

Semantic Technologies an overview

Marko Grobelnik

Marko.Grobelnik@ijs.si

Jozef Stefan Institute

Ljubljana, Slovenia

Outline

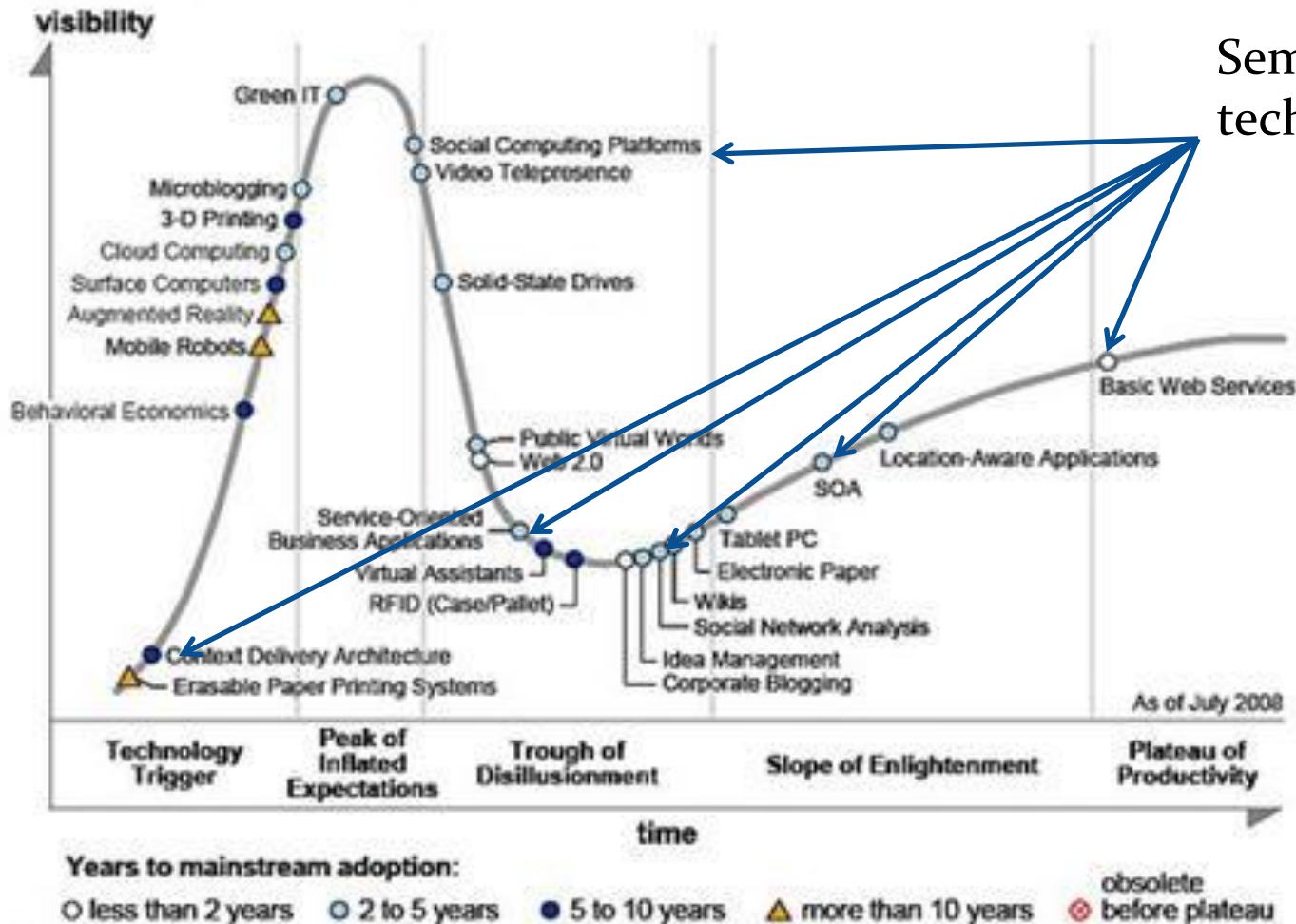
- Quick introduction
 - ...what are semantic technologies?
- Gartner's hype curve
- Semantic Web Technology stack
- Web X.X
- Examples
 - Dealing with legacy relational databases
 - Dealing with legacy software
 - Contextualized search
 - Identifying news reporting bias
 - Common sense reasoning

What are semantic technologies?

- Semantic technologies are interdisciplinary set of technologies with the main goal to make information **interoperable**
- What are the three main “buzzwords”?
 - Semantic Web
 - Semantic Web Services
 - Web2.0
- ...and related ones:
 - W3C, Social computing, Ontologies, ... and many more

