



**Hasso
Plattner
Institut**

IT Systems Engineering | Universität Potsdam

Data Consolidation in Three Steps

Milano, May 9, 2008

Felix Naumann

The HPI – Hasso Plattner Institut

2

- Founded in 1998 as a Public Private Partnership
- Hasso Plattner, co-founder of SAP, endowed over 200 Mio. Euro.
- Adjoined with the University of Potsdam
 - Capital of Brandenburg, bordering Berlin
- 400 students – Bachelor, Master, and PhD



The Information Systems Department

3

project **ViQTOR**



Paul Führung



Patricia Hobro

DQ Assessment



Prof. Felix Naumann

Data Fusion



Jens Bleiholder



Karsten Draba

project **fusem**

Information Integration

Information Quality

project **HumMer**

Duplicate Detection



Data Cleaning

Melanie Weis & Sascha Szott

Peer Data Management Systems



Armin Roth

project **System P**

Matching

Service-Oriented Systems

Data Integration for Life Science Data Sources

project **Aladin**



Alexander Albrecht

project **XClean**

ETL Process Management

Ontologies



Mohammed AbuJarour



Frank Kaufer

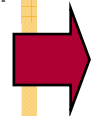


Jana Bauckmann

Data Profiling for Schema Management

Overview

4



- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
- Summary



Data Fusion in Three Steps: Resolving Inconsistencies at Schema-, Tuple-, and Value-level

Felix Naumann¹ Alexander Bilke² Jens Bleiholder¹ Melanie Weis¹
¹ Humboldt-Universität zu Berlin
Unter den Linden 6, 10099 Berlin, Germany
{naumann|bleiho|mweis}@informatik.hu-berlin.de
² Technische Universität Berlin
Strasse des 17. Juni 135, 10623 Berlin, Germany
bilke@cs.tu-berlin.de

