# Towards a Statistically Semantic Web

## Gerhard Weikum

weikum@mpi-sb.mpg.de

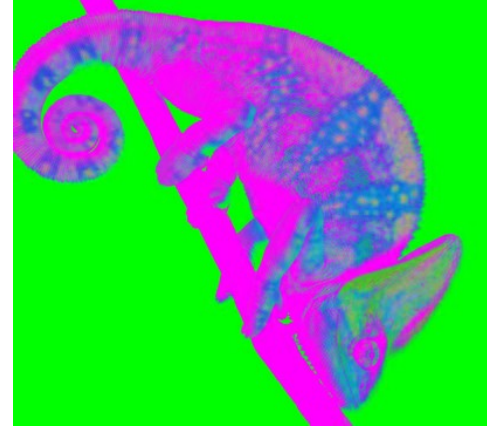http://www.mpi-sb.mpg.de/~weikum/

# Outline

- Motivation and Challenges

- Search                    **(XML, Ontologies)**

- Speed                     **(Top-k Query Processing)**

- Self-Organization    **(P2P, Collaborative Search)**

# Chameleon Words in Computer Science



fragment

page

object

segment

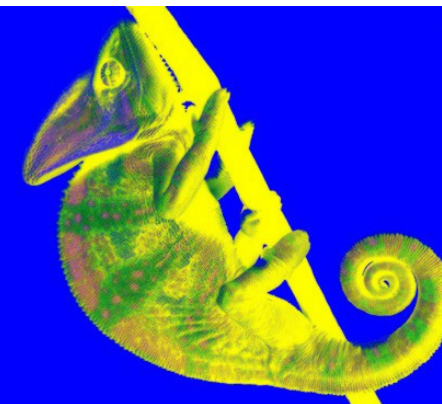performance

ontology

failure

semantics

session

node

service

peer

# Opinions on the Feasibility of the Semantic Web, Universal Data Integration, and Comprehensive Knowledge Bases

**Tim Berners-Lee**: „The Semantic Web is an extension of the current Web in which information is given meaning.“

**Jeff Ullman**: „There is n

**This talk:**

**Names + Statistics → „Semantics“**

**Alon Halevy**: „Structure

**Noam Chomsky**: „Whether there is also a semantics of natural language ... seems to me an open question. Pragmatics must be a central component of linguistic theory.“

**George Lakoff**: „When we have multiple ways of understanding, or 'framing,' a situation, then knowledge, like truth, becomes relative to that understanding.“
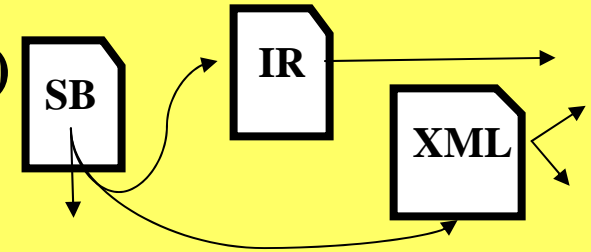
**Confucius**: „Knowledge is to know the extent of one's ignorance.“
孔夫子

# A Few Challenging Queries
## (on Web / Deep Web / Intranet / Personal Info)

- **Which professors from Saarbruecken (SB) are teaching IR and have research projects on XML?**



- **Which gene expression data from Barrett tissue in the esophagus exhibit high levels of gene A01g?**

- **Which drama has a scene in which a woman makes a prophecy to a Scottish nobleman that he will become king?**

- **Who was the woman from Paris that I met at the PC meeting where Paolo Atzeni was PC Chair?**

- **Are there any published theorems that are equivalent to or subsume my latest mathematical conjecture?**

# What if the Semantic Web Existed and All Information Were in XML?



**Professor**
- *Name:* Gerhard Weikum
- *Address* ... *City*: SB
- *Teaching:*

**Lecturer**
- *Name:* Ralf Schenkel
- *Address:* Max-Planck Institute for CS, Germany
- *Teaching:*
- *Interests:* Semistructured Data, IR

**Course**
- *Title:* IR
- *Description:* Information retrieval ...
- *Syllab...*
- *Boo...*

**Seminar**
- *Contents:* Ranked Search ...
- *Literature*

**Project**
- *Title:* Intelligent ...arch ...ML ...Data
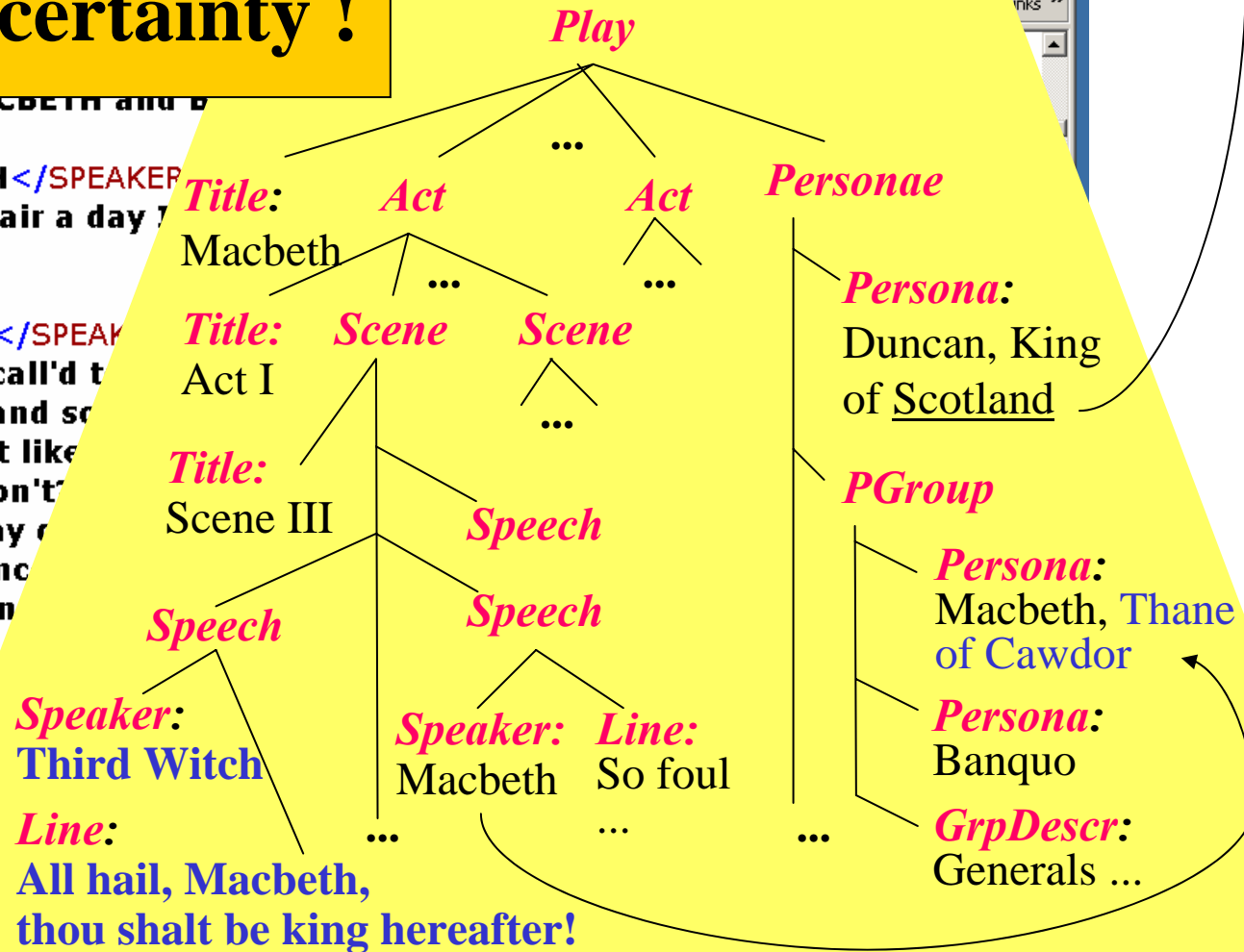- *Sponsor:* German Science Foundation

