

# HumMer and the three steps of information integration



Modena, 26th November 2005

Felix Naumann

Jens Bleiholder

[naumann@informatik.hu-berlin.de](mailto:naumann@informatik.hu-berlin.de)

Humboldt-Universität zu Berlin



## Humboldt-Universität zu Berlin



## Humboldt-Universität zu Berlin



- Wilhelm and Alexander von Humboldt
- Unity of Teaching and Research
- 29 Nobel-prize winners
  - Mommsen, Hertz, Koch, Hahn, Planck, Einstein,...
- 38,000 students, (1100 computer sciences)
- 560 (tenured) professors
  - 11 + 10 in computer sciences



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

3

## Research Group „Information Integration“



- Head: Felix Naumann (naumann@informatik.hu-berlin.de)
- Researchers
  - Jens Bleiholder (bleiho@informatik.hu-berlin.de)
    - Data fusion in relational data
  - Melanie Weis (mweis@informatik.hu-berlin.de)
    - Object identification in XML Data
- Affiliated
  - Armin Roth (aroth@informatik.hu-berlin.de)
    - Data quality in Peer-Data-Management-Systems
  - Alexander Bilke (bilke@cs.tu-berlin.de)
    - Schema Matching
- Research topics around II
  - Object identification
  - Data fusion
  - Optimization
  - Visualization



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

4

## Further Research Topics



- System P
  - Peer Data Management (PDMS) [BTW'05]
- Bioinformatics
  - BioFast [SIGMOD-Record'04, WebDB'05, ...]
    - Links and Paths through Life Sciences Sources
    - With UMD
  - Aladin [CIDR'05]
    - Almost Automatic Data Integration
    - With Ulf Leser (HU Bioinformatics)
- Information Quality
  - ICIQ community and conference (MIT)
- Cooperations with IBM
  - Clio (IBM Almaden)
  - DB2 Wrapper for Search Engines (IBM SVL)
  - DB2 community (IBM Böblingen)

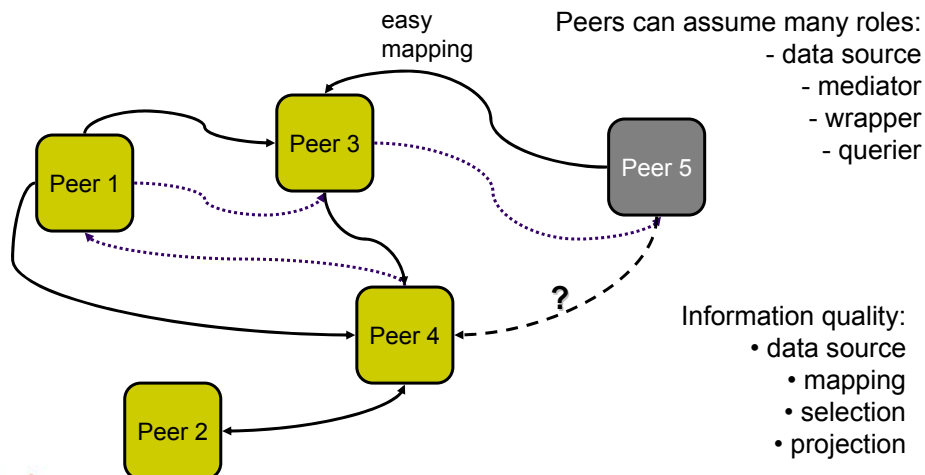


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

5

## Information Quality for Peer-Data-Management (IQ & PDMS)

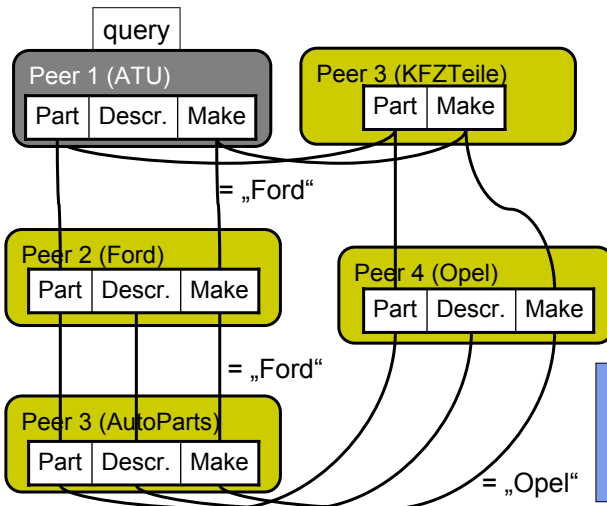


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

6

# PDMS: Incomplete and Selective Mappings



Problem:  
Cumulated selections  
• implicit in schemata  
• explicit in mappings  
• Point selections and range selections

Problem:  
Cumulated projections  
• in schemata  
• in mappings



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

7

# XML related



- XStruct
  - Automatic XML Schema Extraction
    - From many XML files
    - From large XML files
- XQuery Assistance
  - Graphical interface to build XQueries
    - Based on a given Schema
    - Currently: Selection and Projection



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

8

Hummer: XQuery-Generator

christoph@mi-data/ClioSchemas/Chocolate/Customers.xsd

parse XML Schema generate Query execute Query

Element	Incl. Tag	Type Information	Name	Predicate
doc		n.a. CT_IsMixed = false, (1-1)	doc	-
customer		n.a. sequence		-
customer		n.a. CT_IsMixed = false, (0-1)	customer	order/amount > 200
name		n.a. sequence	name	-
address		ST (1-1)	address	-
order		n.a. CT_IsMixed = false, (0-1)	order	-
amount		n.a. sequence	amount	-
date		ST (1-1)	date	-
item		n.a. CT_IsMixed = false, (1-1)	item	-
name		n.a. sequence	name	-
quantity		ST (1-1)	quantity	-
payment		n.a. CT_IsMixed = false, (0-1)	payment	-

```

<result>
{
  for $i55 in doc("Customers.xml")/doc*
  where
    (local-name($i55) = 'customer' and $i55/order/amount > 200)
  return
    if (local-name($i55) = 'customer') then
      element { node-name($i55) }
      {
        for $i56 in $i55/*
        where (local-name($i56) = 'name' or local-name($i56) = 'order')
        return
          (
            if (local-name($i56) = 'name') then
              $i56/text()
            else 0,
            if (local-name($i56) = 'order') then
              for $i57 in $i56/item*
              where local-name($i57) = 'name'
              return
                if (local-name($i57) = 'name') then
                  $i57
                else 0
            else 0
          )
      }
}

```

```

<result>
<customer> Jason Henderson
  <name>Brazilian Coffee </name>
</customer>
<customer> Dean Sura
  <name>Klipsch K3C-50 SportClip Portable Stereophone </name>
  <name>Kodak Max NIMH Rechargeable AA Batteries </name>
  <name>The Complete Monty Python's Flying Circus </name>
  <name>Samsung SC-D67 Digital Video Camcorder </name>
</customer>
</result>

```

9

## ALmost Automatic Data INtegration – ALADIN (with Ulf Leser)

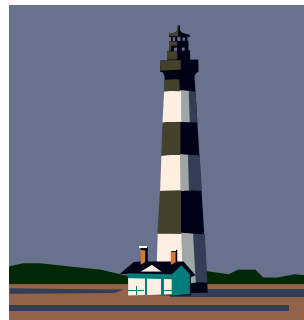
- Observations
  - Data sources have only one "type" of object
    - Genes, proteins, ...
  - Objects have an ID (accession number)
  - Objects have annotations
    - Semi-structured, nested
  - Databases heavily cross-reference each other
    - But cannot keep with the pace of data production
  - Important annotations are often free text
    - which can be compared using text mining
- 5 steps to integration
- Source-specific
  - Download source, parse, import into RDBMS
  - Guess primary objects
  - Guess (hierarchically structured) annotation
- Across data sources
  - Guess cross-references
  - Guess duplicates



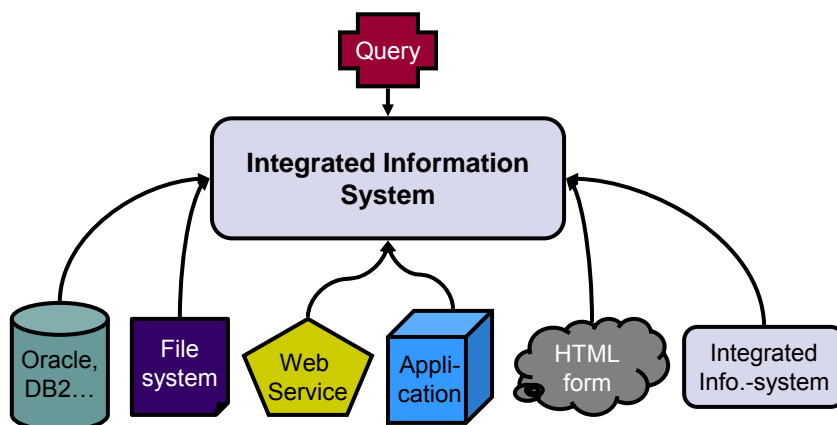
# Overview



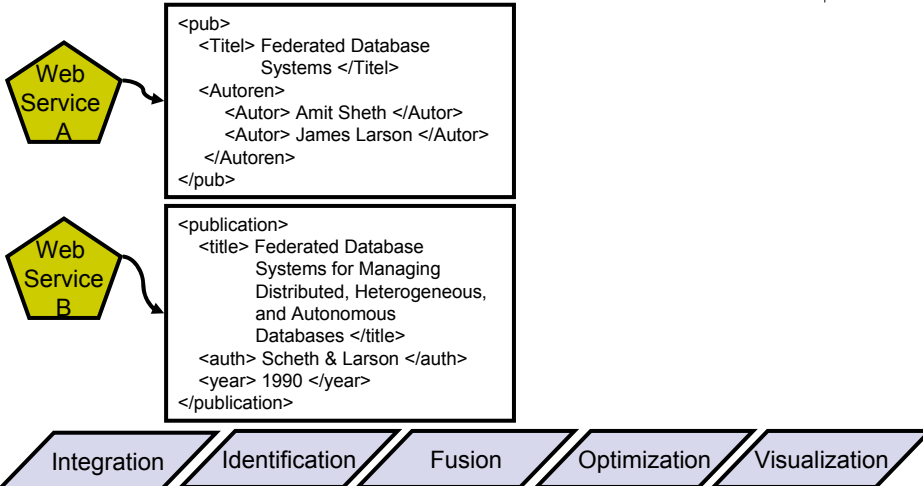
- ➔ ● Introduction and Motivation
- Information Integration in Three Steps
  1. Schema Mapping
    - And Schema Matching
  2. Duplicate Detection
  3. Data Fusion
- HumMer Architecture
- HumMer Demo