Where Semantic technologies fit into Gartner's hype-cycle

Figure 1. Hype Cycle for Emerging Technologies, 2008



Semantic technologies

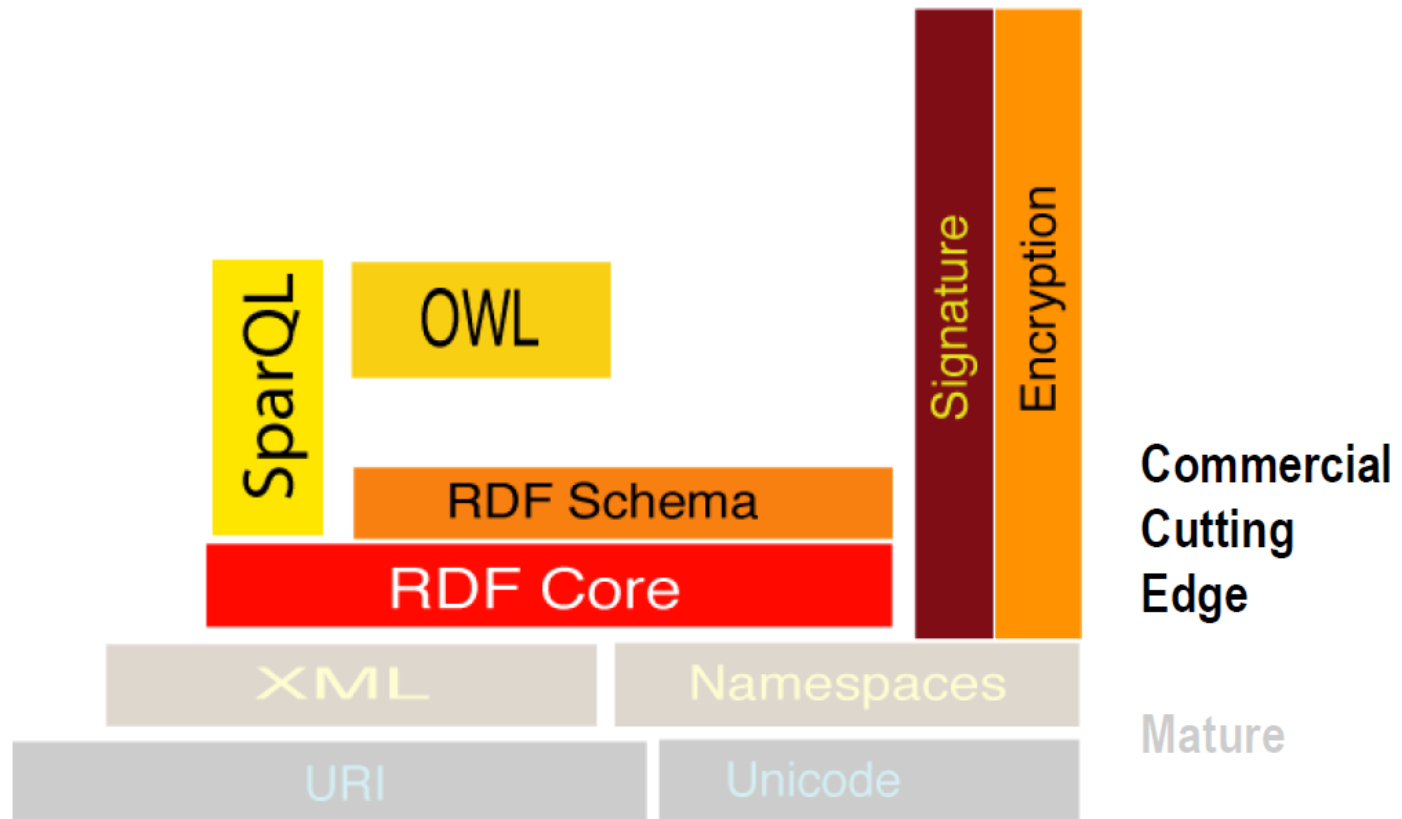
Semantic Web Technology stack

The Semantic Web in 2008



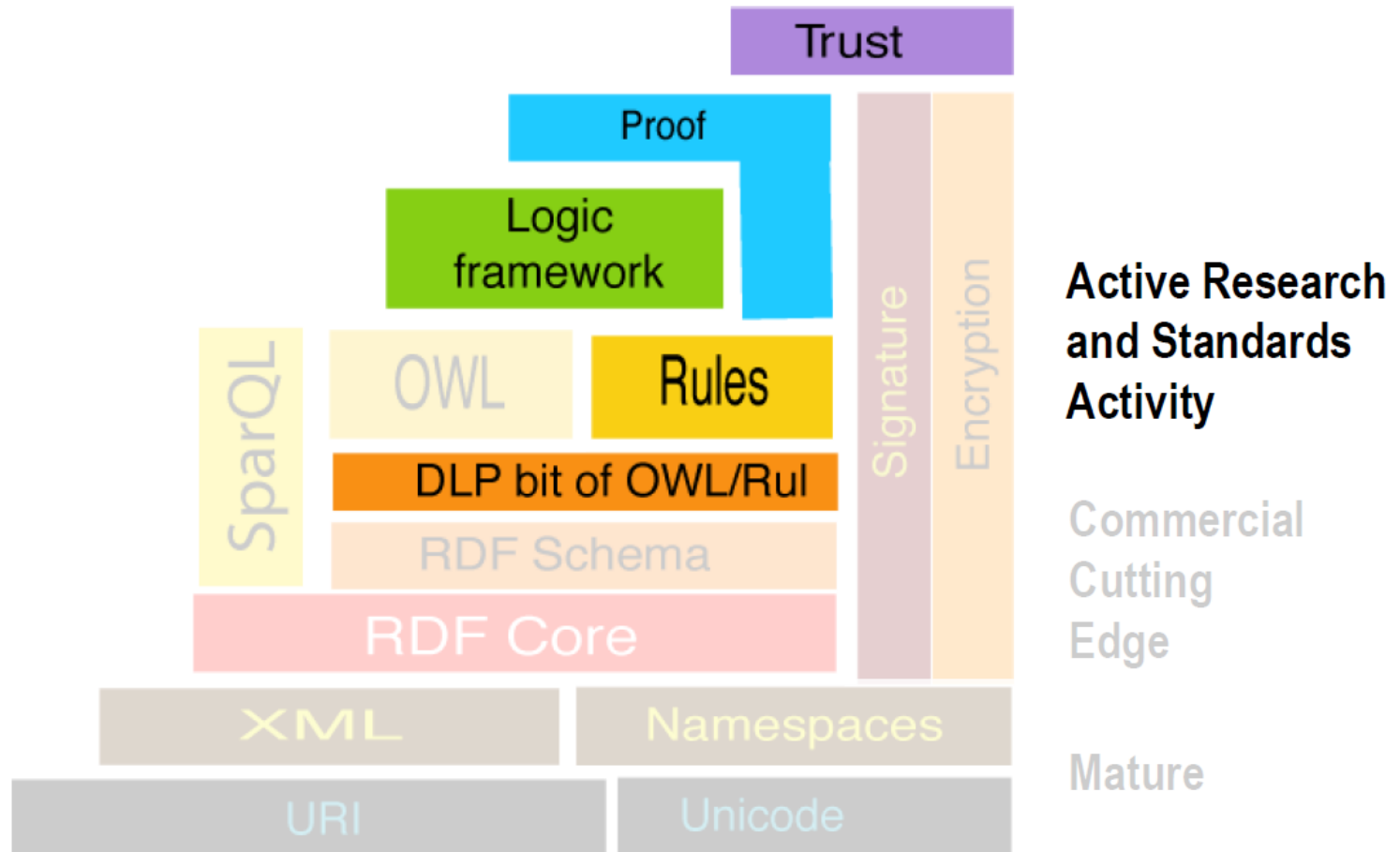
“The Famous Semantic Web Technology Stack”

The Semantic Web in 2008



“The Famous Semantic Web Technology Stack”

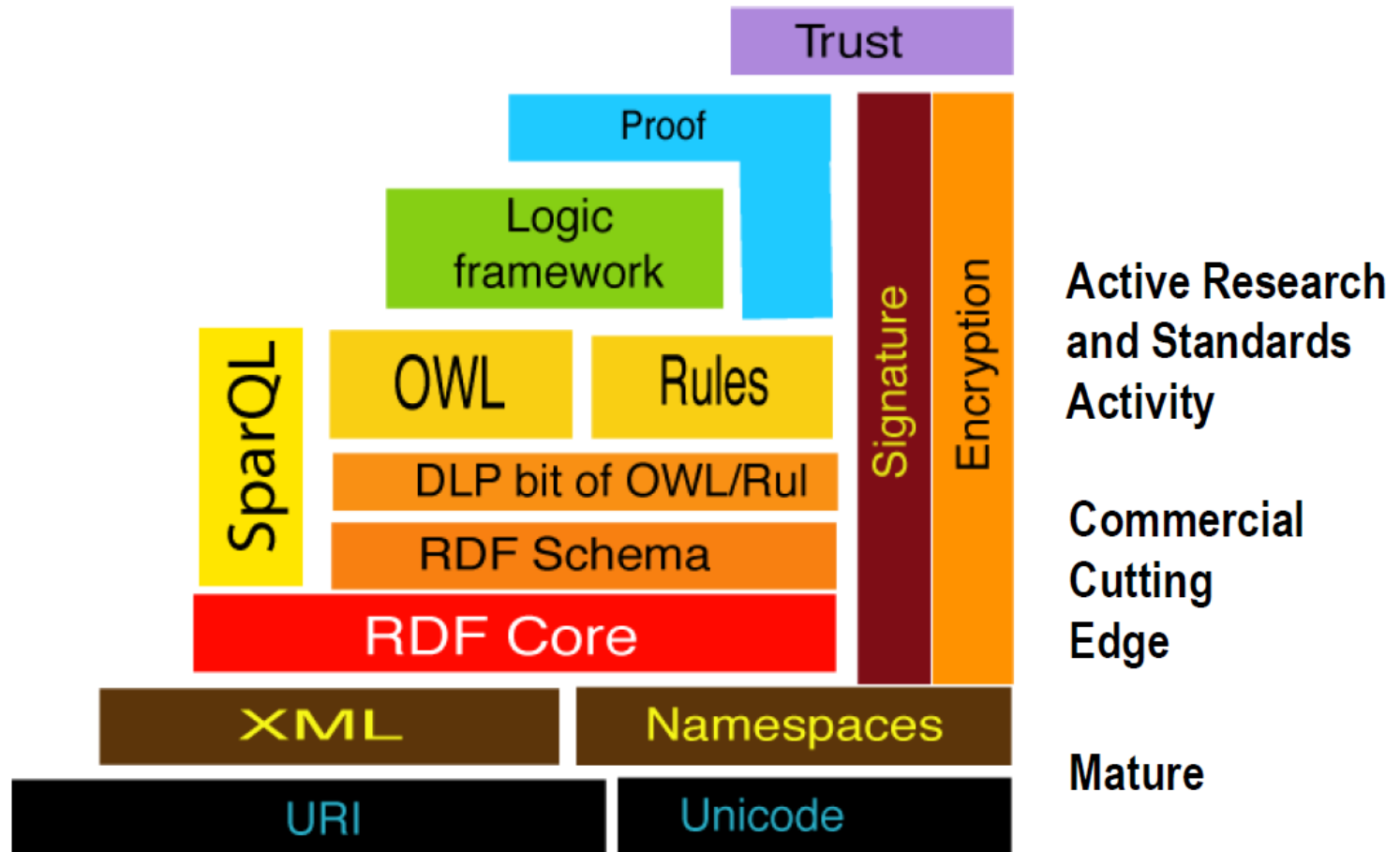
The Semantic Web in 2008



“The Famous Semantic Web Technology Stack”