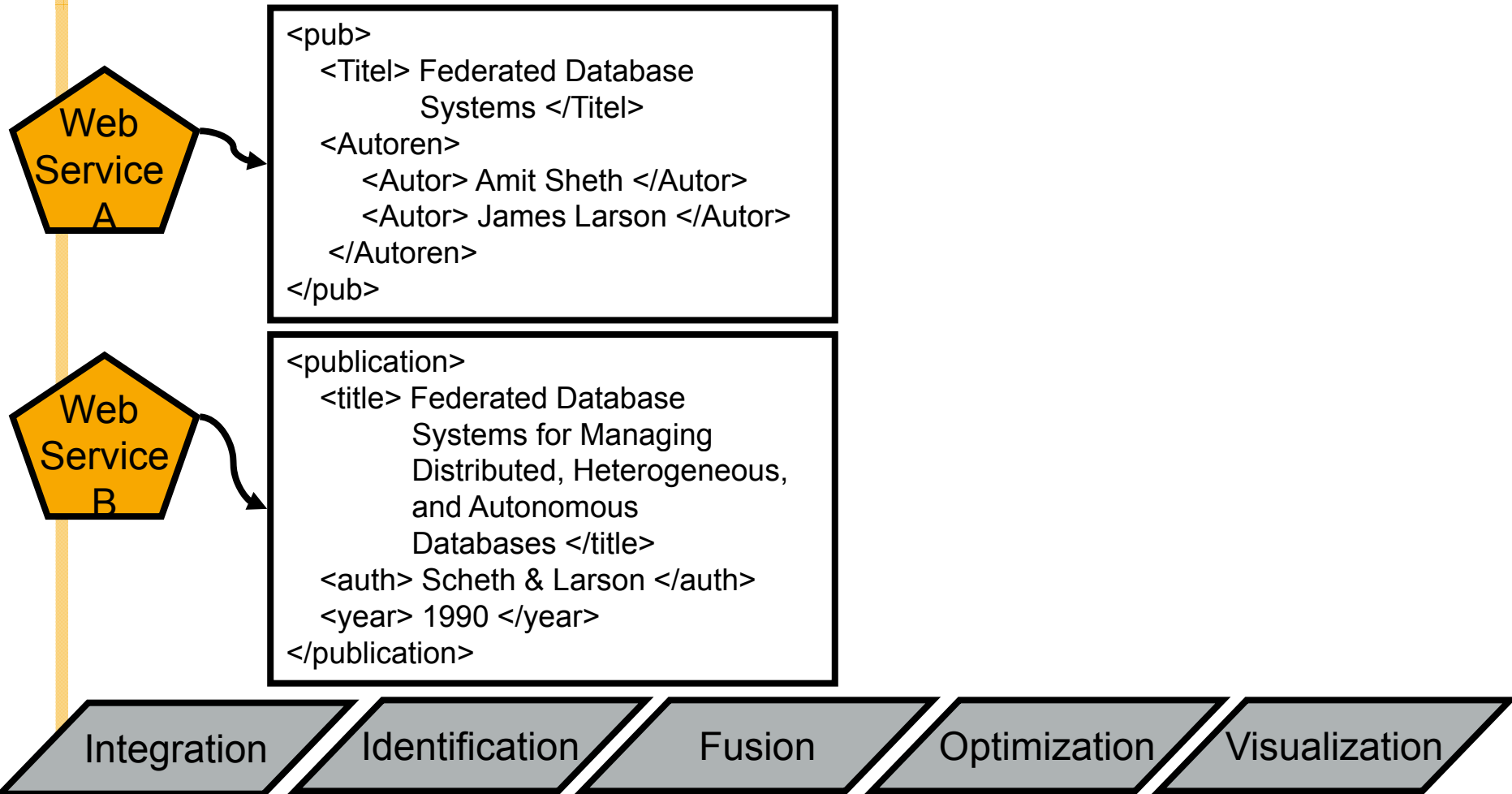
Data Fusion

JENS BLEIHOLDER and FELIX NAUMANN
Hasso-Plattner-Institut, Potsdam, Germany



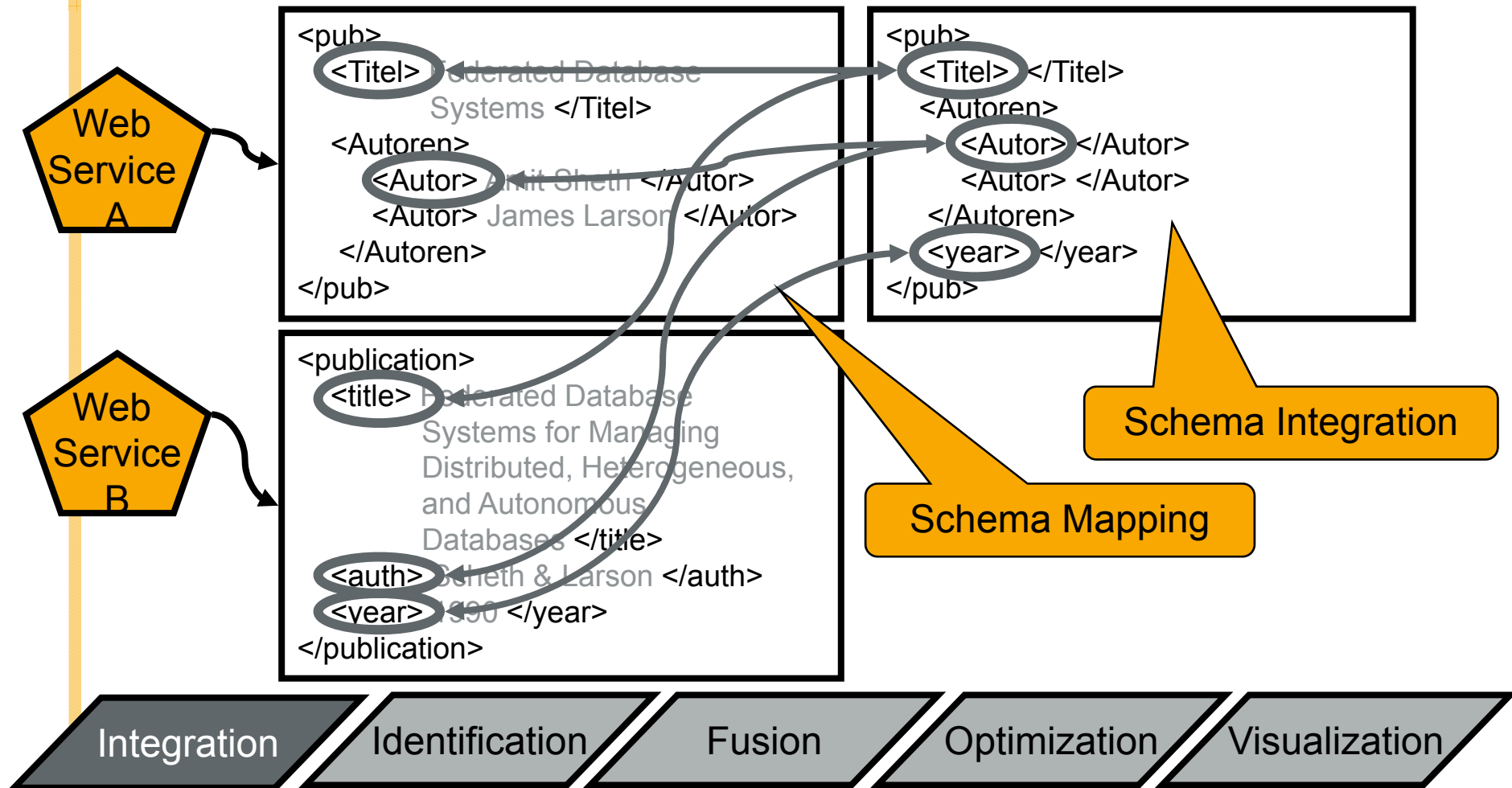
Information Integration

5



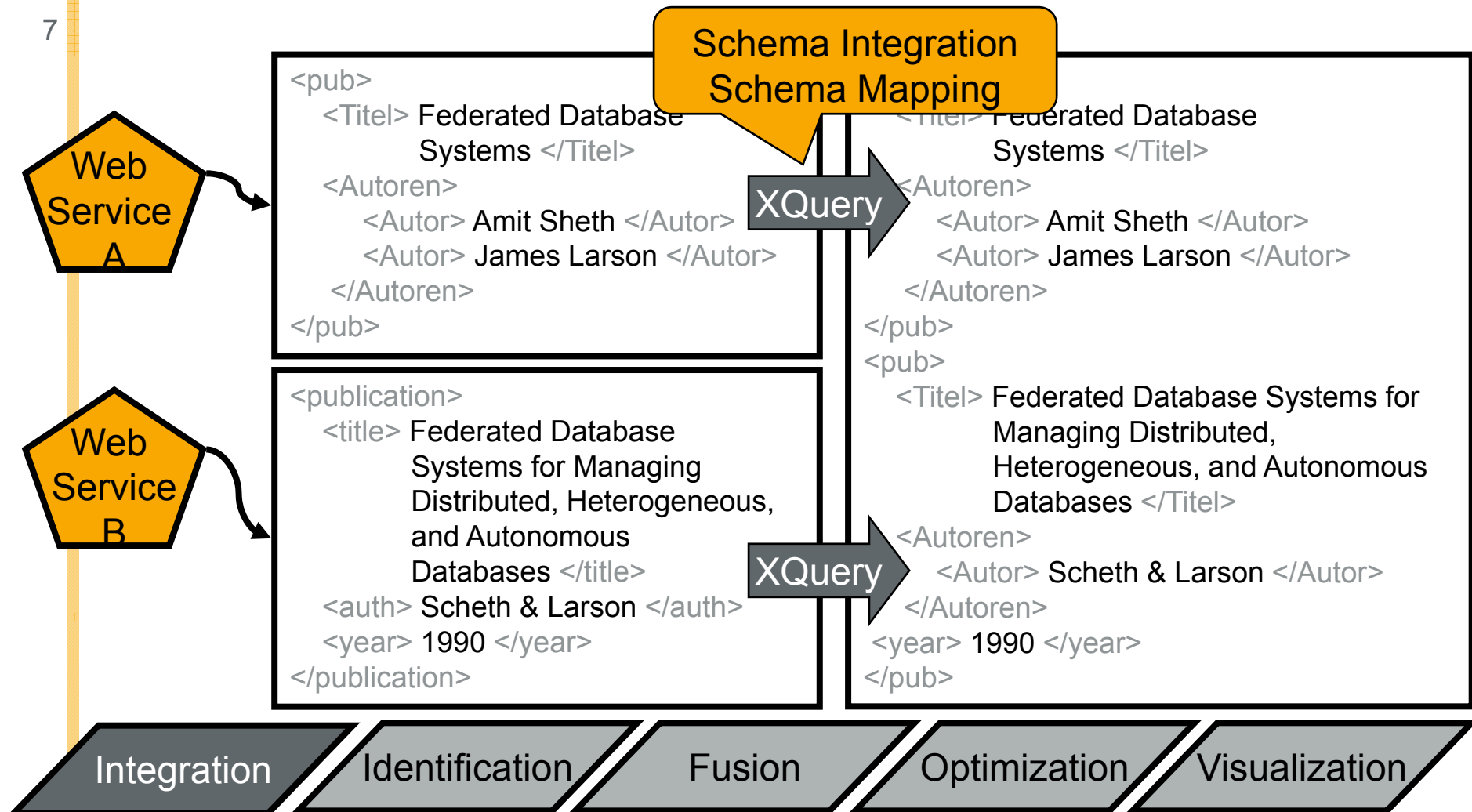
Information Integration

6



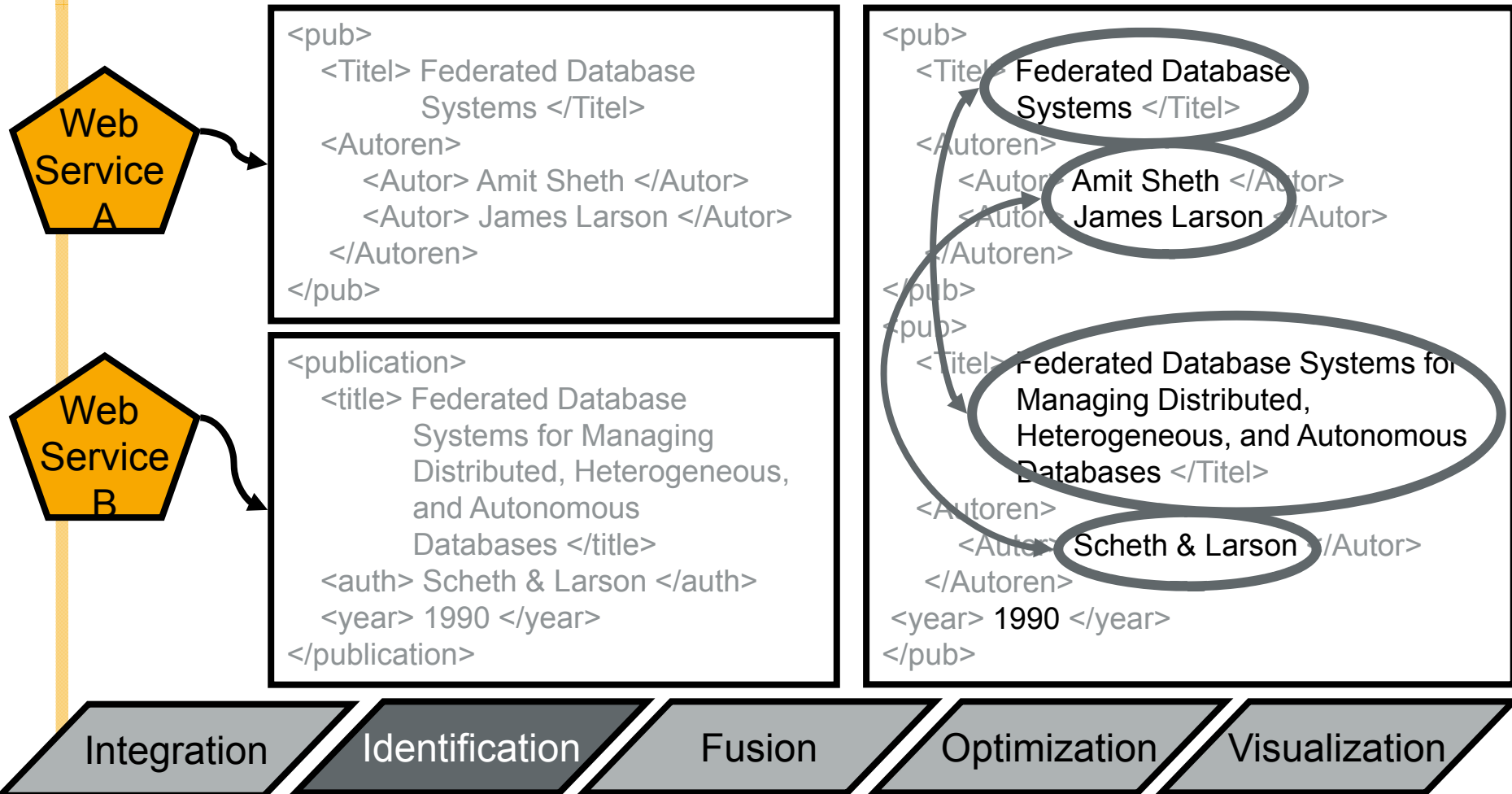
Information Integration

7



Information Integration

8



Information Integration

9

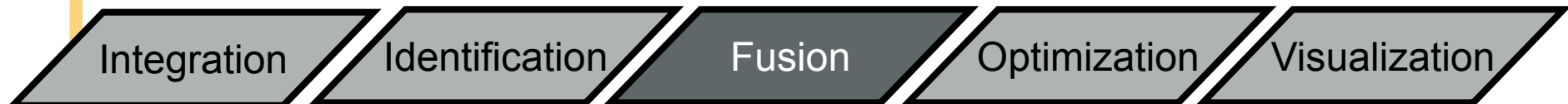


```
<pub>
  <Titel> Federated Database
    Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```



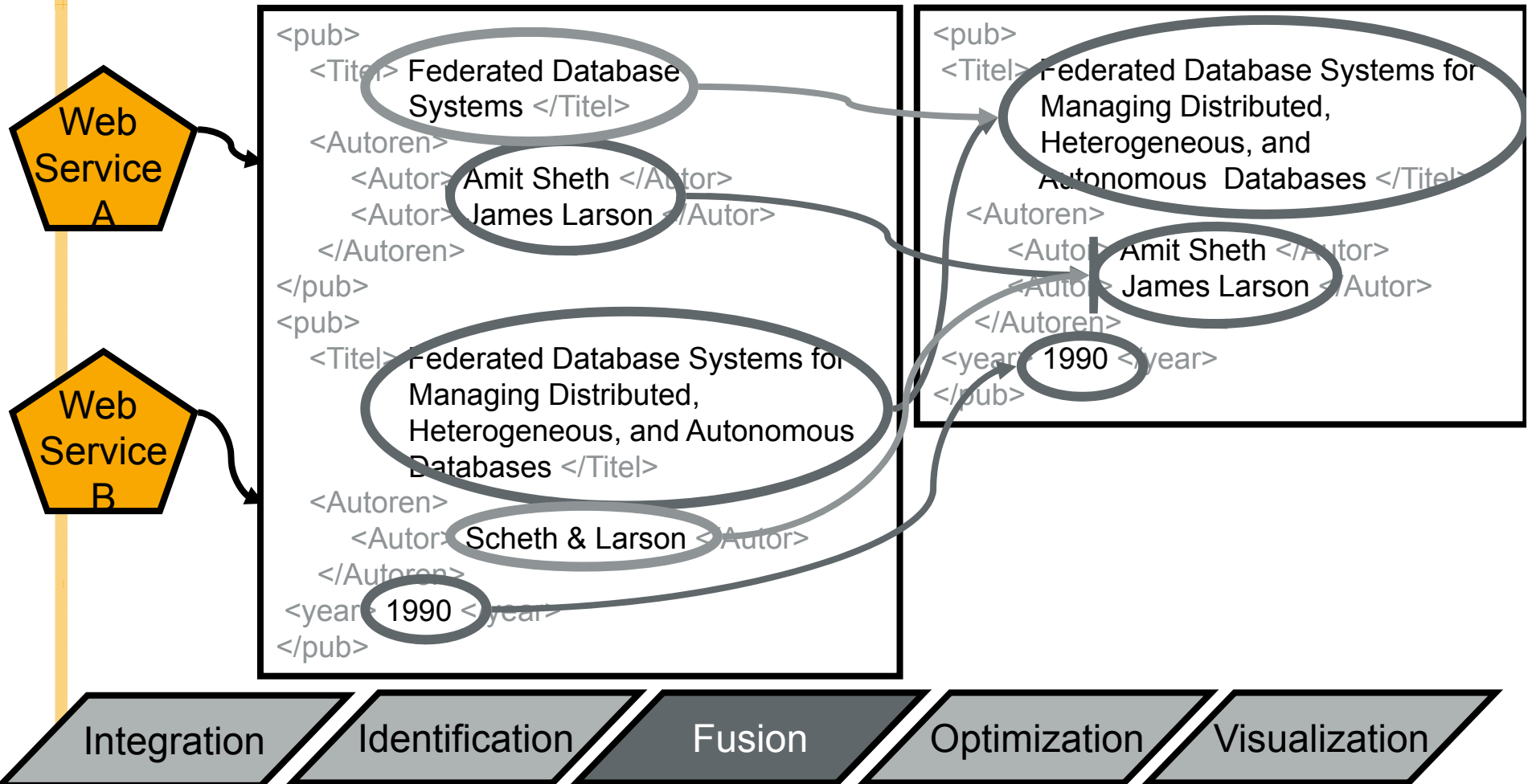
```
<publication>
  <title> Federated Database
    Systems for Managing
    Distributed, Heterogeneous,
    and Autonomous
    Databases </title>
  <auth> Scheth & Larson </auth>
  <year> 1990 </year>
</publication>
```

