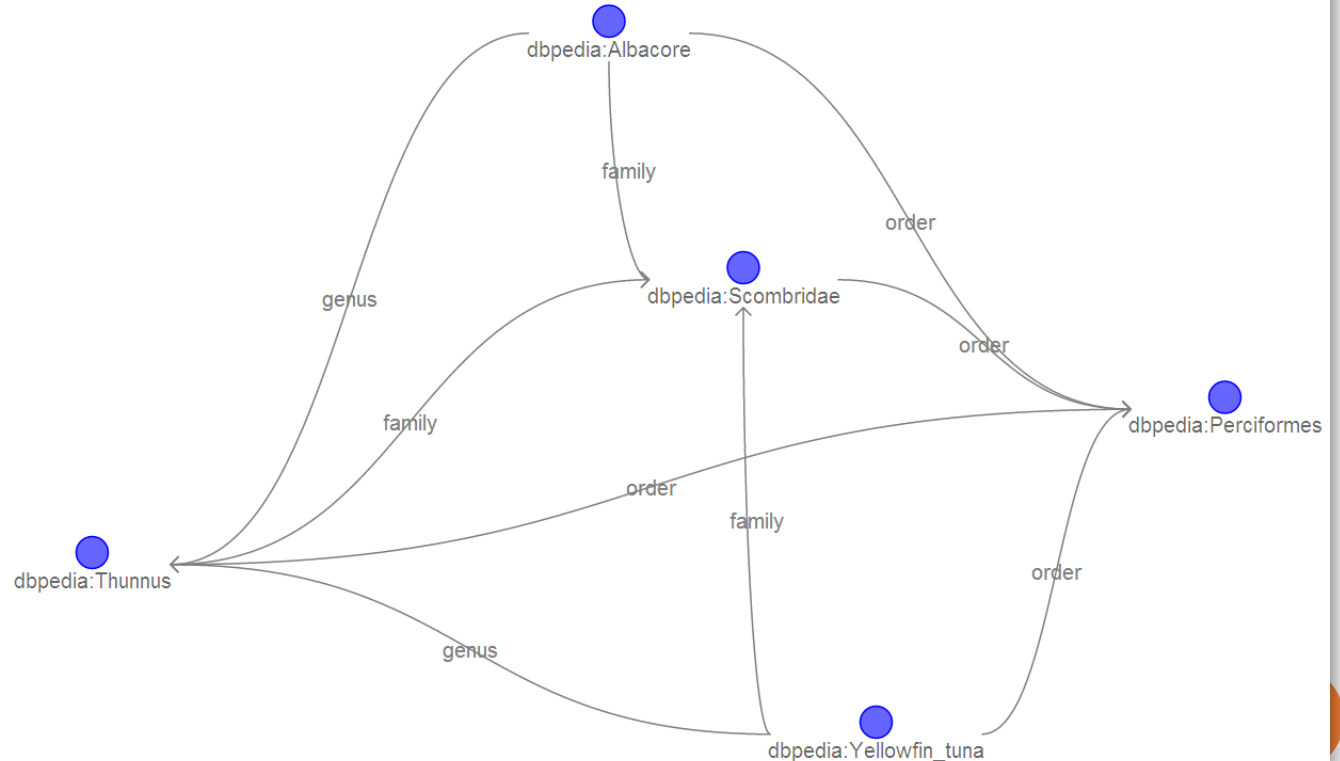
# PROTOTYPE (2014)

<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

# HIPPALUS DEMO

- With Firefox **version 8+** try <http://www.ics.forth.gr/isl/Hippalus>

**CONTACT PERSON:  
YANNIS TZITZIKAS  
([HTTP://WWW.ICS.FORTH.GR/~TZITZIK](http://www.ics.forth.gr/~tzitzik))**

52