**Book**
- *Title:* Statistical Language Models
- ...

## Challenge: Diversity !

# What if the Semantic Web Existed and All Information Were in XML?

**Challenge: Uncertainty !**

C:\jsdk2.1\webpages\macbeth.xml - Microsoft Internet Explorer von T-Online

Datei   Bearbeiten   Ansicht   Favoriten   Extras   ?
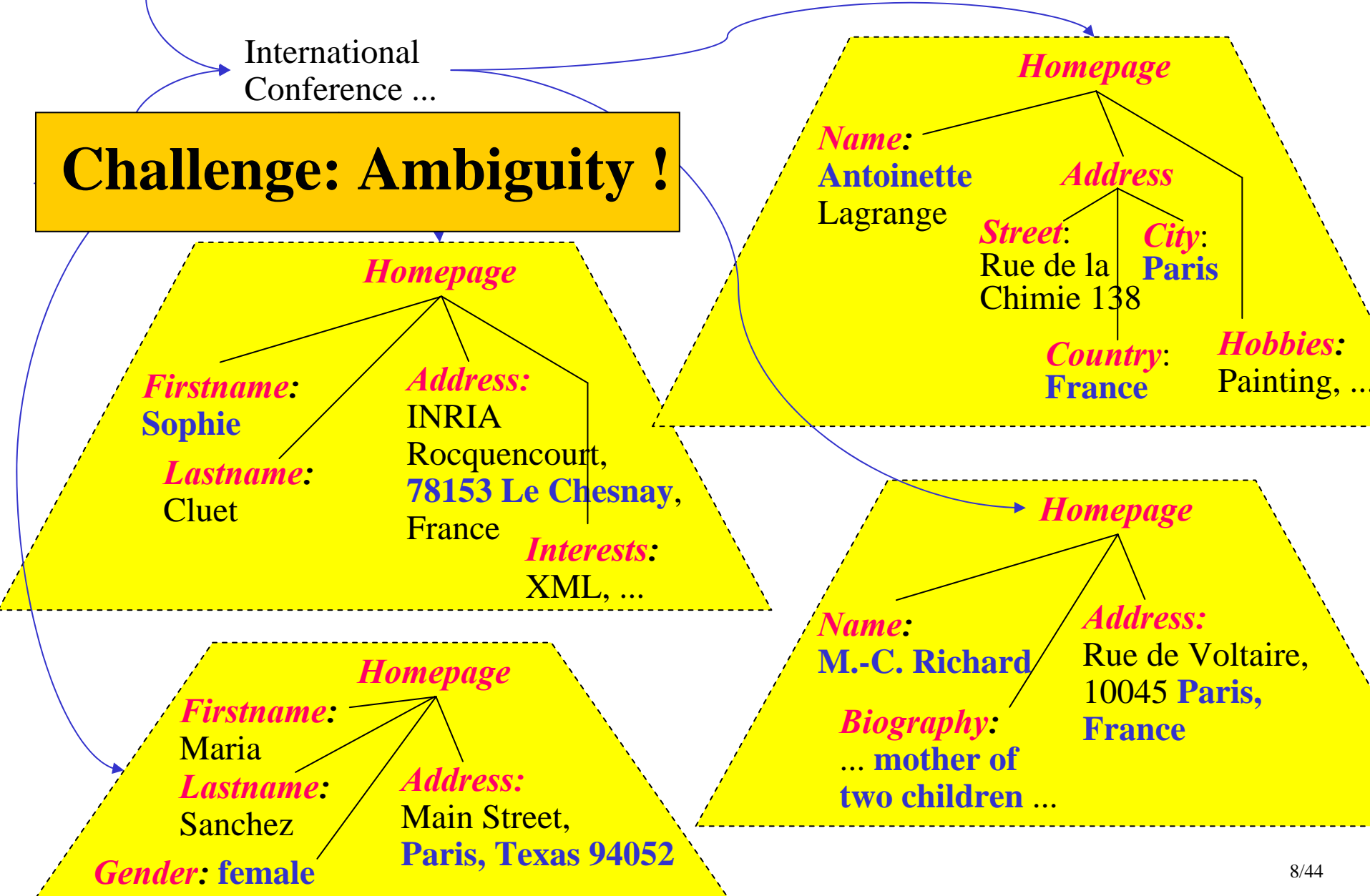
```
<STAGEDIR>Enter MACBETH and B
- <SPEECH>
   <SPEAKER>MACBETH</SPEAKER
   <LINE>So foul and fair a day
 </SPEECH>
- <SPEECH>
   <SPEAKER>BANQUO</SPEAK
   <LINE>How far is't call'd t
   <LINE>So wither'd and so
   <LINE>That look not like
   <LINE>And yet are on't
   <LINE>That man may
   <LINE>By each at onc
   <LINE>Upon her skin
   <LINE>And yet your
   <LINE>That you are
 </SPEECH>
- <SPEECH>
   <SPEAKER>MACB
   <LINE>Speak, if
 </SPEECH>
- <SPEECH>
   <SPEAKER>First Witch</SPEAKER>
   <LINE>All hail, Macbeth! hail to thee, thane of Glamis!</LINE>
```

Fertig                                                             Arbeitsplatz

Tree diagram:

*Play*
- *Title:* Macbeth
- *Act*
  - *Title:* Act I
    - *Scene*
      - *Title:* Scene III
        - *Speech*
          - *Speaker:* Third Witch
          - *Line:* All hail, Macbeth, thou shalt be king hereafter!
        - *Speech*
          - *Speaker:* Macbeth
          - *Line:* So foul
- *Act*
  - *Scene*
- *Personae*
  - *Persona:* Duncan, King of Scotland
  - *PGroup*
    - *Persona:* Macbeth, Thane of Cawdor
    - *Persona:* Banquo
    - *GrpDescr:* Generals ...

# What if the Semantic Web Existed and All Information Were in XML?

International Conference ...

**Challenge: Ambiguity !**

**Homepage**
- *Name*: **Antoinette** Lagrange
- *Address*:
  - *Street*: Rue de la Chimie 138
  - *City*: **Paris**
  - *Country*: **France**
- *Hobbies*: Painting, ...

**Homepage**
- *Firstname*: **Sophie**
- *Lastname*: Cluet
- *Address*: INRIA Rocquencourt, **78153 Le Chesnay**, France
- *Interests*: XML, ...

**Homepage**
- *Name*: **M.-C. Richard**
- *Address*: Rue de Voltaire, 10045 **Paris, France**
- *Biography*: ... **mother of two children** ...

**Homepage**
- *Firstname*: Maria
- *Lastname*: Sanchez
- *Gender*: **female**
- *Address*: Main Street, **Paris, Texas 94052**

# Observations and Challenges (1)

*Observation:*
Despite all structure, tags, and „semantic" metadata, information will exhibit *diversity*, *ambiguity*, *uncertainty*