```
<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>
```



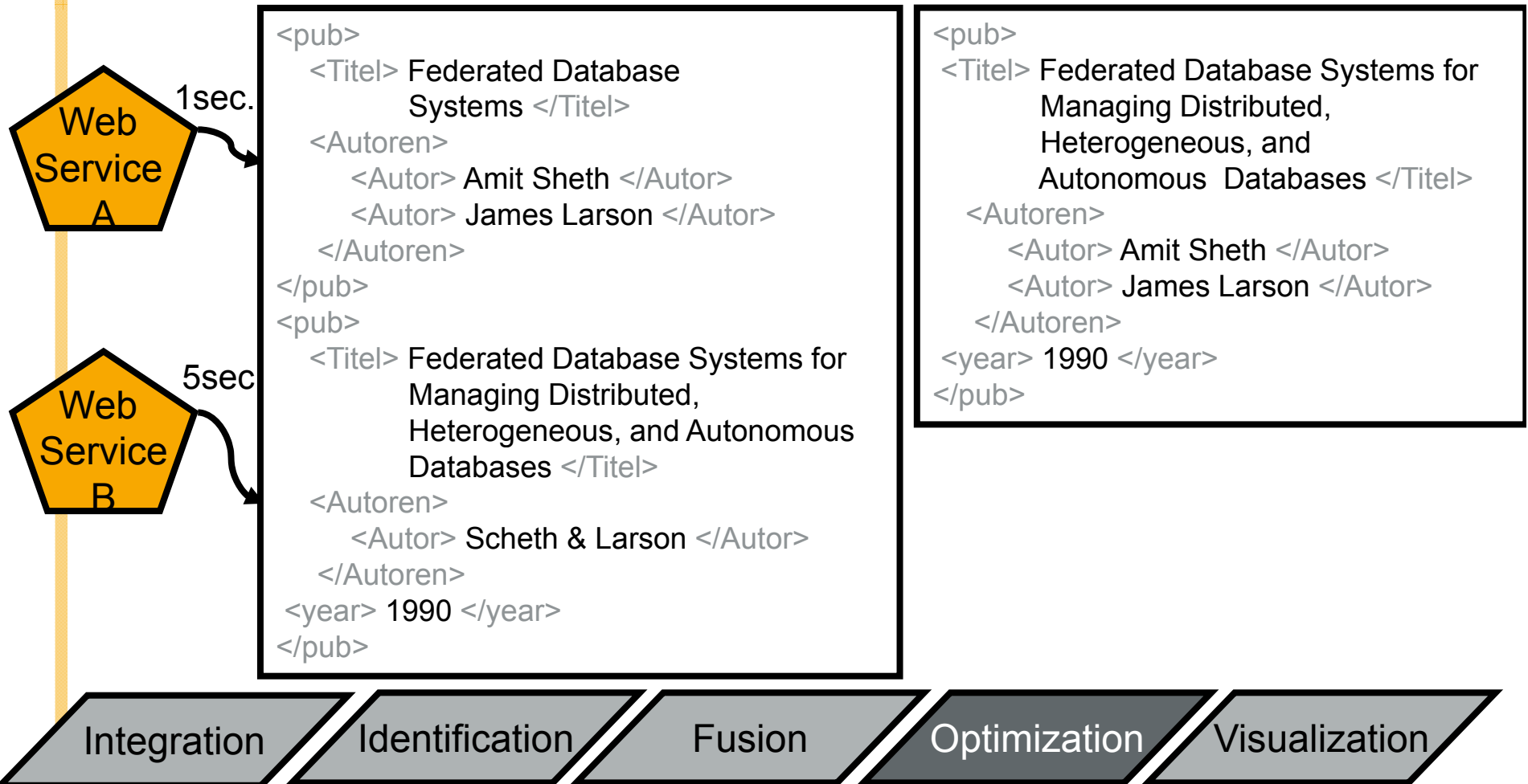
Information Integration

10



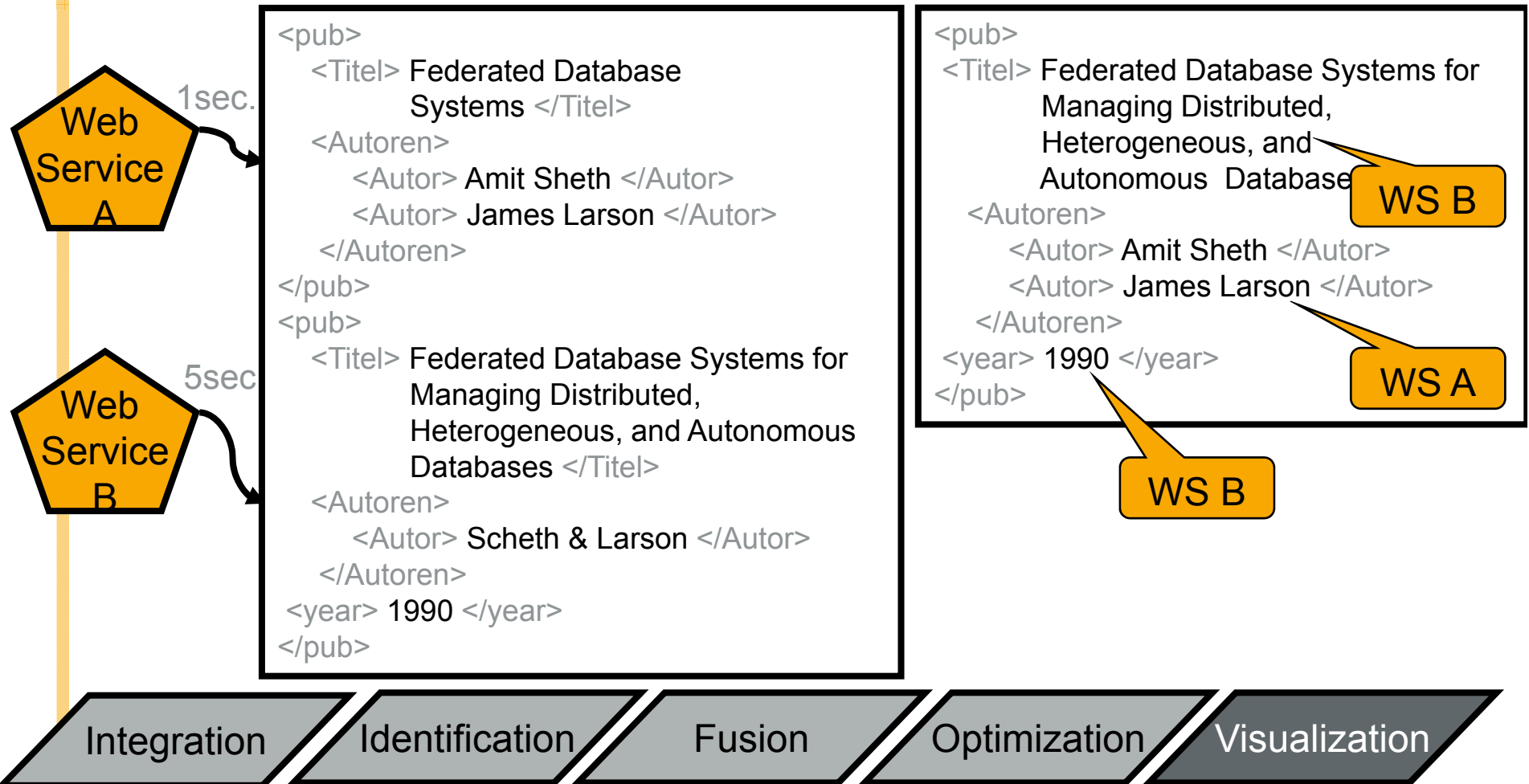
Information Integration

11



Information Integration

12



Overview

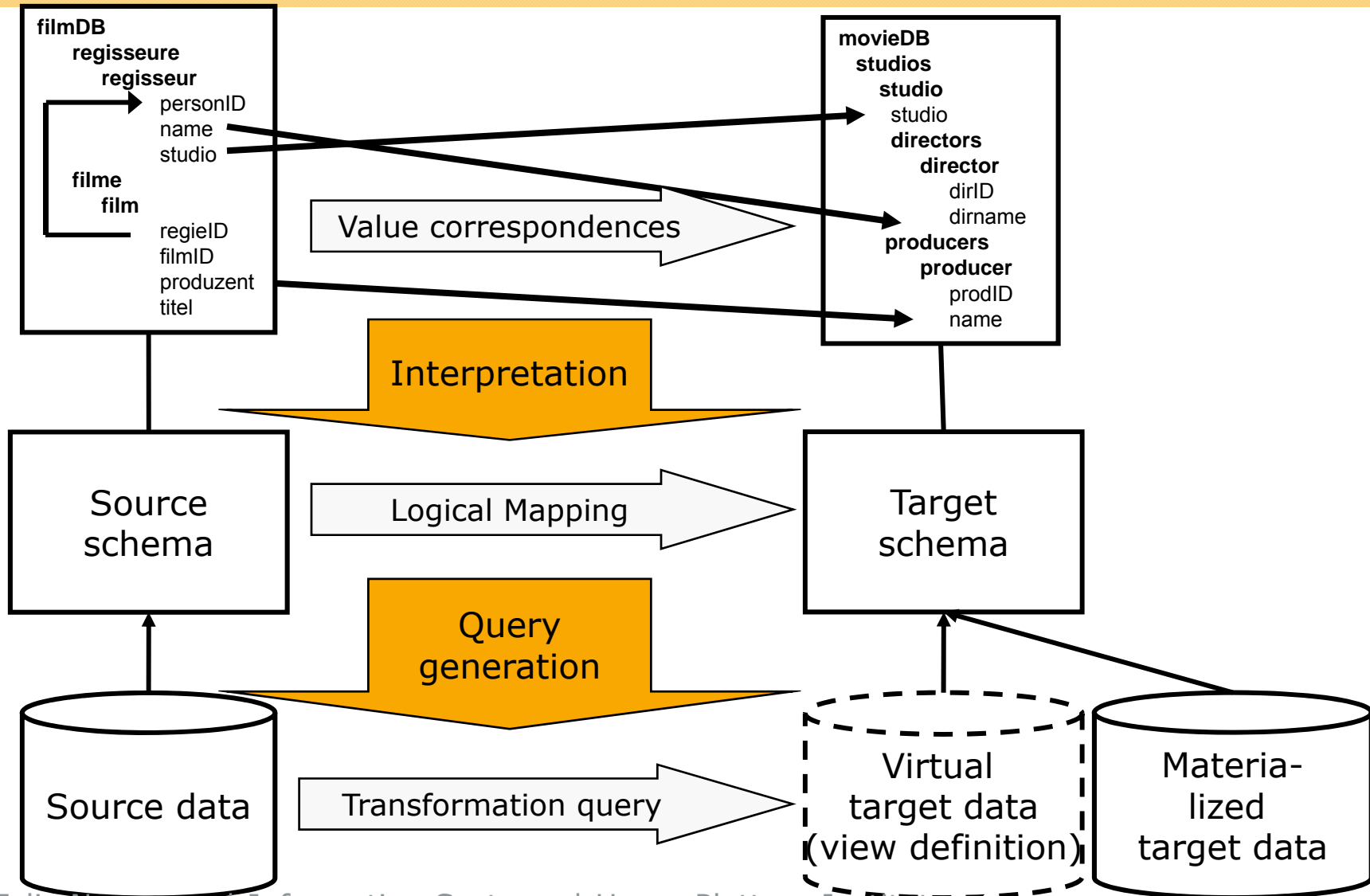
13

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
- Summary



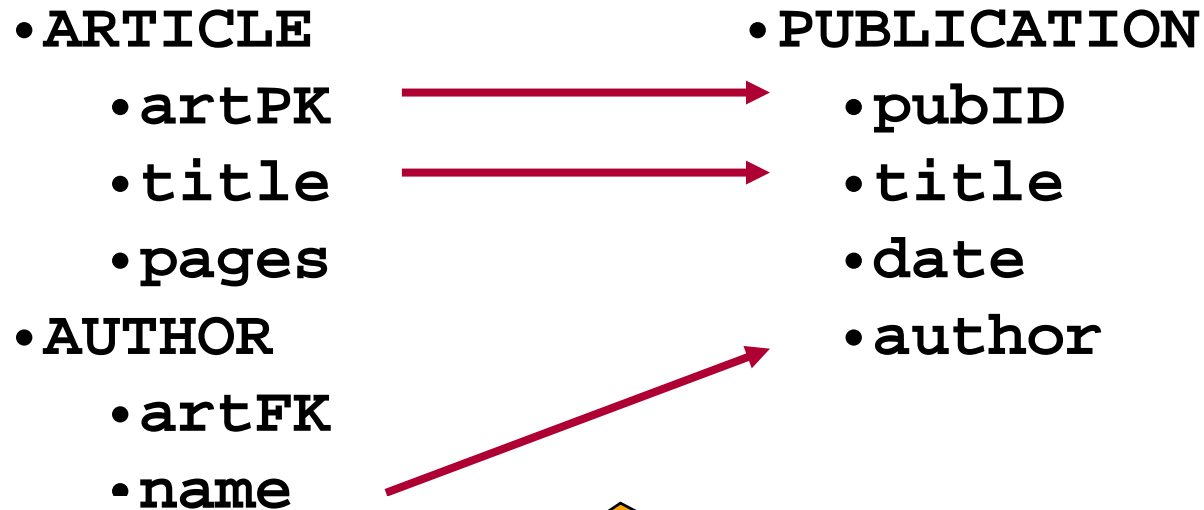
Schema Mapping in Context

14



Schema Mapping Example

15

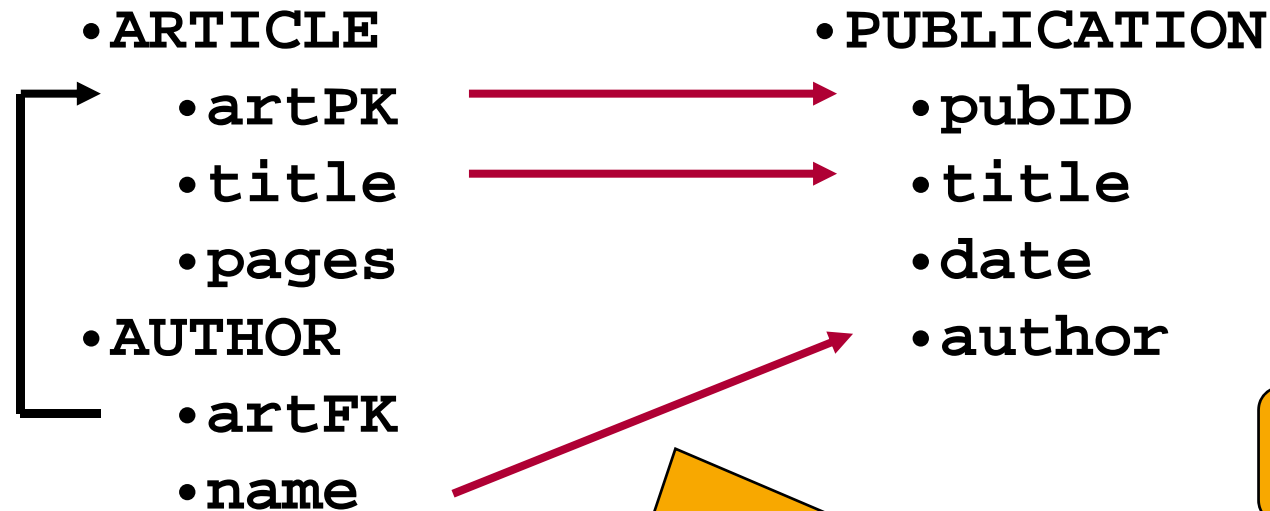