# Integrated Information Systems



# Information Integration

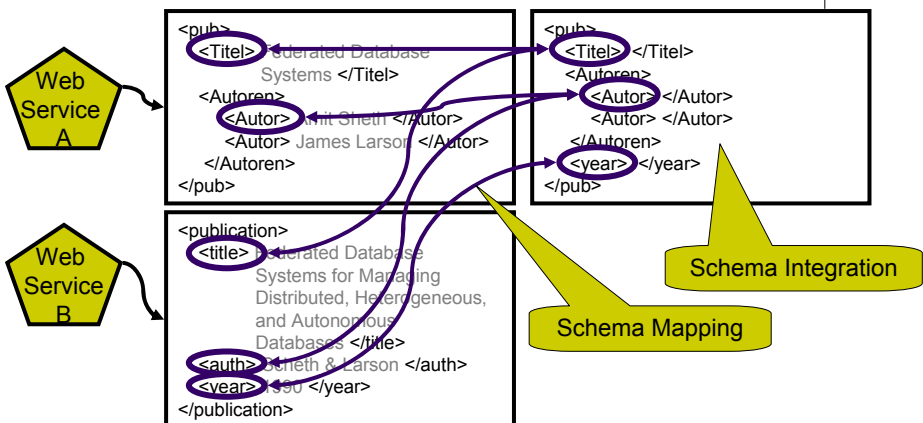


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

13

# Information Integration

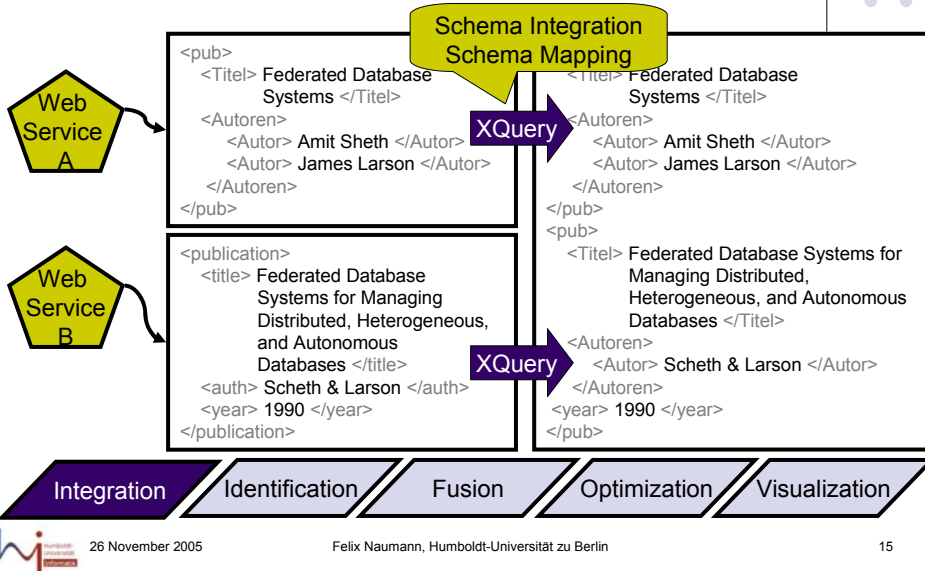


26 November 2005

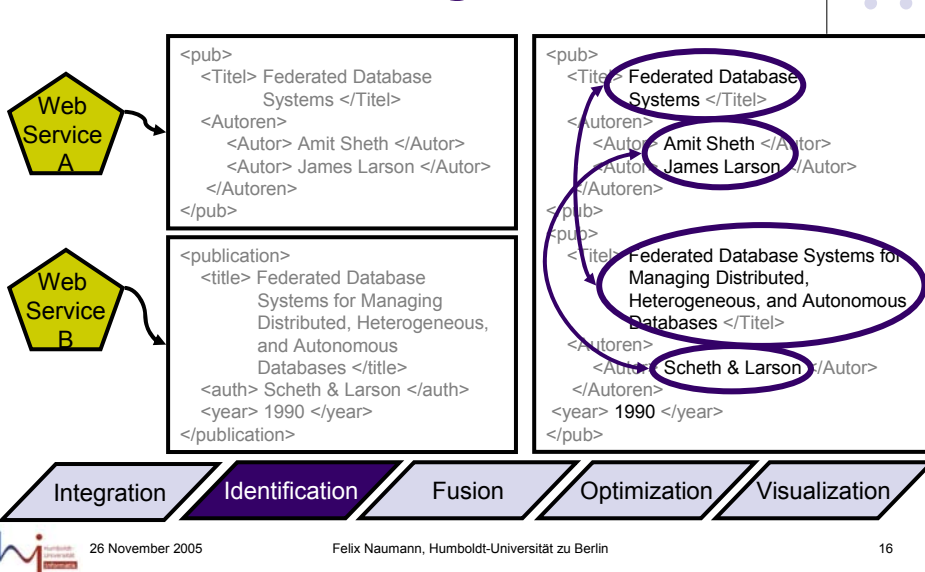
Felix Naumann, Humboldt-Universität zu Berlin

14

# Information Integration



# Information Integration





# Information Integration



Web Service A

```
<pub>
<Titel> Federated Database
Systems </Titel>
<Autoren>
<Autor> Amit Sheth </Autor>
<Autor> James Larson </Autor>
</Autoren>
</pub>
```

Web Service B

```
<publication>
<title> Federated Database
Systems for Managing
Distributed, Heterogeneous,
and Autonomous
Databases </title>
<auth> Scheth & Larson </auth>
<year> 1990 </year>
</publication>
```

```
<pub>
<Titel> Federated Database
Systems </Titel>
<Autoren>
<Autor> Amit Sheth </Autor>
<Autor> James Larson </Autor>
</Autoren>
</pub>
<pub>
<Titel> Federated Database Systems for
Managing Distributed,
Heterogeneous, and Autonomous
Databases </Titel>
<Autoren>
<Autor> Scheth & Larson </Autor>
</Autoren>
<year> 1990 </year>
</pub>
```

Integration

Identification

Fusion

Optimization

Visualization



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

17

# Information Integration



Web Service A

```
<pub>
<Titel> Federated Database
Systems </Titel>
<Autoren>
<Autor> Amit Sheth </Autor>
<Autor> James Larson </Autor>
</Autoren>
</pub>
```

Web Service B

```
<pub>
<Titel> Federated Database Systems for
Managing Distributed,
Heterogeneous, and Autonomous
Databases </Titel>
<Autoren>
<Autor> Scheth & Larson </Autor>
</Autoren>
<year> 1990 </year>
</pub>
```

```
<pub>
<Titel> Federated Database Systems for
Managing Distributed,
Heterogeneous, and
Autonomous Databases </Titel>
<Autoren>
<Autor> Amit Sheth </Autor>
<Autor> James Larson </Autor>
</Autoren>
<year> 1990 </year>
</pub>
```

Integration

Identification

Fusion

Optimization

Visualization

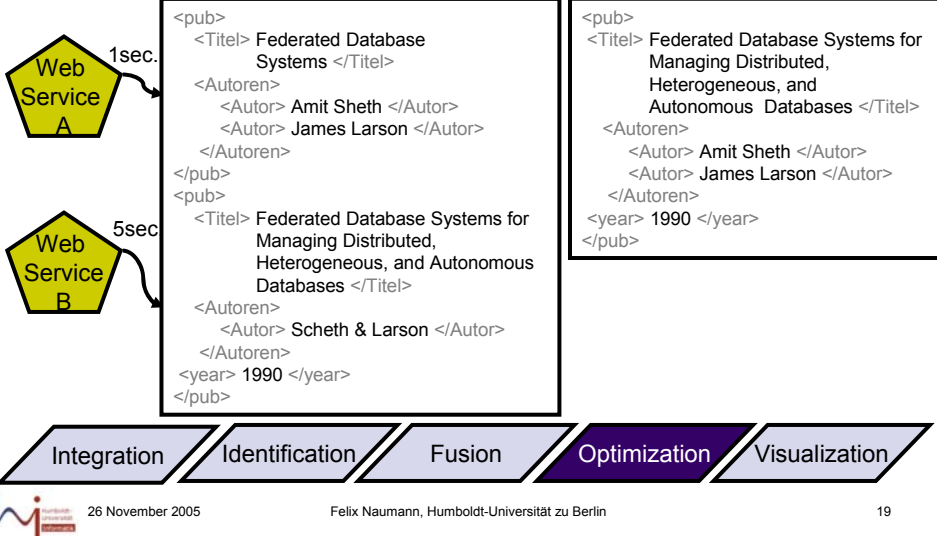


26 November 2005

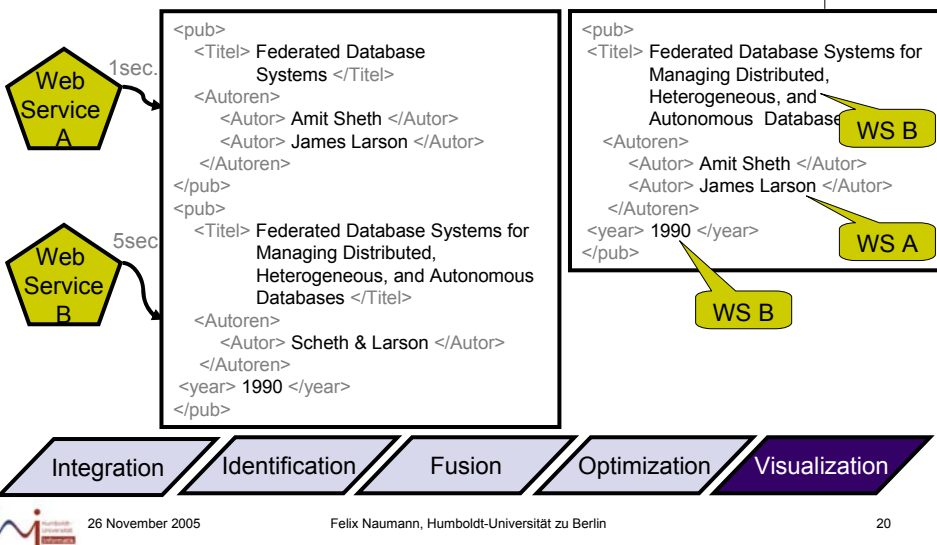
Felix Naumann, Humboldt-Universität zu Berlin

