AIFB

# Semantic Web Technologies I

Lehrveranstaltung im WS12/13

PD Dr. Sebastian Rudolph

Dr. Duc Thanh Tran

M.Sc. Anees ul Mehdi

AIFB

# Agenda

- Introduction
- Semantic Web data
  - The RDF data model
  - Publishing RDF
  - Last Lecture: crawling and indexing RDF data
- Query processing / matching
  - Last lecture: selected problems in structured query (SPARQL) processing
  - Here: big picture of querying with structured and keyword queries
- Ranking
- Result presentation

# Why Semantic Search? I.

- "We are at the beginning of search." (Marissa Mayer)
    - Solved large classes of queries, e.g. navigational
    - Heavy investment in computational power
    - Remaining queries are hard, not solvable by brute force, and require a deep understanding of the world and human cognition
- Background knowledge and metadata can help to address poorly solved queries

# Poorly solved information needs

- Ambiguous searches
  - paris hilton
- Long tail queries
  - george bush (and I mean the beer bre
- Multimedia search
  - paris hilton sexy
- Imprecise or overly precise searches
  - jim hendler
  - pictures of strong adventures people
- Precise searches for descriptions
  - countries in africa
  - 32 year old computer scientist living in barcelona
  - reliable digital camera under 300 dollars

> Many of these queries would not be asked by users, who learned over time what search technology can and can not do.

# Example: multiple interpretations

# Why Semantic Search? II.

AIFB

- The Semantic Web is now a reality
  - Large amounts of data published in RDF
  - Heterogeneous data of varying quality
  - Users who are not skilled in writing complex queries (e.g. SPARQL) and may not be experts in the domain
- Searching data instead or in addition to searching documents
  - Direct answers
  - Novel search tasks

# Example: direct answers in search



Information from the Knowledge Graph

# Document retrieval and data retrieval

AIFB

- Information Retrieval (IR) support the retrieval of documents (document retrieval)
  - Representation based on lightweight syntax-centric models
  - Work well for topical search
  - Not so well for more complex information needs
  - Web scale
- Database (DB) and Knowledge-based Systems (KB) deliver more precise answers (data retrieval)
  - More expressive models
  - Allow for complex queries
  - Retrieve concrete answers that precisely match queries
  - Not just matching and filtering, but also joins
  - Limitations in scalability

# Combination of document and data retrieval

- Documents with metadata
  - Metadata may be embedded inside the document
  - *I'm looking for **documents** that mention countries in Africa.*

- Data retrieval
  - Structured data, but searchable text fields
  - *I'm looking for **directors**, who have directed movies where the synopsis mentions dinosaurs.*

# Semantic Search

AIFB

- Target (combination of) document and data retrieval

- Semantic search is a retrieval paradigm that
  - Exploits the structure/semantics of the data or explicit background knowledge to understand user intent and the meaning of content
  - Incorporates the intent of the query and the meaning of content into the search process (**semantic models**)

- Wide range of semantic search systems
  - Employ different semantic models, possibly at **different steps** of the search process and in order to support **different tasks**

# Semantic Search systems

For **data** / **document** retrieval, semantic search systems might combine a range of techniques, ranging from statistics-based **IR methods** for **ranking**, **database methods** for efficient **indexing** and **query processing**, up to complex **reasoning** techniques for making inferences!

# Repetition: Information Workbench

**AIFB**

- Addressing the **lifecycle of interacting** with the Web of Data
  - Integration of data sources
  - Content generation by the end user
  - **Search and Exploration**
  - **Visualization**
  - Publishing

- Integrated management of **heterogeneous data sources**
  - Structured and unstructured
  - Published and user-generated
  - Static and dynamic
  - Open domain

# Data Sources in the Application

AIFB

- Entire English Wikipedia

- Data from Linked Open Data
    - DBpedia
    - YAGO
    - …

- Data from Data.gov (US Government)
    - E.g. live data about earthquakes

- Many more

# Semantic Search

**AIFB**

- **Hybrid Search**:  Structured queries combined with keywords across structured and unstructured data sources

- **Query interpretation:** Translation of keywords into hybrid queries

- **Keyword search/query interpretation** combined with **faceted search**: iterative refinement process based on keywords and operations on facets

# Search, Refinement and Navigation

**Semantic Web TECHNOLOGIES**

**AIFB** ⬡

## Search Results

**Keywords**

queen single | Search | fluid operations

**Click on one of the suggestions to initiate translation! (can take a few seconds)**

queen single
Set searchfield to "queen single"

**A** (*queen*) is a Single

**B** writer **A** (*queen*)
**B** is a Single

**A** is a Single
**A** producer **B** (*queen*)

**Query Translations**

| RESULT COLUMN 1 | Initial Query *See Entire Query* |
|---|---|
| **producer** | **?sx1** |
| | A Kind of Magic (song) |
| ⊞ Range: All Values (43) | Another One Bites the Dust |
| | Back Chat |
| | Bicycle Race |
| **type** | Body Language (song) |
| | Calling All Girls |
| | Crazy Little Thing Called Love |
| ⊞ Range: All Values (43) | Fat Bottomed Girls |
| | Good Old-Fashioned Lover Boy |
| | Hammer to Fall |
| **writer** | Heaven for Everyone |
| | I Want to Break Free |
| | It's Late |
| ⊟ Range: All Values (42) | It's a Hard Life |
| ⊟ Musical Artist (42) | Keep Yourself Alive |
| Brian May (13) | Killer Queen |
| Frank Musker (1) | Las Palabras de Amor |
| Freddie Mercury (14) | Liar (Queen song) |
| John Deacon (7) | Long Away |
| Roger Meddows-Taylor (7) | Mustapha |

queen single.php
queen singled
queen singler
queen singlerpt
queen singles
queen singles-1997-2007
queen singles/2002/03/04/the
queen singles/2002/03/25/new

**Term Completions**

Queen (band)
Queen (band)
Queen (band)
Queen (band)
Queen (band)
Queen (band)
Queen (band)
Queen (band)

