

Intelligente und effiziente Suche auf semistrukturierten Daten

Gerhard Weikum

weikum@mpi-sb.mpg.de

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

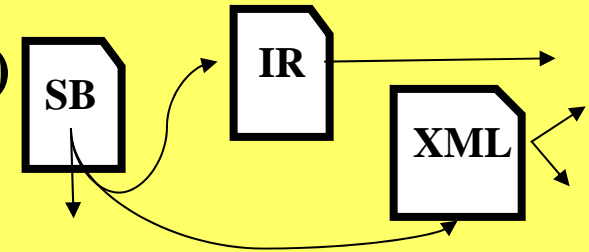
Outline

- Motivation and Challenges
- XXL & XXL-light: IR on XML Data
- Role of Ontologies
- Efficient Evaluation of Top-k Queries
- Ongoing and Future Work

A Few Challenging Queries

(on Web / Deep Web / Intranet / Personal Info)

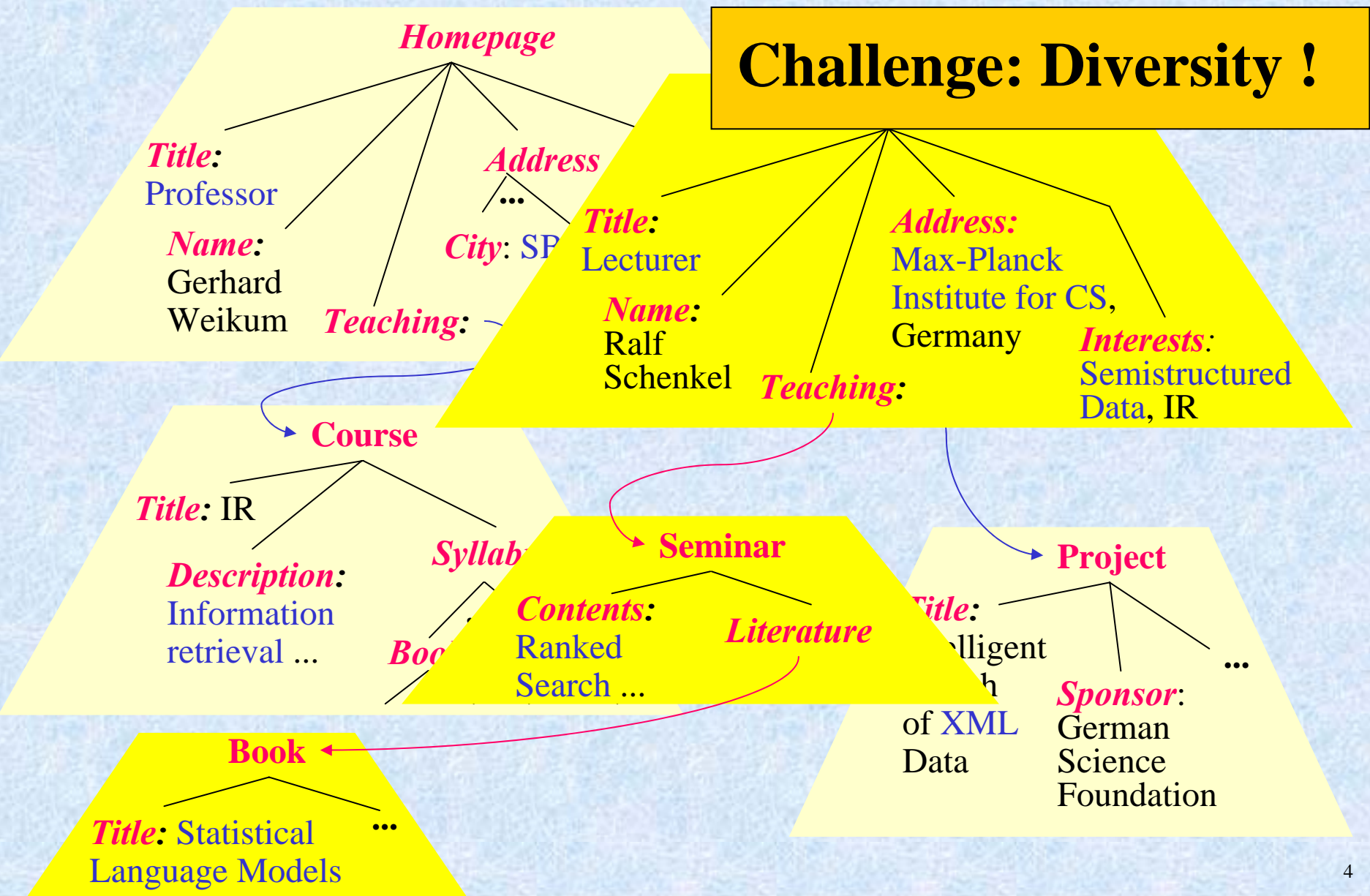
- Which professors from Saarbruecken (SB) are teaching IR and have research projects on XML?



- Which gene expression data from Barrett tissue in the esophagus exhibit high levels of gene A01g?
- What are the most important results on large deviation theory?
- Which drama has a scene in which a woman makes a prophecy to a Scottish nobleman that he will become king?
- Who was the French woman 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?

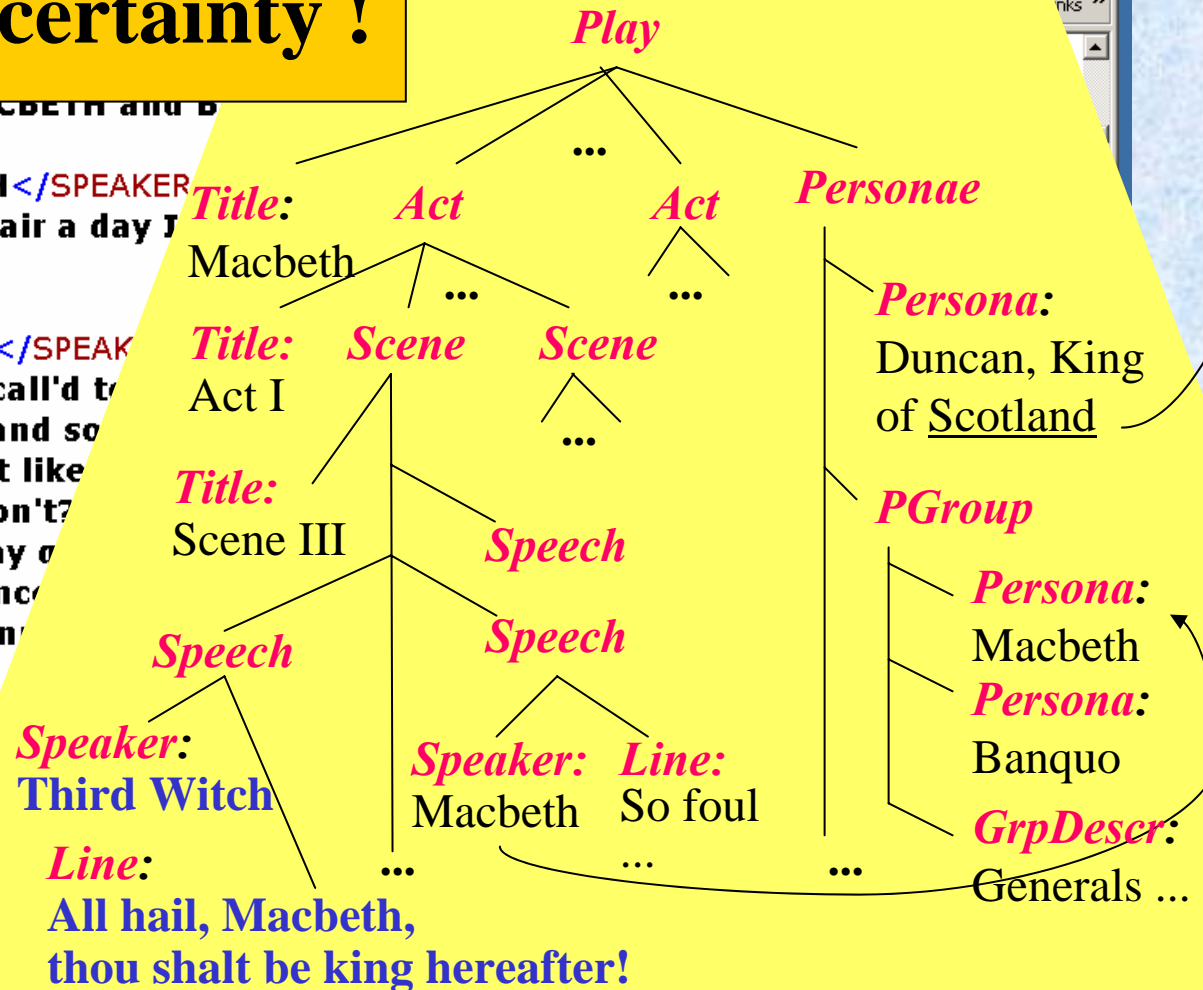
Challenge: Diversity !