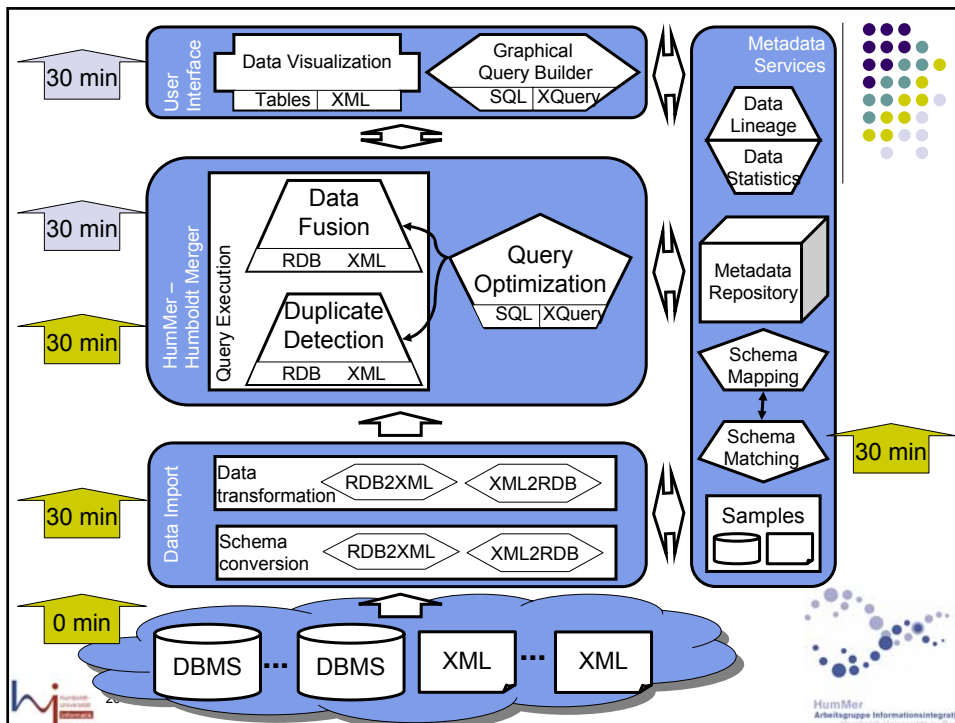
18

# Information Integration



# Information Integration





# The HumboldtMerger Logo

**HumMer**

Humboldt-Universität zu Berlin

26 November 2005      Felix Naumann, Humboldt-Universität zu Berlin      22

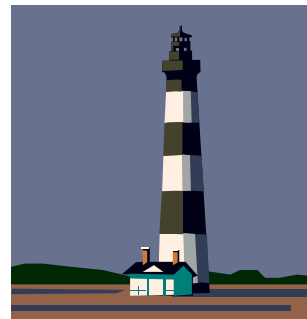
# HumMer The Humboldt Merger



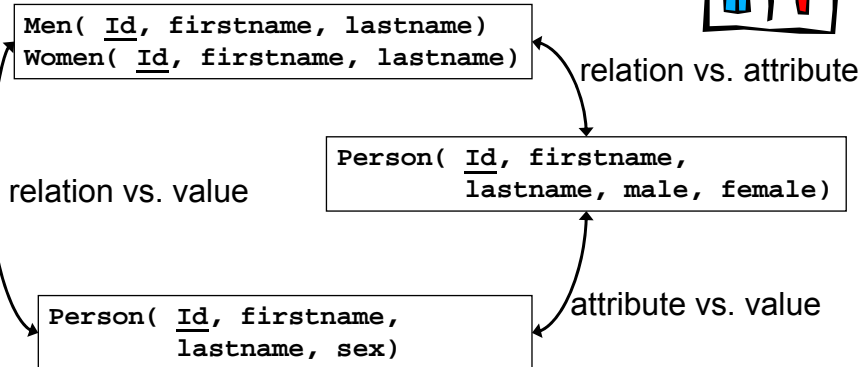
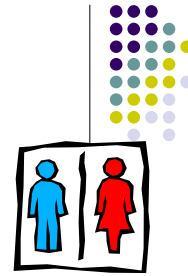
## Overview



- Introduction and Motivation
- Information Integration in Three Steps
  - ➔ 1. Schema Mapping
    - And Schema Matching
  - 2. Duplicate Detection
  - 3. Data Fusion
- HumMer Architecture
- HumMer Demo



## Schematic Heterogeneity



26 November 2005

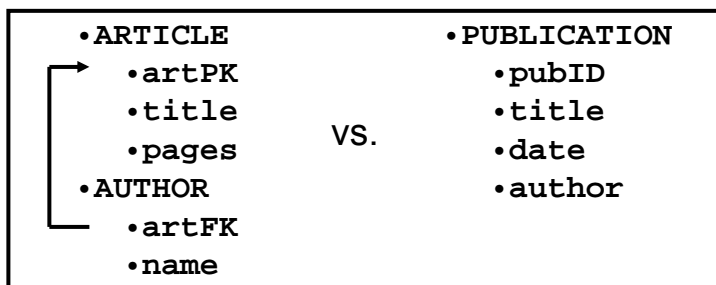
Felix Naumann, Humboldt-Universität zu Berlin

25

## Schematic Heterogeneity – Example



- Normalized vs. Denormalized
  - 1:1 associations between values are represented differently
    - In same tuple
    - Key-foreign key relationship



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

26

## Schematic Heterogeneity – Example



- Nested vs. Flat
  - 1:n associations represented differently
    - As nested elements
    - Key-foreign key relationship
    - Redundant and denormalized

•ARTICLE		•PUBLICATION
•artPK		•pubID
•title	VS.	•title
•pages		•author
•AUTHOR		
•name		



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

27

## Schematic Heterogeneity - Solutions



- Two alternative problems
  1. Access multiple data sets uniformly
    - At schema level:
      - Schema languages (SchemaSQL, MSQL, FRAQL)
      - Schema mapping (Clio, Rondo, Tools)
    - At data level: virtual integration
  2. Integrate multiple data sets into a new database
    - At schema level: schema integration
    - At data level: materialized integration, ETL, DWH



26 November 2005

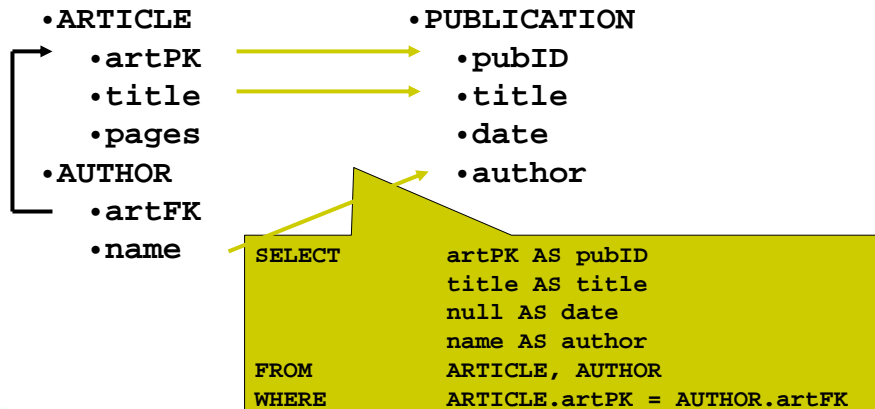
Felix Naumann, Humboldt-Universität zu Berlin

28

# Schema Heterogeneity – Solutions



## • Schema Mapping

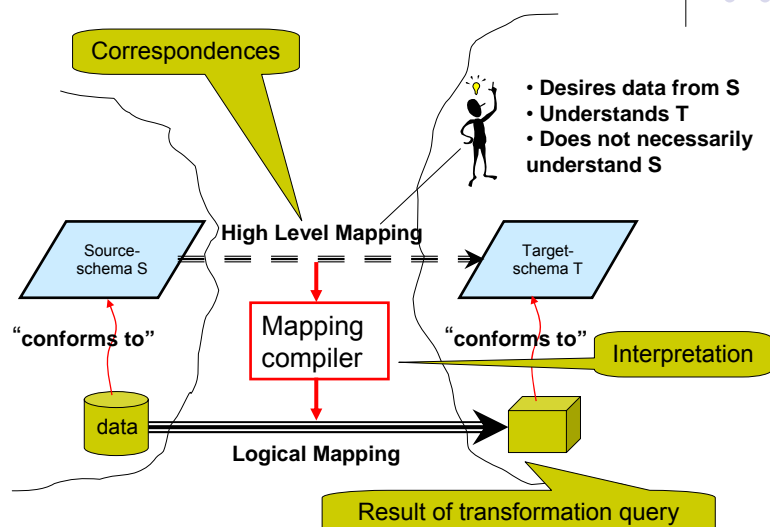


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

29

# Schema Mapping in Context



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

30

# Motivation: Why is Schema Mapping useful?



- Data transformation between heterogeneous schemata
  - Old but omnipresent problem
  - Usually experts formulate complex query or programs
    - Time intensive
    - Expert for domain, for schemata and for query (language)
    - XML makes everything even more difficult
      - XML Schema, XQuery
- Idea: Automation
  - Given: Two schemata and a high-level mapping between them
  - Wanted: Query for data transformation
  - Later: Schema Matching = semi-automatic generation of high-level mapping



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

31

# Schema Mapping in XQuery



1. Schema Matching & Correspondences
2. Schema Mapping
3. Mapping Interpretation
4. Data Transformation

The screenshot shows the XQuery editor interface. On the left, the 'Source Schemas' pane displays a tree view of the 'expenseDB: Reco' schema, including elements like 'company', 'grant', and 'project'. The main editor window shows an XQuery query with the following structure:

```

xquery
  $x0L1/project/text() = $x1L1/name/text() AND
  $x2L1/cid/text() = $x0L1/cid/text() AND
  $x2L1/city/text() = $x2L1/city/text()
RETURN
  <organization>
    <cid> $x0L1/cid/text() </cid>,
    <cname> $x2L1/cname/text() </cname>,
    distinct (
      FOR
        $x0L2 IN $doc/expenseDB/grant,
        $x1L2 IN $doc/expenseDB/project,
        $x2L2 IN $doc/expenseDB/company
      WHERE
        $x0L2/project/text() = $x1L2/name/text() AND
        $x2L2/cid/text() = $x0L2/cid/text() AND
        $x2L1/cname/text() = $x2L2/cname/text() AND
        $x2L1/city/text() = $x2L2/city/text() AND
        $x0L1/cid/text() = $x0L2/cid/text()
      RETURN
        <funding>
          <gid> $x0L2/gid/text() </gid>,
          <proj> $x0L2/project/text() </proj>,
          <faiid> "SK267*", $x0L2/project/text(), ' ,
        </funding> )
    </organization> ),
  distinct (
    FOR
      $x0L1 IN $doc/expenseDB/grant,
      $x1L1 IN $doc/expenseDB/project,
      $x2L1 IN $doc/expenseDB/company
    WHERE
  
```



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin



## Motivation: Why is Schema Mapping Difficult?



- Generation of „correct“ query under constraints of
  - source and target schemata,
  - the mappings,
  - and the users intention: semantics!
- Guarantee that transformed data adheres to target schema
  - Flat or nested
  - Integrity constraints
- Efficient data transformation
  - For materialized integration
  - For virtual integration