**Facets**

# Result Inspection, Analysis and Browsing

# Semantic Web data

# Data on the Web

- Data on the Web is not directly accessible
  - Most web pages are generated from databases, but formatted for human consumption
  - APIs offer limited views over data
- Two solutions
  - Extraction using **Information Extraction** (IE) techniques
    - ✦ Out of scope for this tutorial
  - Relying on publishers to expose structured data using standard **Semantic Web** formats

# Semantic Web

AIFB ⬛

- Sharing data across the Web
  - Standard data model
    - ✦ RDF
  - A number of syntaxes (file formats)
    - ✦ RDF/XML, RDFa
  - Powerful, logic-based languages for schemas
    - ✦ OWL, RIF
  - Query languages and protocols
    - ✦ HTTP, SPARQL

# Publishing RDF

**AIFB**

- Interlinked RDF documents (Linked Data)
  - Each document describes a single resource with URIs pointing to related resources
  - Common RDF file formats are RDF/XML and Turtle
  - Mostly implemented as a wrapper around a database or Web service
- Embedding RDF inside HTML
  - RDFa, microdata
- SPARQL endpoints
  - Triple stores are databases for managing RDF data
  - SPARQL is a standard protocol and query language for accessing triple stores using HTTP

# Example ontologies: schema.org

- Agreement on a shared set of schemas for common types of web content
  - Bing, Google, and Yahoo! as initial supporters
  - Similar in intent to sitemaps.org (2006)
    - ✦ Use a single format to communicate the same information to all three search engines
- Support for microdata
- schema.org covers areas of interest to all search engines
  - Business listings (local), creative works (video), recipes, reviews
  - User defined extensions
- Each search engine continues to develop its products

# Example: Facebook's Open Graph Protocol

- The 'Like' button provides publishers with a way to promote their content on Facebook and build communities
  - Shows up in profiles and news feed
  - Site owners can later reach users who have liked an object
  - Facebook Graph API allows 3rd party developers to

# Example: Facebook's Open Graph Protocol

- RDF vocabulary to be used in conjunction with RDFa
  - Simplify the work of developers by restricting the freedom in RDFa
- Activities, Businesses, Groups, Organizations, People, Places, Products and Entertainment
- Only HTML <head> accepted

```
<html xmlns:og="http://opengraphprotocol.org/schema/">
<head>
<title>The Rock (1996)</title>
<meta property="og:title" content="The Rock" />
<meta property="og:type" content="movie" />
<meta property="og:url"
content="http://www.imdb.com/title/tt0117500/" />
<meta property="og:image" content="http://ia.media-
imdb.com/images/rock.jpg" /> …
</head> ...
```

# Current state of metadata on the Web

AIFB ▣

- ✦ 31% of webpages, 5% of domains contain some metadata
  - Analysis of the Bing Crawl (US crawl, January, 2012)
  - RDFa is most common format
    - ✦ By URL: 25% RDFa, 7% microdata, 9% microformat
    - ✦ By eTLD (PLD): 4% RDFa, 0.3% microdata, 5.4% microformat
  - Adoption is stronger among large publishers
    - ✦ Especially for RDFa and microdata
- See also
  - P. Mika, T. Potter. Metadata Statistics for a Large Web Corpus, LDOW 2012
  - H.Mühleisen, C.Bizer.Web Data Commons - Extracting Structured Data from Two Large Web Corpora, LDOW 2012

# Exponential growth in RDFa data



Another five-fold increase between October 2010 and January, 2012

Five-fold increase between March, 2009 and October, 2010

Sep, 2008
Mar, 2009
Oct, 2010

3.5
3.0
2.5
2.0
1.5
1.0
0.5
0.0

RDFa    eRDF    tag    hcard    adr    hatom    xfn    geo    hreview

Percentage of URLs with embedded metadata in various formats

# Query Processing / Matching

# Structure

- Taxonomy of search approaches
- Query processing / matching techniques for Semantic Search
- Types of semantic data
- Formalisms for querying semantic data
- Approaches
  - General task: hybrid graph pattern matching
  - Matching keyword query against text
  - Matching structured query against structured data
  - Matching keyword query against structured data
  - Matching structured query against text (a hybrid case)
- Main tasks, challenges and opportunities