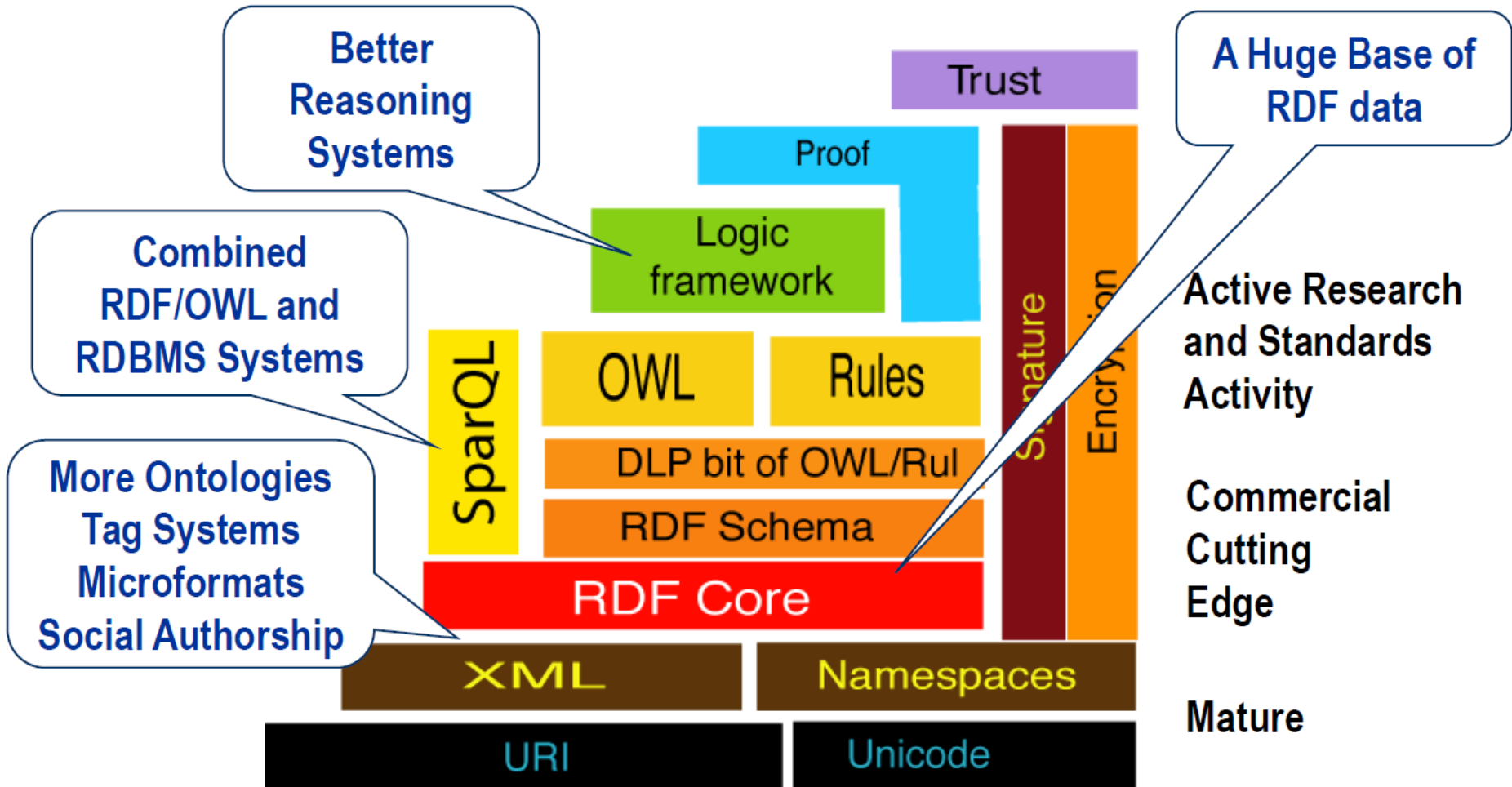
The Semantic Web in 2008



“The Famous Semantic Web Technology Stack”



Completing the Semantic Web Picture



Other Technologies Impact the Semantic Web



The beautiful world of Web X.X

The beautiful world of Web X.X versions

(...a trial to put all of them on one slide)

	Description	Technologies
Web 1.0	Static HTML pages (web as we first learned it)	HTML, HTTP
Web 1.5	Dynamic HTML content (web as we know it)	Client side (JavaScript, DHTML, Flash, ...), server side (CGI, PHP, Perl, ASP/.NET, JSP, ...)
Web 2.0	Integration on all levels, collaboration, sharing vocabularies (web as it is being sold)	weblogs, social bookmarking, social tagging, wikis, podcasts, RSS feeds, many-to-many publishing, web services, ... URI, XML, RDF, OWL, ...
Web 3.0	...adding meaning to semantics - AI dream revival (web as we would need it)	Closest area of a research would be “common sense reasoning” and the “Cyc system” (http://www.nytimes.com/2006/11/12/business/12web.html?ref=business)

Web 2.0 –is there any new quality?

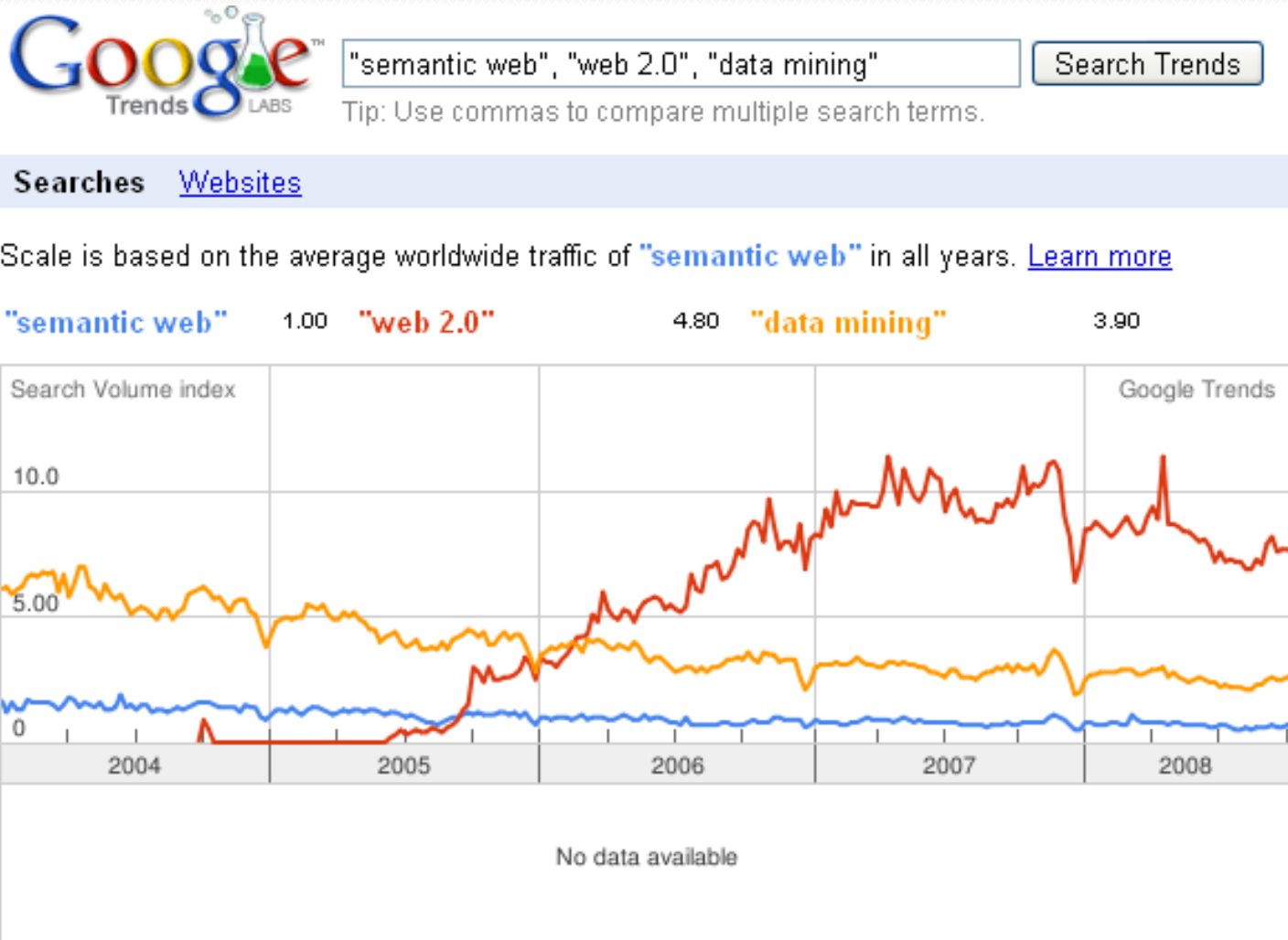
- IMHO, with “Web 2.0” the Web community became **really aware** of the importance of the global collaborative work
 - ...next step in globalization of the Web
 - **Bottom-up** “social networking” seems to nicely complement the traditional **top-down** schema design approaches