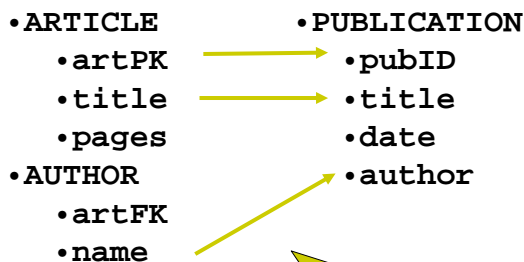


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

33

## Schema Mapping Example



```
SELECT artPK AS pubID UNION SELECT null AS pubID
      title AS title      null AS title
      null AS date        null AS date
      null AS author      name AS author
FROM ARTICLE              FROM AUTHOR
```

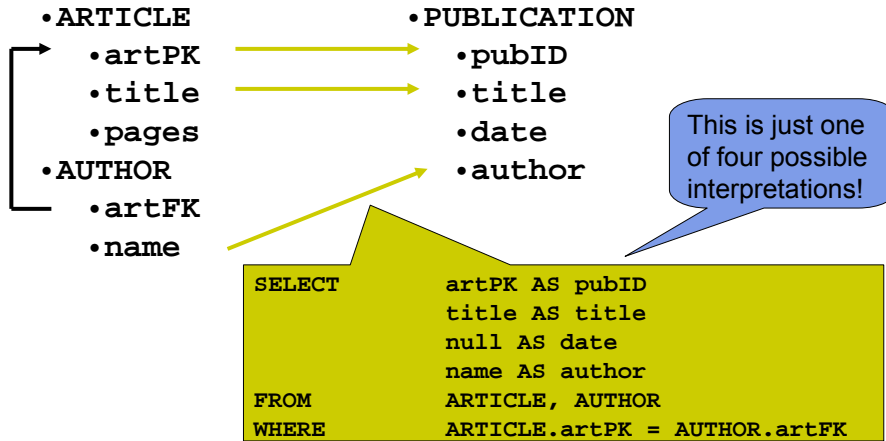


26 November 2005

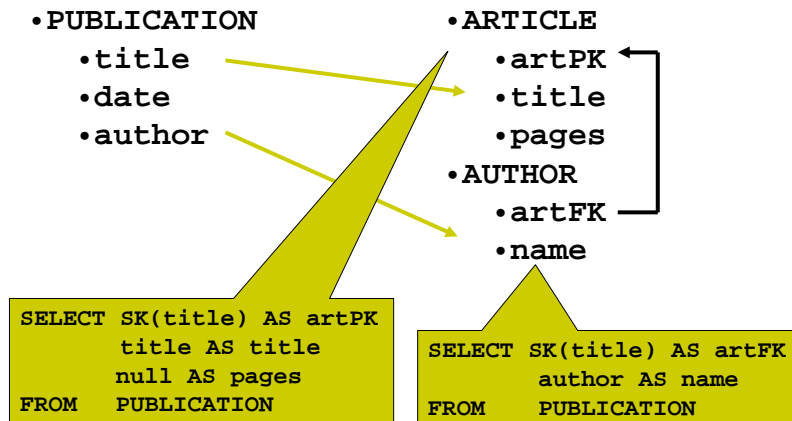
Felix Naumann, Humboldt-Universität zu Berlin

34

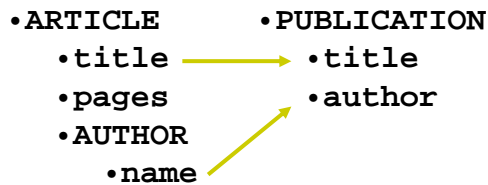
# Schema Mapping Example



# Schema Mapping Example



## Schema Mapping Example



```
LET $doc0 := document("articlcs.xml") RETURN
<dblp> { distinct-values (
  FOR $x0 IN $doc0/authorDB/ARTICLE, $x1 IN $x0/AUTHOR
  RETURN
    <publication>
      <title> { $x0/title/text() } </title>
      <author> { $x1/name/text() } </author>
    </publication> )
} </dblp>
```

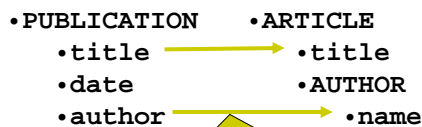


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

37

## Schema Mapping Example



```
LET $doc0 := document("publication.xml")
RETURN
<articles> { distinct-values (
  FOR $x0 IN $doc0/dblp/publication RETURN
    <ARTICLE>
      <title> { $x0/title/text() } </title>
      { distinct-values (
        FOR $x0L1 IN $doc0/dblp/publications
        WHERE $x0/title/text() = $x0L1/title/text()
        RETURN
          <AUTHOR>
            <name> { $x0L1/author/text() } </name>
          </AUTHOR> ) }
    </ARTICLE> ) } </articles>
```



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

38

## Schema Mapping Tools



- Research
  - Clio
  - Rondo
  - MOMIS
- Industrial
  - WSAD
  - BizTalk
  - MapForce
  - ...
  - ETL
- Comparing tools
  - Data Models
  - Query Languages
  - Interface
  - Interpretation!



26 November 2005

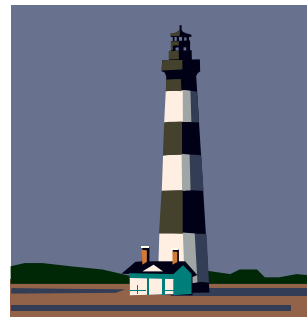
Felix Naumann, Humboldt-Universität zu Berlin

39

## Overview



- Introduction and Motivation
- Information Integration in Three Steps
  1. Schema Mapping
    - Optional: Mapping Interpretation
    - And Schema Matching
  2. Duplicate Detection
  3. Data Fusion
- HumMer Architecture
- HumMer Demo



26 November 2005

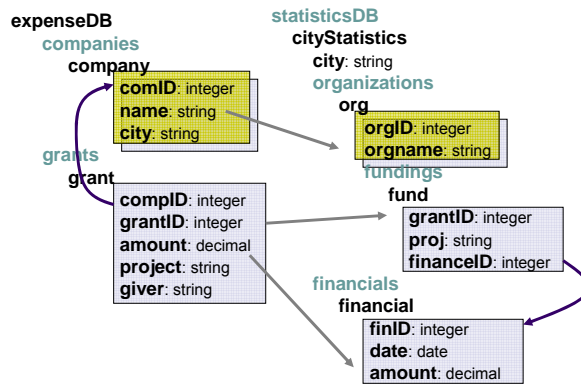
Felix Naumann, Humboldt-Universität zu Berlin

40

# Mapping – Clio's Algorithm [FHP+02]



- Three steps
  1. Discovery of intra-schema associations
  2. Discovery of inter-Schema logical mappings
  3. Query generation



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

41

# Discovery of Associations



- Step 1
  - Intra-schema associations between schema elements
  - Relational views contain maximal groups of associated elements
  - Each view represents a unique category of data of the data source
  - Independent of mapping (but restricted to mapped elements)



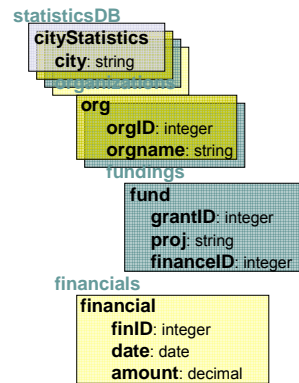
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

Source: [FHP+02]

## Discovery of Associations

- Start: All „primary“ paths
  - Associations in schema without integrity constraints
- For relational schemas
  - Each relation is one primary path
- For nested schemas
  - Attributes of one level
  - Attributes of nested levels



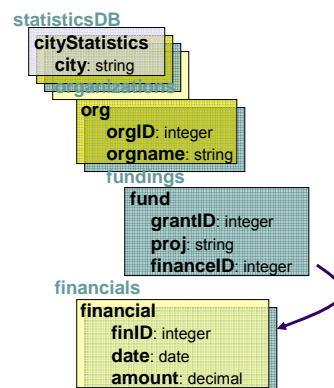
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

43

## Discovery of Associations

- Include key/foreign key constraints
- Logical relation
  - Extend each primary path by chasing constraints

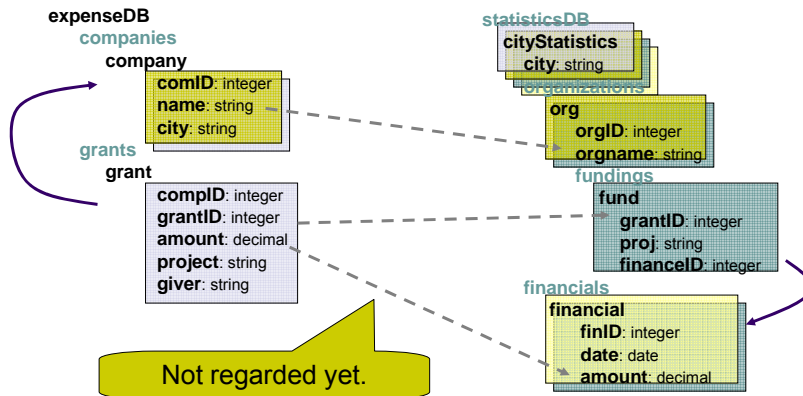


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

44

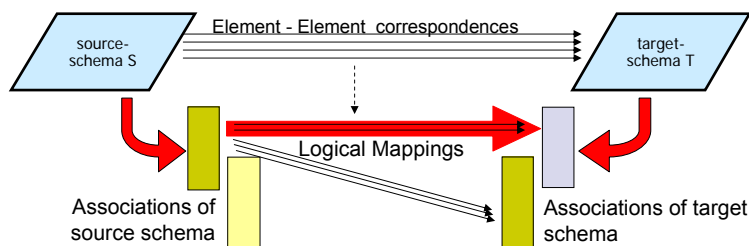
# Discovery of Associations



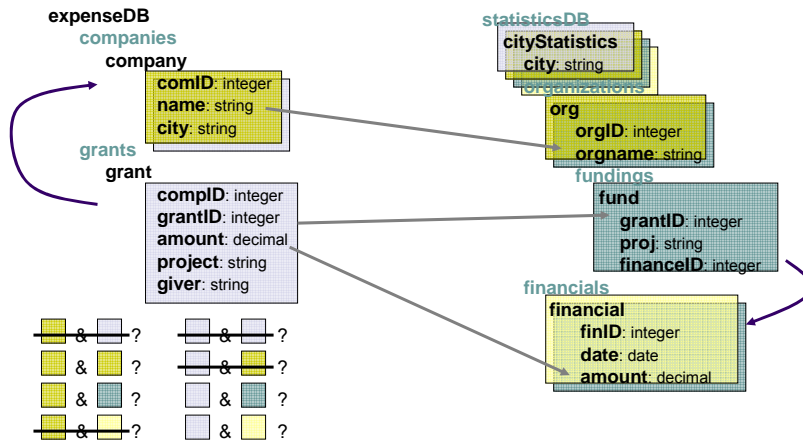
# Discovery of Logical Mappings



- Step 2
  - Discover logical mappings between source and target schema
  - Regard all combinations of associations of source schema and associations of target schema
    - Under consideration of all correspondences