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

Challenge: Uncertainty !

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



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

International Conference ...

Challenge: Ambiguity !

Homepage

Firstname:

Sophie

Lastname:

Cluet

Address:

INRIA
Rocquencourt,
78153 Le Chesnay,
France

Interests:

XML, ...

Homepage

Firstname:

Maria

Lastname:

Sanchez

Gender: female

Address:

Main Street,
Paris, Texas 94052

Homepage

Name:

Antoinette
Lagrange

Address

Street:

Rue de la
Chimie 138

City:

Paris

Country:

France

Hobbies:

Painting, ...

Homepage

Name:

M.-C. Richard

Address:

Rue de Voltaire,
10045 Paris,
France

Biography:

... mother of
two children ...

Our Research Agenda

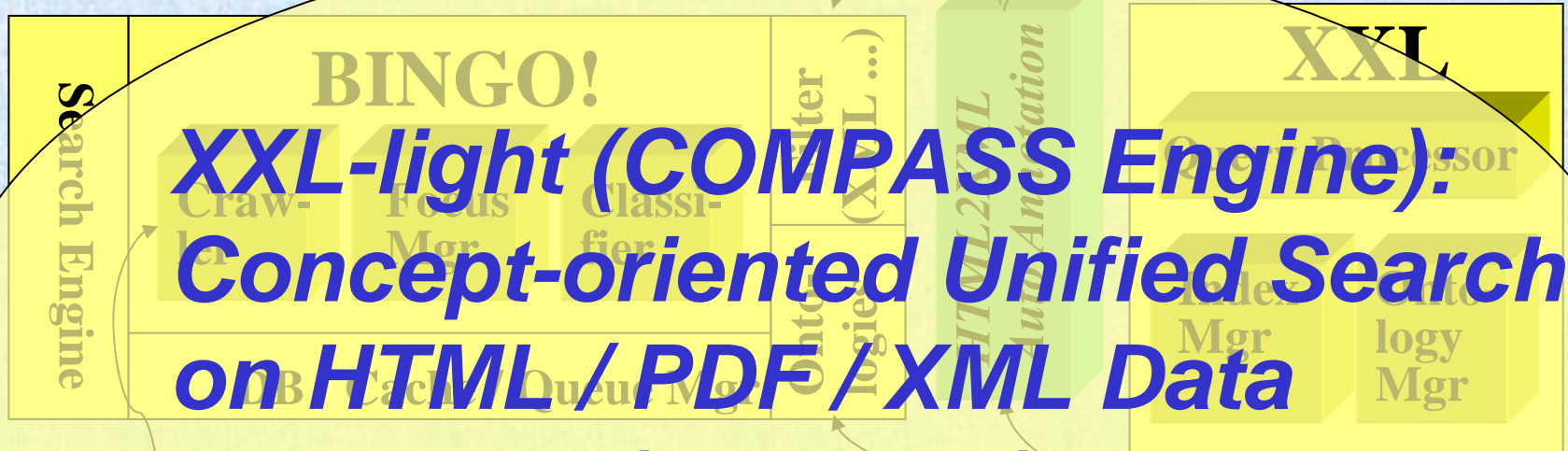
Web Applications

Intranet Applications

Chamber of Crafts
Portal

...

Materialography
Portal



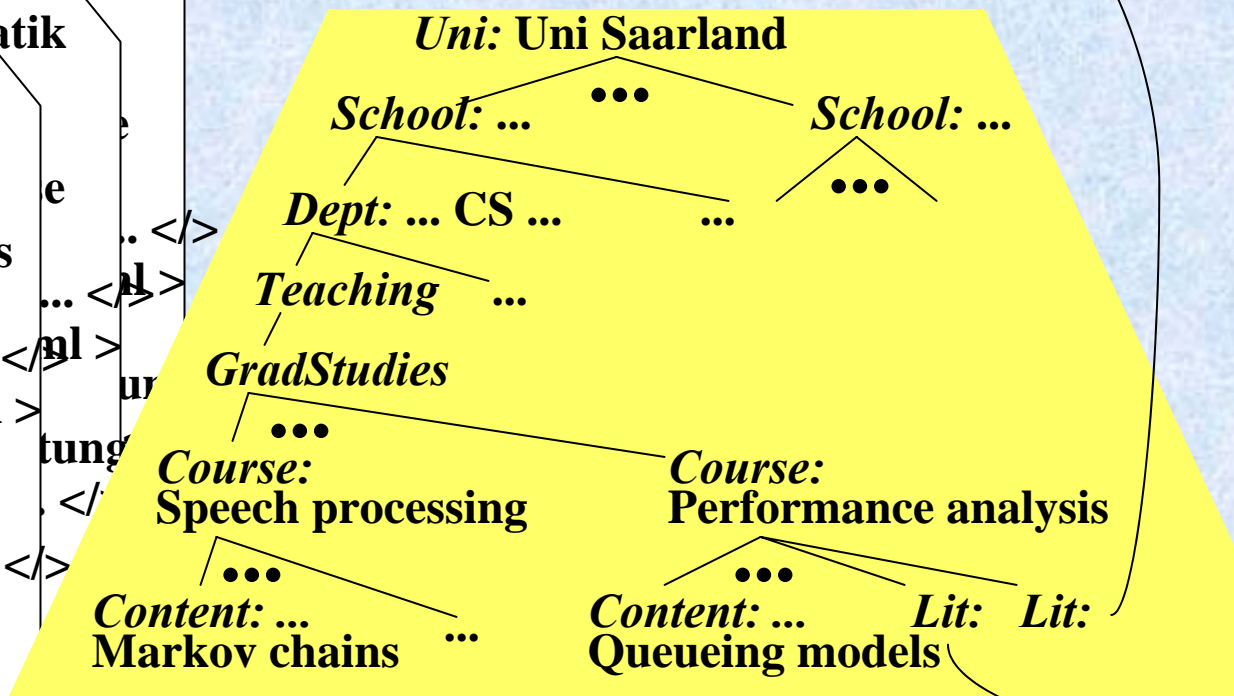
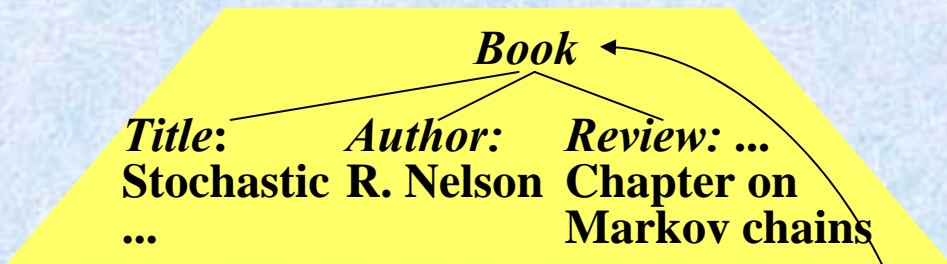
Outline