Visualization of Web 2.0 typical vocabulary
(http://en.wikipedia.org/wiki/Image:Webzo_en.png)

Web 2.0 – the current hype

Google search volume of “Web 2.0” vs. “semantic web” vs. “data mining”



...scale and dynamics of Web 2.0

- Per minute, there are:
 - 100 edits in Wikipedia (144K/day)
 - 200 tags in del.icio.us (288K/day)
 - 270 image uploads to flickr (388K/day)
 - 1100 blog entries (1.6M/day)

What about Web 4.0? 😊

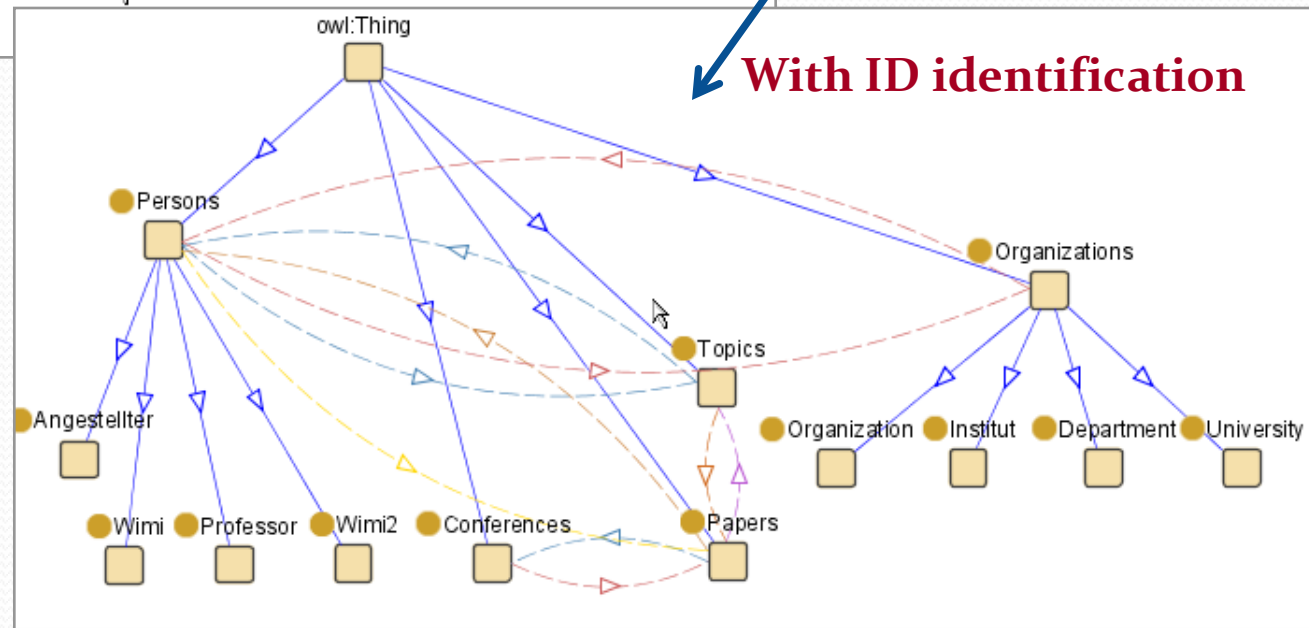
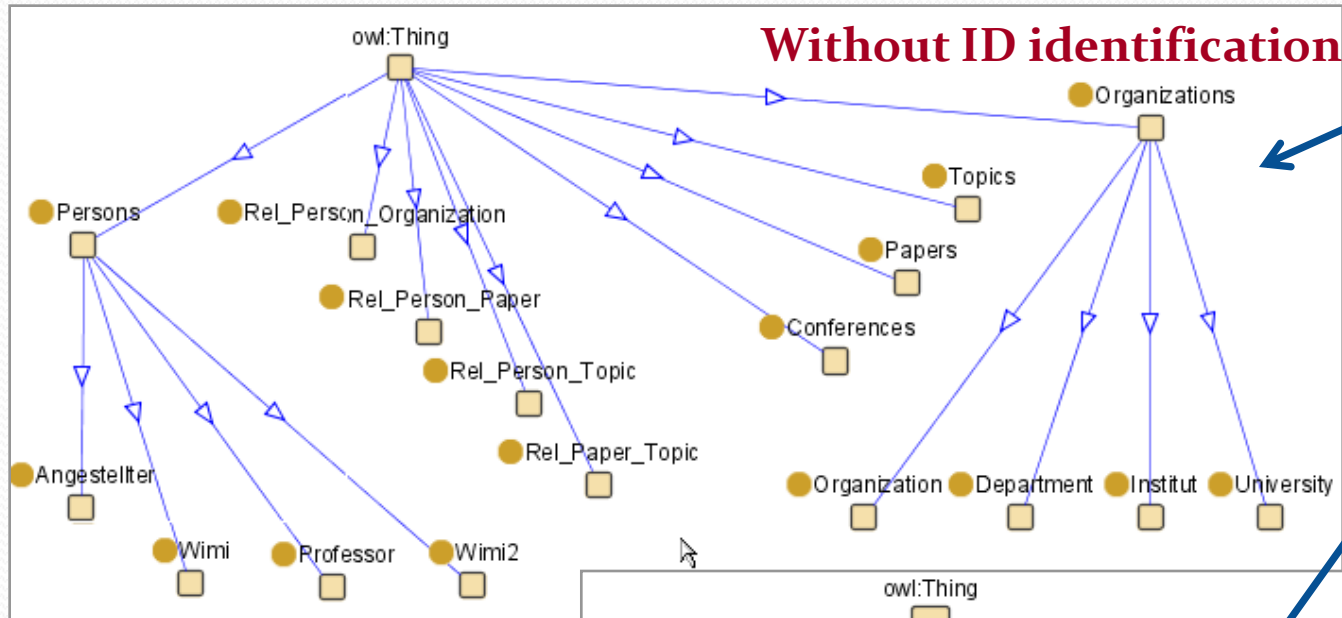
- Citation from some Intel blog:
 - *“...Web 4.0 is the impending state at which all information converges into a great ball of benevolent self-aware light, and solves every problem from world peace to ...”* http://blogs.intel.com/it/2006/11/web_40_a_new_hype.html
- Ultimate stage in web development...
 - ...will prevent Web 5.0 to happen since everything will be resolved already by Web 4.0.