# Taxonomy of search approaches

**AIFB** ⬡

- The search problem
  - A collection of resources, called *data*
  - Information needs expressed as *queries*
  - Search is the task of **efficiently computing results** from data that are **relevant** to queries
- **Document** data retrieval vs. **structured data** retrieval
  - Differences in query and data representation and matching
  - Efficiently retrieve structured data that exactly match formal information needs expressed as structured queries
  - Effectively rank textual results that match ambiguous NL / keyword queries to a certain degree (notions of relevance)
- Semantic search: **ranked** retrieval of document and structured data (given **ambiguous** queries / data)
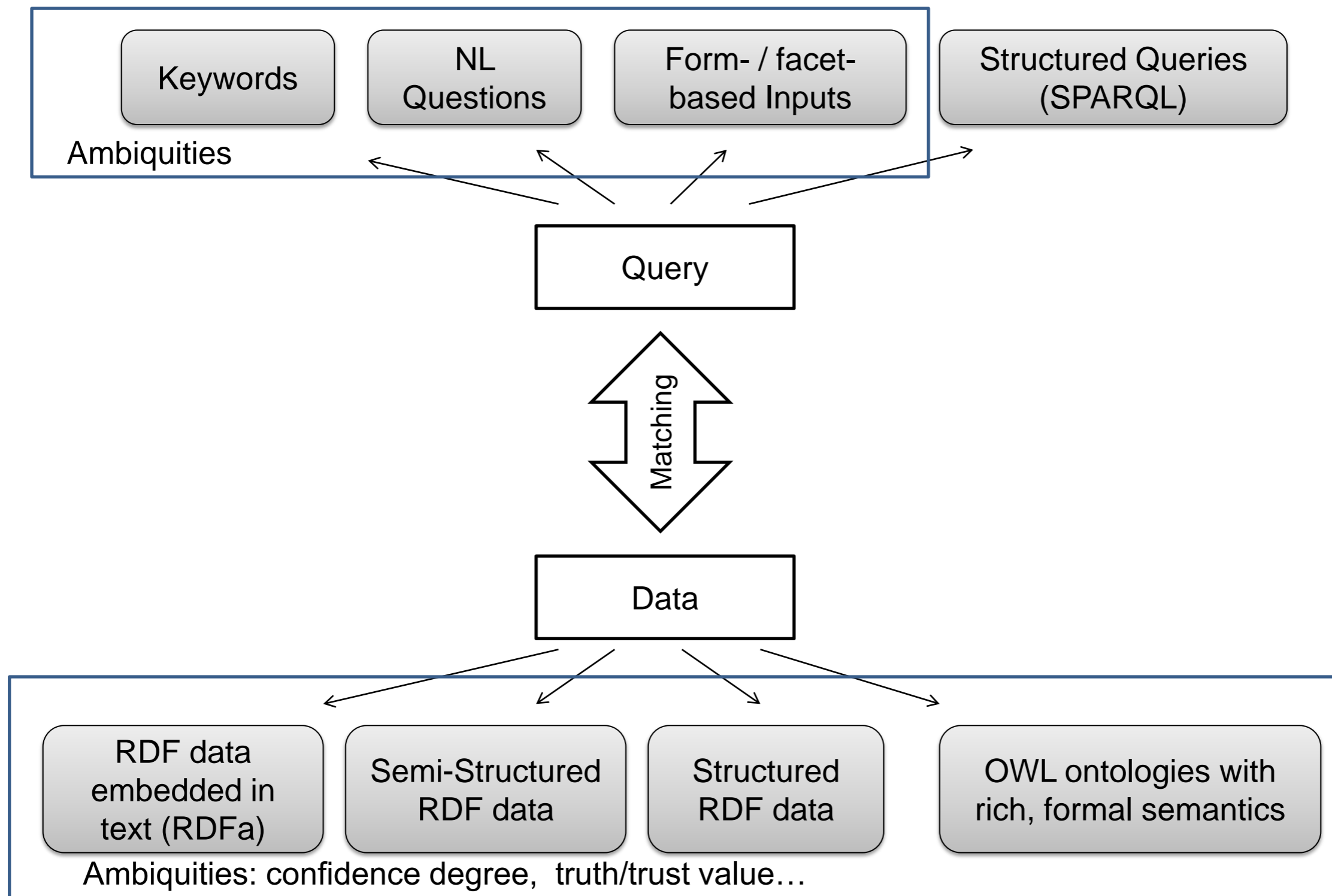
# Query processing for Semantic Search (1)

AIFB

- Resources represented by semantic data ranging from
  - Structured data with well defined schemas
  - Semi-structured data with incomplete or no schemas
  - Data that largely comprise text
  - Hybrid / embedded data
- Information needs of varying complexity, captured using different formalisms and querying paradigms
  - Natural language texts and keywords
  - Form-based inputs
  - Formal structured queries

  (*Search is end-user oriented paradigm, requires "natural", intuitive querying interfaces*)

- Semantic search: efficiently computing results (**query processing)** from data that are relevant to queries (**ranking**)

# Query processing for Semantic Search (2)

**AIFB** ⬡

Keywords | NL Questions | Form- / facet-based Inputs | Structured Queries (SPARQL)

Ambiquities

**Query**

**Matching**

**Data**

RDF data embedded in text (RDFa) | Semi-Structured RDF data | Structured RDF data | OWL ontologies with rich, formal semantics