```

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
  
```


Schematic heterogeneity – solutions

16



Further interpretations?

```

SELECT      artPK AS pubID
            title AS title
            null AS date
            name AS author
FROM        ARTICLE, AUTHOR
WHERE      ARTICLE.artPK = AUTHOR.artFK
    
```

Schema Matching – Motivation

17

Schemata are

- large
- complex
- foreign
- confusing
- different language
- cryptic

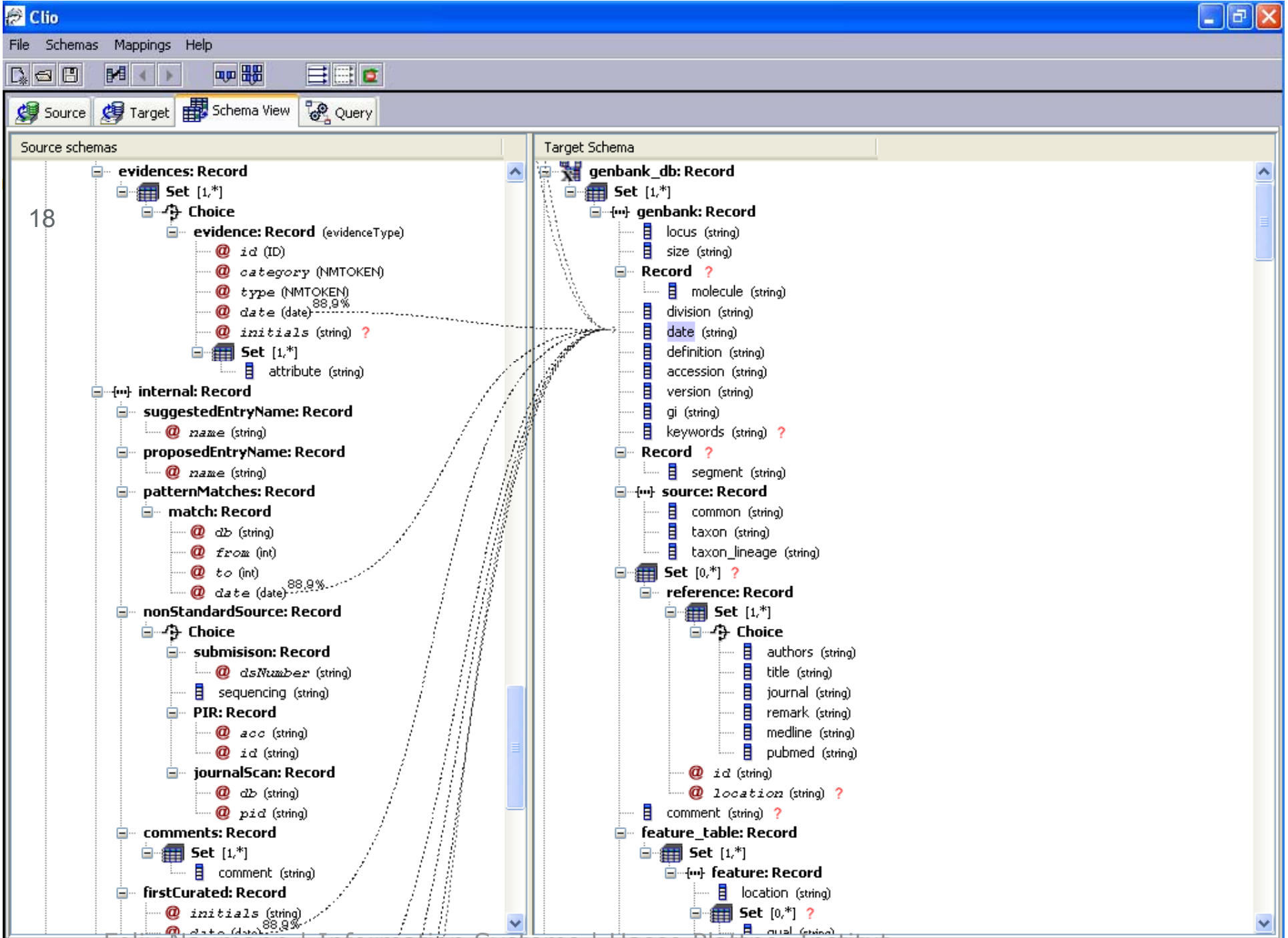
> 100 tables, many attributes

Deep Nesting
Foreign keys
XML Schema

Unknown synonyms

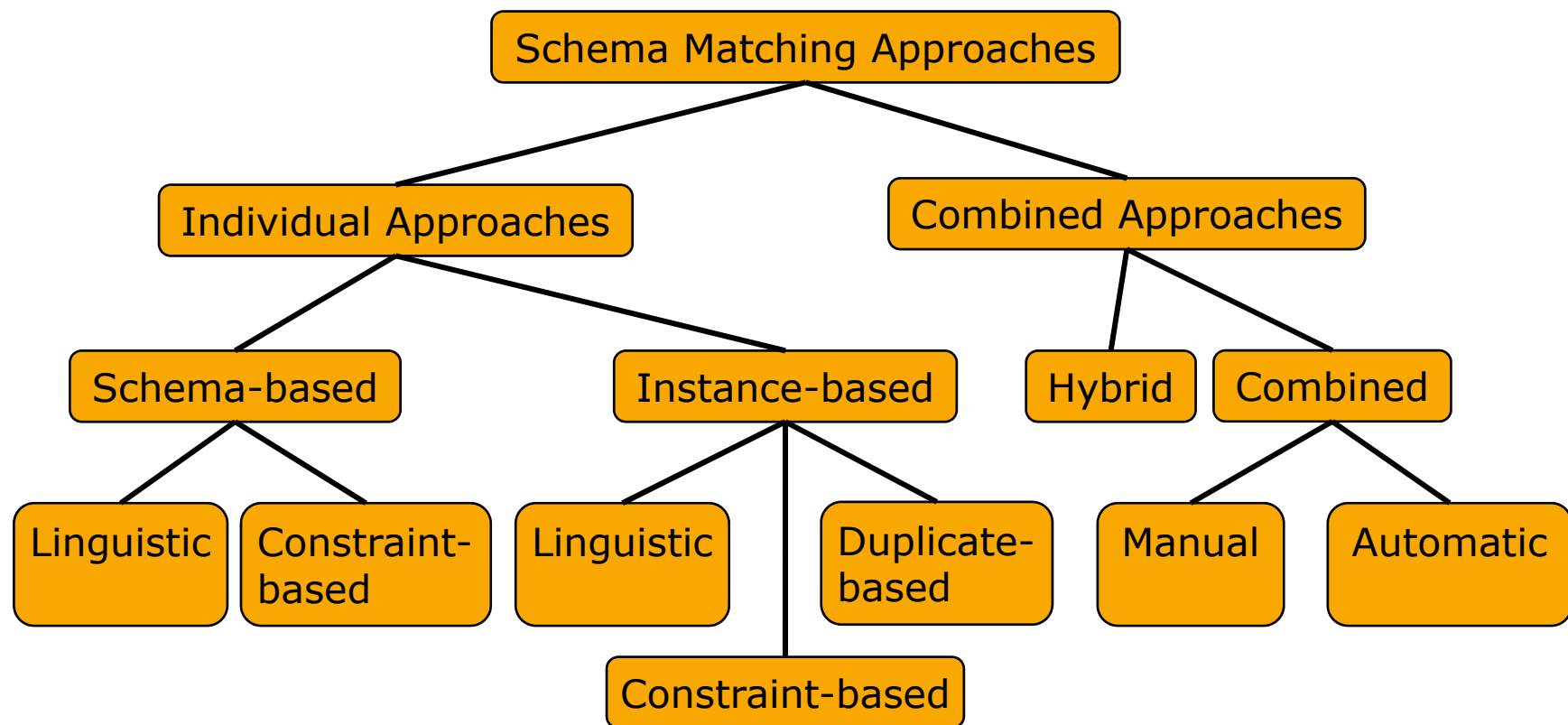
Unknown homonyms

$|\text{attribute name}| \leq 8$
 $|\text{table name}| \leq 8$



Schema Matching Classification [RB01]

19



Instance-based Schema Matching

20

Instance-based Schema Matching:

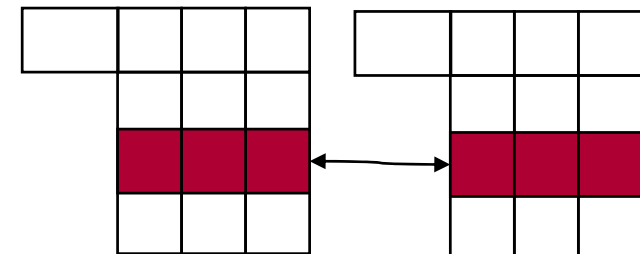
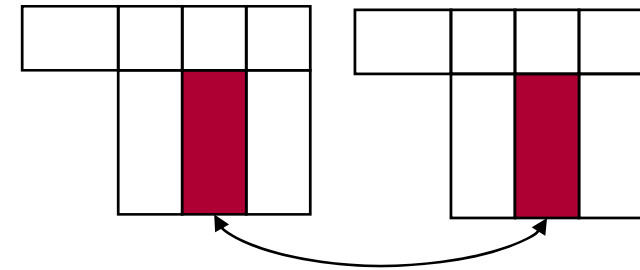
- Correspondences based on similar data values or their properties

Conventional solution: Vertical

- Comparison of columns
- = Attribute classification

Our solution: Horizontal

- Comparison of rows
- = Duplicate detection (despite missing attribute correspondences)

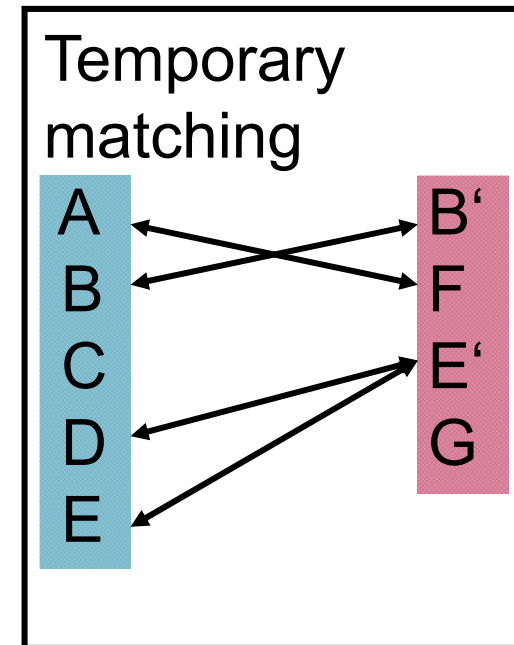
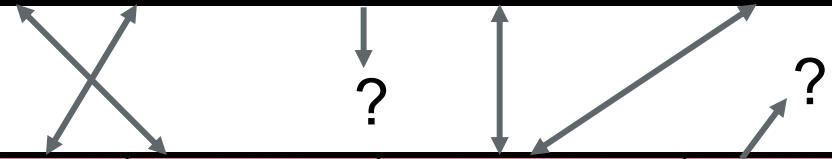


DUMAS Matcher

21

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

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

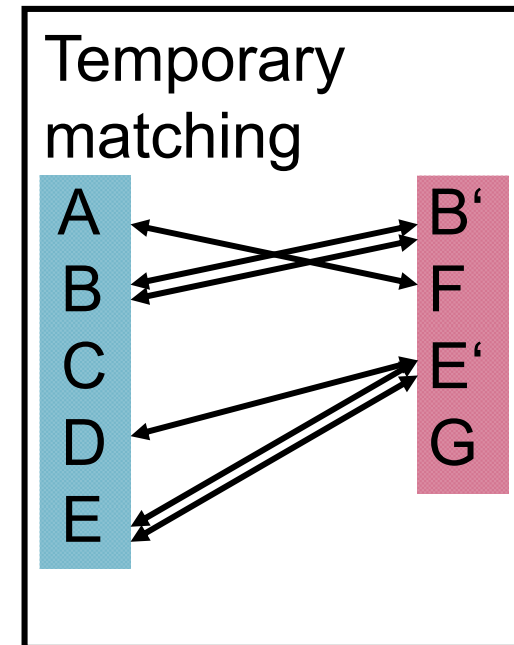
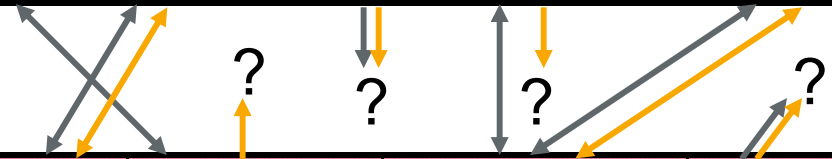


DUMAS Matcher

22

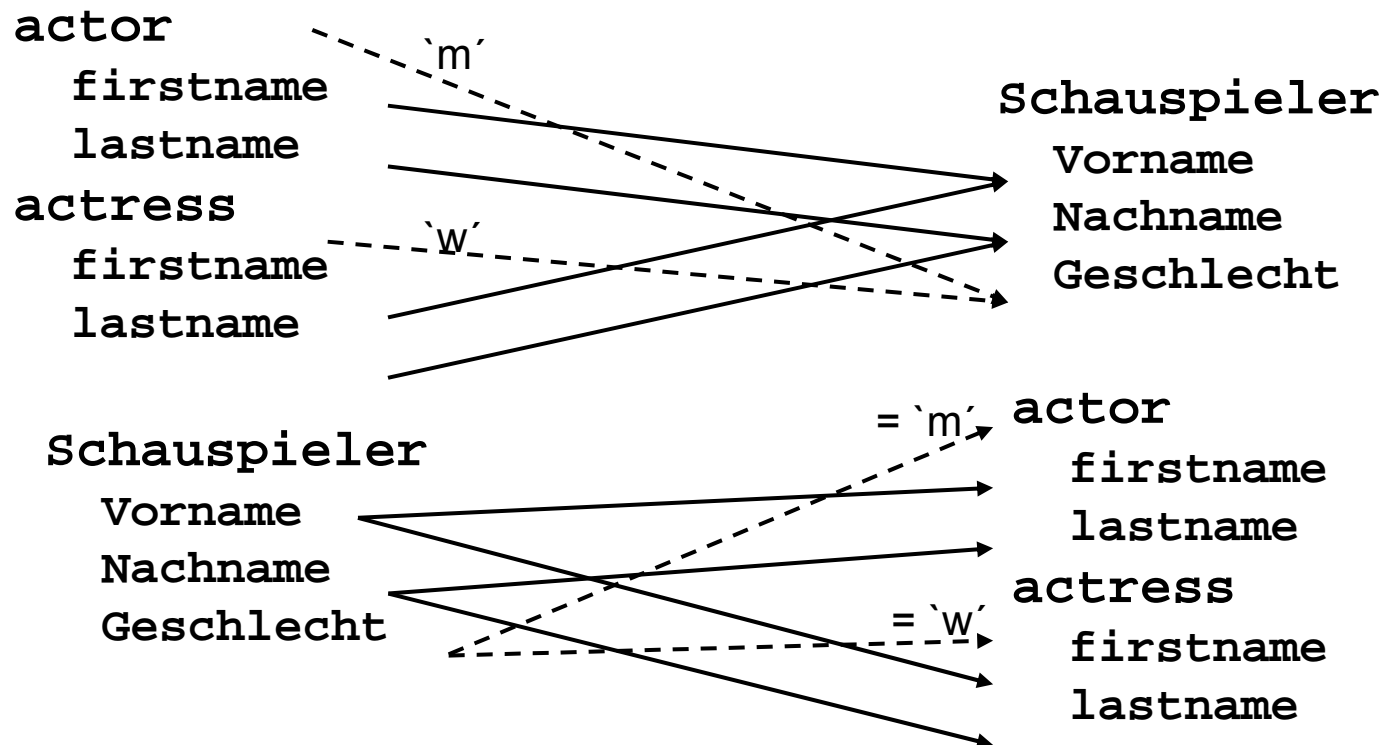
A	B	C	D	E
Max	Michel	m	601- 4839204	601- 4839204
Sam	Adams	m	541- 8127100	541- 8121164

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



Schema Matching – High-level Matching

23



Schema Matching – Extensions

24

n:1 und 1:n matches

- Many combinations
- Many functions
- Parsing
- iMap

Matching in complex schemata

- Find mapping, not only correspondences
- Unions and joins

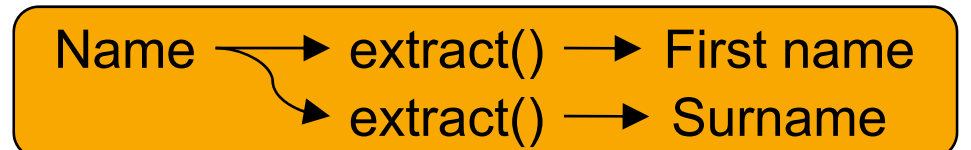
Global Matching

- Match Table and Schema, not just Attributes

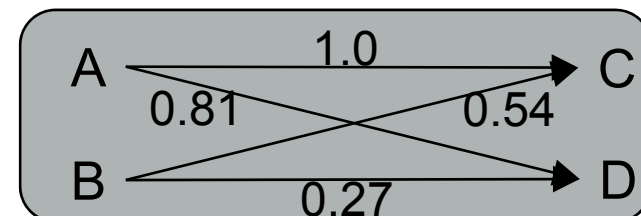
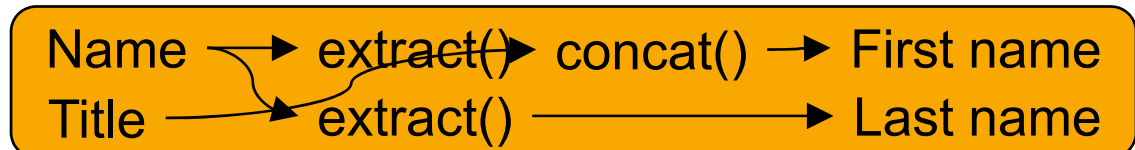
n:1 Matching



1:n Matching



m:n matching