*Implication:*
Information search faces IR dilemma – drown in results or almost empty result – and thus needs *ranked retrieval*
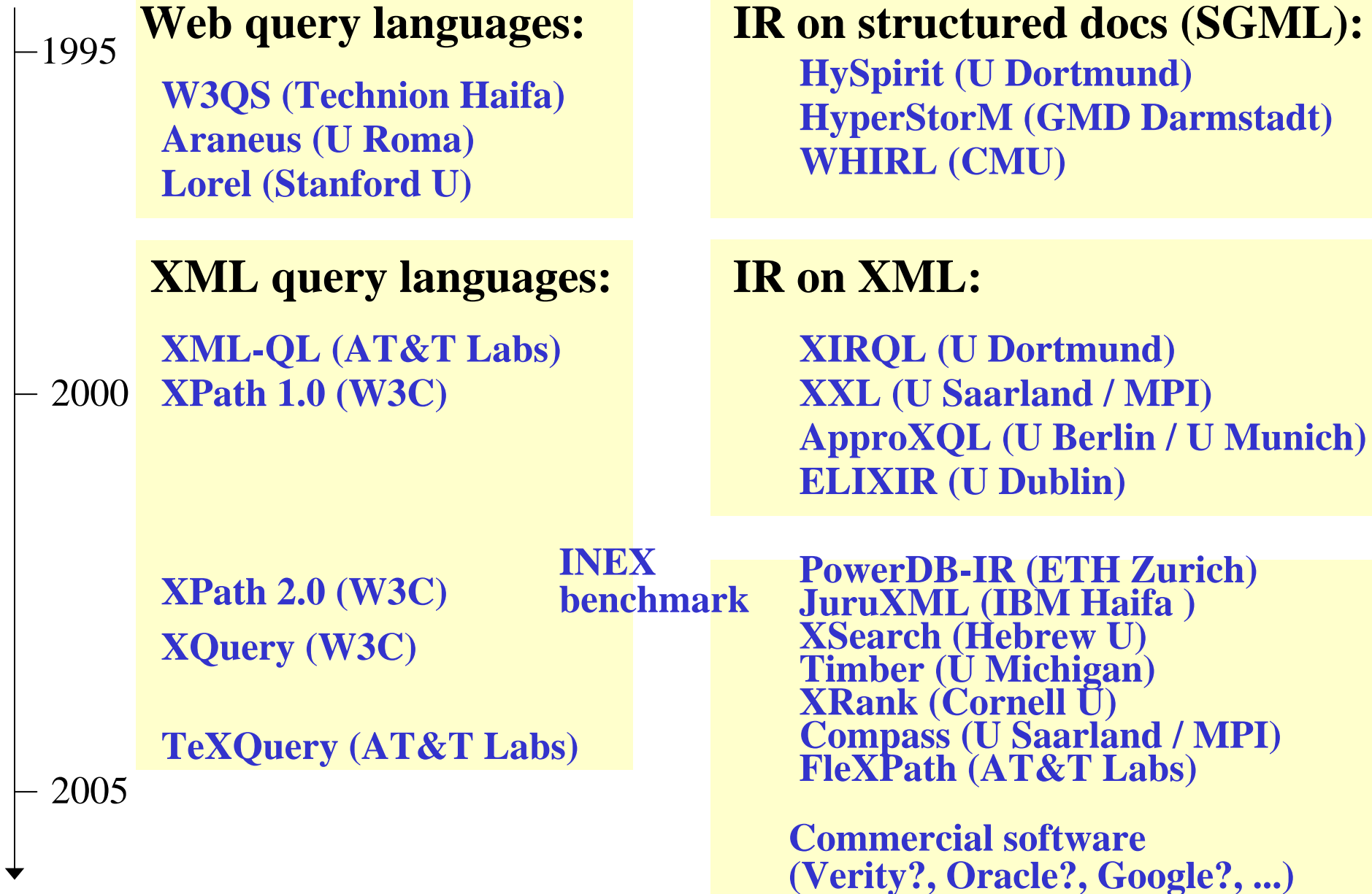
*Challenge:*
Combine the best of *precise querying* from DB world with *vague search* and *relevance assessment* from IR, Web & learning communities
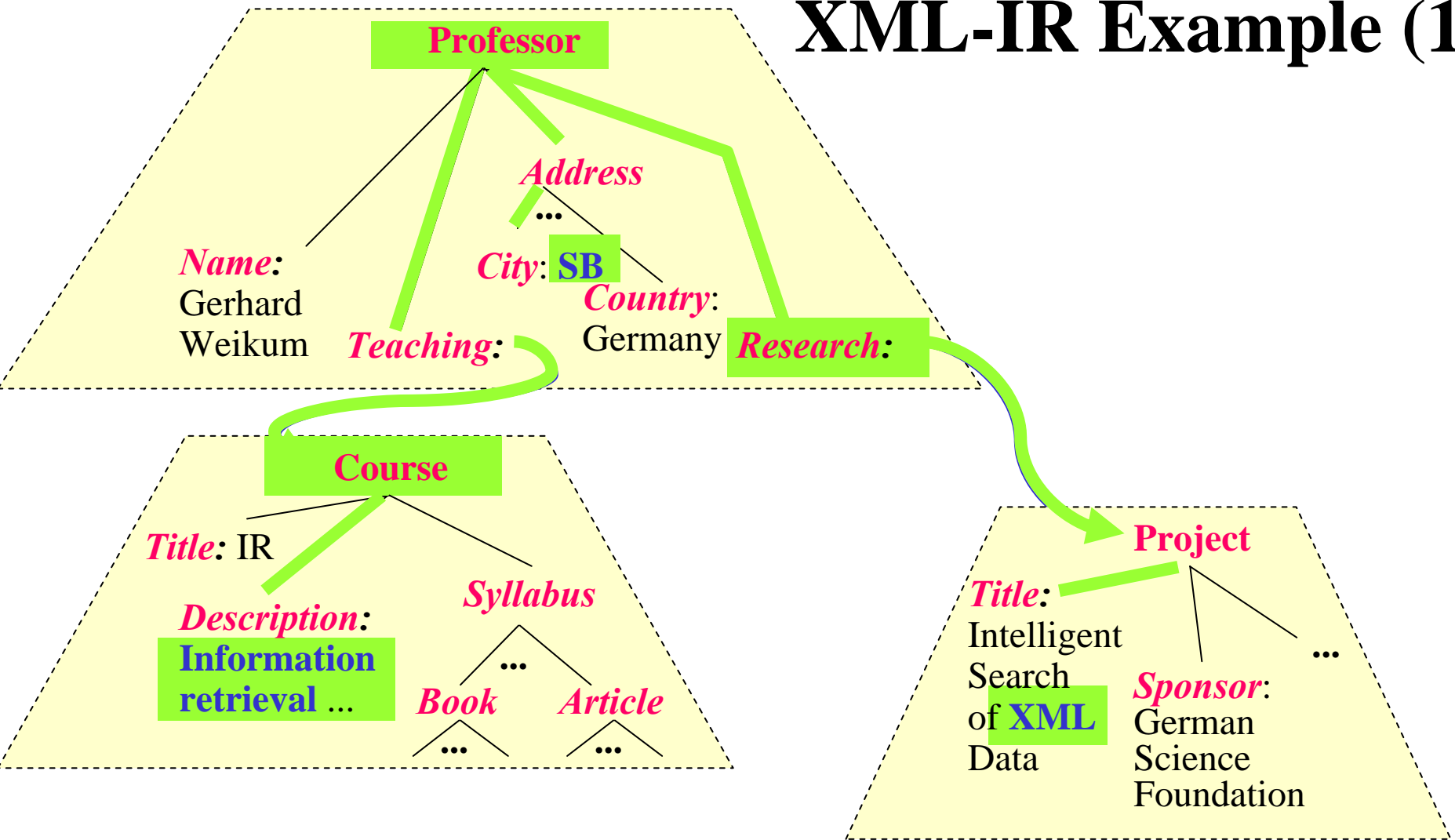[ and expressive *logical inferences* from AI ]