Ambiquities: confidence degree, truth/trust value…

# Query processing for Semantic Search (3)

Textual Data

Unstructured Query

Semantic Search target different group of users, information needs, and types of data. Query processing for semantic search is **hybrid combination of techniques**!
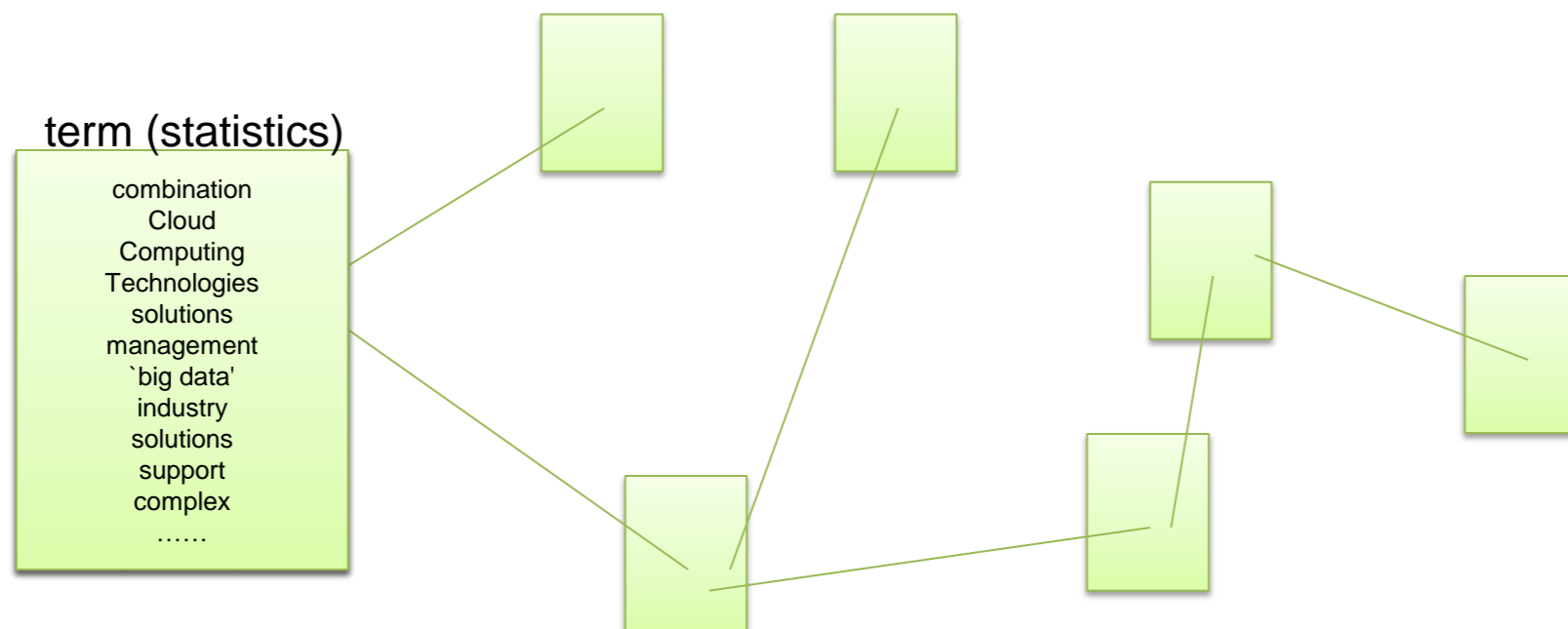
Structured Query
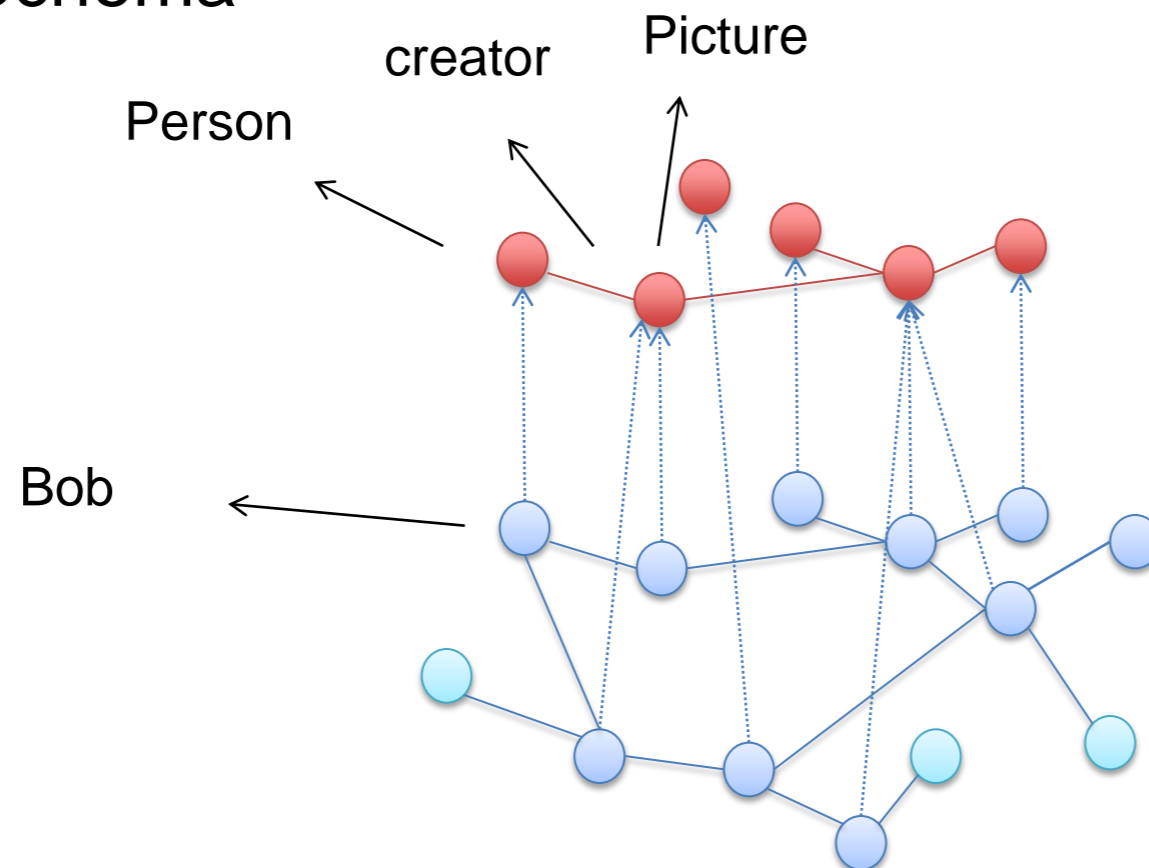
Structured Data

# Types of data models (1)

- Textual
  - **Bag-of-words**
  - Represent documents, text in structured data,…, real-world objects (captured as structured data)
  - Lacks "structure"
    - in text, e.g. linguistic structure, hyperlinks, (positional information)
    - Structure in structured data representation

term (statistics)