Overview

25

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
- Summary



Duplicate Detection

26

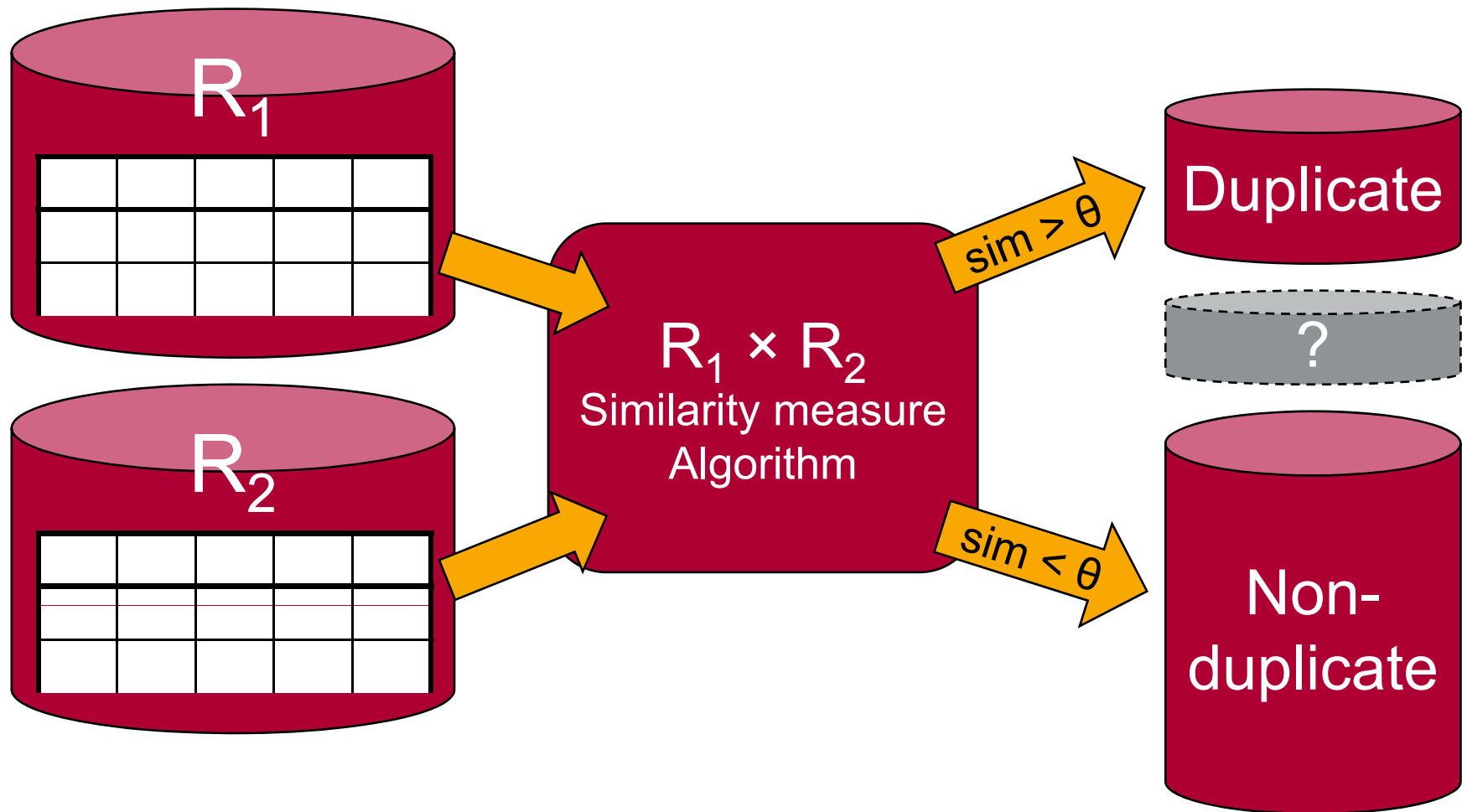
Duplicate detection is the discovery of multiple representations of the same real-world object.

- Problem 1: Representations are not identical.
 - *Fuzzy duplicates*
- Solution: Similarity measures
 - Value- and record-comparisons
 - Domain-dependent or domain-independent

- Problem 2: Data sets are large.
 - Quadratic complexity: Comparison of every pair of records.
- Solution: Algorithms
 - E.g., avoid comparisons by partitioning.

Duplicate Detection

27



Motivation

28

- Possible effects

- Example: Portfolio Management Offers
- Credit maximum not detected
- Too low inventory levels
- No quantity discount for multiple orders
- Total revenue of preferred customers unknown
- Multiple mailings of same catalog to same household

Customer	Revenue
BMW	20.000
BaMoWe	5.000.000
Bayerische Motorenwerke	300.000
...	...

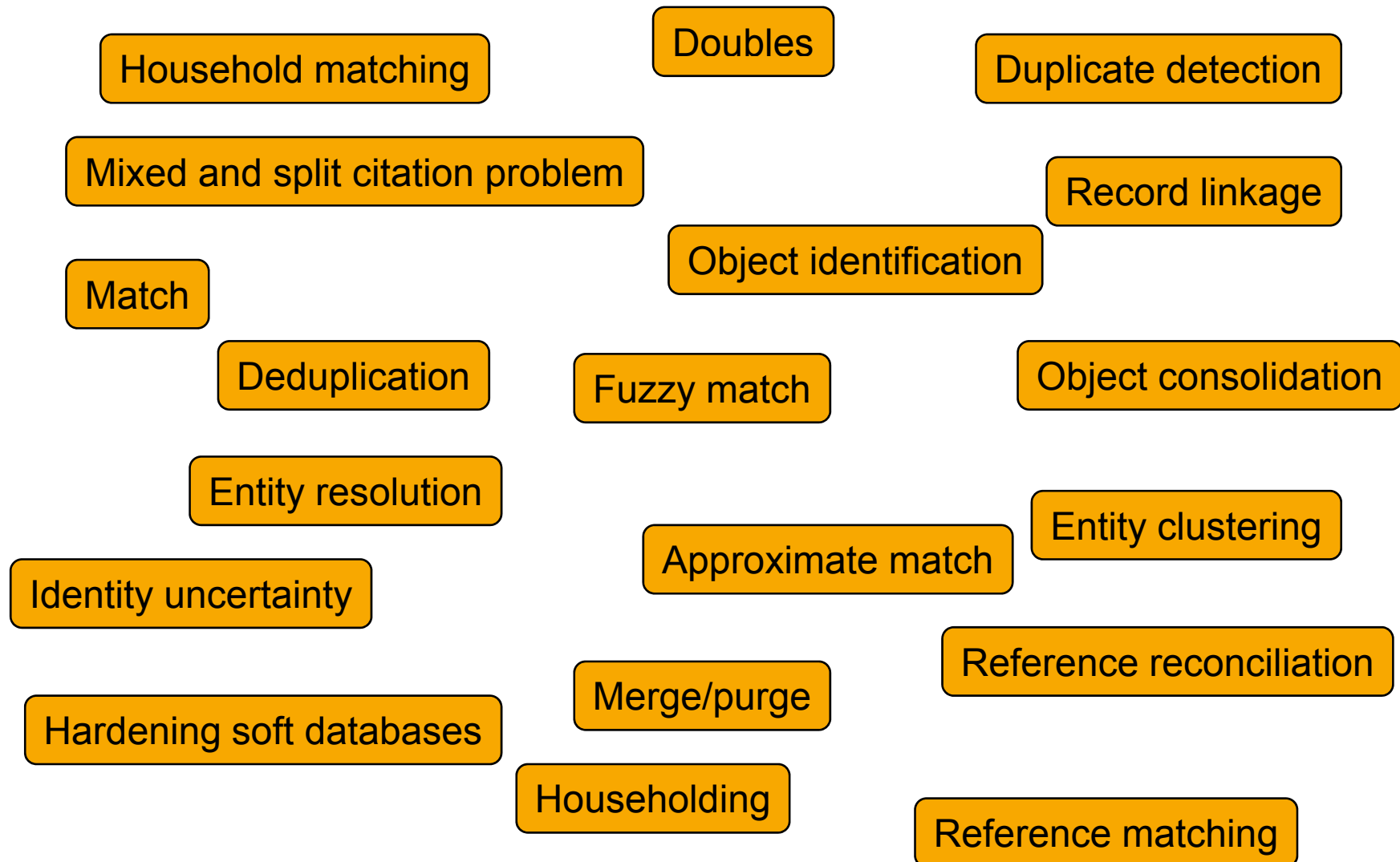
- General problems

- Additional, unnecessary IT expenses
- Low customer satisfaction
- Potentials and dangers not detected
- Poor quality financial data



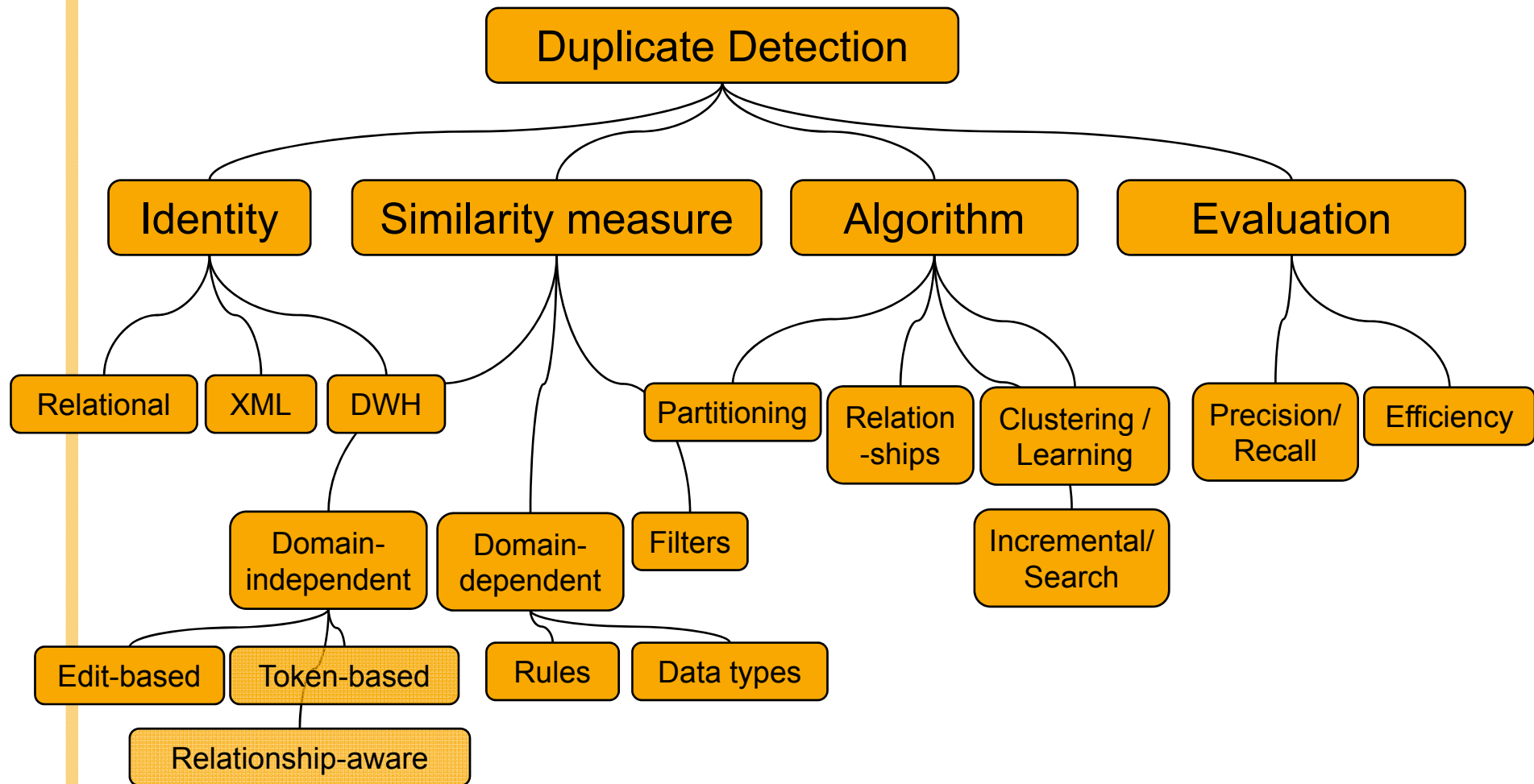
Ironically, “Duplicate Detection” has many Duplicates

29



Duplicate Detection – Research

30



Token-based Similarity Measures

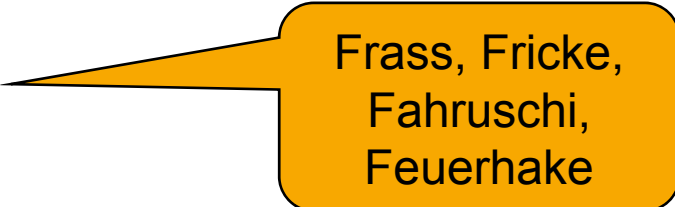
31

- Tokens
 - Words / Terms
 - n-grams
- Jaccard
 - $|\{\text{common tokens}\}| / |\{\text{all tokens}\}|$
- TFIDF [Cohen et al. 2003]
 - Term frequency: *tf*
 - Inverse document frequency: *idf*
 - TFIDF: $\log (tf+1) \times \log (idf)$
 - Common words have low weight
 - Similarity measure: Cosine similarity of term vectors weighted by TFIDF
- And many more
[Koudas Srivastava 2005]

Edit-based Similarity Measures

32

- Jaro [Jaro 1989] / Jaro-Winkler [Winkler 1999]
 - Common letters within $\frac{1}{2}$ string length
 - Transposed letters
- Edit-distance / Levenshtein-distance [Levenshtein 1965]
 - Minimum number of edits from one word to the other
 - Domain-specific costing
 - Dynamic Programming
- Soundex
 - 4-letter code for each word