Examples & Demos

Example: mining data models in legacy databases

- Data models in relational databases are often not designed properly
 - ...especially after many patches and many people being involved
- In the next example we show how a large relational database (~500 tables) from DESSAULT (airplane producer) was corrected with semiautomatic system

Example: finding hidden foreign-key relationships in large relational databases



Example: legacy software mining

- Software is as any other data source possible domain for analysis
- ...in the following example we are mining large legacy software package GATE written in Java and present some alternative views

Software Data Sources

- Structured
 - Code samples
 - Web service usage logs
 - Source code
 - DB schemas ...
- Unstructured
 - Web pages
 - User's/Reference manual
 - Tutorials, lectures, forums, newsgroups, etc.
 - Source code comments
 - DB content ...

A Typical Java Class

Class
comment

```
/** The format of Documents. Subclasses of DocumentFormat know about
 * particular MIME types and how to unpack the information in any
 * markup or formatting they contain into GATE annotations. Each MIME
 * type has its own subclass of DocumentFormat, e.g. XmlDocumentFormat,
 * RtfDocumentFormat, MpegDocumentFormat. These classes register themselves
 * with a static index residing here when they are constructed. Static
 * getDocumentFormat methods can then be used to get the appropriate
 * format class for a particular document.
 */
public abstract class DocumentFormat
extends AbstractLanguageResource implements LanguageResource{
```

```
/** The MIME type of this format. */
private MimeType mimeType = null;
```

Field comment

```
/**
 * Find a DocumentFormat implementation that deals with a particular
 * MIME type, given that type.
 * @param aGateDocument this document will receive as a feature
 * the associated Mime Type. The name of the feature is
 * MimeType and its value is in the format type/subtype
 * @param mimeType the mime type that is given as input
 */
static public DocumentFormat getDocumentFormat(gate.Document aGateDocument,
MimeType mimeType){
```

```
    } // getDocumentFormat(aGateDocument, MimeType)

} // class DocumentFormat
```

Creating a Document Network

DocumentFormat.class

```
/** The format of Documents. Subclasses of DocumentFormat know about
 * particular MIME types and how to unpack the information in any
 * markup or formatting they contain into GATE annotations. Each MIME
 * type has its own subclass of DocumentFormat, e.g. XmlDocumentFormat,
 * RtfDocumentFormat, MpegDocumentFormat. These classes register themselves
 * with a static index residing here when they are constructed. Static
 * getDocumentFormat methods can then be used to get the appropriate
 * format class for a particular document.
 */
public abstract class DocumentFormat
extends AbstractLanguageResource implements LanguageResource{

    /** The MIME type of this format. */
    private MimeType mimeType = null;

    /**
     * Find a DocumentFormat implementation that deals with a particular
     * MIME type, given that type.
     * @param aGateDocument this document will receive as a feature
     * the associated Mime Type. The name of the feature is
     * MimeType and its value is in the format type/subtype
     * @param mimeType the mime type that is given as input
     */
    static public DocumentFormat getDocumentFormat(gate.Document aGateDocument,
                                                    MimeType mimeType){

    } // getDocumentFormat(aGateDocument, mimeType)

} // class DocumentFormat
```