combination
Cloud
Computing
Technologies
solutions
management
`big data'
industry
solutions
support
complex
……
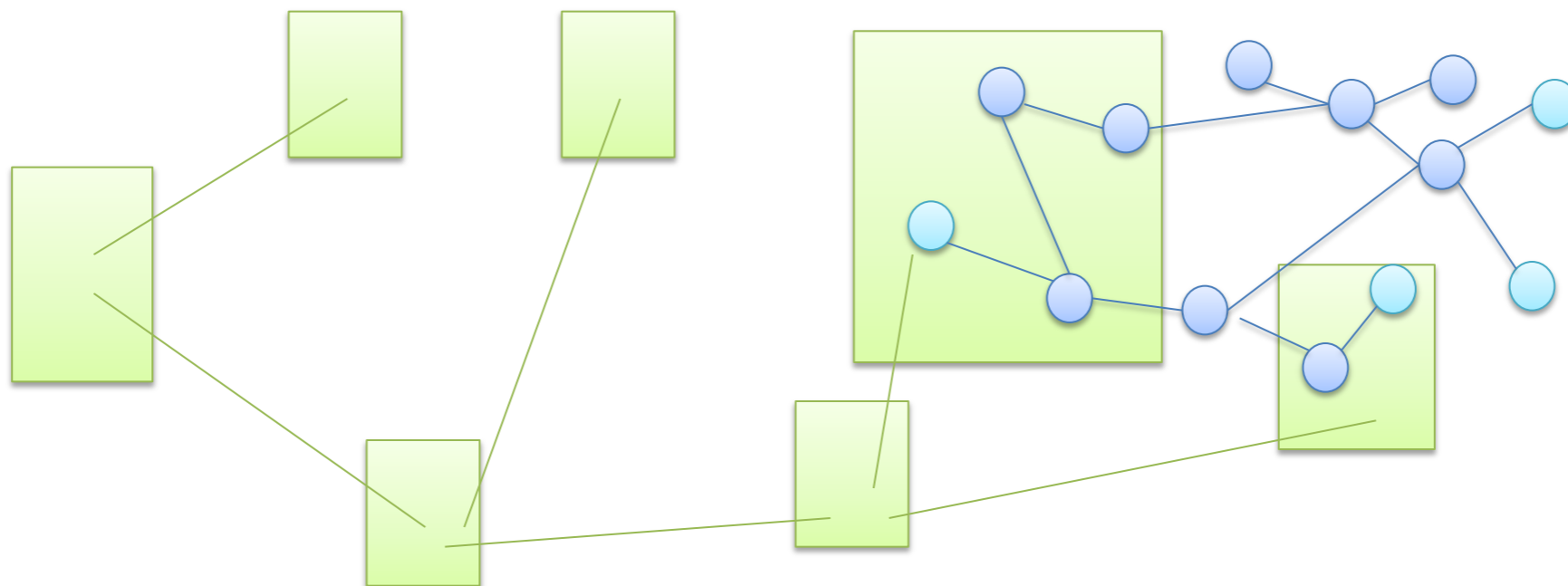
# Types of data models (2)

- Textual
- Structured
  - **Resource Description Framework (RDF)**
  - Represent real-world objects, services, applications, …. documents
  - Resource attribute values and relationships between resources
  - Schema

creator     Picture

Person

Bob

# Types of data models (3)

- Textual
- Structured
- **Hybrid**
  - RDF data embedded in text (RDFa)

# Types of data models – RDFa (1)

AIFB ⬢

```
…
<div about="/alice/posts/trouble_with_bob">
    <h2 property="dc:title">The trouble with Bob</h2>
    <h3 property="dc:creator">Alice</h3>

        Bob is a good friend of mine. We went to the same university, and
        also shared an apartment in Berlin in 2008. The trouble with Bob is        that he takes
much better photos than I do:

    <div about="http://example.com/bob/photos/sunset.jpg">
      <img src="http://example.com/bob/photos/sunset.jpg" />
      <span property="dc:title">Beautiful Sunset</span>
      by <span property="dc:creator">Bob</span>.
    </div>
</div>
…
```

adopted from : http://www.w3.org/TR/xhtml-rdfa-primer/

# Types of semantic data – RDFa (2)

**AIFB** ▢

Bob is a good friend of mine. ← content

We went to the same university, and also shared an apartment in Berlin in 2008. The trouble with Bob is that he takes much better photos than I do:

content

`<http://example.com/alice/posts/trouble_with_bob>`

dc:creator

dc:title

"The Trouble with Bob"   "Alice"

`<http://example.com/bob/photos/sunset.jpg>`

dc:creator

dc:title

"Beautiful Sunset"   "Bob"

adopted from : http://www.w3.org/TR/xhtml-rdfa-primer/

# Types of semantic data - conclusion

Semantic data in general can be conceived as a **graph** with **text** and **structured data** items as nodes, and edges represent different types of relationships including explicit **semantic relationships** and vaguely specified ones such as **hyperlinks**!

# Formalisms for querying semantic data (1)

**Example information need**
*"Information about a friend of Alice, who shared an apartment with her in Berlin and knows someone working at KIT."*

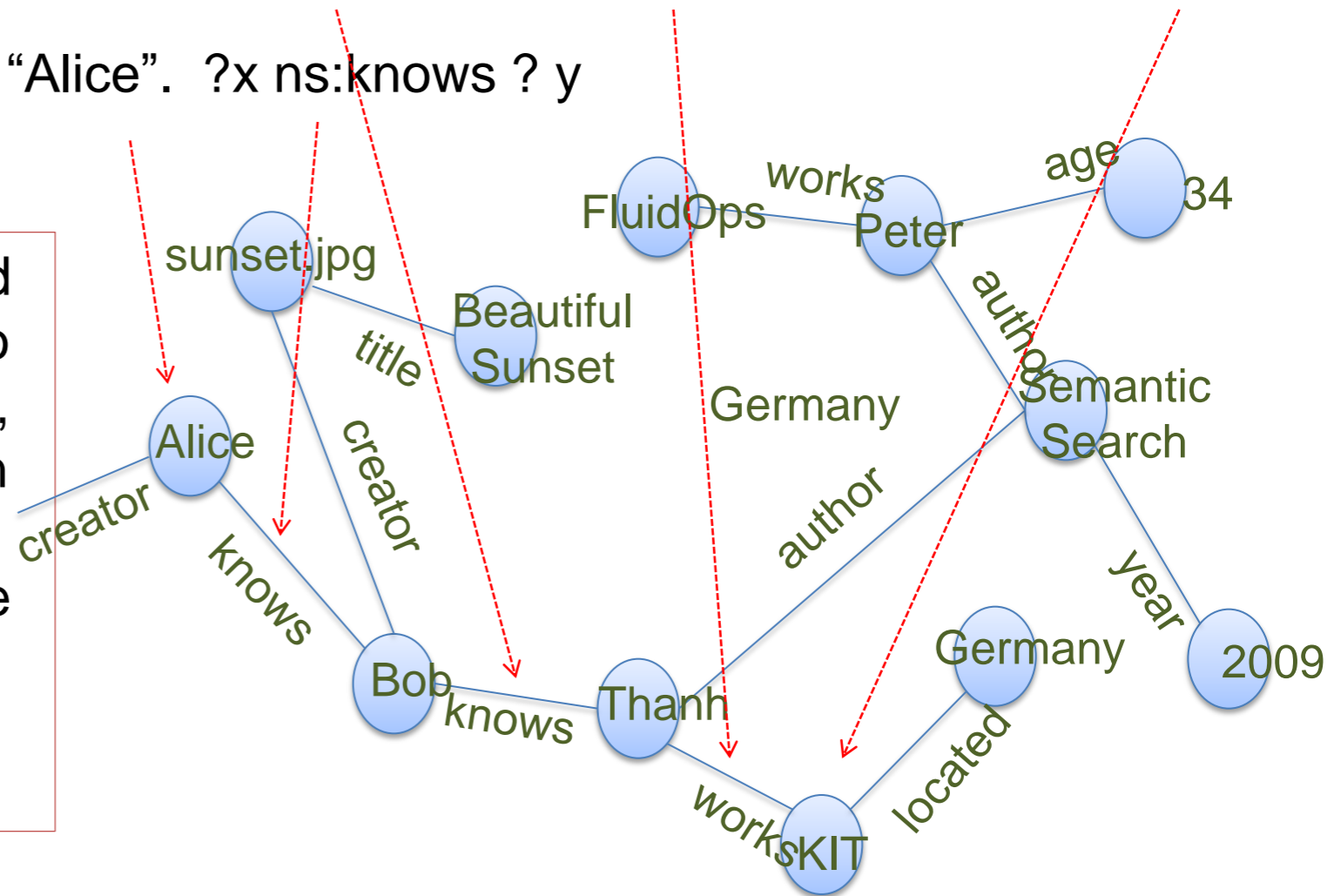- Unstructured queries
- Fully-structured queries
- Hybrid queries: unstructured + structured