 - `SOUNDEX('Farwick ')` = F620
- ...



Frass, Fricke,
Fahruschi,
Feuerhake

Domain-dependent Similarity Measures

33

■ Data Types

- Special similarity for dates
- Special similarity for numerical attributes
- ...

■ Rules

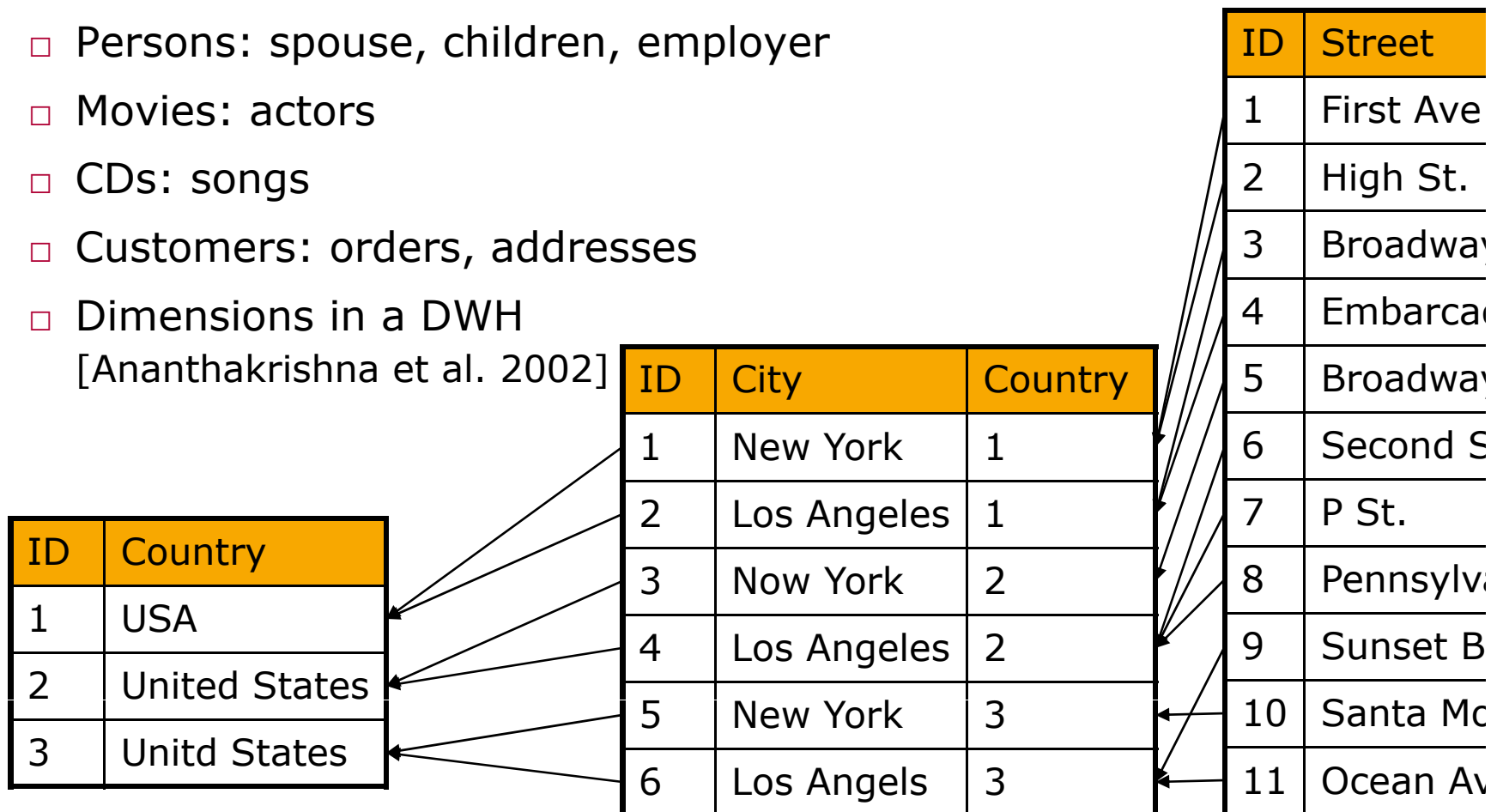
- [Hernandez Stolfo 1998], [Lee et al. 2000]
- Given two records, *r1* and *r2*.
IF last name of *r1* = last name of *r2*,
AND first names differ slightly,
AND address of *r1* = address of *r2*
THEN *r1* is equivalent to *r2*.

Relationship-aware Similarity Measures

34

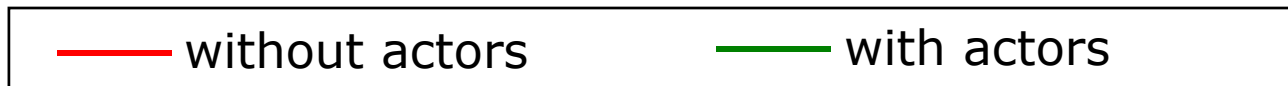
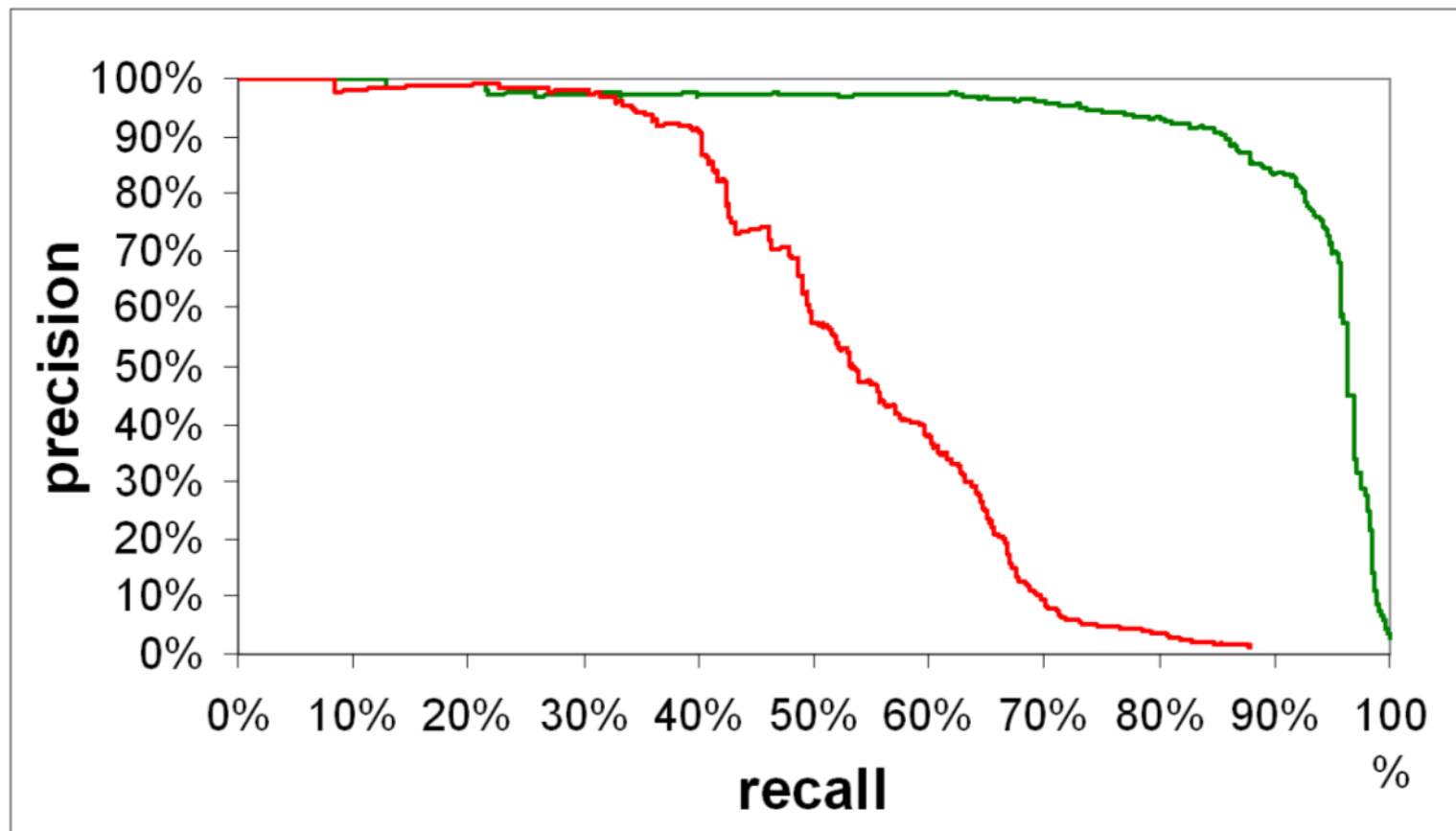
- Idea: Not only values of the records, but values of related records are relevant for similarity.

- Persons: spouse, children, employer
- Movies: actors
- CDs: songs
- Customers: orders, addresses
- Dimensions in a DWH
[Ananthakrishna et al. 2002]



Relationship-aware Similarity Measures – Evaluation

35



Partitioning / Blocking

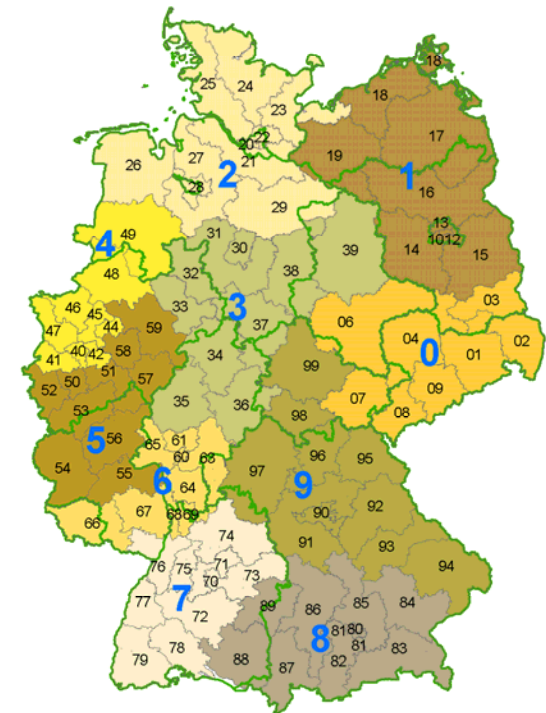
36

- Partition the records (horizontally) and compare pairs of records only within a partition.

- Partitioning by first two zip-digits
 - ◇ Ca. 100 partitions in Germany
 - ◇ Ca. 100 customers per partition
 - ◇ => 495.000 comparisons
- Partition by first letter of surname
- ...

- Idea: Partition multiple times by different criteria.

- Then apply transitive closure on discovered duplicates.



Source: wikipedia.de

Sorted Neighborhood

[Hernandez Stolfo 1998]

37

- Idea
 - Sort tuples so that similar tuples are close to each other.
 - Only compare tuples within a small neighborhood (window).
- 1. Generate key
 - E.g.: SSN+“first 3 letters of name” + ...
- 2. Sort by key
 - Similar tuples end up close to each other.
- 3. Slide window over sorted tuples
 - Compare all pairs of tuples within window.
- Problems
 - Choice of key
 - Choice of window size
- Complexity: At least 3 passes over data
 - Sorting!

Overview

38

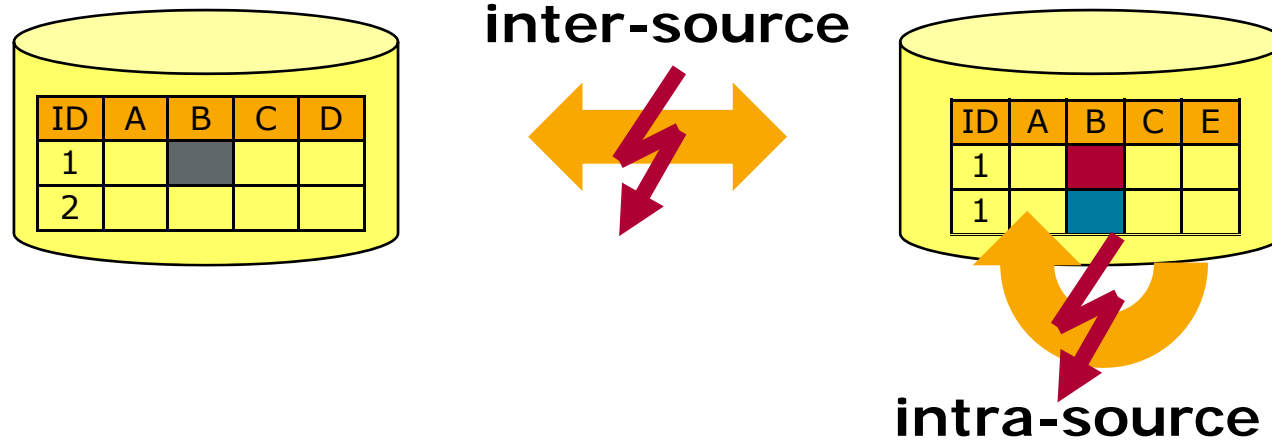
- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
 - ➔ □ Data Conflicts
 - Relational Operators
 - Conflict Resolution
 - Tools
- Summary



Data conflicts