# Outline

✓ Motivation and Challenges

- Search **(XML, Ontologies)**

- Speed **(Top-k Query Processing)**

- Self-Organization **(P2P Collaborative Search)**

# XML-IR: History and Related Work

**1995**

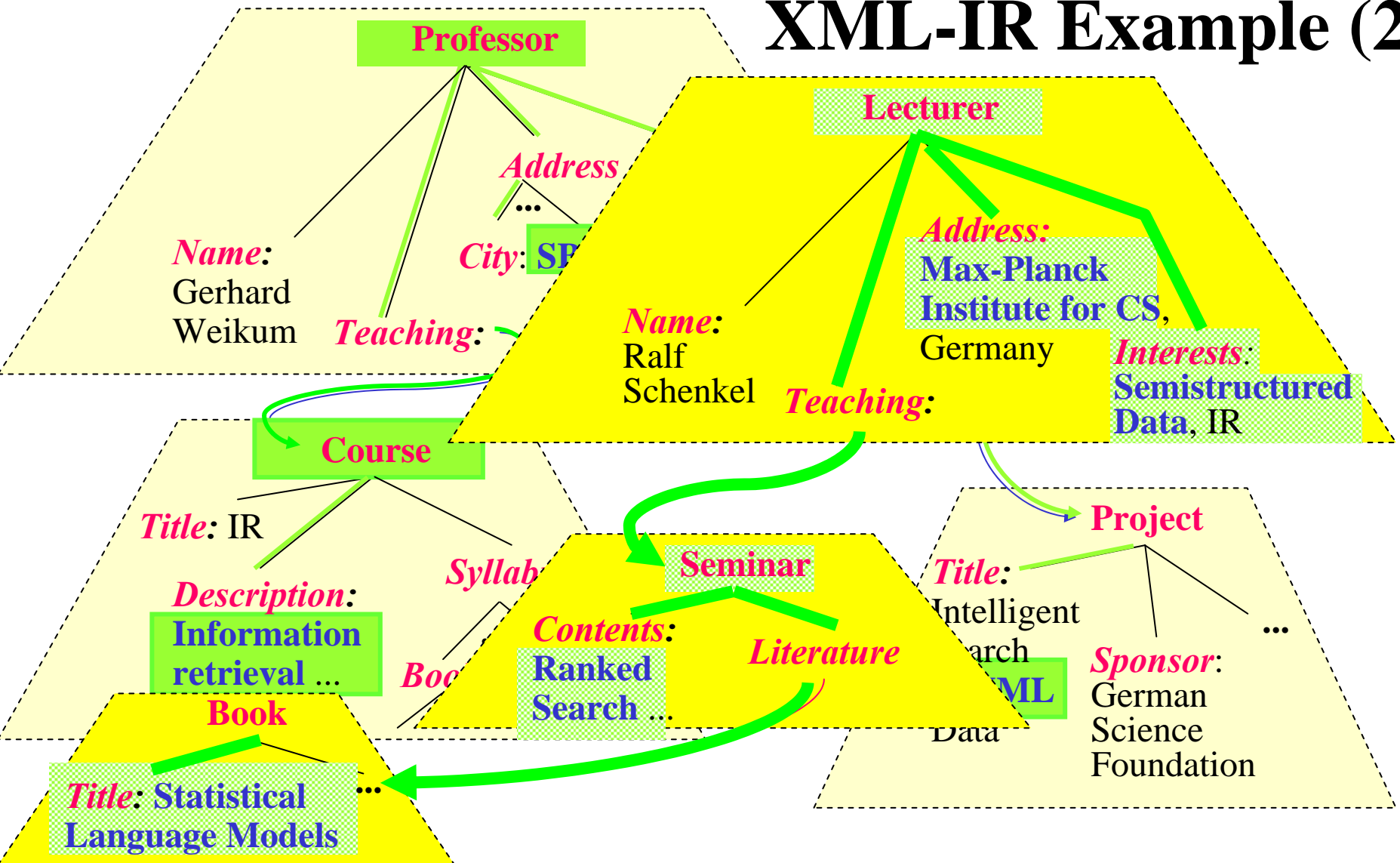**Web query languages:**

**W3QS (Technion Haifa)**
**Araneus (U Roma)**
**Lorel (Stanford U)**

**IR on structured docs (SGML):**

**HySpirit (U Dortmund)**
**HyperStorM (GMD Darmstadt)**
**WHIRL (CMU)**

**XML query languages:**

**XML-QL (AT&T Labs)**
**XPath 1.0 (W3C)**

**2000**

**IR on XML:**

**XIRQL (U Dortmund)**
**XXL (U Saarland / MPI)**
**ApproXQL (U Berlin / U Munich)**
**ELIXIR (U Dublin)**

**INEX benchmark**

**PowerDB-IR (ETH Zurich)**
**JuruXML (IBM Haifa )**
**XSearch (Hebrew U)**
**Timber (U Michigan)**
**XRank (Cornell U)**
**Compass (U Saarland / MPI)**
**FleXPath (AT&T Labs)**

**XPath 2.0 (W3C)**

**XQuery (W3C)**

**TeXQuery (AT&T Labs)**

**2005**

**Commercial software (Verity?, Oracle?, Google?, ...)**

# XML-IR Example (1)

**Professor**

*Address*
...

*Name*:
Gerhard
Weikum

*City*: **SB**

*Country*:
Germany

*Research*:

*Teaching*:

**Course**

*Title*: IR

*Description*:
**Information retrieval ...**

*Syllabus*
...
*Book*  *Article*
...  ...

**Project**

*Title*:
Intelligent
Search
of **XML**
Data

*Sponsor*:
German
Science
Foundation

...

Select *P, C, R* From *Index*
Where *Professor* As *P* And *P* = „ *Saarbruecken*"
And *P//* *Course* = „ *IR*" As *C* And *P//* *Research* = „ *XML*" As *R*

# XML-IR Example (2)

**Professor**

*Name:* Gerhard Weikum

*Address*
...
*City*: SB

*Teaching:* ~

**Lecturer**

*Name:* Ralf Schenkel

*Address:* Max-Planck Institute for CS, Germany

*Interests*: Semistructured Data, IR

*Teaching:*

**Course**

*Title:* IR

*Description:* Information retrieval ...

*Syllab...*

*Boo...*

**Seminar**

*Contents:* Ranked Search ...

*Literature*

**Project**

*Title:* Intelligent Search XML Data

*Sponsor*: German Science Foundation

**Book**

*Title:* Statistical Language Models
...

Select *P, C, R* From *Index*
Where *~Professor* As *P* And *P* = „*~Saarbruecken"*
And *P//~Course* = „*~IR"* As *C* And *P//~Research* = „*~XML"* As *R*

# XML-IR Concepts

**applicable to both XML and HTML data graphs**

**Where clause: conjunction of restricted** *path expressions*
                **with binding of variab**

*Elementary conditions* **on names and c**

<div style="background:#9EF5C9;">

*Query result:*
• **query is a pattern**
  **with relaxable conditions**
• **results are approximate**
  **matches to query**
  **with similarity scores**

</div>

Select *P, C, R* From *Index*
Where *~Professor* As *P*
And *P = „Saarbruecken"*
And *P//~Course = „Information Retrieval"* As *C*
And *P//~Research = „~XML"* As *R*

**„Semantic"** *similarity conditions* **on names and contents**
        *~Research = „~XML"*

**Relevance scoring based on**
        **tf*idf similarity of contents,**
        **ontological similarity of names,**
        **aggregation of local scores into global scores**

# XML-IR Scoring Model

**Homepage**

*Title*: Professor

*Name*: Gerhard Weikum

*Address*

*City*: SB

*Teaching*:

*Research*:

**Course**

*Title*: IR

*Description*: ... IR ... XML ... ... IR ... DB ... IR

**Project**

...

*Description*: ... XML ... XSLT ...

**Course**

*Title*: Ranked Search

*Syllabus*

*Book* ...

**Book**

*Title*: Statistical Language Models

*Contents*: ... IR ... ... IR ...

**local score** for elementary condition: based on tf*idf-style statistics for node or node context with score propagation

**global score** for query:
$\sum$ local scores  * compactness

**compactness** of result:
max{$\sum$ node & edge weights | graph connecting matching nodes}
$\rightarrow$ generalized MST (related to Steiner trees)

# XML-IR Scoring Model



**local score** for
elementary condition:
based on tf*idf-style statistics
for node or node context
with score propagation

**global score** for query:
∑ local scores * compactness

**compactness** of result:
max{∑ node & edge weights |
    graph connecting
    matching nodes}
→ generalized MST
  (related to Steiner trees)

# XML-IR Scoring Model

**Homepage**

*Title*: Professor

*Name*: Gerhard Weikum

*Address*

*City*: SB

*Teaching*:

*Research*:

**Course**

*Title*: IR

*Description*: ... IR ... XML ... ... IR ... DB ... IR

**Project**

*Description*: ... XML ... XSLT ...

**Course**

*Title*: Ranked Search

*Syllabus*

*Book* ...

**Book**

*Title*: Statistical Language Models

*Contents*: ... IR ... ... IR ...

**local score** for

**Efficient score computation: heuristics work; advanced algorithms is open issue**

statistics

ext

on

**global score** for query:
$\sum$ local scores * compactness

**compactness** of result:
max{$\sum$ node & edge weights | graph connecting matching nodes}
$\rightarrow$ generalized MST (related to Steiner trees)

# On Thesauri and Ontologies

**Taxonomy**: classification of concepts into groups (and trees of groups)

**Thesaurus**: repository („treasure") of synonyms
(and other relationships between words and concepts)

**Ontology**: metaphysical study of the nature of being & existence

**Ontology (new definition)**: structured repository of knowledge
with a description of concepts and relationships,
possibly in the form of description logics formula

## Reasoning on Ontologies and Thesauri:

Professor $\subseteq$ Lecturer $\cap$ $\exists$ hasStaff.Secretary
Teaching $\supseteq$ Cour
Professor $\subseteq$ Acad
Academician $\subseteq$ H
Human $\subseteq$ Carniv
...

**poor man's ontology:
pragmatic, rich, efficient**

$\rightarrow$ logical inferences
with sub-FOL calculus

$\rightarrow$ transitive closures,
shortest paths, etc.
along generalizations

# Example WordNet

File   History   Options   Help

Search Word: woman | Redisplay Overview

Searches for woman:   Noun | Senses:

1 of 4 senses of woman

Sense 1
woman, adult female -- (an adult female person (as opposed to a man); "the woman kept house while the man hunted")
   => Eve -- ((Old Testament) Adam's wife in Judeo-Christian mythology: the first woman and mother of the human race; God created Eve from Adam's rib and placed Ada
   => black woman -- (a woman who is Black)
   => white woman -- (a woman who is White)
   => yellow woman -- (offensive term for an Oriental woman)

**woman, adult female – (an adult female person)**
   **=> amazon, virago – (a large strong and aggressive woman)**
   **=> donna -- (an Italian woman of rank)**
   **=> geisha, geisha girl -- (...)**
   **=> lady (a polite name for any woman)**