# Formalisms for querying semantic data (2)

AIFB

**Example information need**
*"Information about a friend of **Alice**, who **shared an apartment with her in Berlin** and knows someone working at KIT."*
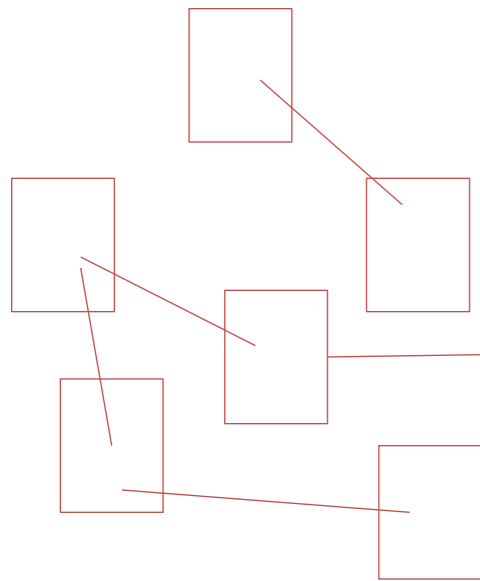
- Unstructured
  - NL
  - **Keywords**

| shared | apartment | Berlin | Alice |

# Formalisms for querying semantic data (3)

AI

**Example information need**

*"Information about **a friend of Alice**, who shared an apartment with her in Berlin and **knows someone working at KIT**."*

- Unstructured

- Fully-structured

  - SPARQL: **BGP**, filter, optional, union, select, construct, ask, describe

    ✦ PREFIX ns: <http://example.org/ns#>
      SELECT ?x
      WHERE { ?x ns:knows ? y. ?y ns:name "Alice".
          ?x ns:knows ?z.  ?z ns: works ?v. ?v ns:name "KIT" }

# Formalisms for querying semantic data (4)

- Fully-structured

- Unstructured

- Hybrid: content and structure constraints

"shared apartment Berlin Alice"

?x ns:knows ? y. ?y ns:name "Alice".
?x ns:knows ?z.  ?z ns: works ?v.
?v ns:name "KIT"

# Formalisms for querying semantic data (5)

- Fully-structured
- Unstructured
- Hybrid: content and structure constraints

"shared apartment Berlin Alice"

?x ns:knows ? y. ?y ns:name "Alice".
?x ns:knows ?z.  ?z ns: works ?v.
?v ns:name "KIT"

# Formalisms for querying semantic data - conclusion

AIFB

Semantic search queries can be conceived as **graph patterns** with nodes referring to **text** and **structured data** items, and edges referring to relationships between these items!

# Processing hybrid graph patterns (1)

AI

**Example information need**

*"Information about a friend of **Alice**, who **shared an apartment with her in Berlin** and knows someone working at KIT."*

apartment shared Berlin Alice

?x ns:knows ?z.  ?z ns: works ?v. ?v ns:name "KIT"

?y ns:name "Alice".  ?x ns:knows ? y

trouble with bob

Bob is a good friend of mine. We went to the same university, and  also shared an apartment in Berlin in 2008. The trouble with Bob is that he takes much better photos than I do:

- **Matching** hybrid graph patterns against data

# Matching keyword query against text

**AIFB** ⬛

- Retrieve documents
  - Inverted list (inverted index)

keyword → {<doc1, pos, score, ...>,

<doc2, pos, score, ...>, ...}

- AND-semantics: top-k join

shared     Berlin     Alice            shared     Berlin           Alice

D1     D1     D1

shared    =    berlin    =    alice

shared

# Matching structured query against structured data

**AIFB** ⬜

- Retrieve data for triple patterns
    - Index on tables
    - Multiple "redundant" indexes to cover different access patterns
- Join (conjunction of triples)
    - Blocking, e.g. linear merge join (required sorted input)
    - Non-blocking, e.g. symmetric hash-join
    - Materialized join indexes

Per1 ns:works **?v**    **?v** ns:name "KIT"
SP-index        PO-index

?x ns:knows ?y. ?x ns:knows ?z.
?z ns: works ?v. ?v ns:name "KIT"