✓ Motivation and Challenges

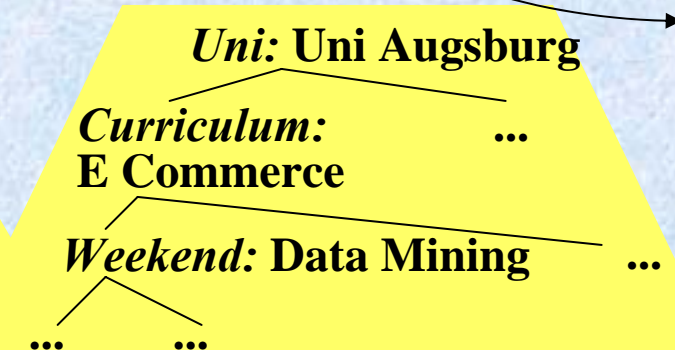
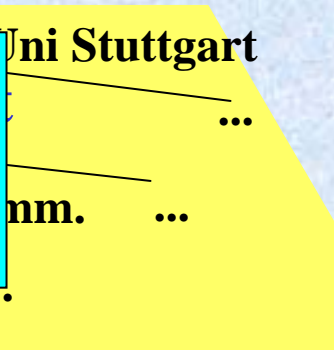
- **XXL & XXL-light: IR on XML Data**
- Role of Ontologies
- Efficient Evaluation of Top-k Queries
- Ongoing and Future Work

XML-IR Example (1)

```
<Uni> ETH Zürich  
...  
<Uni> Uni Stuttgart  
...  
<Uni> Uni Saarland  
<School> Math & Engineering  
<Dept> CS  
<Teaching> ...  
<GradStudies>  
<Course> Performance analysis  
<Lecturer> ...  
<Content> Queueing models ..  
<Lit href=springer/nelson.xml >  
<Lit href=... >  
<Course> Speech processing  
<Content> ... Markov chains...  
</Course>  
...  
</Teaching> .. </Dept> .. </School>  
</Uni>
```



**Semistructured data:
elements, attributes, links
organized as labeled graph**



XML-IR Example (2)

www.allunis.de/unis.xml

Regular expressions
over path labels
+ Logical conditions
over element contents

Uni: Uni Stuttgart

Inst: CS

Course: Mobile comm.

Prerequisites:

... Markov processes

Uni: Uni Augsburg

Curriculum:
E Commerce

Weekend: Data Mining

Outline: ...
statistical methods
for classification ...

Tit
Sto
...

Uni: Uni Saarland

School: ...

School: ...

Dept: ... CS

Teaching

GradStudies

Course:
Speech processing

Content: ...
Markov chains

Course:
Performance analysis

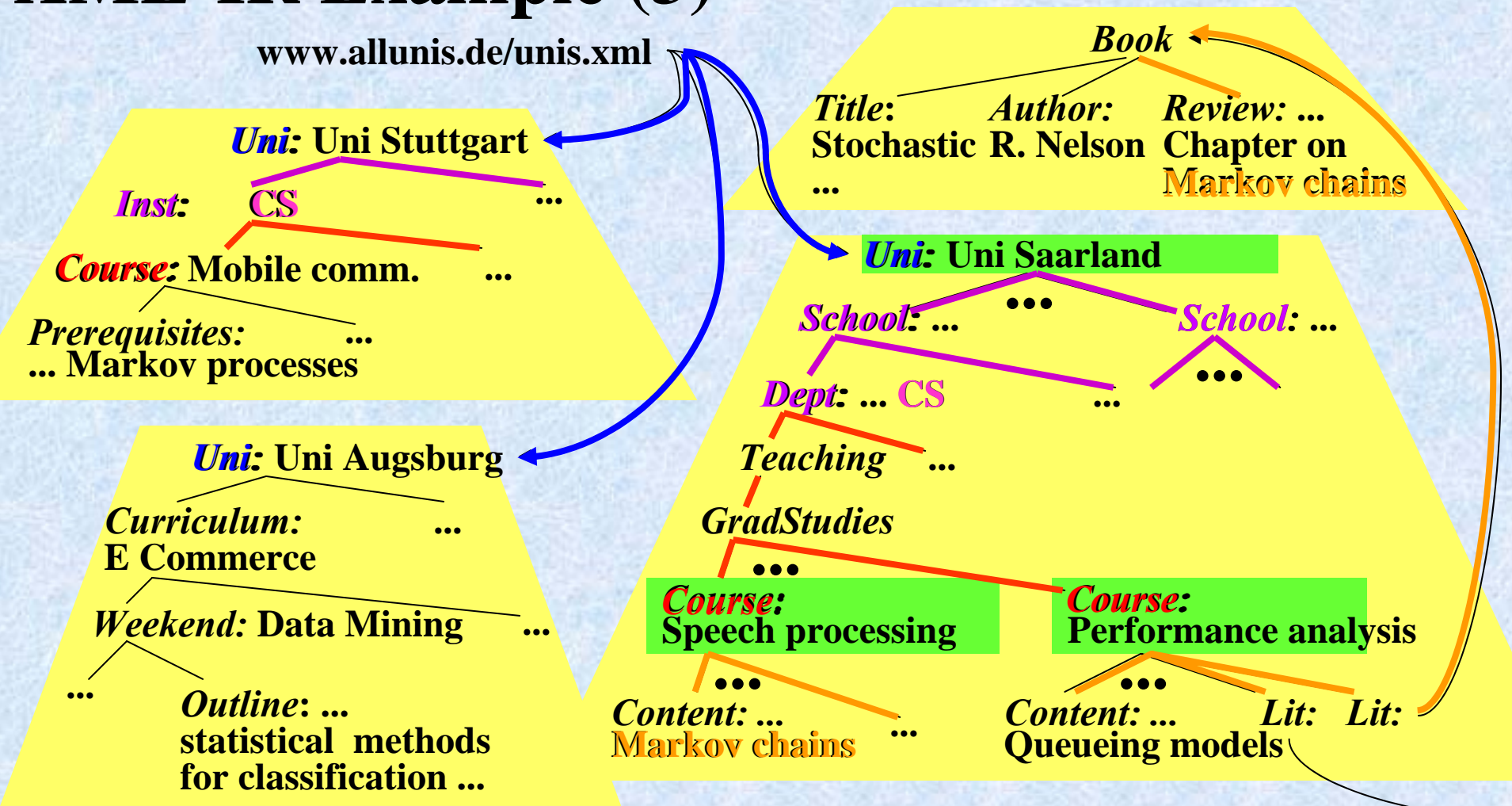
Content: ...
Queueing models

Lit: Lit:

Select U, C From www.allunis.de/unis.xml Where Uni As U
And U.#.School?.#.(Inst | Dept)+ As D And D Like „%CS%“
And D.#.Course As C And C.# Like „%Markov chain%“

XML-IR Example (3)

www.allunis.de/unis.xml



Select **U, C** From www.allunis.de/unis.xml Where **Uni** As **U**
 And **U.#.School?.#.(Inst | Dept)+** As **D** And **D Like „%CS%“**
 And **D.#.Course** As **C** And **C.# Like „%Markov chain%“**

XML-IR Example (4)

www.allunis.de/unis.xml

Uni: Uni Stuttgart

Inst: CS

Course: Mobile comm.

Prerequisites:
... Markov processes

Uni: Uni Augsburg

Curriculum:
E Commerce

Weekend: Data Mining

Outline: ...
statistical methods
for classification ...

Book

Title: Stochastic ...
Author: R. Nelson
Review: ... Chapter on
Markov chains

Uni: Uni Saarland

School: ...

Dept: ... CS

Teaching ...