# Discovery of Logical Mappings

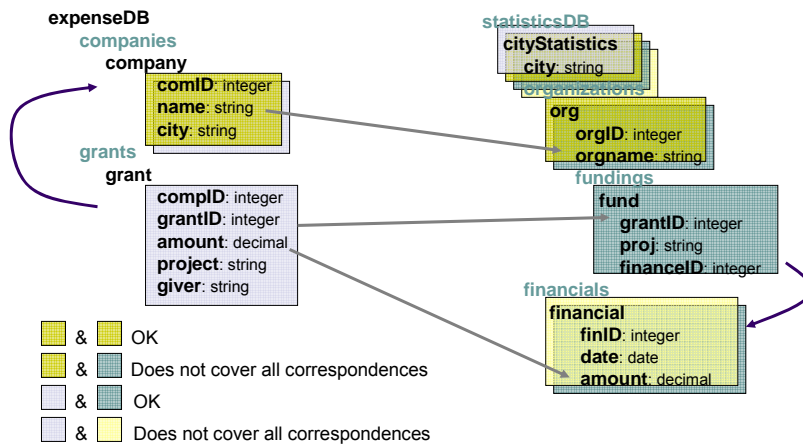


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

47

# Discovery of Logical Mappings



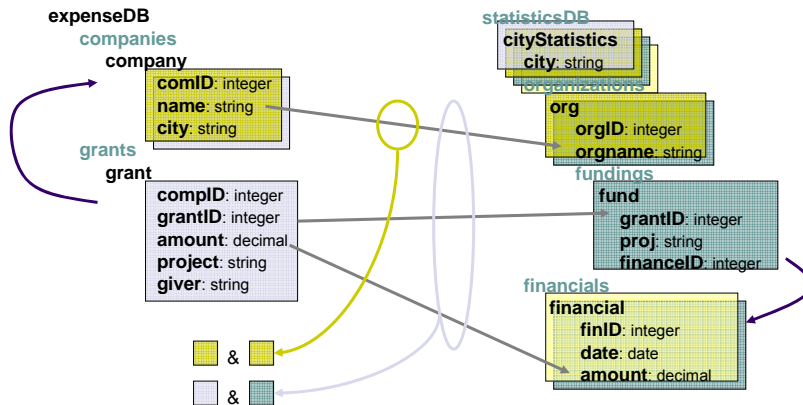
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

48



# Discovery of Logical Mappings

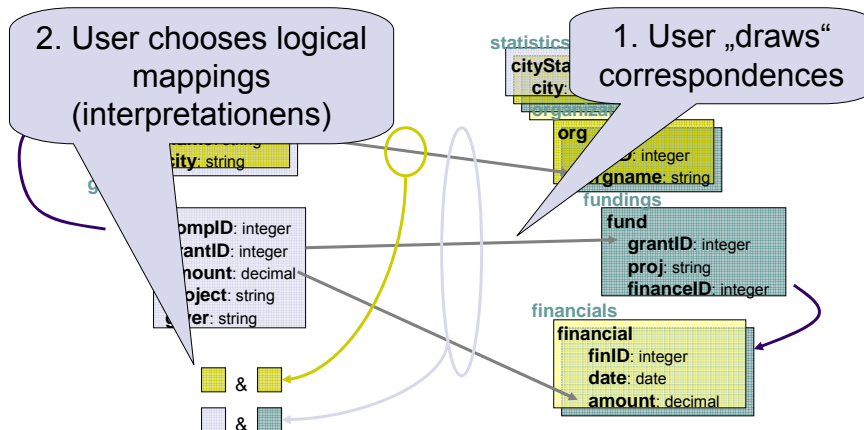


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

49

# User-Input or Expert-Input



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

50

## Final Step: Query Generation



- Implementationspecific
- Depending on data models
  - Relational → Relational: SQL
  - Relational → XML: SQLX (nesting and tagging of result)
  - XML → Relational: XQuery or XSLT (remove *tags*)
  - XML → XML: XQuery or XSLT
- In Clio
  - Generate proprietary rules
  - Then: Translate rules into query of specific language
- Not here...



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

51

## Overview



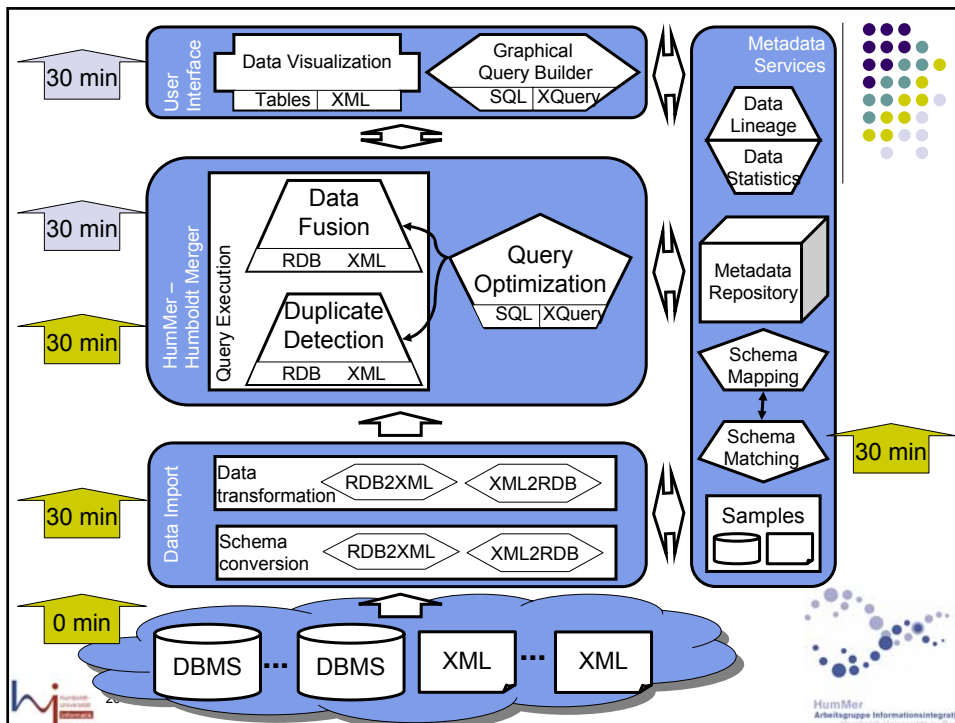
- Introduction and Motivation
- Information Integration in Three Steps
  1. Schema Mapping
    - And Schema Matching
  2. Duplicate Detection
  3. Data Fusion
- HumMer Architecture
- HumMer Demo



26 November 2005

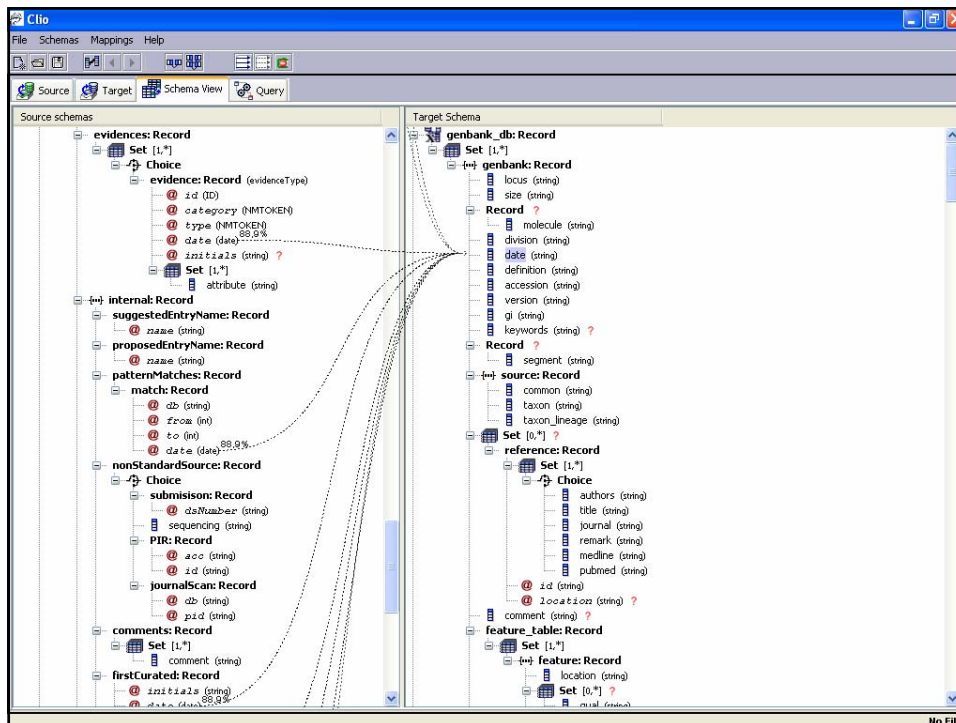
Felix Naumann, Humboldt-Universität zu Berlin

52



## Schema Matching – Motivation

- Schemata are
  - large
  - complex
  - foreign
  - confusing
  - different language
  - cryptic
- > 100 tables, many attributes
- Deep Nesting
- Foreign keys
- XML Schema
- Unknown synonyms
- Unknown homonyms
- |attribute name| ≤ 8
- |table name| ≤ 8



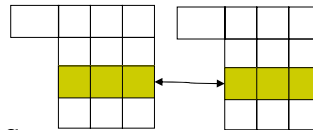
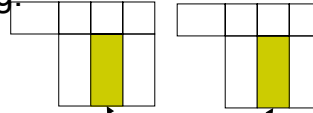
## Schema Matching Algorithms

- Core algorithm of Schema Matching:
  - Compare attribute pairs from two schemata.
  - Choose most similar pairs using similarity measure.
  - Semi-automation: Conformation of suggestions by expert
- Similarity measures based on
  - Schema elements (label-based)
    - And other metadata, such as ontologies
  - Data (instance-based)
  - Structure (structure-based)
  - Combination, Meta-Matcher

# Instance-based Schema Matching



- Instance-based Schema Matching:
  - Correspondences based on similar data values or their properties
- Conventional solution: Vertical
  - Comparison of columns
  - = Attribute classification
  - [ICDE'02]
- Our solution: Horizontal
  - Comparison of rows
  - = Duplicate detection (despite missing attribute correspondences)
  - [ICDE'05]



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

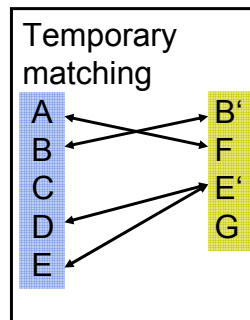
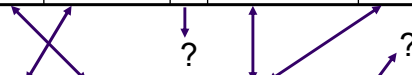
57

# Duplicate-driven Schema Matching



A	B	C	D	E
Max	Michel	m	601- 4839204	601- 4839204
...	...	...	...	...