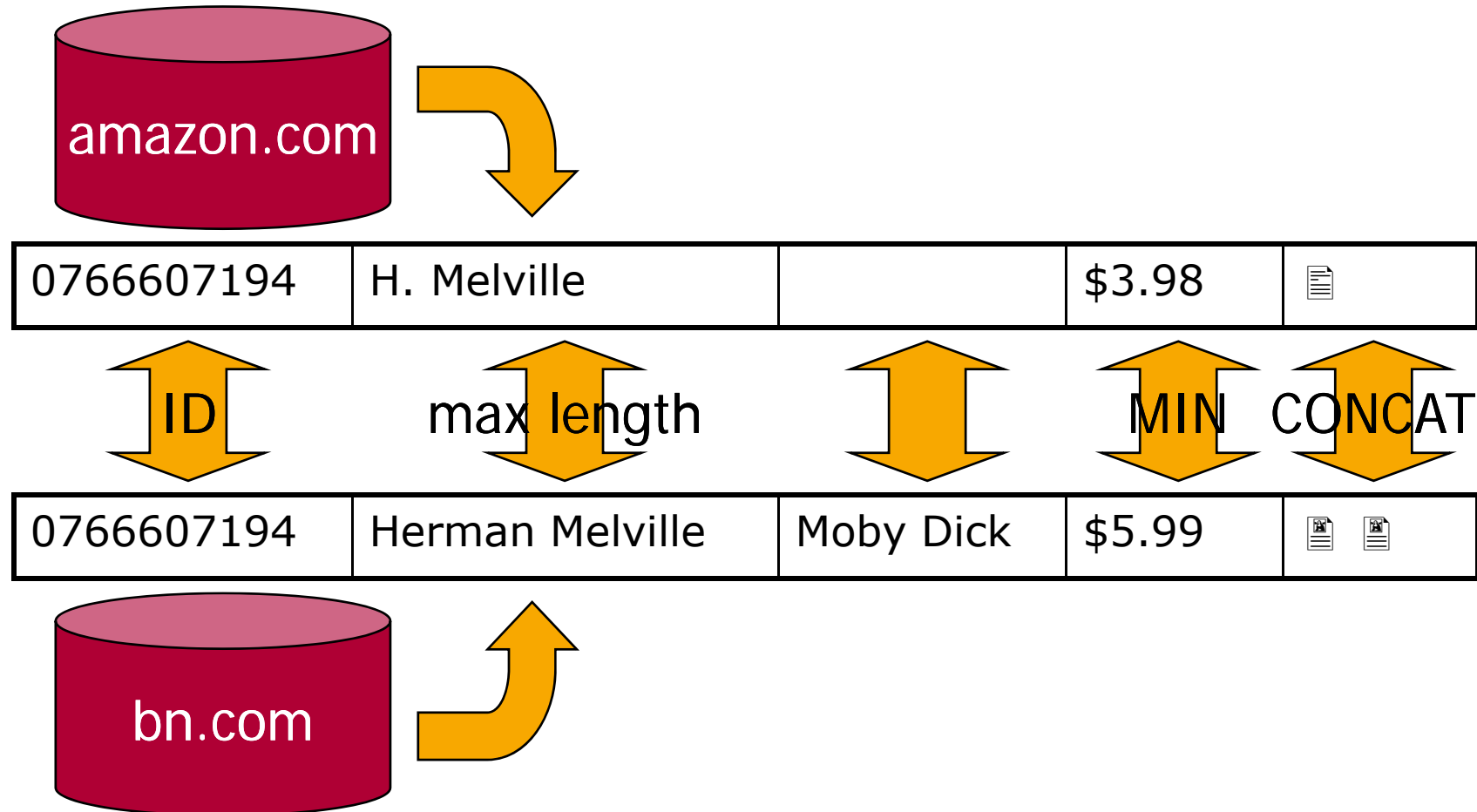
39

Two **duplicates** have different values for **semantically same attribute**.



Data Fusion

40



Data conflicts – origins

41

- No integrity or consistency checks
- Redundant schemata
- Typos, transmission errors, incorrect calculations
- Variants
 - Kantstr. / Kantstrasse / Kant Str. / Kant Strasse
 - Kolmogorov / Kolmogoroff / Kolmogorow
- Typical confusion (OCR)
 - U<->V, 0<->o, 1<->l, etc.
- Obsolete values
 - Different update frequencies, forgotten update

Within information system

Data conflicts – origins

42

- Locally consistent but globally inconsistent
- Duplicates
- Different data types
- Local spelling variations and conventions

Across information systems

Examples for errors

43

■ Addresses

- Str. → Straße, Ch. → Chaussee, etc.
- R.-Breitscheid-Str. 72 a → Rudolf-Breitscheid.-Str. 72A
- 128 spellings for Frankfurt am Main
 - ◇ Frankfurt a.M., Frankfurt/M Frankfurt, Frankfurt a. Main, ...

■ Names

- Dr. Ing. h.c. F. Porsche AG
- Hewlett-Packard Development Company, L.P.



■ Numerical data

- 10.000 € = 10T EURO = 10k EUR = 10.000,00€ = 10,000.- €

■ Phone numbers

■ Birth dates

Numerical data - consequences

44

Southwest
NEWSGROUP

Published on Chanhassen Villager (<http://www.chanvillager.com>)

Property mistakenly valued at \$189 million

By rcrw

Created 12/03/2007 - 4:46pm

Property mistakenly valued at \$189 million results in tax adjustments in county

An \$18,900 Waconia property that was mistakenly valued at \$189 million is "throwing a wrench" into property tax statements and the Carver County budget. County officials issued a press release Monday detailing the problem that came to light last week.

An error was identified in the estimated market valuations used to calculate Pay 2008 Proposed Property Taxes, according to the release. The County Assessor's Office placed an incorrect estimated market value on a parcel located in the city of Waconia, apparently resulting in extra zeroes being added to the value.

The mistake results in an imbalance in the amount of property taxes the county was expecting to collect. The mistake added about \$900,000 in expected revenue, according to County Administrator David Hemze.