   **...**
   **=> wife – (a married woman, a man's partner in marriage)**

   **=> witch – (a being, usually female, imagined to
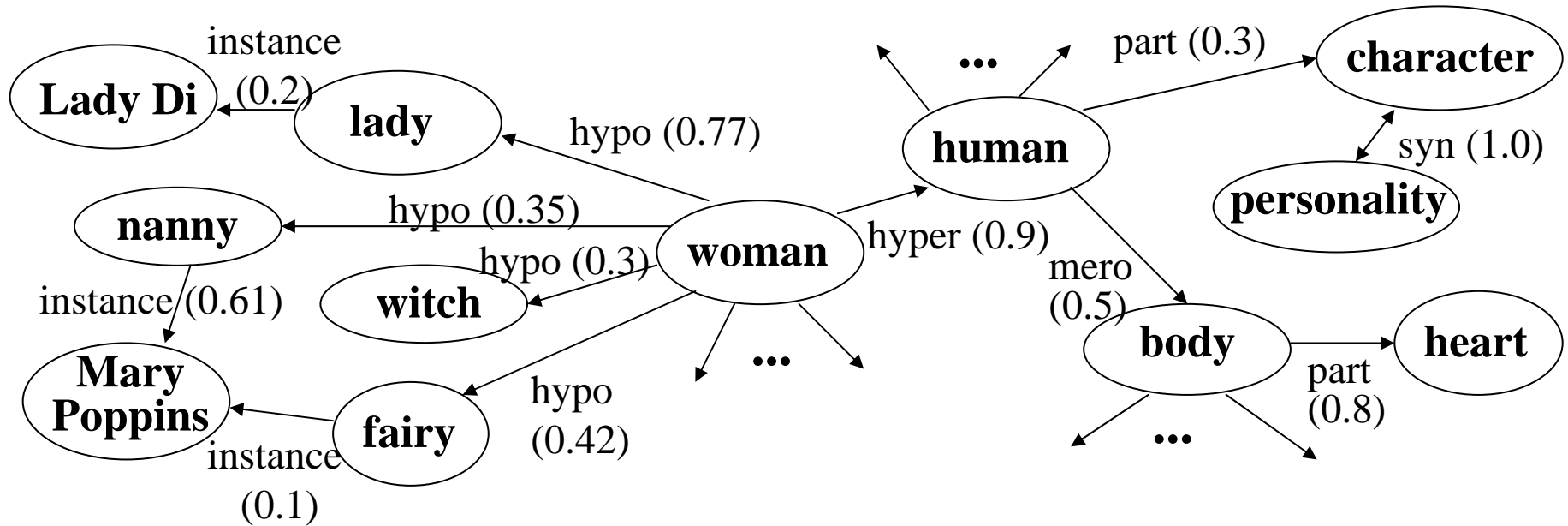                 have special powers derived from the devil)**

=> maenad -- (an unnaturally frenzied or distraught woman)
=> matron, head nurse -- (a woman in charge of nursing in a medical institution)

"Hyponyms (...is a kind of this), brief" search for noun "woman"

# Ontology Graph

An ontology graph is a directed graph with concepts (and their descriptions) as nodes and semantic relationships as edges (e.g., hypernyms).



**Weighted edges capture strength of relationships**
**→ key for identifying closely related concepts**

# Statistics for Weighted Ontological Relations

Gather statistics from large corpus or by (focused) Web crawl

*Various correlation measures for sim(c1, c2):*

**Dice coefficient:**
$$\frac{2\left|\{docs\ with\ c1\}\cap\{docs\ with\ c2\}\right|}{\left|\{docs\ with\ c1\}\right|\ +\ \left|\{docs\ with\ c2\}\right|}$$

**Jaccard coefficient:**
$$\frac{\left|\{docs\ with\ c1\}\cap\{docs\ with\ c2\}\right|}{\left|\{docs\ with\ c1\}\right|\ +\ \left|\{docs\ with\ c2\}\right|-\left|\{docs\ with\ c1\ and\ c2\}\right|}$$

**Conditional probabilites:**
$$P[\,doc\ has\ c1\mid doc\ has\ c2\,]$$

**Transitive similarity:**

$$sim^*(c1,\ cn)\ =\ max\{\prod_{i=1..n-1}sim(c_i,c_{i+1})\mid all\ paths\ from\ c1\ to\ cn\}$$

compute by (adaptation of) Dijkstra's shortest-path algorithm

# Benefits from Ontology Service

Ontology service accessible via SOAP or RMI
Ontology filled with WordNet, geo gazetteer,
                                      focused crawl results, extracted tables & forms

useful for:

- Threshold-based query expansion

- Query keyword disambiguation

- Support for automatic tagging of HTML
  and enhanced XML tags

- Mapping of concept-value query conditions
  onto Deep-Web portals

# Query Expansion

***Threshold-based query expansion:***

substitute ~w by $(c_1 \mid ... \mid c_k)$ with all $c_i$ for which $\mathrm{sim}(w, c_i) \geq \delta$

*„Old hat" in IR; highly disputed for danger of topic dilution*

***Approach to careful expansion:***
- determine phrases from query or best initial query results
  (e.g., forming 3-grams and looking up ontology/thesaurus entries)
- if uniquely mapped to one concept
  then expand with synonyms and weighted hyponyms

Problem: choice of threshold $\delta \;\rightarrow\;$ see Top-k QP

# Query Expansion Example

From TREC 2004 Robust Track:

**Title: International Organized Crime**

**Description:** Identify organizations that participate in international criminal activity, the activity, and, if possible, collaborating organizations and the countries involved.

**Query = {international[0.145|1.00],**
**~META[1.00|1.00][{gangdom[1.00|1.00], gangland[0.742|1.00],**
**organ[0.213|1.00] & crime[0...**
**mafia[0.154|1.00], "sicilian[0...**
**!black[0.066|1.00] & hand[0...**
**organ[0.213|1.00], crime[0.31...**
**columbian[0.686|0.20], cartel...**

**Let us take, for example, the case of Medellin cartel's boss Pablo Escobar. Will the fact that he was eliminated change anything at all? No, it may perhaps have a psychological effect on other drug dealers but,** ...

**... for organizing the illicit export of metals and import of arms. It is extremely difficult for the law-enforcement organs to investigate and stamp out corruption among leading officials.** ...

**A parliamentary commission accused Swiss prosecutors today of doing little to stop drug and money-laundering international networks from pumping billions of dollars through Swiss companies.** ...
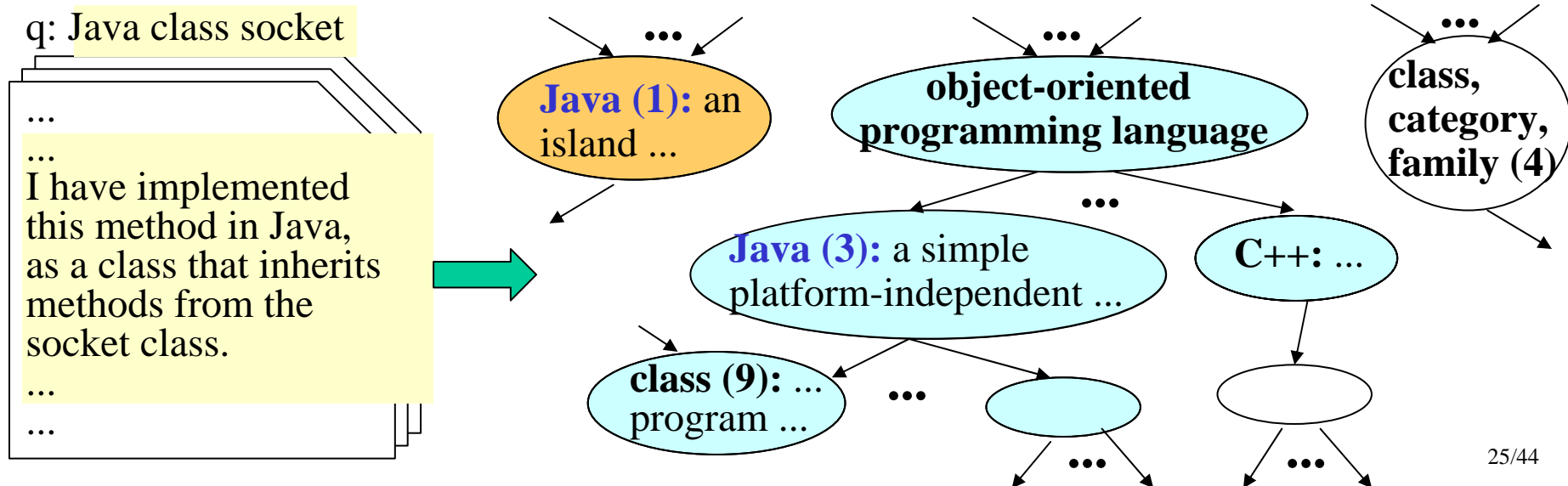
**Results:**

1. Interpol Chief on Fight Against...
2. Economic Counterintelligen...
3. Dresden Conference Views...
4. Report on Drug, Weapons S...
5. SWITZERLAND CALLED SOFT ON CRIME
...