GradStudies

Course:
Speech processing

Content: ...
Markov chains

Course:
Performance analysis

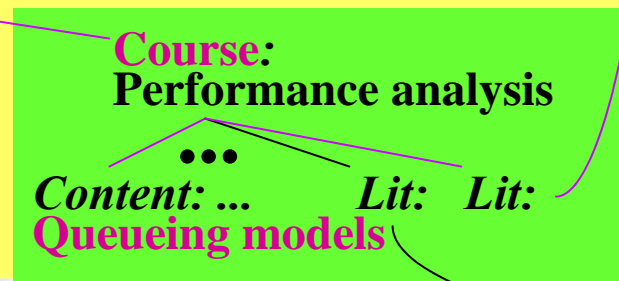
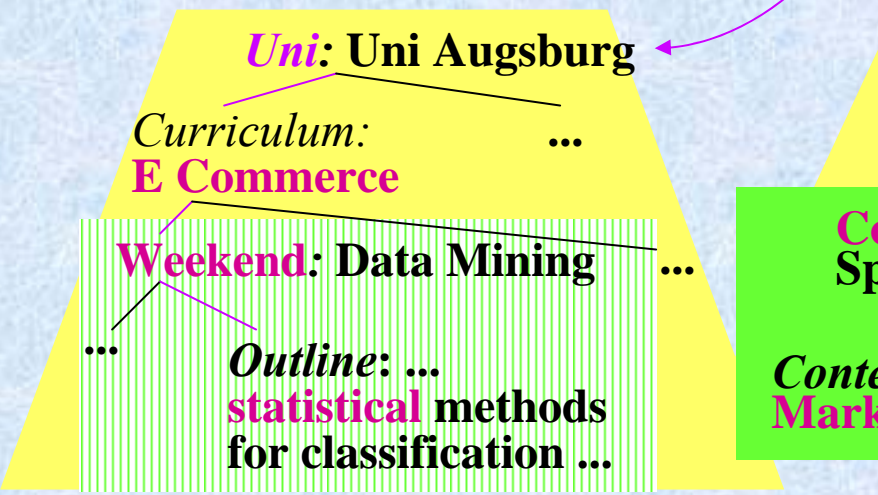
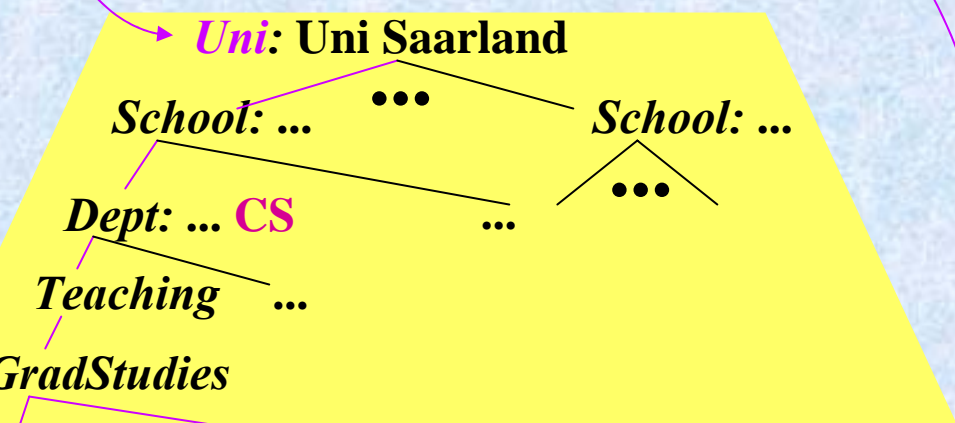
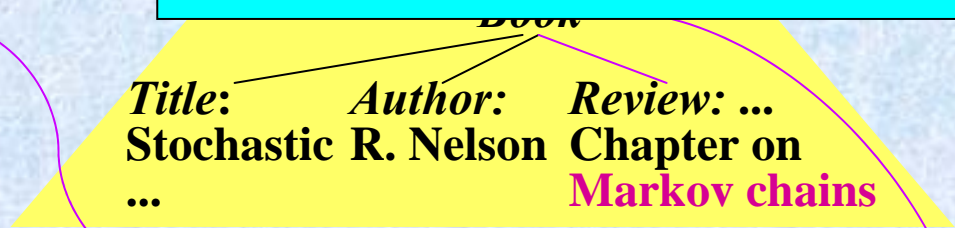
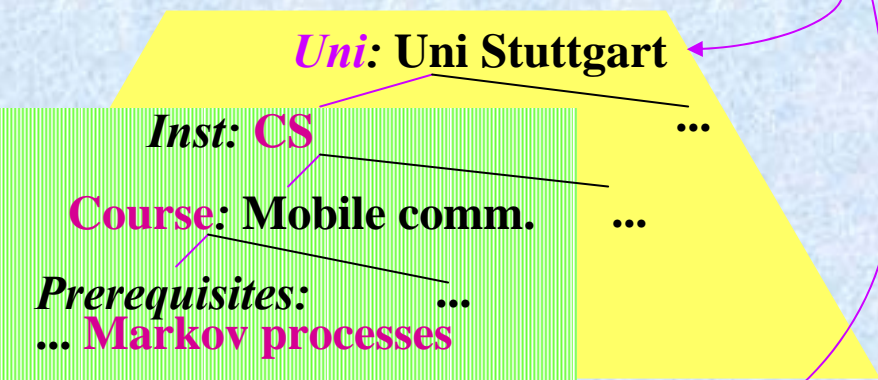
Content: ...
Queueing models

Select U, C From www.allunis.de/unis.xml Where Uni As U
And U.# As D And D ~ „CS“
And D.#.~Course As C And C.# ~ „Markov chain“

XML-IR Example (5)

Result ranking of XML data based on semantic similarity

www.allunis.de/unis.xml



Select U, C From www.allunis.de/unis.xml Where Uni As U And U.# As D And D ~ „CS“ And D.#.~Course As C And C.# ~ „Markov chain“

XML-IR Concepts

Example COMPASS

(Concept-Oriented Multi-Format Portability)
simple, extensible core language – application

Where clause: conjunction of restricted conditions
with binding of variables

Query Semantics:

- query is a pattern with relaxable conditions
- results are approximate matches to query with similarity scores

Elementary conditions on names and contents

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*

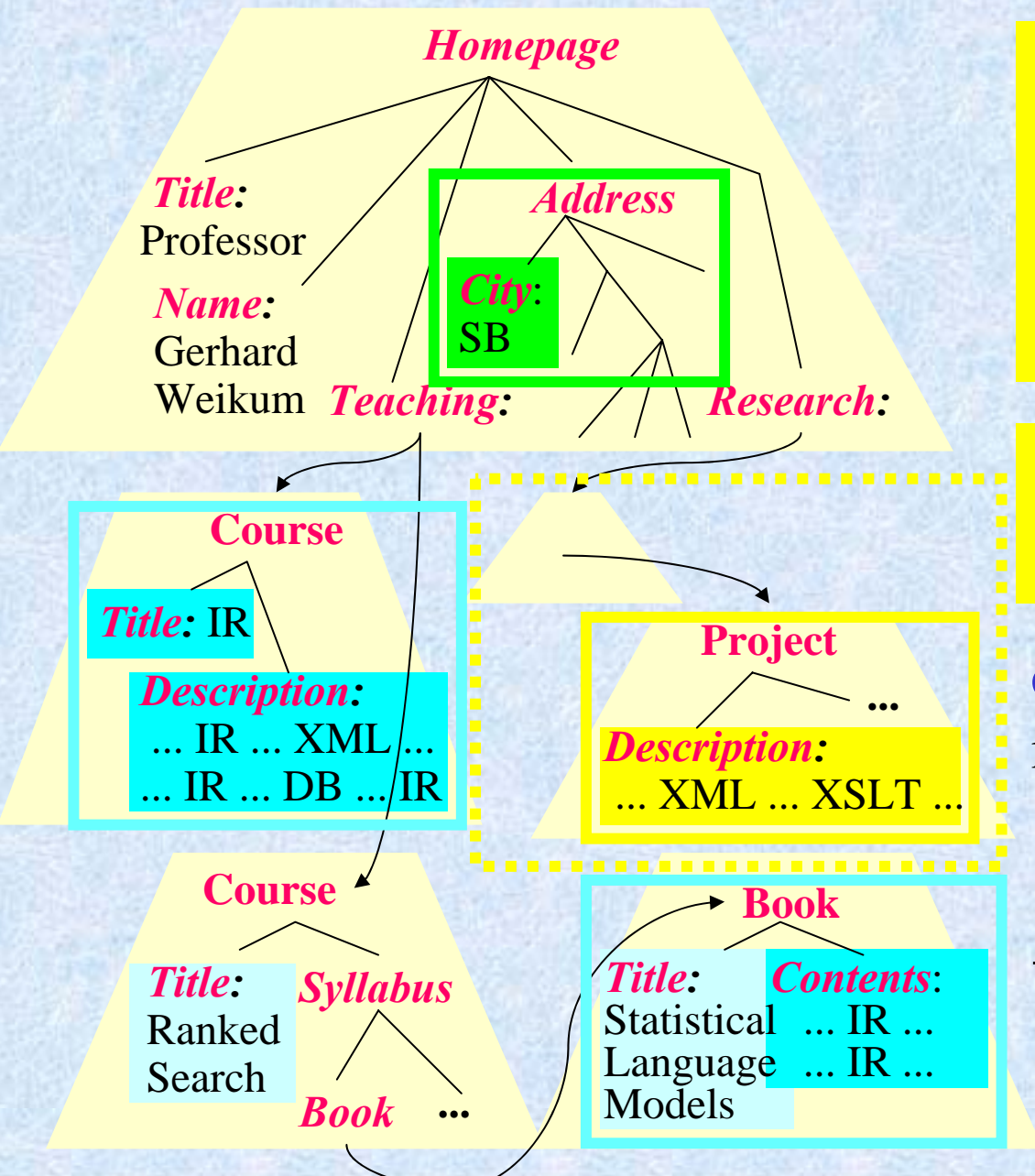
Semantic similarity conditions on names and contents