Per1 ns:works **Ins1 Ins1** ns:name KIT

Per1 ns:works Ins1 Ins1 ns:name KIT

# Matching keyword query against structured data

**AIFB** ⬛

- Retrieve keyword elements
  - Using inverted index
    
    keyword → {<el1, score, ...>, <el2, score, ...>,…}
- Exploration / "Join"
  - Data indexes for triple lookup
  - Materialized index (paths up to graphs)
  - Top-k Steiner tree search, top-k subgraph exploration

Alice  Bob  KIT

Alice ns:knows **Bob**  **Inst1** ns:name KIT

  **Bob** ns:works **Inst1**

Alice  Bob       KIT

# Matching structured query against text

- Based on offline IE (offline see Peter's slides)
- Based on online IE, i.e., "retrieve " is as follows
  - Derive keywords to retrieve relevant documents
  - On-the-fly information extraction, i.e., phrase pattern matching  "X name Y"
  - Retrieve extracted data for structured part
  - Retrieve documents for derived text patterns, e.g. sequence, windows, reg. exp.

?x ns:knows ?y. ?x ns:knows ?z.
?z ns: works ?v. ?v ns:name "KIT"

knows

name    KIT

# Matching structured query against text

AIFB

- Index

  - Inverted index for document retrieval and pattern matching

  - Join index → inverted index for storing materialized joins between keywords

  - Neighborhood indexes for phrase patterns

?x ns:knows ?y. ?x ns:knows ?z.
?z ns: works ?v. ?v ns:name "KIT"

# Query processing – main tasks

Query

Matching

Data

- Retrieval
  - Documents , data elements, triples, paths, graphs
  - Inverted index,…, but also other (B+ tree)
  - Index documents, triples, materialized paths
- Join
  - Different join implementations, efficiency depends on availability of indexes
  - Non-blocking join good for early result reporting and for "unpredictable" Linked Data / data streams scenario

AIFB

# Ranking

# Structure

AIFB

- Problem definition
- Types of ambiguities
- Ranking paradigms
- Model construction
  - Content-based
  - Structure-based

# Ranking – problem definition

Query

Matching

Data

- Ambiguities arise when representation is incomplete / imprecise
- Ambiguities at the level of
    - elements (**content ambiguity**)
    - structure between elements (**structure ambiguity**)

Due to ambiguities in the representation of the information needs and the underlying resources, the results cannot be guaranteed to exactly match the query. Ranking is the problem of determining the **degree of matching** using some notions of **relevance**.

# Content ambiguity

apartment shared Berlin Alice

?x ns:knows ?z.  ?z ns: works ?v. ?v ns:name "KIT"

?y ns:name "Alice".  ?x ns:knows ? y

trouble with bob

Bob is a good friend of mine. We went to the same university, and  also shared an  apartment in Berlin in 2008. The trouble with Bob is that he takes much better photos than I do:

sunset.jpg

Beautiful Sunset

title

creator

Alice

creator

knows

Bob

knows

Thanh

FluidOps

works

Peter

age

34

author

Germany

author

Semantic Search

year

2009

Germany

located

works

KIT

What is meant by "Berlin" in the query?
What is meant by "Berlin" in the data?
A city with the name Berlin?  a person?

What is meant by "KIT" in the query?
What is meant by "KIT" in the data?
A research group?  a university? a location?

# Structure ambiguity

apartment shared Berlin Alice

?x ns:knows ?z. ?z ns: works ?v. ?v ns:name "KIT"

?y ns:name "Alice". ?x ns:knows ? y

trouble with bob

Bob is a good friend of mine. We went to the same university, and also shared an apartment in Berlin in 2008. The trouble with Bob is that he takes much better photos than I do:

sunset.jpg

FluidOps    works    Peter    age    34

Beautiful Sunset

title    Germany    Semantic Search

creator    author

Alice

creator

knows    author

year    2009

Bob    Thanh    Germany

knows    located

works KIT

What is the connection between "Berlin" and "Alice"?
Friend? Co-worker?

What is meant by "works"?
Works at? employed?

# Ambiguity

- Recall: query processing is matching at the level of syntax and semantics

- Ambiguities arise when data or query allow for **multiple interpretations**, i.e. multiple matches
  - **Syntactic**, e.g. works vs. works at
  - **Semantic**, e.g. works vs. employ

- "**Aboutness**", i.e., contain some elements which represent the correct interpretation
  - Ambiguities arise when matching elements of **different granularities**
  - Does *i* contains the interpretation for *j,* given some part(s) of *i* (syntactically/semantically) match *j*
  - E.g. Berlin vs. "…we went to the same university, and also, we shared an apartment in Berlin in 2008…"

- Strictly speaking, ranking is performed after syntactic / semantic matching is done!

# Features: What to use to deal with ambiguities?

**AIFB** ⬛

> What is meant by "Berlin"? What is the connection between "Berlin" and "Alice"?

- ## Content features
  - **Frequencies** of terms: $d$ more likely to be "about" a query term $k$ when $d$ more often, mentions $k$ *(probabilistic IR)*
  - **Co-occurrences**: *terms K that often co-occur form a contextual interpretation, i.e., topics (cluster hypothesis)*

- ## Structure features
  - Consider relevance at level of fields
  - Linked-based popularity

# Ranking paradigms

- Explicit relevance model
  - Foundation: **probability ranking principle**
  - Ranking results by the posterior probability (odds) of being observed in the relevant class:
  - P(w|R) varies in different approaches, e.g., binary independence model, 2-poisson model, **relevance model**

$$\frac{P(D|R)}{P(D/N)}$$

$$P(D \mid R) = \prod_{w \in D} P(w \mid R) \prod_{w \notin D} (1 - P(w \mid N))$$

$$P(w \mid R) \approx P(w \mid q_1, ..., q_k) = \sum_{m \in M} P(m) P(w \mid m) \sum_{i=1}^{k} P(q_k \mid m)$$

# Ranking paradigms

- No explicit notion of relevance: similarity between the query and the document model
    - Vector space model (cosine similarity)
    - Language models (KL divergence)

$$Sim(q,d) = Cos((w_{1,d},...,w_{t,d}),(w_{1,q},...,w_{k,q}))$$

$$Sim(q,d) = -KL(\theta_q \| \theta_d) = -\sum_{t \in V} P(t|\theta_q) \log(\frac{P(t|\theta_q)}{P(t|\theta_d)})$$

# Model construction

- How to obtain

  - Relevance models?