DocumentFormat

Creating a Document Network

DocumentFormat.class

```

/** The Format of Documents. Subclasses of DocumentFormat know about
 * particular MIME types and how to unpack the information in any
 * markup or formatting they contain into GATE annotations. Each MIME
 * type has its own subclass of DocumentFormat, e.g. XmlDocumentFormat,
 * RtfDocumentFormat, MpegDocumentFormat. These classes register themselves
 * with a static index residing here when they are constructed. Static
 * getDocumentFormat methods can then be used to get the appropriate
 * format class for a particular document.
 */
public abstract class DocumentFormat
extends AbstractLanguageResource implements LanguageResource{

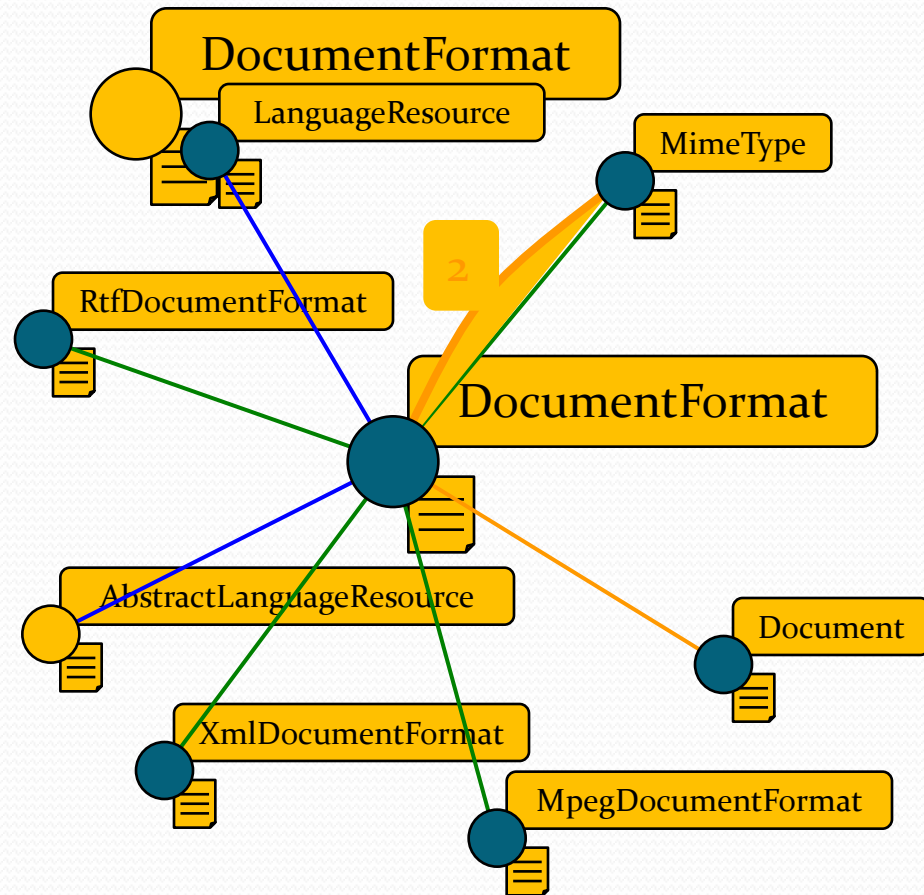
    /** The MIME type of this format. */
    private MimeType mimeType = null;

    /**
     * Find a DocumentFormat implementation that deals with a particular
     * MIME type, given that type.
     * @param aGateDocument this document will receive as a feature
     *                       the associated Mime Type. The name of the feature is
     *                       MimeType and its value is in the format type/subtype
     * @param mimeType the mime type that is given as input
     */
    static public DocumentFormat getDocumentFormat(gate.Document aGateDocument,
                                                    MimeType mimeType){

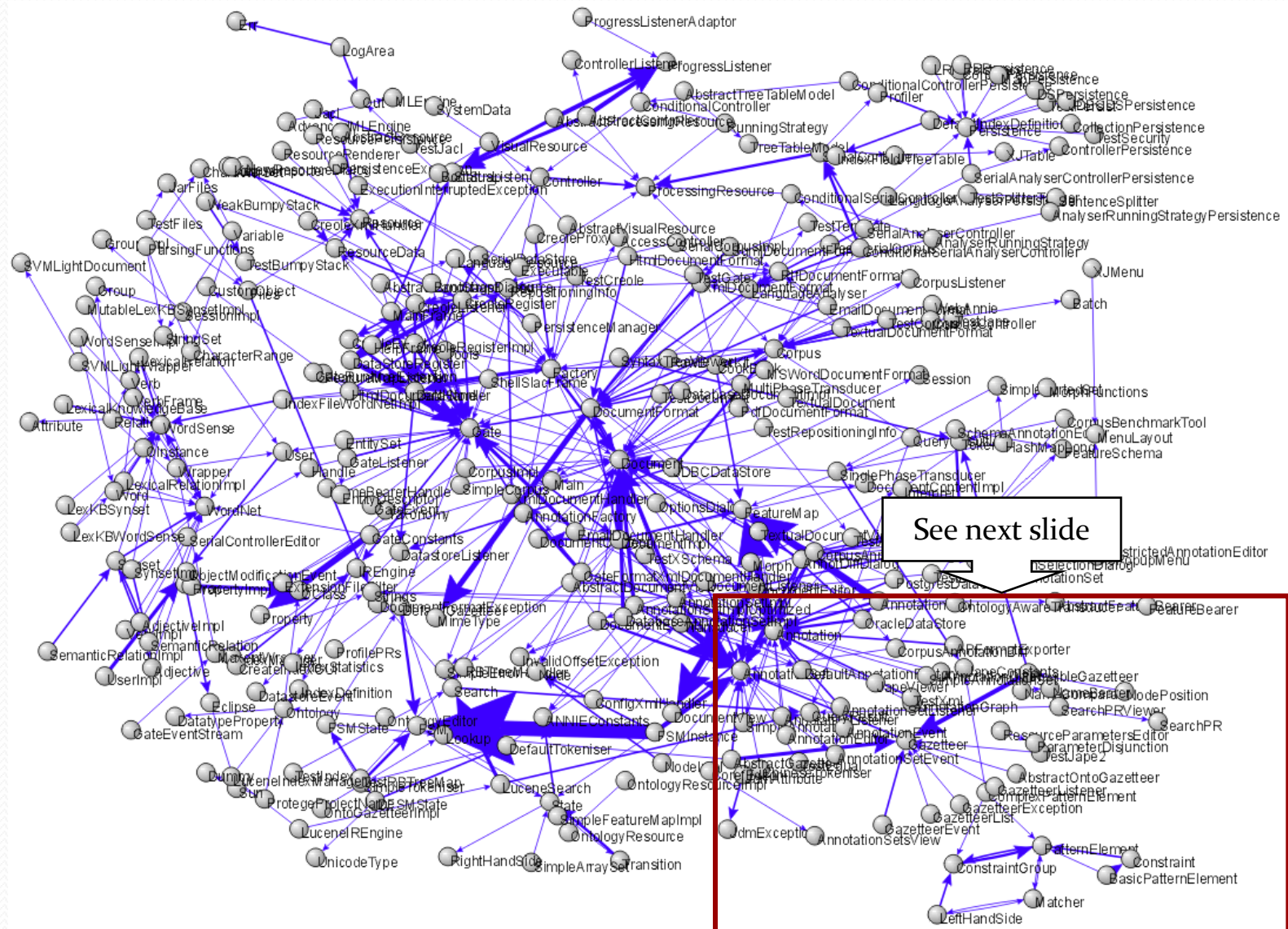
        } // getDocumentFormat(aGateDocument, MimeType)

    } // class DocumentFormat

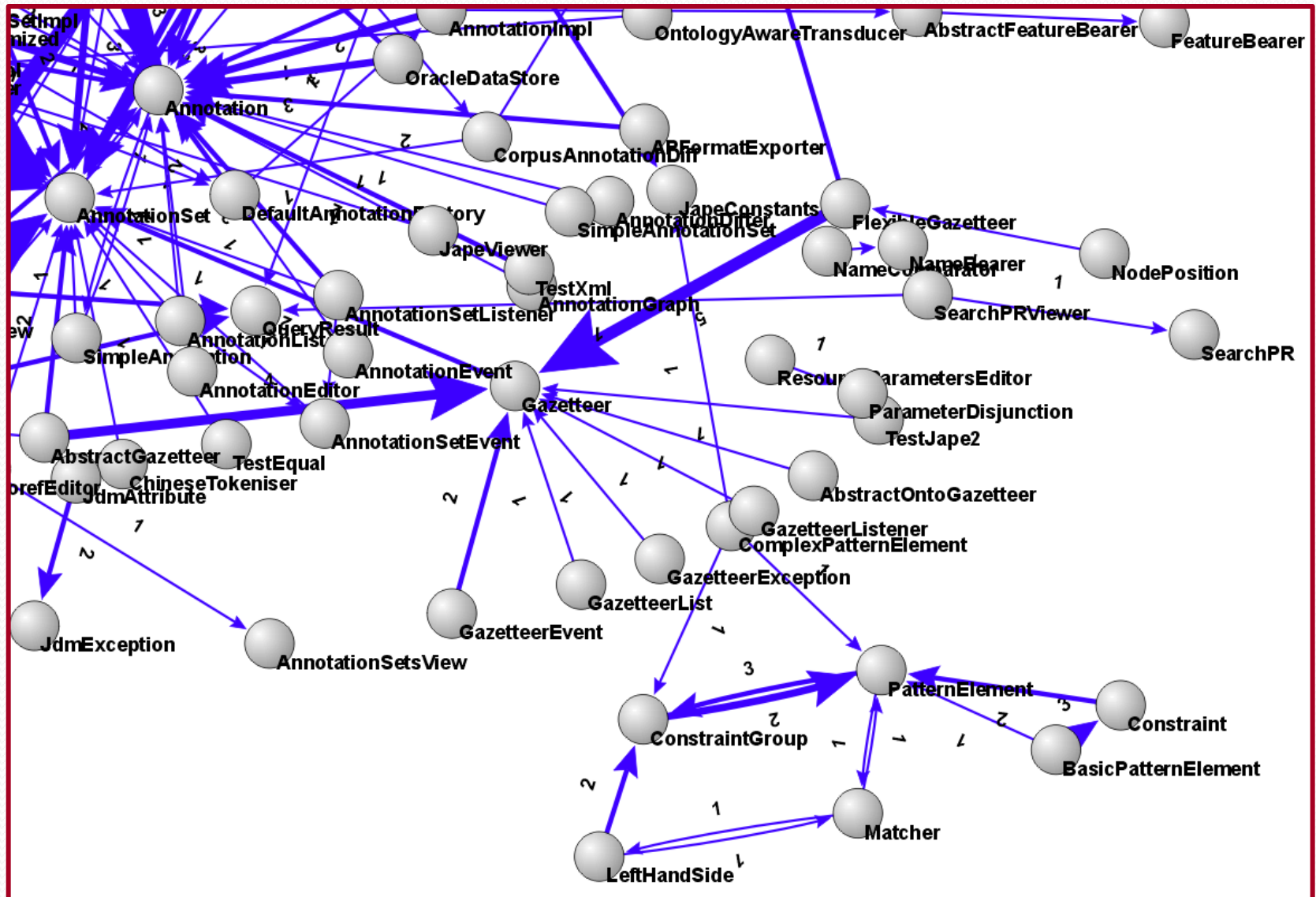
```



GATE Comment Reference Network



GATE Comment Reference Network



FileAbout

Concepts

NewMoveDelete

root

annotation, event, document

exception, impl, words

test, indexed, property

Concept properties

DetailsSuggestionsRelations

Suggestk-MeansQueryAdd

No. suggestions: 3Docs: ☒ All ☐ Unused

KeywordsNo. ... [%]

Ontology details

Ontology visualizationConcept's documentsConcept Visualization

Map propertiesZoom modeSelect mode

Annotation

Data storage

GUI

WordNet

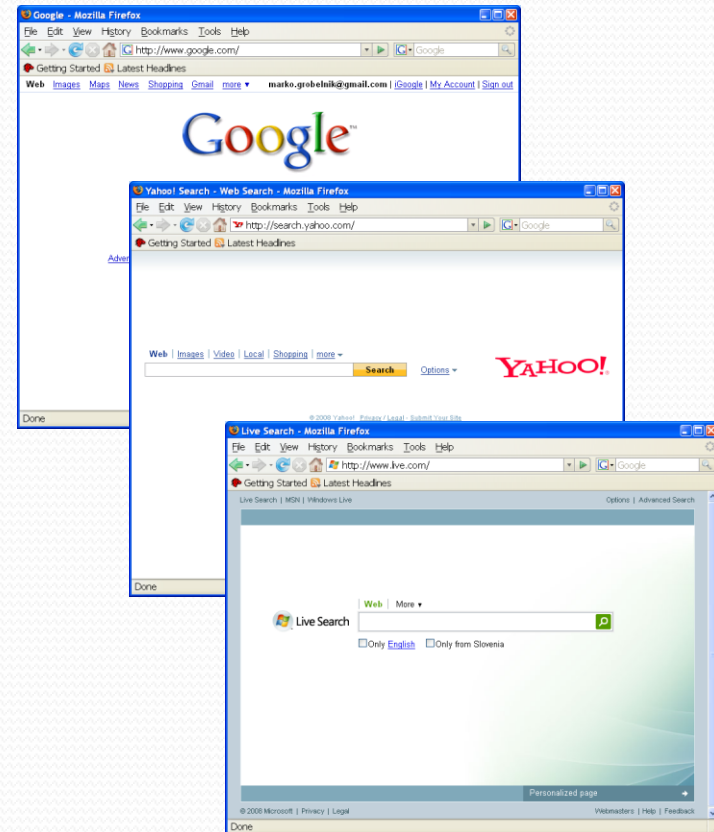
Exceptions

Test

OntoGen news: [OntoGen won Best Demo Award at ESWC 2006](#)