B'	F	E'	G
Michel	maxm	601- 4839204	UNIX
...	...	...	...



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

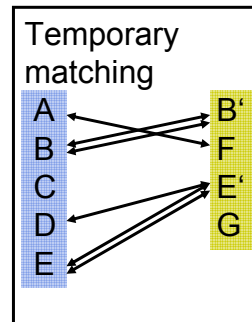
58

# Duplicate-driven Schema Matching



A	B	C	D	E
Max	Michel	m	601- 4839204	601- 4839204
Sam	Adams	m	541- 8127100	541- 8121164

B'	F	E'	G
Michel	maxm	601- 4839204	UNIX
Adams	beer	541- 8127164	WinXP



- Assumptions
  - There is data in both DBs.
  - There are (at least a few) duplicates in both DBs.
  - Equal or similar values reflect same semantics of attributes.



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

59

# Duplicate-driven Schema Matching



- Duplicate detection
  - Goal: Find the top-k duplicates.
  - Problems
    - Correspondences unknown
    - Possibly small intensional overlap
- Schema Matching
  - Goal: Derive attribute correspondences from attribute values.
  - Problems
    - Attribute values only similar, not equal
    - Synonyms and homonyms in values



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

60

# Duplicate Detection in Unaligned Tables



- Cosine measure with TFIDF weights

- Tuple as vector of term weights

- Term weights is high if term appears often in tuple (TF) and it appears in only few tuples (IDF).
 
$$w'(r,t) = \log(tf_{r,t} + 1) \cdot \log\left(\frac{N}{df_t} + 1\right)$$

- Tuple similarity is the cosine of the angle of both (normalized) vectors.

$$tsim(r,s) = \sum_{t \in r \cap s} w(r,t) \cdot w(s,t)$$

123	Max Michel	max@michel.com	601- 4839204	601- 4839204
1	Max Michel	601- 4839204	max@michel.com	

„1“	0	0.06
„123“	0.21	0
„Max“	0.3	0.28
„601“	0.14	0.16
⋮	⋮	⋮



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

61

# Duplicate Detection



- Algorithm

- Difference to regular duplicate detection:
  - Only top K duplicates necessary  $\Rightarrow$  no threshold
  - Comparisons for top 10 duplicates in 5,000 tuples
    - Naive: 25,000,000
    - Inverted index: 340.333
    - WHIRL (A\*): 8.833

- Synthetic person data

- Different intensional overlaps
- Different degree of duplication

- Berlin real estate offers

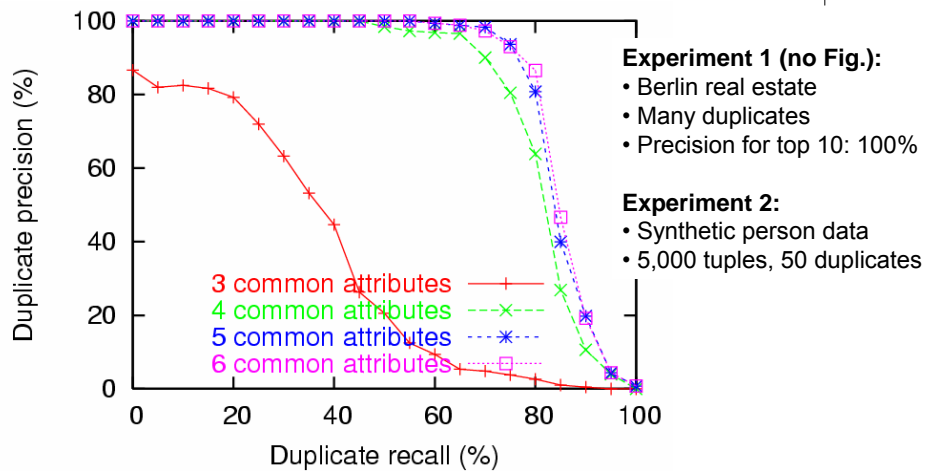


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

62

## Duplicate detection



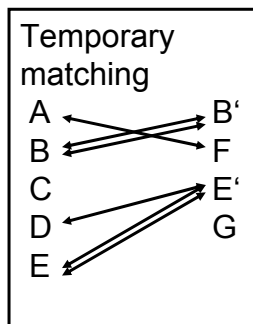
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

63

## Schema Matching

- Given the top K duplicates.
- Find a (global) matching
  - I.e., each attribute has 0 or 1 correspondence.



Formalized by similarity matrix

Average SoftTFIDF

	A	B	C	D	E
B'	0.22	0.92	0.07	0	0
F	0.60	0.60	0.07	0	0
E'	0	0	0	0.58	0.64
G	0	0.07	0	0.07	0.02



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

64



## 2. Schema Matching – Similarity Measure



- Fieldwise Comparison
  - Compare attribute values of record pairs.  
Similarity measure: Edit-distance vs. SoftTFIDF
- Edit distance
  - Minimal number of edit operations (substitute, insert, delete); several variations exist
- SoftTFIDF
  - 'Soft' variation of TFIDF that also considers similar terms

$$tfidf(r, s) = \sum w(r, t) \cdot w(s, t)$$

$$soft - tfidf(r, s) = \sum_{t \in r \cap s} w(r, t) \cdot w(s, t') \cdot sim(t, t')$$

$$t \in CLOSURE(\theta, r, s)$$



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

65

## Graph Matching



- Given:
  - Similarity matrix
  - = weighted bipartite graph
- Find:
  - Maximal weight matching
  - Alternative: Matching with *stable marriage* property
- Outlook
  - Produce mappings and not just correspondences

	A	B	C	D	E
B'	0.22	0.92	0.07	0	0
F	0.60	0.60	0.07	0	0
E'	0	0	0	0.58	0.64
G	0	0.07	0	0.07	0.02

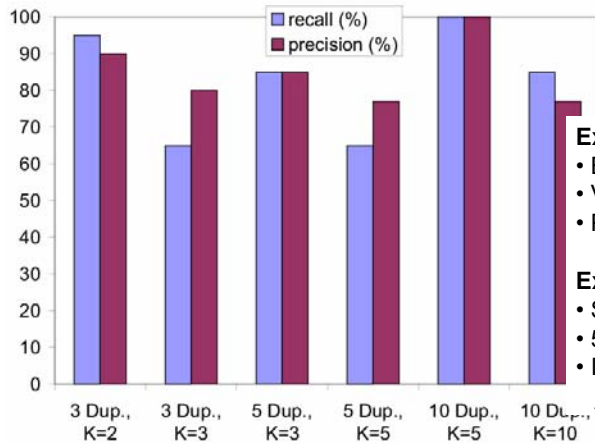


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

66

# Schema Matching



### Experiment 1 (ohne Abb.):

- Berliner Wohnungsmarkt
- Viele Duplikate
- Precision und recall immer 100%

### Experiment 2:

- Synthetische Personendaten
- 5,000 Tupel
- Intensionale Überlappung: jeweils 4 von 8



26 November 2005

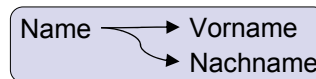
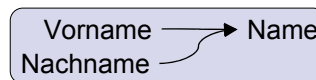
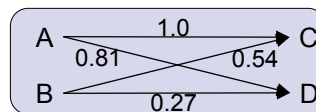
Felix Naumann, Humboldt-Universität zu Berlin

67

# Schema Matching – Extensions



- Globales Matching
  - Matche Tabellen und Schemata, nicht nur Attribute
  - Stable Marriage bzw. Maximum Weighted Matching
- n:1 und 1:n Matches
  - Viele Kombinationsmöglichkeiten
  - Viele Funktionen denkbar
- Matching in komplexen Schemata
  - Ziel: Finde Mapping, nicht Korrespondenzen



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

68

## Coffee break?



- After the break
  - Duplicate Detection
  - Data Fusion
  - Demo



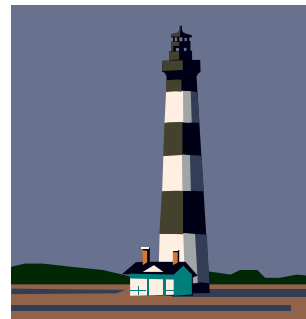
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

69

## Overview

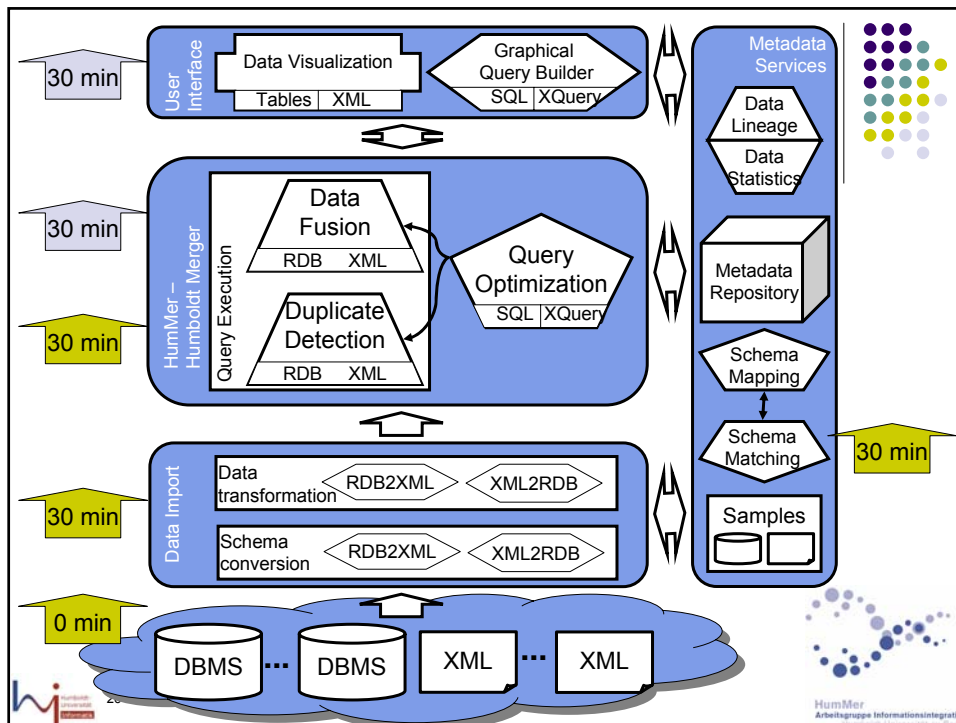
- Introduction and Motivation
- Information Integration in Three Steps
  1. Schema Mapping
    - And Schema Matching
  2. Duplicate Detection
  3. Data Fusion
- HumMer Architecture
- HumMer Demo



26 November 2005

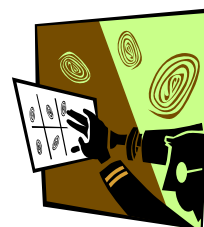
Felix Naumann, Humboldt-Universität zu Berlin