~research = „~XML“

Relevance scoring based on

tf*idf similarity of contents,
ontological similarity of names,
probabilistic combination of conditions

XML-IR Scoring Model

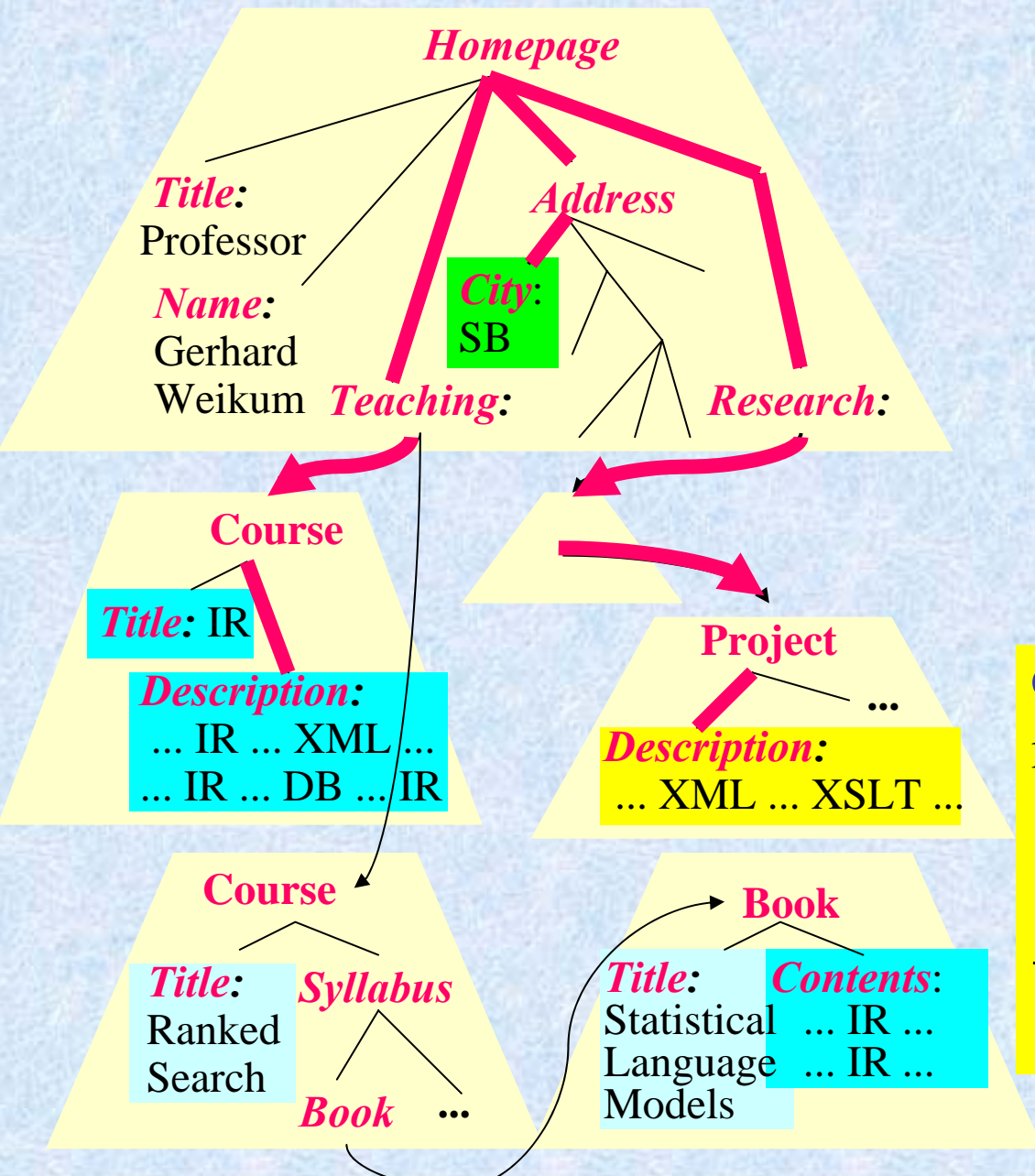


local score for elementary condition:
based on tf*idf-style statistics for node or node context with score propagation

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

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

XML-IR Scoring Model

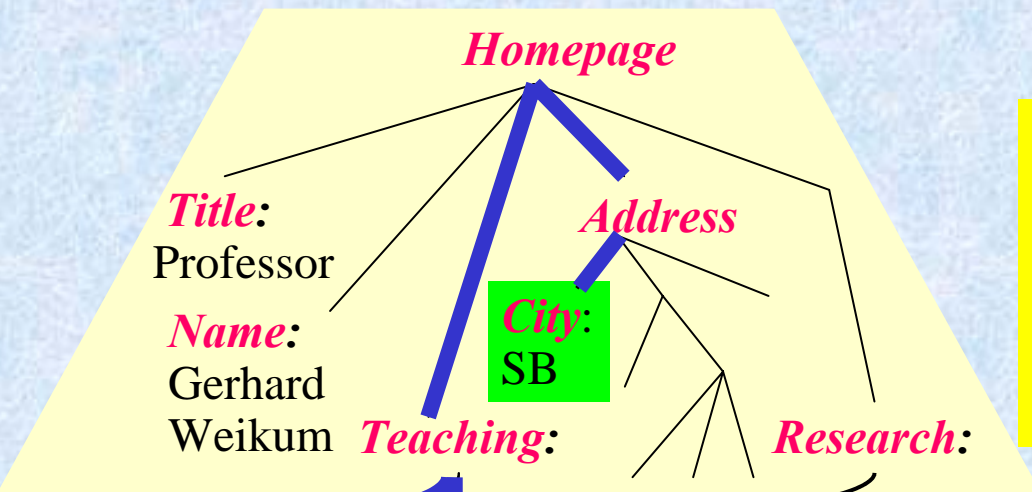


local score for elementary condition:
based on tf*idf-style statistics for node or node context with score propagation