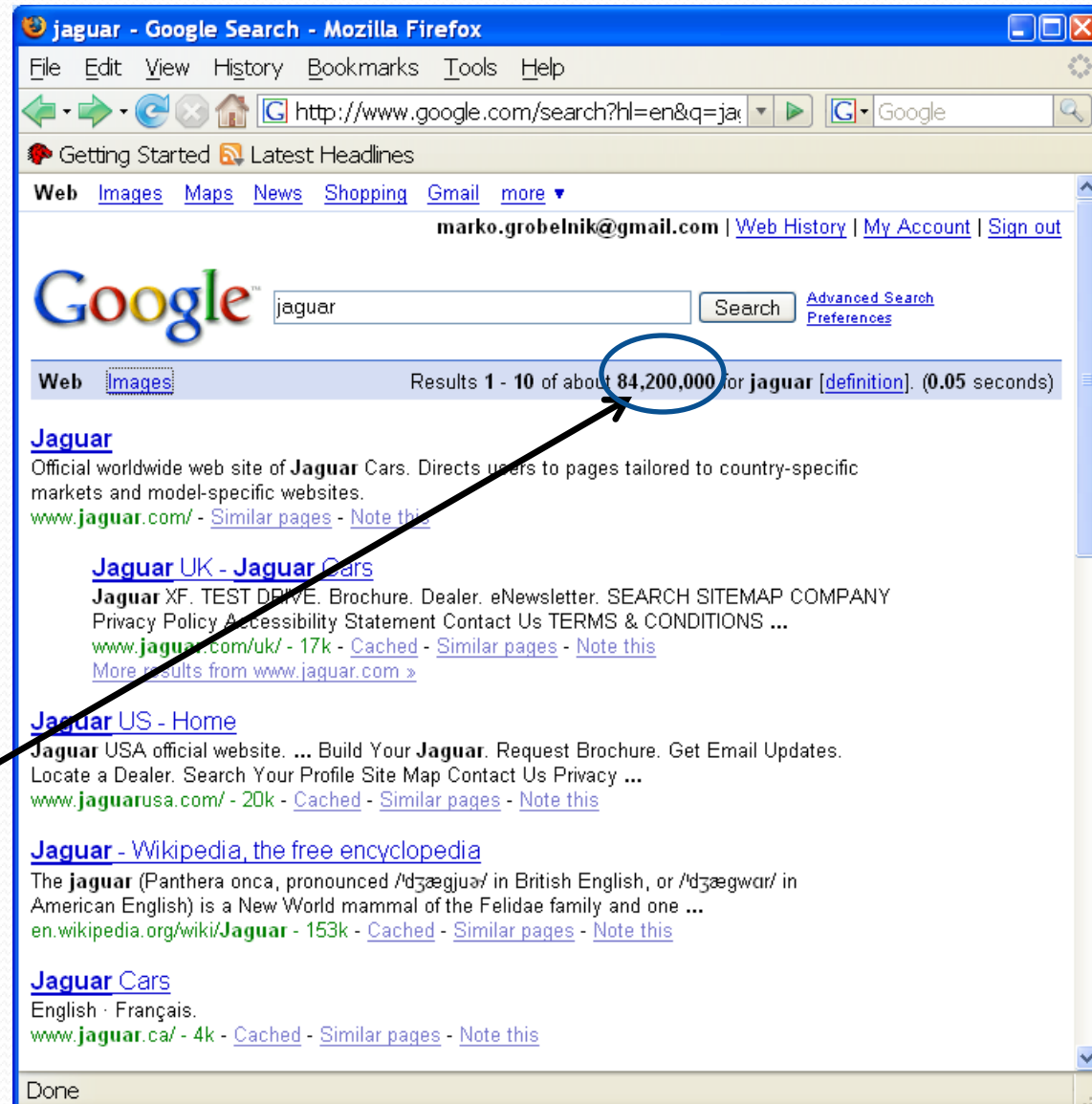
Example: contextualized search

- What are the most common tasks where we manipulate texts in everyday life?
 - “Internet search”!
- ...but – how smart is search technology today?
 - ...not too smart!
 - It is sophisticated, but not smart



Example: searching for “Jaguar”

- Query “jaguar” has many meanings...
- ...but the first page of search engines doesn't provide us with many answers
- ...there are 84M more results



Context sensitive search with

<http://searchpoint.ijs.si>

The screenshot shows a Windows Internet Explorer browser window titled "jaguar - SearchPoint - Windows Internet Explorer". The address bar displays "http://searchpoint.ijs.si/Result.aspx". The search bar contains the query "jaguar". Below the search bar are three buttons: "Search via topics", "Search via query to ontology", and "Search via hits to ontology". The search results are listed on the left, with the top result being "(9) Jaguar" from Big Cats Online. To the right of the search results is a conceptual map. The map is a network of nodes connected by lines. The central node is "Top". Other nodes include "Mammalia", "Vehicles", "Shopping", "Sports", "Games", "Console Platforms", "Aviation", "Society", "Recreation", "Models", "Enthusiasts", "Aircraft", "Parts and Accessories", and "NFL". A red dot is placed on the "Vehicles" node, and a black arrow points from the "Vehicles" node in the map to the "Vehicles" node in the search results.

Query

Conceptual map

Search Point

Dynamic contextual ranking based on the search point

Example:

Detecting News Reporting Bias

- The task:
 - Given a news story, are we able to say from which news source it came?
- We compared **CNN** and **Aljazeera** reports about the same events from the war in Iraq
 - ...300 aligned articles describing the same story from both sources
- The same topics are expressed in both sources with the following keywords:
 - CNN with:
 - **Insurgents**, Troops, Baghdad, Iran, **Militant**, Police, **Suicide**, **Terrorist**, United, National, Hussein, **Alleged**, Israeli, Syria, Terrorism...
 - Aljazeera with:
 - Attacks, Claims, **Rebels**, Withdrawing, Report, **Fighters**, President, **Resistance**, Occupation, Injured, Army, Demanded, Hit, Muslim, ...

Semantics, Knowledge and Common sense reasoning

Towards text understanding...