# Keyword-to-Concept Mapping and Word Sense Disambiguation (1)

Example: „**Java class socket**" vs. „**Java** beach snorkeling"
Which concept should „Java" be mapped to for query expansion?

*Note: unlike in LSI or pLSI, concepts are explicit, not latent!*

***Approach for query keyword disambiguation:***
- form **contexts con(w) and con($c_i$)**
  for keyword w and potential target concepts $c_i \in \{c_1, ..., c_k\}$
- bag-of-words similarity **sim(con(w), con(c))** based on cos or KL diff
- choose concept $\text{argmax}_c \{\text{sim}(\text{con}(w), \text{con}(c))\}$

q: Java class socket

...
...
I have implemented this method in Java, as a class that inherits methods from the socket class.
...
...

**Java (1):** an island ...

**object-oriented programming language**

class, category, family (4)

**Java (3):** a simple platform-independent ...
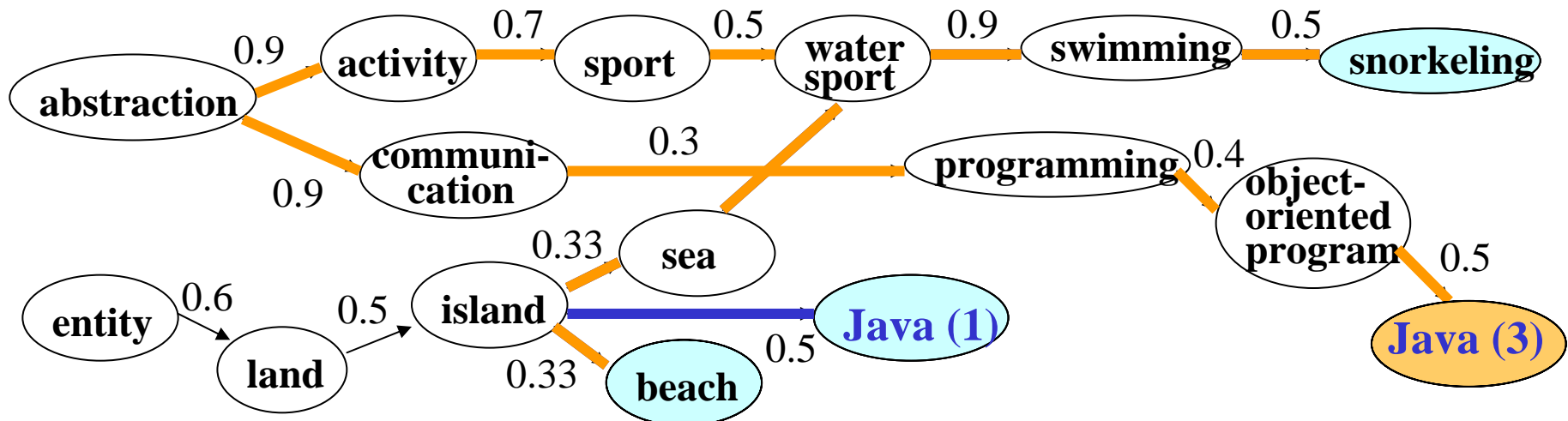
**C++:** ...

**class (9):** ... program ...

# Keyword-to-Concept Mapping and Word Sense Disambiguation (2)

Example: „**Java** class socket" vs. „**Java beach snorkeling**"
Which concept should „Java" be mapped to for query expansion?

*Alternative approach for query keyword disambiguation:*
• consider potential target concepts for all query keywords together
• choose **concepts ($c_1$, ..., $c_m$) for words ($w_1$, ..., $w_m$)** according to
  **sim * compactness** where
    • **sim** ~ aggregation of $con(w_i)$-to-$con(c_i)$ similarities,
    • **compactness** ~ weight of MST for $\{c_1, ..., c_m\}$

# Observations and Challenges (2)

*Observation:*

Explicit ontologies/thesauri and statistical models
need to be combined for ranked retrieval of richly annotated
but highly heterogeneous XML data

*Challenges:*

- Develop full-fledged statistical language model
  for XML subgraph scoring
- Constructing statistically quantified ontologies
  from rich sources (WordNet, Wikipedia, bookmarks, etc.)
- Combine uncertainty of automatic tagging and
  query mapping with query result ranking
  in a comprehensive probabilistic algebra
- Efficient query processing and optimization

# Outline

✓ Motivation and Challenges

✓ Search             **(XML, Ontologies)**

- Speed             **(Top-k Query Processing)**

- Self-Organization      **(P2P Collaborative Search)**

# Top-k Query Processing with Scoring

q: algorithm
 performance
 z-transform

B+ tree on terms

| algorithm | ... | performance | ... | z-transform |
|---|---|---|---|---|
| **17: 0.3** | | **12: 0.5** | | **11: 0.6** |
| **44: 0.4** | | **14: 0.4** | | **17: 0.1** |
| **52: 0.1** | | **28: 0.1** | | **28: 0.7** |
| **53: 0.8** | | **44: 0.2** | | ⋮ |
| **55: 0.6** | | **51: 0.6** | | |
| ⋮ | | **52: 0.3** | | |
| | | ⋮ | | |

index lists with
(DocId, tf*idf)
sorted by DocId

Google:
> 10 mio. terms
> 4 bio. docs
> 2 TB index

<u>Given:</u> query $q = t_1\ t_2\ ...\ t_z$ with z (conjunctive) keywords
 similarity scoring function score(q,d) for docs $d \in D$, e.g.: $\vec{q} \cdot \vec{d}$
<u>Find:</u> top k results with regard to score(q,d) = aggr{$s_i(d)$} (e.g.: $\Sigma_{i \in q}\ s_i(d)$)

*Naive QP algorithm:*
 candidate-docs := $\varnothing$;
 for i=1 to z do {
 candidate-docs := candidate-docs $\cup$ index-lookup(ti) };
 for each dj $\in$ candidate-docs do {compute score(q,dj)};
 sort candidate-docs by score(q,dj) descending;

# TA (Fagin'01; Güntzer/Kießling/Balke; Nepal et al.)

scan all lists $L_i$ (i=1..m) in parallel:
   consider $d_j$ at position $pos_i$ in Li;
   $high_i := s_i(d_j)$;
   if $d_j \notin$ top-k then {
      look up $s_\nu(d_j)$ in all lists $L_\nu$ with $\nu \neq i$; // random access
      compute $s(d_j) := aggr \{s_\nu(dj) \mid \nu=1..m\}$;
      if $s(d_j) >$ min score among top-k then
        add $d_j$ to top-k and remove min-score d from top-k; };
if **min score among top-k** $\geq$ **aggr {$high_\nu \mid \nu=1..m$}** then exit;

> *but random accesses*
> *are expensive !*
> *→ TA-sorted*
> *→ Prob-sorted*

m=3
aggr: sum
k=2

| L1 | L2 | L3 |
|---|---|---|
| f: 0.5 | a: 0.55 | h: 0.35 |
| b: 0.4 | b: 0.2 | d: 0.35 |
| c: 0.35 | f: 0.2 | b: 0.2 |
| a: 0.3 | g: 0.2 | a: 0.1 |
| h: 0.1 | c: 0.1 | c: 0.05 |
| d: 0.1 | | f: 0.05 |

**top-k:**

~~f: 0.75~~

a: 0.95

b: 0.8

*applicable to XML data:*
*course = „~ Internet"  and  ~topic = „performance"*

# TA-Sorted

scan index lists in parallel:
consider $d_j$ at position $pos_i$ in $L_i$;
$E(d_j) := E(d_j) \cup \{i\}$; $high_i := s_i(q,d_j)$;
bestscore($d_j$) := aggr$\{x_1, ..., x_m)$
   with $x_i := s_i(q,d_j)$ for $i \in E(d_j)$, $high_i$ for $i \notin E(d_j)$;
worstscore(dj) := aggr$\{x_1, ..., x_m)$
   with $x_i := si(q,d_j)$ for $i \in E(d_j)$, 0 for $i \notin E(d_j)$;
top-k := k docs with largest worstscore;
if min worstscore among top-k $\geq$
   max bestscore{d | d not in top-k} then exit;

m=3
aggr: sum
k=2

| | | |
|---|---|---|
| f: 0.5 | a: 0.55 | h: 0.35 |
| b: 0.4 | b: 0.2 | d: 0.35 |
| c: 0.35 | f: 0.2 | b: 0.2 |
| a: 0.3 | g: 0.2 | a: 0.1 |
| h: 0.1 | c: 0.1 | c: 0.05 |
| d: 0.1 | | f: 0.05 |