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

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

XML-IR Scoring Model



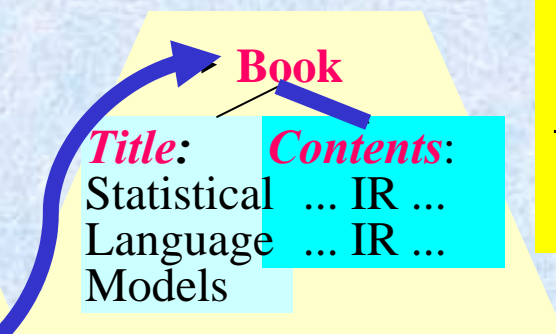
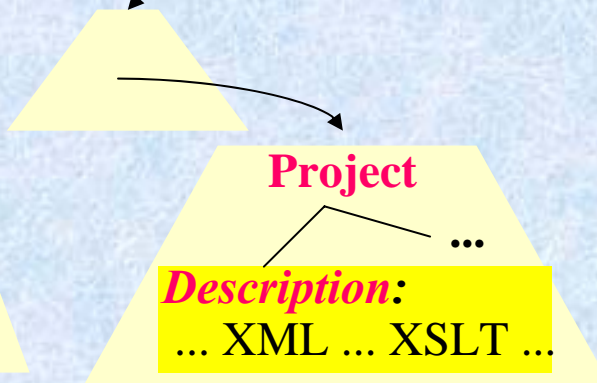
local score for

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

global score for query:

\sum local scores * compactness

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



Outline

- ✓ Motivation and Challenges
- ✓ XXL & XXL-light: IR on XML Data
- Role of Ontologies
- Efficient Evaluation of Top-k Queries
- Ongoing and Future Work

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

XML schemas, DTDs, namespaces:

syntactic conventions and standardized naming (plus typing info)

Gazetteer: (geographical) dictionary of names

Reasoning on Ontologies and Thesauri:

Professor \subseteq Lecturer $\cap \exists$ hasStaff.Secretary

Teaching \supseteq Cou

Professor \subseteq Acad

Academician \subseteq H

Human \subseteq Carniv

...

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

→ logical inferences
with sub-FOL calculus

→ transitive closures,
shortest paths, etc.
along generalizations

Example WordNet

WordNet 1.6 Browser

File History Options Help

Search Word: 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)

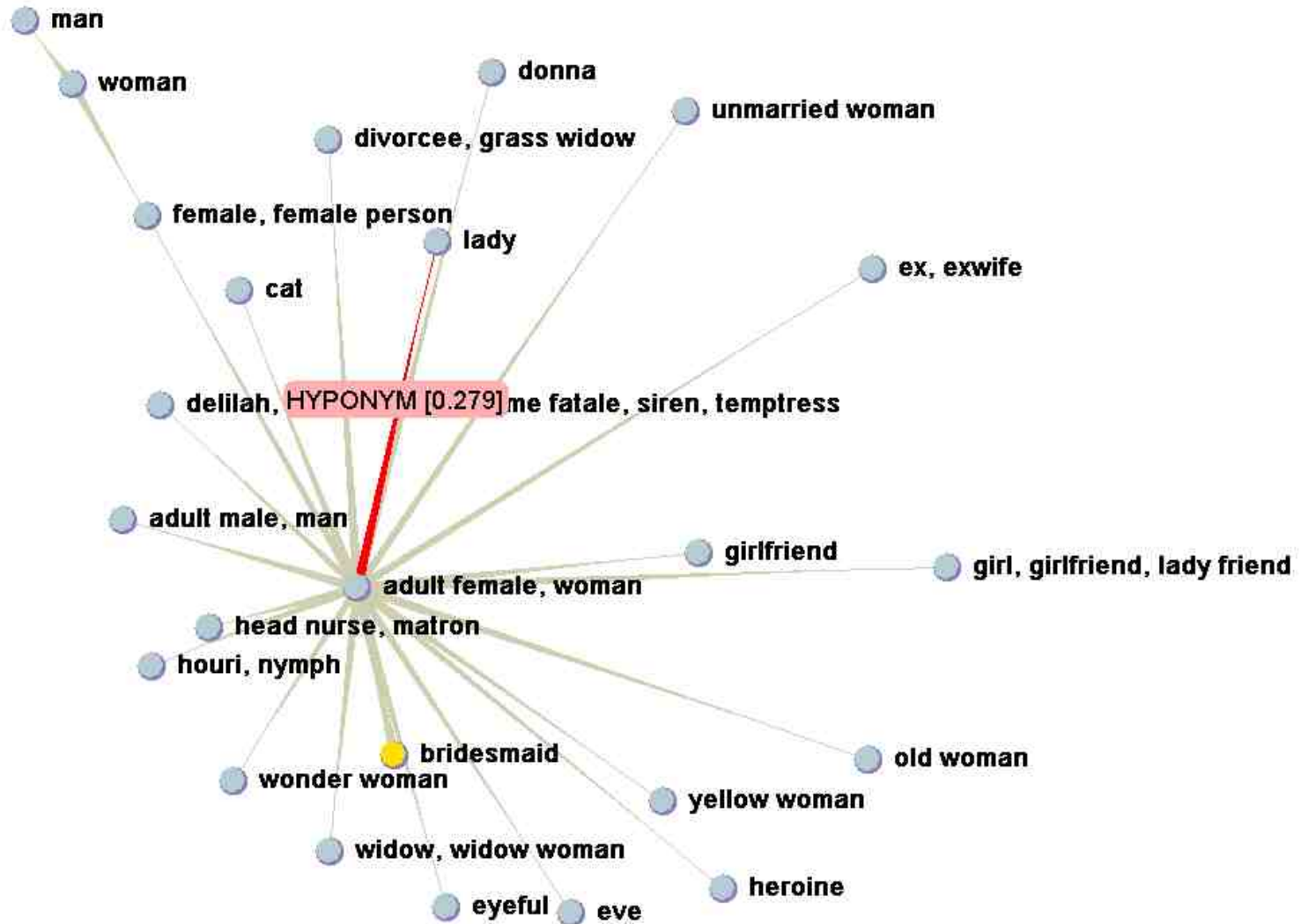
woman, adult female – (an adult female person)

- => 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 Visualization

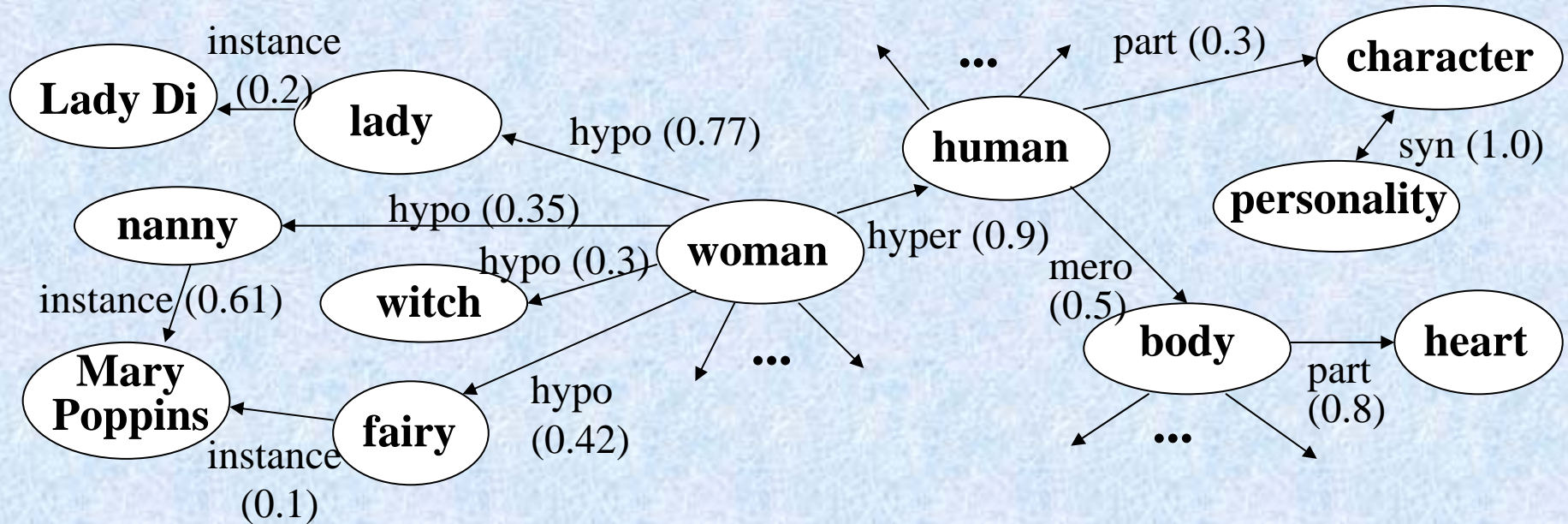


SEARCH ONTOLOGY:

FONT SIZE:

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 $\text{sim}(c1, c2)$:

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

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

Conditional probabilities:
$$P[\text{doc has } c1 \mid \text{doc has } c2]$$

Transitive similarity:

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

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

usefor 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 $\sim w$ by $(c_1 \mid \dots \mid c_k)$ with all c_i for which $\text{sim}(w, c_i) \geq \delta$

„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

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[0.145|1.00], "mafia[0.154|1.00], "sicilian[0.066|1.00], "black[0.066|1.00] & hand[0.066|1.00], organ[0.213|1.00], crime[0.311|1.00], 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.