70



## Duplicate detection

- Detection of multiple same or similar representations of a real-world object
- Also
  - Object Identification
  - Data Cleansing
  - Record Linkage
  - Reference Reconciliation
  - ...
- Domain-specific algorithms
  - Address data („mailman similarity“)
  - Genome data
- Sometimes there are IDs
  - ISBN, SSN, URL

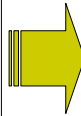


# Duplicate detection

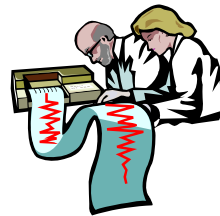
- Definition of an Object
  - Tuple
  - XML element



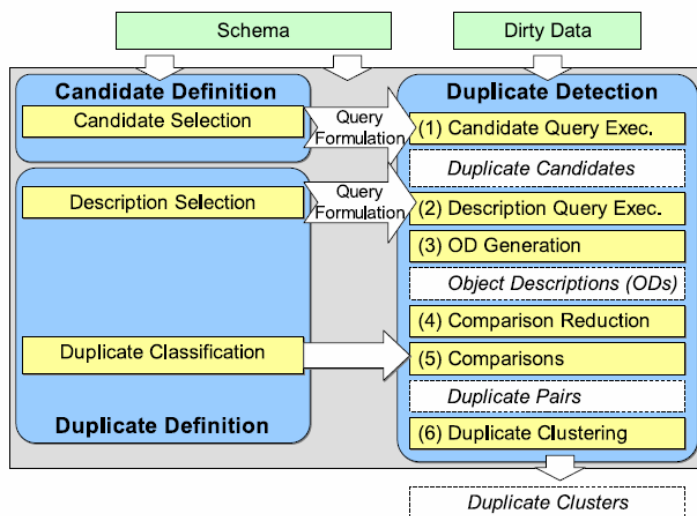
- Similarity measure
  - Domain independent
    - Edit-distance
    - TFIDF-based
    - Tree-distance
  - Domain dependent
    - Rules



- Algorithm
  - Acceleration of similarity comparisons
    - Filter
  - Avoidance of similarity comparisons
    - Partitioning



# The General Framework



# Definition of an XML object



```

- <author>
  <name>Bernd Amann</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
</author>
- <author>
  <name>Sophie Cluet</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Views in a Large Scale XML Repository.</title>
  <year>2001</year>
</publication>
</author>
  
```

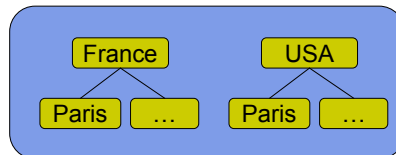
- Compare <author>
  - Include sub-elements (<publication>)?
  - How deep?
- Compare <publication>
  - Only in the scope of the parent element?
  - Schema, or Data?
- In short: What is an object?



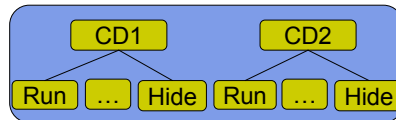
# Top-down, bottom-up, and through the middle



**Top-down [SIGMOD'05]**  
 - Only compare if parents are equal or duplicates  
 - Improved efficiency



**Bottom-up [EDBT'06]**  
 - Similar Children => Duplicate  
 - Improved effectiveness



**Through the middle [ICDE'06]**  
 - Begin with most promising pairs  
 - Best of both worlds

**Further techniques**  
 - Object filter  
 - Edit distance filter  
 - Transitivity



# DELPHI – Data Warehouse Duplicates



ID	country
1	USA
2	United States
3	Unitd States

ID	City	Country_ID
1	New York	1
2	Los Angeles	1
3	Now York	2
4	Los Angeles	2
5	New York	3
6	Los Angels	3

String similarity  
→ 2 ≈ 3

Common children  
→ 1 ≈ 2 ≈ 3



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

77

# Similarity measure



- Token-Set  $TS(e)$ 
  - Data of XML element  $e$  as tokens
- Children-Set  $CS(e)$ 
  - Data of children elements of  $e$  as strings
- $D(e) = TS(e) \cup CS(e)$

$$softidf(s) = |\{e \mid d_{edit}(s, s') \leq t_{edit}, s' \in e\}|$$

Soft document frequency

$$softidf_D(S) = \sum_{s \in S} \log \frac{|D|}{softidf_D(s)}$$

Soft inverse doc. frequency

$$sim(e_1, e_2) = \frac{softidf(D(e_1) \cap D(e_2))}{softidf((D(e_1) \cup D(e_2)) \setminus D(e_1) \cap D(e_2))}$$

Importance of common tokens

Importance of non-common tokens



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

## Algorithm



1. Create data structure
  - XQueries to extract relevant elements
    - Elements of one type at one level
    - plus descendants
2. Acceleration of similarity comparisons
  - Edit-distance filter
3. Avoidance of similarity comparisons
  - Element-similarity-filter
  - Connected components
4. Similarity comparisons
  - Among remaining elements
  - $sim(e_1, e_2) \geq t_{dup}$

Supported by graph-based data structure:

Similarity of tokens are edges between token-nodes

Similarity of elements are edges between element-nodes



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

79

## Edit-distance Filter



- Goal:
  - Accelerate calculation of  $d_{edit}(s_1, s_2) \leq t_{edit}$
  - $d_{edit}(s_1, s_2)$ : Minimum number of edit-operations  $s_1 \leftrightarrow s_2$
- 1. Length filter
  - $|l(s_1) - l(s_2)| \leq d_{edit}(s_1, s_2)$
- 2. Triangular inequation property
  - $|d_{edit}(s_1, s_2) + d_{edit}(s_2, s_3)| \geq d_{edit}(s_1, s_3)$
  - $|d_{edit}(s_1, s_2) - d_{edit}(s_2, s_3)| \leq d_{edit}(s_1, s_3)$
- 3. Bag filter
  - $d_{bag}(s_1, s_2) = \max\{|S_1 - S_2|, |S_2 - S_1|\} \leq d_{edit}(s_1, s_2)$



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

80



## Element-similarity Filter



- Goal
  - Avoid calculation of  $sim(e_1, e_2) \geq t_{dup}$
- Element-similarity filter

$$sim(e_1, e_2) = \frac{softidf(D(e_1) \cap D(e_2))}{softidf((D(e_1) \cup D(e_2)) \setminus D(e_1) \cap D(e_2))}$$
$$\leq \frac{softidf(D(e_1) \cap (Dok - D(e_1)))}{softidf(D(e_1) \setminus (D(e_1) \cap (Dok - D(e_1))))}$$

- Intuition
  - Similarity to an element  $\leq$  similarity to all elements
- Question
  - What do I save?



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

81

## Similarity comparisons



1.  $sim(e_1, e_2) \leq t_{dup}$  for remaining element pairs
2. Transitive closure of all duplicates
  - Duplicate-cluster
3. Repeat for next level of hierarchy
  - But: Only within clusters



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

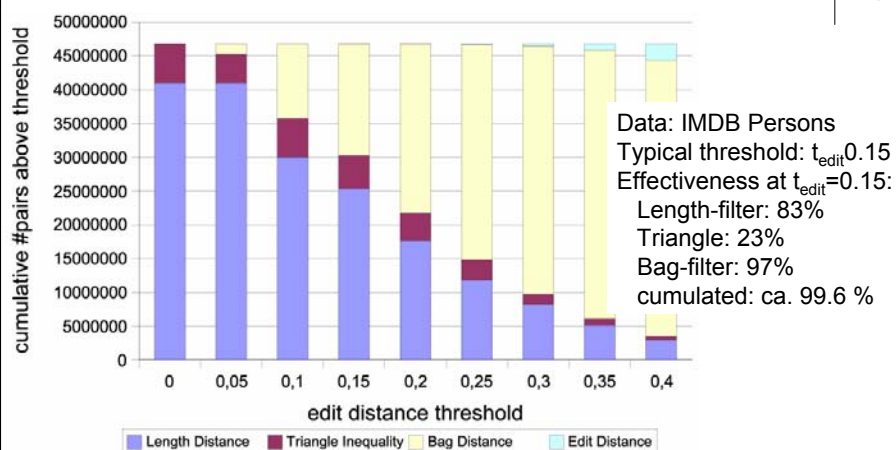
82

# Experiments

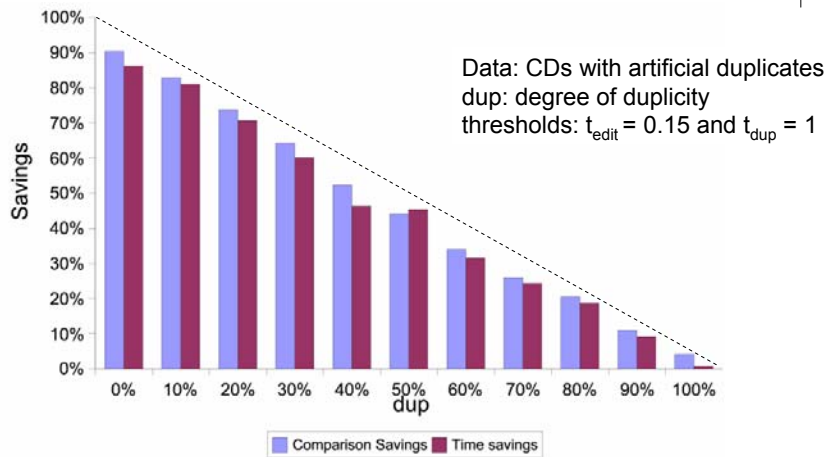
- Data
  - IMDB: 54.000 persons
  - freedb.org: 10.000 CDs
  - mondialDB: 260 countries
  - Synthetic duplicates
- Effectiveness of edit-distance filter
- Effectiveness of element-similarity filter
- Precision & recall of duplicate detection



# Edit-distance Filter



## Element-similarity Filter

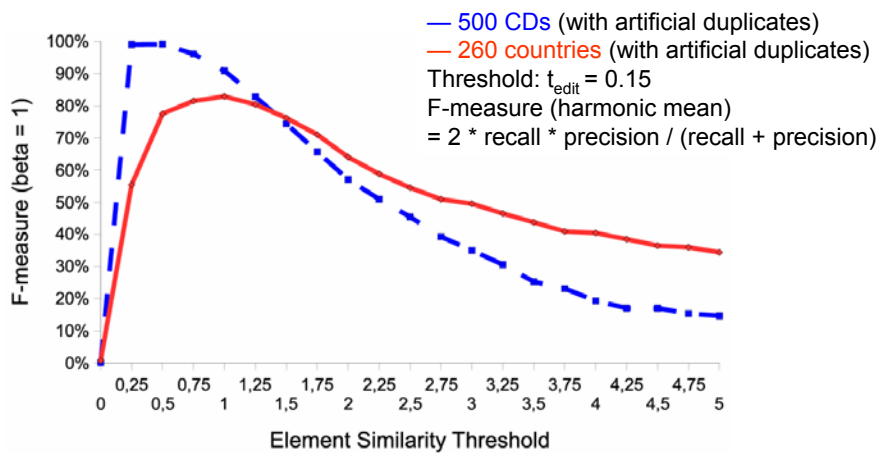


26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

85

## Precision & Recall



26 November 2005

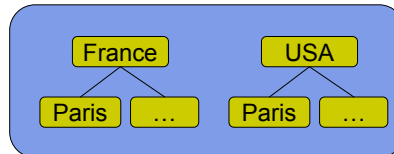
Felix Naumann, Humboldt-Universität zu Berlin

86

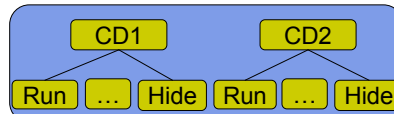
# Top-down, bottom-up, and through the middle



Top-down [SIGMOD'05]  
- Only compare if parents are equal or duplicates  
- Improved efficiency



Bottom-up [EDBT'06]  
- Similar Children  $\Rightarrow$  Duplicate  
- Improved effectiveness



Through the middle [ICDE'06]  
- Begin with most promising pairs  
- Best of both worlds

Further techniques  
- Object filter  
- Edit distance filter  
- Transitivity



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

87

# The Bottom-Up Approach



- Idea: Apapt Sorted-Neighborhood algorithm
  - Construct sorted neighborhood for leaves.
  - Construct sorted neighborhood for intermediate nodes with many similar children
- Implemented on top of relational DBMS
- Not as efficient as top down.
- But more effective
  - More comparisons
  - But not too many (windowing!)