**top-k:**
a: 0.95
b: 0.8

**candidates:**
f: 0.7 + ? $\leq$ 0.7 + 0.1
h: 0.45 + ? $\leq$ 0.45 + 0.2
c: 0.35 + ? $\leq$ 0.35 + 0.3
d: 0.35 + ? $\leq$ 0.35 + 0.3
g: 0.2 + ? $\leq$ 0.2 + 0.4

# Top-k Queries with Probabilistic Guarantees

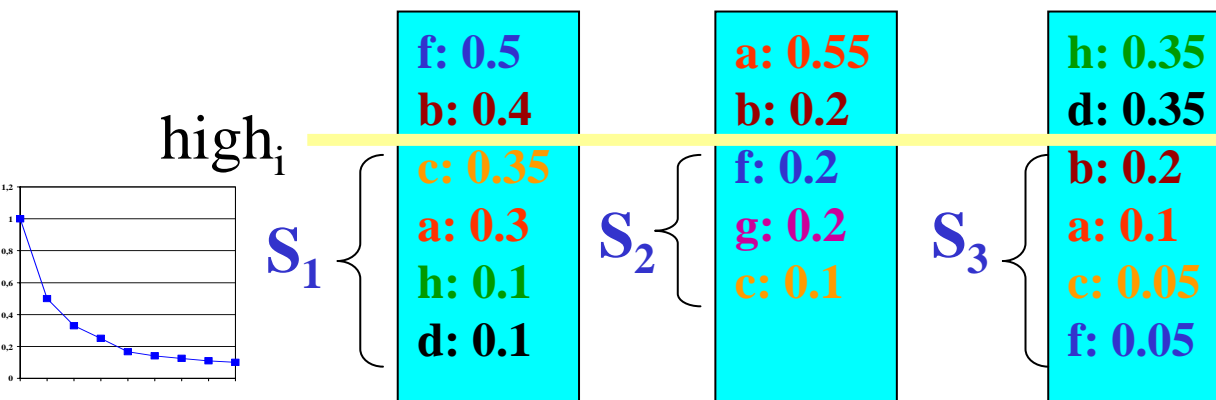**TA family of algorithms based on invariant (with sum as aggr)**

$$\underbrace{\sum_{i \in E(d)} s_i(d)}_{\textbf{worstscore(d)}} \leq \; s(d) \; \leq \underbrace{\sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} high_i}_{\textbf{bestscore(d)}}$$

**Relaxed into probabilistic invariant**

$$p(d) := P[s(d) > \delta] = P[\sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} S_i > threshold]$$

$$= P[\sum_{i \notin E(d)} S_i > threshold - \sum_{i \in E(d)} s_i(d)] =: P[\sum_{i \notin E(d)} S_i > \delta'] \leq \varepsilon$$

**where the RV $S_i$ has some (postulated and/or estimated) distribution in the interval $(0, high_i]$**

high$_i$



| $S_1$ | $S_2$ | $S_3$ |
|---|---|---|
| f: 0.5 | a: 0.55 | h: 0.35 |
| b: 0.4 | b: 0.2 | d: 0.35 |
| c: 0.35 | f: 0.2 | b: 0.2 |
| a: 0.3 | g: 0.2 | a: 0.1 |
| h: 0.1 | c: 0.1 | c: 0.05 |
| d: 0.1 | | f: 0.05 |

- *Discard candidates with p(d) ≤ ε*
- **Exit index scan when candidate list empty**

# Probabilistic Threshold Test

cand doc d
with
$2 \notin E(d)$,
$3 \notin E(d)$

$f_2(x)$  $\oplus$  $f_3(x)$  $\rightarrow$  *Convolution* $(f_2(x), f_3(x))$

1   $high_2$   0       1   $high_3$   0       2   $\delta(d)$   0

- **postulating *uniform or Zipf* score distribution in [0, $high_i$]**
  - **compute convolution using LSTs**
  - **use Chernoff-Hoeffding tail bounds or generalized bounds for correlated dimensions (Siegel 1995)**
- **fitting *Poisson* distribution (or Poisson mixture)**
  - **over equidistant values:** $P[d = v_j] = e^{-\alpha_i} \dfrac{\alpha_i^{\,j-1}}{(j-1)!}$
  - **easy and exact convolution**
- **distribution approximated by *histograms*** *engineering-wise*
  - **precomputed for each dimension** *histograms work best!*
  - **dynamic convolution at query-execution time**

**with *independent* Si's or with *correlated* Si's**

# Performance Results for .Gov Queries

*on .GOV corpus from TREC-12 Web track:*
**1.25 Mio. docs (html, pdf, etc.)**

**50 keyword queries, e.g.:**
- *„Lewis Clark expedition",*
- *„juvenile delinquency",*
- *„legalization Marihuana",*
- *„air bag safety reducing injuries death facts"*

*speedup by factor 10 at high precision/recall (relative to TA-sorted);*

*aggressive queue mgt. even yields factor 100 at 30-50 % prec./recall*

|  | **TA-sorted** | **Prob-sorted (smart)** |
|---|---|---|
| **#sorted accesses** | **2,263,652** | **527,980** |
| **elapsed time [s]** | **148.7** | **15.9** |
| **max queue size** | **10849** | **400** |
| **relative recall** | **1** | **0.69** |
| **rank distance** | **0** | **39.5** |
| **score error** | **0** | **0.031** |

# .Gov Expanded Queries

*on .GOV corpus with query expansion based on WordNet synonyms:*
**50 keyword queries, e.g.:**
- **„juvenile delinquency** *youth minor crime law jurisdiction offense prevention*",
- **„legalization marijuana** *cannabis drug soft leaves plant smoked chewed euphoric abuse substance possession control pot grass dope weed smoke*"

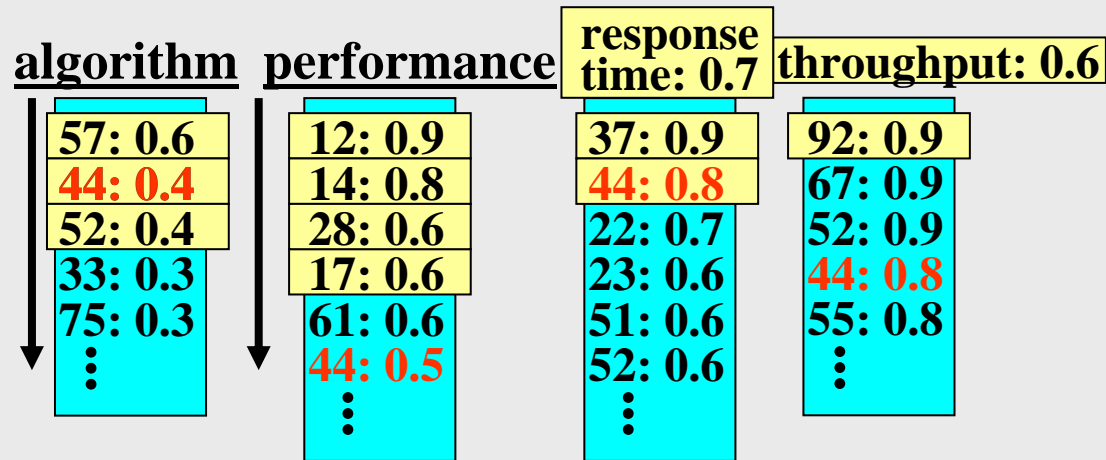|                    | TA-sorted   | Prob-sorted (smart) |
|--------------------|-------------|---------------------|
| #sorted accesses   | 22,403,490  | 18,287,636          |
| elapsed time [s]   | 7908        | 1066                |
| max queue size     | 70896       | 400                 |
| relative recall    | 1           | 0.88                |
| rank distance      | 0           | 14.5                |
| score error        | 0           | 0.035               |

# Handling Ontology-Based Query Expansions
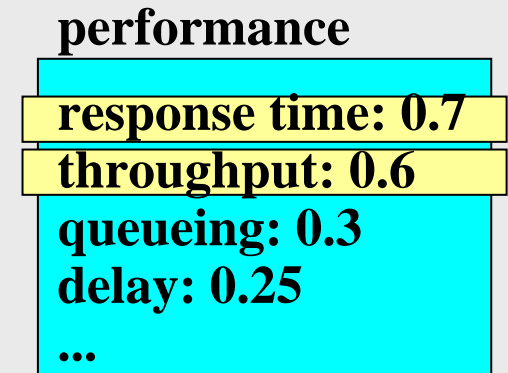
**consider expandable query** *„algorithm and ~performance"*
**with score** $\Sigma_{i \in q} \{ \max_{j \in onto(i)} \{ sim(i,j)*sj(d)) \} \}$

**dynamic query expansion with**
**incremental on-demand merging of additional index lists**

B+ tree index on terms

| algorithm | performance | response time: 0.7 | throughput: 0.6 |
|---|---|---|---|
| 57: 0.6 | 12: 0.9 | 37: 0.9 | 92: 0.9 |
| 44: 0.4 | 14: 0.8 | 44: 0.8 | 67: 0.9 |
| 52: 0.4 | 28: 0.6 | 22: 0.7 | 52: 0.9 |
| 33: 0.3 | 17: 0.6 | 23: 0.6 | 44: 0.8 |
| 75: 0.3 | 61: 0.6 | 51: 0.6 | 55: 0.8 |
| ⋮ | 44: 0.5 | 52: 0.6 | ⋮ |
| | ⋮ | ⋮ | |

ontology / meta-index

performance

response time: 0.7
throughput: 0.6
queueing: 0.3
delay: 0.25
...