135530 sorted accesses in 11.073

Results:

1. Interpol Chief on Fight Against ...
 2. Economic Counterintelligence ...
 3. Supreme Procuratorate Wor ...
 4. Crime and Punishment in the ...
 5. SWITZERLAND CALLED ...
- 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.

...

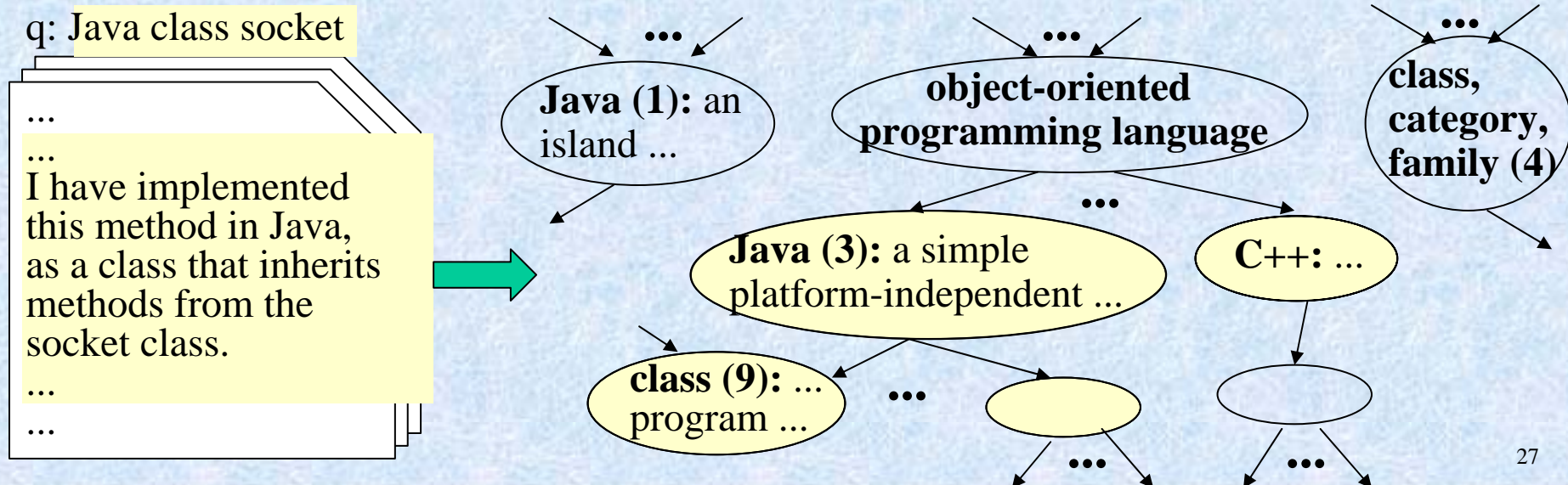
Keyword-to-Concept Mapping and Word Sense Disambiguation

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



What About Deep Web and Web Services?

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

~sheetmusic = „~flute“ → instrument = (flute | piccolo | recorder)
→ category = reeds
→ style = (classical | jazz | folk)

digital sheet music | music books | power search | wish list

power search find your favorite sheet mu

search for:

available in: *Applies to c

notation type: all easy play piano/vocal TAB

keyword:

title/song:

artist/composer:

instrument:

style:

scoring:

difficulty: (for digital sheet music only)

lyrics: (for digital sheet music only)

show items per page:

```
<element name="WSF_Form0Select0_Enum">
  <simpleType>
    <restriction base="string">
</restriction>
</simpleType>
</element>
<simpleType name="WSF_Form0Select1_Enum">
  <restriction base="string">
    <enumeration value="Alternative"/>
    <enumeration value="Blues"/>
    <enumeration value="Children's"/>
    <enumeration value="Classical"/>
    <enumeration value="Country"/>
```

Observations:

- Deep Web has > 500 000 hidden databases with > 500 billion ($5 \cdot 10^{11}$) dynamic pages
- High „redundancy“ among query forms → enables exploitation of statistics

What About Deep Web and Web Services?

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

~sheetmusic = „~flute“ → instrument = (flute | piccolo | recorder)
→ category = reeds
→ style = (classical | jazz | folk)

digital sheet music | music books | power search | wish list

power search find your favorite sheet music

```
<element name="WSF_Form0Select0_Enum">  
<simpleType>  
<restriction base="string">
```

Approach:

- crawl Web forms, tables, WSDL
- mine for multivariate correlations among N-tuples of (concept, value) pairs, examples:
 - ((instrument, flute), (category, reeds))
 - ((make, Audi), (model, A4), (color, red))
 - ((title, *), (author, <person name>))
- map query to target form/WSDL for max likelihood (taking into account ontological similarities)

Enum">

go search now!

Outline

- ✓ Motivation and Challenges
- ✓ XXL & XXL-light: IR on XML Data
- ✓ Role of Ontologies
- Efficient Evaluation of Top-k Queries
- Ongoing and Future Work

Top-k Query Processing with Scoring

B+ tree on terms

algorithm

...

performance

...

z-transform

17: 0.3
44: 0.4
52: 0.1
53: 0.8
55: 0.6
⋮

12: 0.5
14: 0.4
28: 0.1
44: 0.2
51: 0.6
52: 0.3
⋮

11: 0.6
17: 0.1
28: 0.7
⋮

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

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

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

Naive QP algorithm:

```

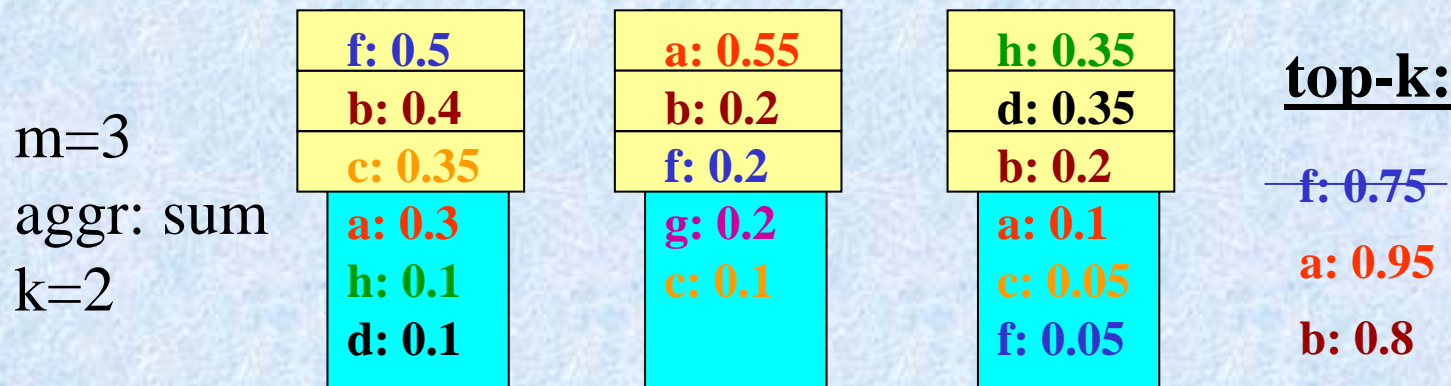
candidate-docs := ∅;
for i=1 to z do {
    candidate-docs := candidate-docs ∪ index-lookup(ti) };
for each dj ∈ candidate-docs do {compute score(q,dj)};
sort candidate-docs by score(q,dj) descending;
    
```