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

88

## A Through-the-Middle Approach



- Intuition:
  - Any detected duplicate pair affects neighbors (parents and children)
  - Rank all pairs of (comparable) objects by their degree of influence on neighbors
  - Compare objects with low influence first
    - Either: Reduce number of recomparisons
    - Or: Improve effectiveness of recomparisons are not allowed.



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

89

## Outlook – Duplicate detection



- Definition of Objects
  - Declarative specification
  - Graphical tool
- Similarity measures
  - Use schema information
- Algorithm
  - Scalability (RAM, CPU)
  - XML DBMS (Tamino)
- Experiments
  - Movie data from five sources
  - Graphical tool for manual duplicate detection



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

90

# Relaxing the object definition



```

- <author>
  <name>Bernd Amann</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <author>
  <name>Sophie Cluet</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Views in a Large Scale XML Repository.</title>
  <year>2001</year>
</publication>
</author>

- <inproceedings key="conf/vldb/AbiteboulAAACHHMMSTV99">
  <author>Serge Abiteboul</author>
  <author>Vincent Aguilera</author>
  <author>Sébastien Ailleret</author>
  <author>Bernd Amann</author>
  <author>Sophie Cluet</author>
  <author>Brendan Hills</author>
  <author>François Martignoli</author>
  <title>XML Repository and Active Views Demonstration.</title>
  <pages>199-216</pages>
  <year>1999</year>
  <booktitle>VLDB</booktitle>
  <url>db/conf/vldb/vldb99.html#AbiteboulAAACHHMMSTV99.html</url>
  <crossref>conf/vldb/99</crossref>
  <ee>db/conf/vldb/AbiteboulAAACHHMMSTV99.html</ee>
  <cdrom>VLDB99/P73.pdf</cdrom>
  <cite>conf/edbt/SantosAD94</cite>
  <cite>www/org/w3/dom</cite>
</inproceedings>
  
```

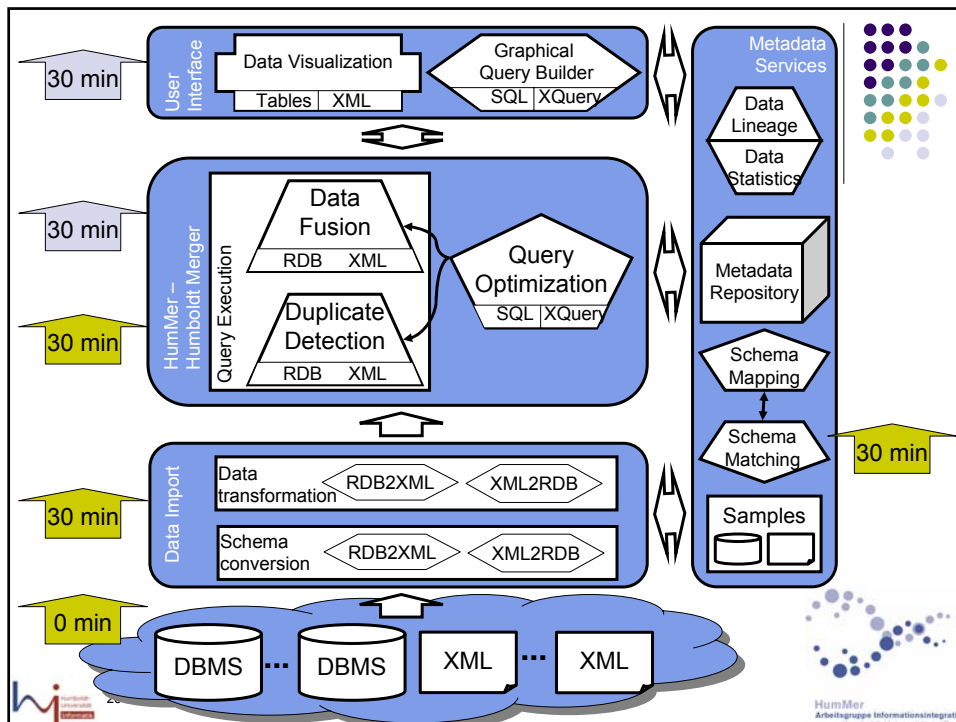
Remember Schema Matching?



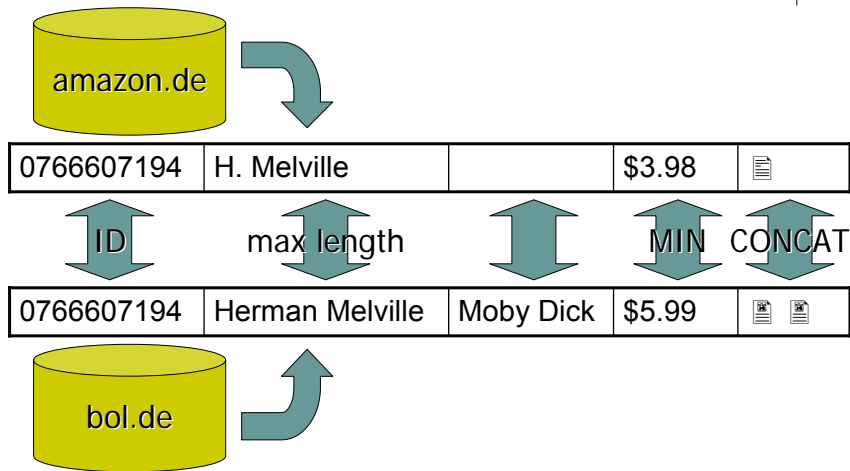
26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

91



# Data Fusion



# Visualization of Integrated Data



- Why Provenance & Where Provenance
- Conflicts en detail and in overview (zoom out)

The screenshot shows the 'Fuse By GUI - Example' interface. It features a table with columns for Title, Year, Director, and Country. The table is color-coded to show different types of conflicts: Duplicate (green), Contradiction (red), Uncertainty (orange), and Unique (blue). Navigation controls (back, next) and a status bar are also visible.

	Title	Year	Director	Country
		CR_MIN	CR_COALESCE	CR_COALESCE
0	Ying Xiong	2002		
1	Serenity	2005	Joss Whedon	USA
2	Metropolis	1927		Germany
3	Donnie Darko	2001	Richard Kelly	
4	Citizen Kane	1941	Orson Welles	USA

Rows: 0:4 5/5

Navigation: go to 0, back, next 100

Statusbar: Duplicate, Contradiction, Uncertainty, Unique



# Thank You for Listening

naumann@informatik.hu-berlin.de

## Information Integration in Three Steps

1. Schema Mapping
  - 1.1 Schema Matching
2. Duplicate Detection
3. Data Fusion



## References

- [SIGMOD-Record'04] BioFast: Challenges in Exploring Linked Life Science Sources. J. Bleiholder, Z. Lacroix, H. Murthy, F. Naumann, L. Raschid, and M-E. Vidal: SIGMOD Record 33(2), June 2004
- [ICDE'05] Schema Matching Using Duplicates. Alexander Bilke, Felix Naumann: Proceedings of the International Conference on Data Engineering (ICDE 05) Tokyo, Japan.
- [WebDB'05] A Data Model and Query Language to Explore Enhanced Links and Paths in Life Science Sources. George Mihaila, Felix Naumann, Louiqa Raschid, Maria-Esther Vidal. In SIGMOD WebDB 2005, Baltimore, MD.
- [ICDE'02] Attribute Classification Using Feature Analysis. Felix Naumann, Ching-Tien Ho, Xuqing Tian, Laura M. Haas, Nimrod Megiddo: ICDE 2002: 271
- [VLDB'05] Automatic Data Fusion with HumMer (demo). Alexander Bilke, Jens Bleiholder, Christoph Böhm, Karsten Draba, Felix Naumann, Melanie Weis. Proceedings of the Int. Conf. on Very Large Databases (VLDB) 2005, Troindheim, Norway
- [SIGMOD'05] Dogmatix Tracks down Duplicates in XML. Melanie Weis and Felix Naumann, (SIGMOD 2005), Baltimore, MD
- [ADBIS'05] Declarative Data Fusion - Syntax, Semantics, and Implementation. Jens Bleiholder and Felix Naumann, (ADBIS 2005), Tallin, Estonia.
- [CIDR'05] (Almost) Hands-Off Information Integration for the Life Sciences. Ulf Leser, Felix Naumann. Proceedings of the Conference in Innovative Database Research (CIDR 2005), Asilomar, CA.
- [BTW'05] Self-Extending Peer Data Management. Ralf Heese, Sven Herschel, Felix Naumann, Armin Roth (BTW 2005), Karlsruhe, Germany.
- [FHP+02] Ron Fagin, Mauricio Hernandez, Lucian Popa, Renee Miller, and Yannis Velegrakis, Translating Web Data, VLDB 2002, Hong Kong, China.
- [EDBT'06] XML Duplicate Detection Using Sorted Neighborhoods. Sven Puhlmann, Melanie Weis, Felix Naumann: *Proceedings of the International Conference on Extending Database Technology (EDBT) 2006*, Munich, Germany.
- [ICDE'06] Detecting Duplicates in Complex XML Data. Melanie Weis, Felix Naumann. *Proceedings of the International Conference on Data Engineering (ICDE) 2006*, Atlanta, GA. poster



26 November 2005

Felix Naumann, Humboldt-Universität zu Berlin

96