  - Weights for query / document terms?

  - Language models for document / queries?

# Content-based model construction

AIFB

- Document statistics, e.g.
  - Term frequency
  - Document length
- Collection statistics, e.g.
  - Inverse document frequency
  - Background language models

- An object is more likely about "Berlin"?
  - When it contains a **relatively** high number of **mentions** of the term "Berlin"
  - When the number of mentions of this term in the overall collection is relatively low

$$w_{t,d} = \frac{tf}{|d|} * idf$$

$$P(t \mid \theta_d) = \lambda \frac{tf}{|d|} + (1-\lambda)P(t \mid C)$$

# Structure-based model construction

- Consider structure of objects during content-based modeling, i.e., to obtain structured content-based model

    - Content-based model for structured objects, documents and for general tuples

$$P(t \mid \theta_d) = \sum_{f \in F_d} \alpha_f P(t \mid \theta_f)$$

- An object is more likely about "Berlin"?
    - When one of its (important) **fields** contains a relatively high number of mentions of the term "Berlin"

# Structure-based model construction

- PageRank
  - Link analysis algorithm
  - Measuring relative importance of nodes
  - Link counts as a vote of support
  - The PageRank of a node recursively depends on the number and PageRank of all nodes that link to it (incoming links)
- ObjectRank
  - Types and semantics of links vary in structured data setting
  - Authority transfer schema graph specifies connection strengths
  - Recursively compute authority transfer data graph

- An object about "Berlin" is more important than one another?
  - When a relatively large number of objects are linked to it

# Taxonomy of ranking approaches

- Explicitly vs. non-explicitly relevance-based
- Content-based ranking
- Structure-based ranking
- Content- and-structure-based ranking

AIFB ⬛

# Result Presentation

# Search interface

- Input and output functionality
  - helping the user to formulate complex queries
  - presenting the results in an intelligent manner
- Semantic Search brings improvements in
  - Query formulation
  - Snippet generation
  - Suggesting related entities
  - Adaptive and interactive presentation
    - ✦ Presentation adapts to the kind of query and results presented
    - ✦ Object results can be actionable, e.g. buy this product
  - Aggregated search
    - ✦ Grouping similar items, summarizing results in various ways
    - ✦ Filtering (facets), possibly across different dimensions
  - Task completion
    - ✦ Help the user to fulfill the task by placing the query in a task context

# Query formulation

- "Snap-to-grid": suggest the most likely interpretation of the query

  - Given the ontology or a summary of the data

## Freebase Suggest

Freebase Suggest is a jQuery plugin that adds Freebase topic autocomplete to search boxes on your site. Start typing text and the widget suggests relevant matches from the millions of topics on Freebase.com or any subset of types like People, Locations or Animals. Topic flyouts help the user select the correct item which is uniquely identified with a Freebase id.

**Try it out:**          Favorite movie director: | Steven Seagal |

Select an item from the list:

**Steven Seagal**                                          Film director

**view more**

**Features:**

- Cross browser - based on jQuery, teste
- 31KB Minified (+ 19KB for jQuery)
- Cross-domain. No proxy servers required thanks to JSONP.
- Hosted on freebaselibs.com
- Free! (The standard Freebase ToS apply.)

**Steven Seagal**

Date of birth: **Apr 10, 1952**

Place of birth: **Lansing**

Religion: **Tibetan Buddhism, Buddhism**

Steven Frederic Seagal (pronounced /sɨ'gɑːl/; born April 10, 1952) is an American action film actor, producer, writer, martial artist, guitarist and a reserve deputy sheriff. A 7th-dan black belt in aikido, Seagal began his adult life as an aikido in...

*Film actor, Film producer, Martial Artist*

**Add to your site**

It's easy to add Freebase Suggeest to your web page. Just include this html in your document head:

# Enhanced results/Rich Snippets

AIFB

- Use mark-up from the webpage to generate search snippets
  - ◆ Originally invented at Yahoo! (SearchMonkey)

# Other result presentation tasks

- Select the most relevant resources within an RDF document
  - Penin et al. Snippet Generation for Semantic Web Search Engines, ASWC 2010
- For each resource, rank the properties to be

# Aggregated search: facets

# Aggregated search: Sig.ma

# Related entities

Related actors and movies

# Adaptive presentation: housing search

# Resources

- Books
  - Ricardo Baeza-Yates and Berthier Ribeiro-Neto. Modern Information Retrieval. ACM Press. 2011
- Survey papers
  - Thanh Tran, Peter Mika. Survey of Semantic Search Approaches. Under submission, 2012.
- Conferences and workshops
  - ISWC, ESWC, WWW, SIGIR, CIKM, SemTech
  - Semantic Search workshop series
  - Exploiting Semantic Annotations in Information Retrieval (ESAIR)
  - Entity-oriented Search (EOS) workshop

# Plan

AIFB

| |
|---|
| XML und URIs |
| Einleitung in RDF |
| RDF Schema |
| Logik – Grundlagen |
| Semantik von RDF(S) |
| SPARQL – Syntax und Intuition |
| Semantik von SPARQL |
| Linked Data |
| Semantic Search |
| **OWL – Syntax und Intuition I** |
| OWL – Syntax und Intuition II |
| OWL – Semantik und Reasoning |
| Konjunktive Anfragen und Regelsprachen |
| Applications |

**Semantic Web**

**TECHNOLOGIES**

**AIFB** ⬜

- Slides erstellt von Thanh Tran, Peter Mika für das Tutorial "**Semantic Search**"

  - https://sites.google.com/site/kimducthanh/activity