“Fagin’s 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 L_i ;
 $high_i := s_i(d_j)$;

*but random accesses
are expensive !*

if $d_j \notin \text{top-k}$ then {
 look up $s_v(d_j)$ in all lists L_v with $v \neq i$; // random access
 compute $s(d_j) := \text{aggr} \{s_v(d_j) \mid v=1..m\}$;
 if $s(d_j) > \text{min score among top-k}$ then
 add d_j to top-k and remove min-score d from top-k; }
threshold := aggr { $high_v \mid v=1..m$ };
if min score among top-k \geq threshold then exit;



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 ;

$E(d_j) := E(d_j) \cup \{i\}$; $high_i := si(q, d_j)$;

$bestscore(d_j) := aggr\{x_1, \dots, x_m\}$

with $x_i := si(q, d_j)$ for $i \in E(d_j)$, $high_i$ for $i \notin E(d_j)$;

$worstscore(d_j) := aggr\{x_1, \dots, x_m\}$

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 $worstscore$;

$threshold := bestscore\{d \mid d \text{ not in } top-k\}$;

if $\min worstscore$ among $top-k \geq threshold$ then exit;

$m=3$

aggr: sum

$k=2$

f: 0.5
b: 0.4
c: 0.35
a: 0.3
h: 0.1
d: 0.1

a: 0.55
b: 0.2
f: 0.2
g: 0.2
c: 0.1

h: 0.35
d: 0.35
b: 0.2
a: 0.1
c: 0.05
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)

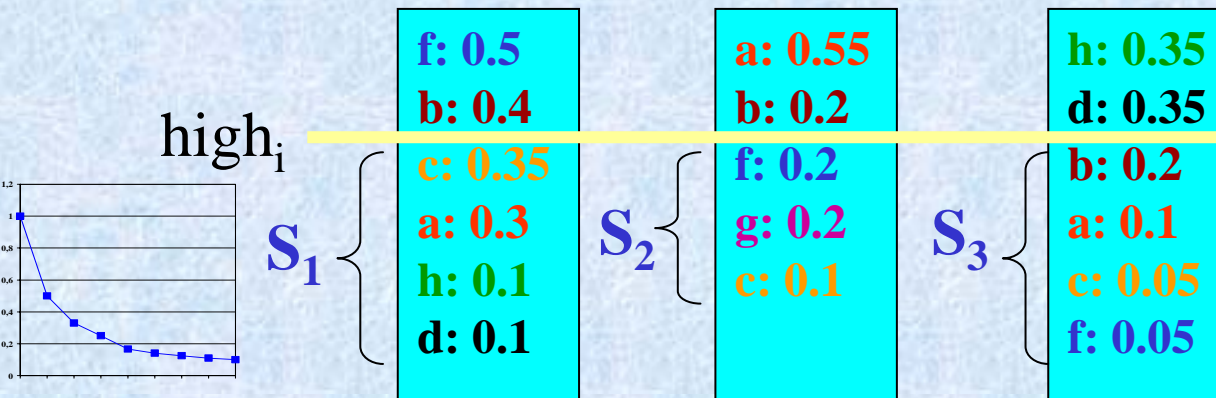
$$\sum_{i \in E(d)} s_i(d) \leq s(d) \leq \sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} high_i$$

Relaxed into probabilistic invariant

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

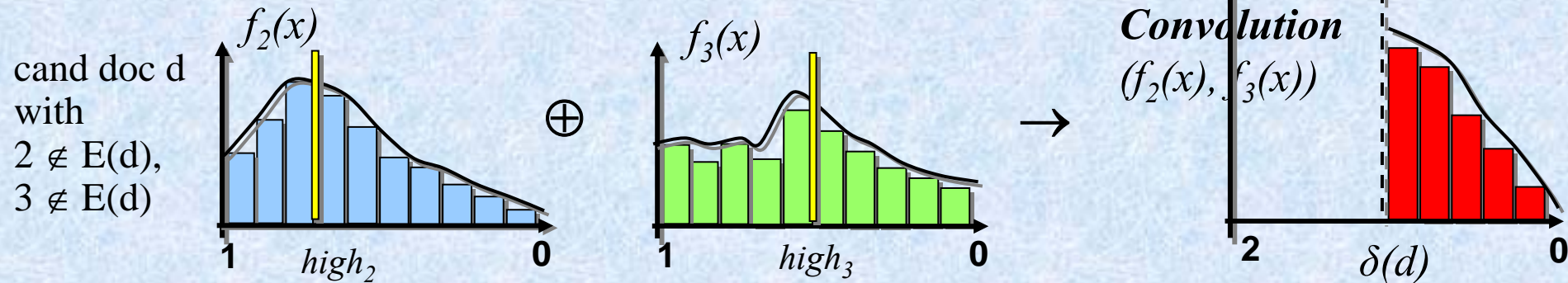
$$= P\left[\sum_{i \notin E(d)} S_i > threshold - \sum_{i \in E(d)} s_i(d)\right] =: P\left[\sum_{i \notin E(d)} S_i > \delta'\right] \leq \varepsilon$$

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



- *Discard candidates with $p(d) \leq \varepsilon$*
- **Exit index scan when candidate list empty**

Probabilistic Threshold Test



- 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: $P[d = v_j] = e^{-\alpha_i} \frac{\alpha_i^{j-1}}{(j-1)!}$
 - 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* S_i 's or with *correlated* S_i 's

Prob-sorted Algorithm (Smart Variant)

Prob-sorted (RebuildPeriod r, QueueBound b):

...

scan all lists L_i ($i=1..m$) in parallel:

...same code as TA-sorted...

// queue management

for all priority queues q for which d is relevant do

insert d into q with priority $\text{bestscore}(d)$;

// periodic clean-up

if step-number mod $r = 0$ then

// rebuild; single bounded queue

if strategy = Smart then

for all queue elements e in q do

update $\text{bestscore}(e)$ with current high_i values;

rebuild bounded queue with best b elements;

if $\text{prob}[\text{top}(q) \text{ can qualify for top-}k] < \epsilon$ then exit;

if all queues are empty then exit;

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

Performance Results for IMDB Queries

on IMDB corpus (Web site: Internet Movie Database):

375 000 movies, 1.2 Mio. persons (html/xml)

20 structured/text queries with Dice-coefficient-based similarities of categorical attributes Genre and Actor, e.g.:

- $Genre \supseteq \{Western\} \wedge Actor \supseteq \{John\ Wayne, Katherine\ Hepburn\} \wedge Description \supseteq \{sheriff, marshall\}$,
- $Genre \supseteq \{Thriller\} \wedge Actor \supseteq \{Arnold\ Schwarzenegger\} \wedge Description \supseteq \{robot\}$

	TA-sorted	Prob-sorted (smart)
#sorted accesses	1,003,650	403,981
elapsed time [s]	201.9	12.7
max queue size	12628	400
relative recall	1	0.75
rank distance	0	126.7
score error	0	0.25

Outline

- ✓ Motivation and Challenges
- ✓ XXL & XXL-light: IR on XML Data
- ✓ Role of Ontologies
- ✓ Efficient Evaluation of Top-k Queries
- Ongoing and Future Work

Exploiting Collective Human Input for Collaborative Web Search

- Beyond Relevance Feedback and Beyond Google -

- href links are human endorsements → PageRank, etc.
- Opportunity: online analysis of human input & behavior may compensate deficiencies of search engine

Typical scenario for 3-keyword user query: a & b & c

→ top 10 results: user clicks on ranks 2, 5, 7

→ top 10 results: u query logs, bookmarks, etc. provide
u
u • human assessments & endorsements
u • correlations among words & concepts
u and among documents

user asks friend for tips

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

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* → absolute need for ranked retrieval
- view *information organization* and *information search* as dual views of the same problem
- combine techniques from *DBS*, *IR*, *CL*, *AI*, and *ML*
- need better *theory* about quality/efficiency *tradeoffs* as well as *large-scale experiments*