- The key element to understand the text is to go beyond characters and words...
 - ...meaning, we need to have knowledge in the form of a “**world model**” where all the facts from text fit,
 - ...we need to be able to deal with **contexts**, and
 - ...we need to be able to **reason**
- Do we have something which would go in this direction?
 - ...there were couple of trials in the last decades
 - ...the only marketable system is **Cyc** from a company **CyCorp** (US and Europe/Slovenia based)
 - New York Times article on Cyc and Web 3.0:
 - <http://www.nytimes.com/2006/11/12/business/12web.html>

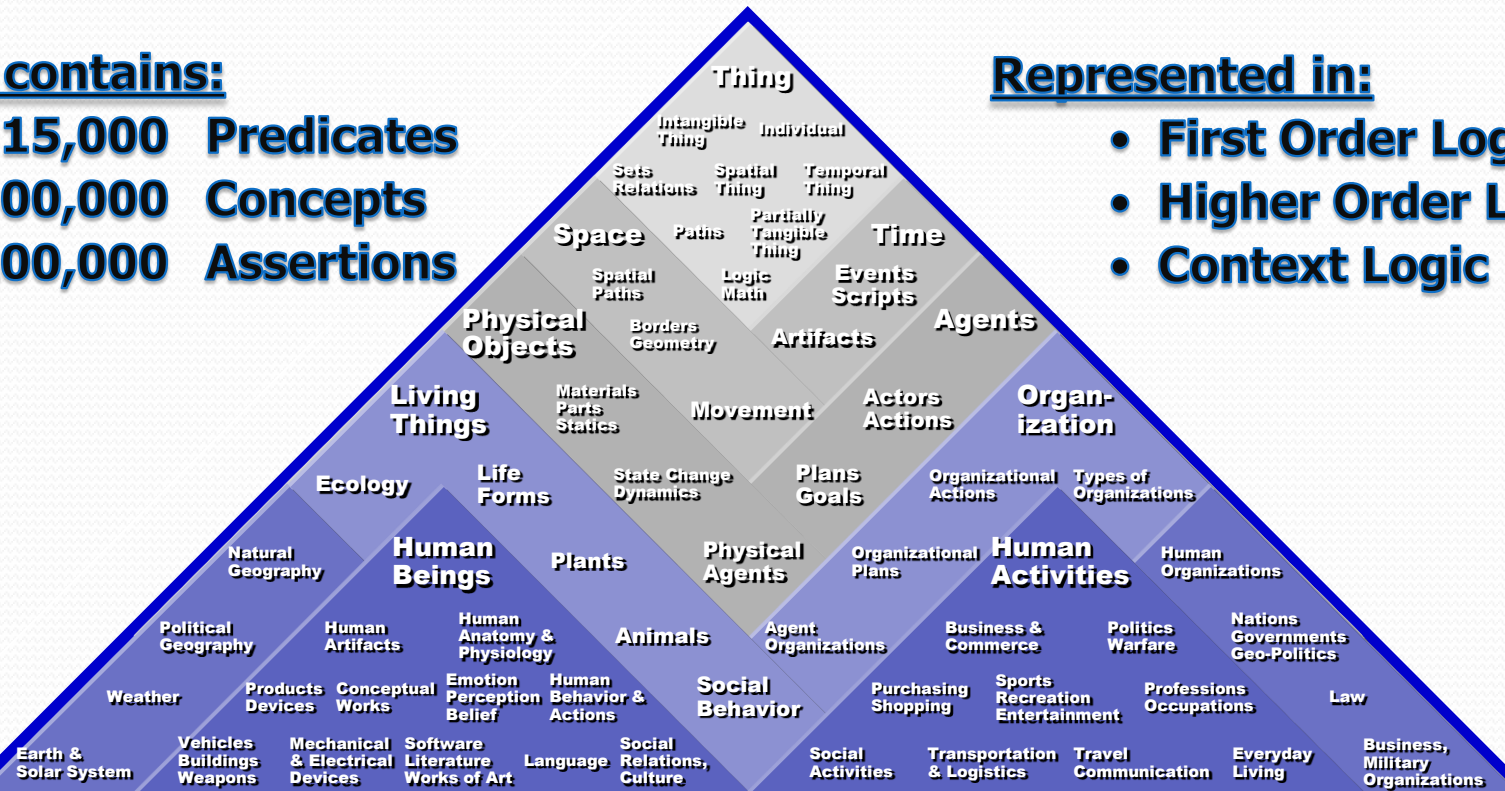
The Cyc Ontology – knowledge about common sense

Cyc contains:

15,000 Predicates
300,000 Concepts
3,200,000 Assertions

Represented in:

- **First Order Logic**
- **Higher Order Logic**
- **Context Logic**



General Knowledge about Various Domains

Specific data, facts, and observations

Cyc's front-end: "Cyc Analytic Environment" – querying (1/2)

Task Info Document Search Concepts Related-to Query Creator Queries

Find Stop

Text query

Query (semi) automatically translated in the First Order Logic

Answers to the query

WHO had a motive for the assassination of Hariri.

Continue
Save
New Tab
Reset

5 answers
Timed out

Allow speculation?

Answers (5)

Answer	Speculation Level	Sources
Bashar al-Assad	No Speculation	W CNN
Syria	Mildly Speculative	CNN
al Qaeda	Moderately Speculative	SAIC CNN
United States, the	No Speculation	2
Israel	No Speculation	2

Justify Fact Sheet Visualize Visualize All


Status: Finished Message: No appropriate visualizations found

Cyc's front-end: "Cyc Analytic Environment" – justification (2/2)


Task Info Document Search Concepts Related-to Query Creator Queries Justification Justification Justification

Proof 1 Save... Copy

▼ Query: Who or what had a motive for the assassination of Hariri?
Answer: al Qaeda
Because:


Since 2000, Lebanon has been responsible for according with Lebanese economic reform.  1


February 14, 2005 was the date of the assassination of Hariri.  2

Rafik Hariri was killed during the assassination of Hariri.  2
Rafik Hariri is an advocate of Lebanese economic reform.
Al Qaeda opposes Lebanese economic reform.

▼ Detailed Justification:
▶ Al Qaeda had a motive for the assassination of Hariri.

▼ External Sources:

1  Gary C. Gambill, "Dossier: Rafiq Hariri", *United States Committee for a Free Lebanon*, July 2001, http://www.meib.org/articles/0107_id1.htm.

2  "Huge blast kills Lebanese ex-PM", *the Cable News Network*, February 14, 2005, <http://www.cnn.com/2005/WORLD/meast/02/14/beirut.explosion.1910/>.

Query & Answer

Justification

Sources for Reasoning and Justification

▼ Options
▼ Options

Further online information

Recorded tutorials, lectures, summer-schools available from <http://videolectures.net>

- Semantic Web:
http://videolectures.net/Top/Computer_Science/Semantic_Web/

The screenshot shows the VideoLectures website in a Mozilla Firefox browser window. The browser's address bar displays <http://videolectures.net/>. The website's header includes the title "VideoLectures - exchange ideas & share knowledge - Mozilla Firefox" and a navigation menu with links: HOME, MOST POPULAR, LATEST LECTURES, CATEGORIES, EVENTS, PEOPLE, INTERVIEWS, TUTORIALS, and CONTACT US. A search bar is located in the top right corner, and a status bar at the bottom indicates "191 events, 3500 authors, 3948 lectures, 5505 videos".

The main content area is titled "FEATURED LECTURES:" and displays a grid of video thumbnails. Each thumbnail includes a video player, a title, a view count, and a duration. The featured lectures are:

- Some Mathematical Tools for Machine Learning** by Chris Burges (2211 views, 02:54:53)
- Introduction to the KDD07 Conference** by Pavel Berkhin (365 views, 00:08:06)
- A service robot named Markovito** by Hector Aviles, Elva Corona-Xelhuantzi, Sergio Cabello, Victor Manuel Jaquez Leal, Enrique Sucar, Eduardo Morales (136 views, 00:04:59)
- Inference and Learning with Networked Data** by Foster Provost (86 views, 01:58:53)
- Shrinkage Estimator for Bayesian Network Parameters** by John Burge (23 views, 00:20:31)

Below the featured lectures, there are three sections:

- RECENT EVENTS:** A link to "more" and a featured event titled "2.993 Spec. Topics in Mechanical Engineering: The Art and Science of Boat Design" by MIT, which is offered during the Independent Activities Period (IAP) in January 2007.
- NEWS:** A section titled "MIT OpenCourseWare Collection" and "MITOPENCOURSEWARE" announcing a collaboration with MIT OpenCourseWare, and a section titled "Cambridge University Engineering Department - Machine Learning seminars" announcing a collaboration with the University of Cambridge.
- CATEGORIES:** A list of categories with their respective counts: Architecture (2), Arts (24), Biology (38), Business (63), Chemistry (12), Computers (15), Computer Science (1357), Economics (6), Education (4), Environment (12), Events (35), History (2), Law (12), Mathematics (71), Medicine (29), Philosophy (7), Physics (4), Psychology (3), and Science (10).