**+ much more efficient than threshold-based expansion**
**+ no threshold tuning**
**+ no topic drift**

# Observations and Challenges (3)

**Observation:**

Approximations with ***statistical guarantees*** are key to obtaining ***Web-scale efficiency***
(e.g., TREC'04 Terabyte benchmark:
ca. 25 Mio. docs, ca. 700 000 terms, 5-50 terms per query)

**Challenges:**

- Efficient consideration of ***correlated dimensions***
- Integrated support for all kinds of XML similarity search: content & ontological sim, ***structural sim***
- ***Scheduling*** of index-scan steps and few random accesses
- Integration of top-k operator into ***physical algebra*** and ***query optimizer*** of XML engine

# Outline

✓ Motivation and Challenges

✓ Search               **(XML, Ontologies)**

✓ Speed               **(Top-k Query Processing)**

- Self-Organization     **(P2P Collaborative Search)**

# P2P for Web Search ?

**Given: overlay networks (often DHTs) à la Chord, CAN, P-Grid**
**How do we exploit this technology for keyword queries?**

**Naive idea: use multidimensional keys as** ~~~
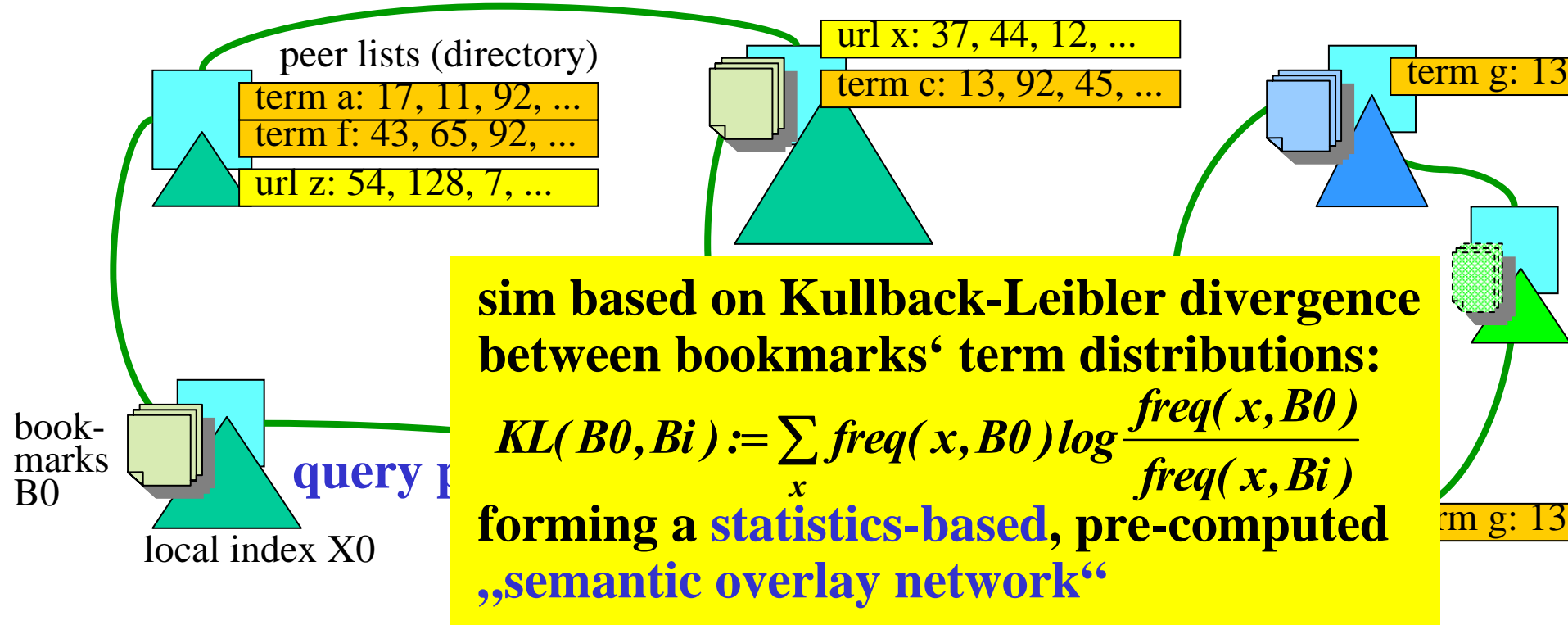                 **or encode doc/query vectors a**

→ **grand challenge for performance at W**
→ **infeasible for very-high-dimensional, v**
→ **no support for ranking (similarity que**
→ **breach with autonomous behavior of i**

**Ongoing projects:**
**PlanetP (Rutgers U)**
**Odissea (Polytec Brooklyn)**
**Pepper (U Duisburg / CMU)**
**Peers (Stanford)**
**Pier (Berkeley)**
**YouServ (IBM)**
**GridVine (EPFL)**
**Minerva (MPI)**
**Evergrow (EU)**
**PeerSDI (Fudan U)**

**Our approach: use DHT's for managing (statistical) metadata only**
                 **with single-dim. keys (PeerId, term, URL)**

# Our Approach to P2P Query Routing

peer lists (directory)

term a: 17, 11, 92, ...
term f: 43, 65, 92, ...
url z: 54, 128, 7, ...

url x: 37, 44, 12, ...
term c: 13, 92, 45, ...

term g: 13

book-marks B0

**query p**

local index X0

**sim based on Kullback-Leibler divergence between bookmarks' term distributions:**

$$KL(B0, Bi) := \sum_x freq(x, B0) \log \frac{freq(x, B0)}{freq(x, Bi)}$$

**forming a statistics-based, pre-computed „semantic overlay network"**

rm g: 13

**peer P0 first executes query locally**

**P2P directory has peer lists for posted terms and bookmark URLs**

**P0 identifies best peers Pi in terms of benefit/cost:**
**( sim (P0, Pi) / overlap (P0, Pi) ) / cost(Pi)**

# Exploiting Collective Human Input
# for Collaborative Web Search
## - Beyond Relevance Feedback -

- href links are human endorsements $\rightarrow$ PageRank, etc.
- <u>Opportunity</u>: online analysis of human input & behavior
  may compensate deficiencies of search engine

**<u>Typical scenario</u> for 3-keyword user query: a & b & c**
$\rightarrow$ **top 10 results: user clicks on ranks 2, 5, 7**

$\rightarrow$ **top 10 results: u**

$\rightarrow$ **top 10 results: u**

**user asks friend for tips**

**query logs, bookmarks, etc. provide**
- **human assessments & endorsements**
- **correlations among words & concepts
  and among documents**

**Challenge: How can we use knowledge about the collective
input of all users in a large community?**

# Observations and Challenges (4)

*Observation:*
*Semantic overlay networks* for P2P Web search
build on statistical similarity among peers

*Challenges:*
- Efficient benefit/cost estimation and efficient computation
  of global measures from local ones (idf, PageRank, KL, ...)
- From bookmark-driven query routing towards
  exploiting query logs and click streams
- Distributed, self-optimizing TA-sorted and Prob-sorted
- Caching, lazy replication, proactive dissemination
- Incentive mechanisms and trust management

# Outline

✓ Motivation and Challenges

✓ Search                    **(XML, Ontologies)**

✓ Speed                     **(Top-k Query Processing)**

✓ Self-Organization         **(P2P Collaborative Search)**

# Concluding Remarks

*long-term goal:*  exploit the Web's potential
for being the world's largest knowledge base

- *XML* and *Semantic Web* are key assets,
  but by themselves not sufficient; we need to
  cope with *diversity*, *incompleteness*, and *uncertainty*
  $\rightarrow$ absolute need for ranked retrieval
  $\rightarrow$ *statistics* is key

- combine techniques from *DBS*, *IR*, *CL*, *AI*, and *ML*

- *P2P* is intriguing paradigm:
  computing power, community input, anti-monopoly

- key issue is *quality/efficiency tradeoffs*