The county is planning to consider recommendations to cut the 2008 budget by \$900,000 so that proposed property taxes will match tax notices sent to residents in November.

"It kind of threw a wrench into everything," said Hemze. "It's unfortunate. It's a mistake and we're concentrating on responding to the mistake and trying to ensure that it doesn't happen again." If the county does not cut the budget by \$900,000, the county portion of property taxes would go up for all properties in the county. The effect would be greatest in Waconia, but Hemze said the average-valued home outside of Waconia would also experience a \$29 increase on top of the number indicated on the November tax notices.

Numerical data - consequences

45

SPIEGEL ONLINE

28. Januar 2008, 11:27 Uhr

FRANKREICH

Telefonkundin erhält Rechnung über 63 Millionen Euro

Als eine Französin aus Lothringen unlängst ihre Telefonrechnung bekam, blieb ihr buchstäblich die Spucke weg: 63 Millionen Euro sollte sie begleichen. Dabei hatte sie ursprünglich nur um Korrektur einer Abrechnung in Höhe von 67 Euro gebeten.

Paris - "Da muss wohl ein Komma verrutscht sein", zitiert "Le Figaro" heute den Vizedirektor der französischen Telefongesellschaft Télé2, Olivier Anstett. Die Kundin aus dem Ort Herserange in der Nähe von Metz hatte sich zunächst über einen ihrer Meinung nach zu hohen Rechnungsbetrag von 67,69 Euro bei der Telefongesellschaft beschwert. Als eine Antwort ausblieb, schickte sie einen zweiten Brief. Daraufhin erhielt sie eine "korrigierte" Rechnung über die Summe 63.280.067,96 Euro.

"Uns bleibt nur, uns bei der Kundin zu entschuldigen und dafür zu sorgen, dass so etwas nie wieder vorkommt", so der lapidare Kommentar des Vizechefs von Télé2.

Data conflict – elimination

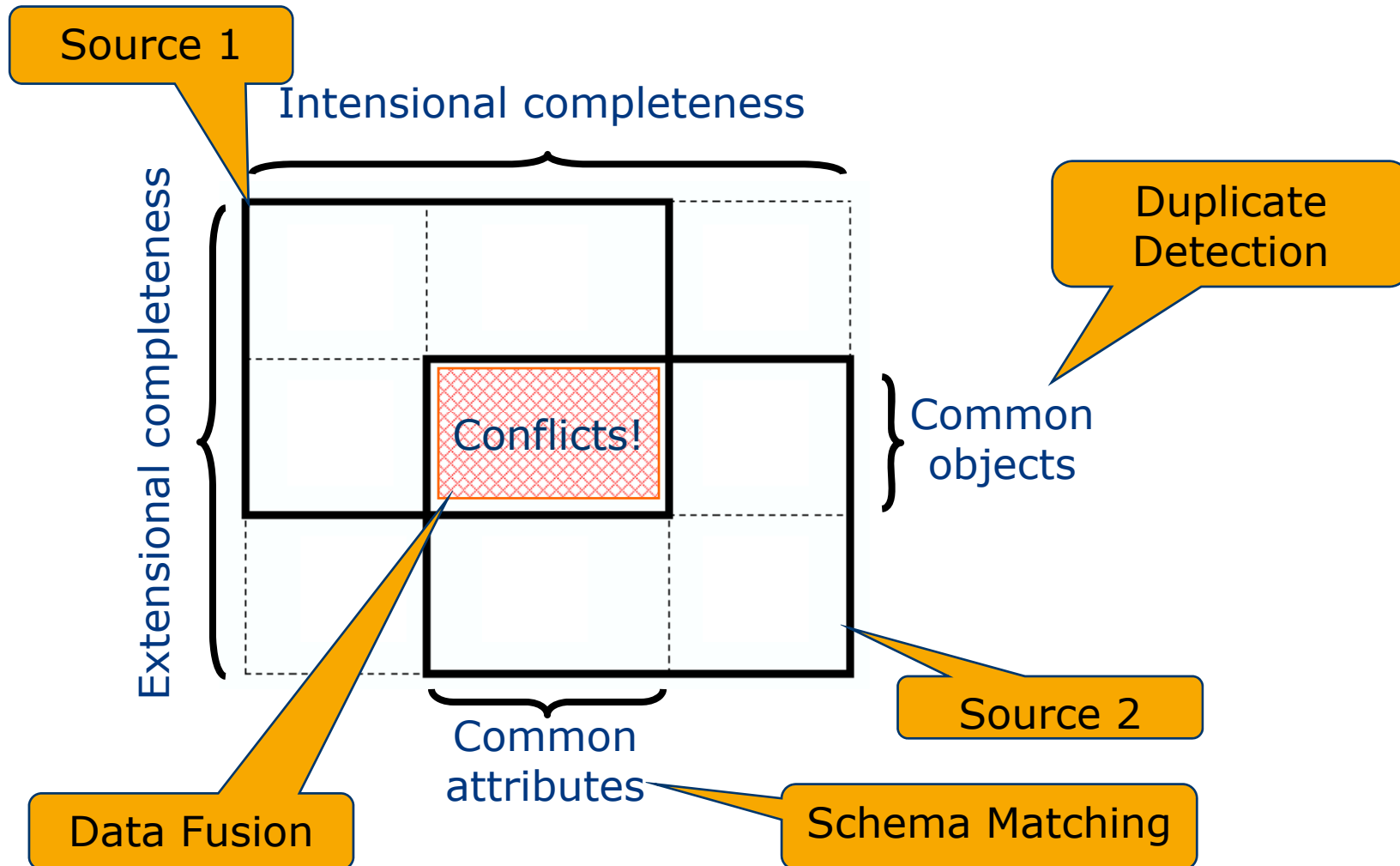
46

- Error correction
 - Reference tables
 - ◇ Cities, countries, products ...
 - Similarity measures
 - ◇ For typos
 - ◇ For language-specific variants (Meier, Mayer,...)
 - Standardization and Transformation
 - Domain-knowledge (meta data)
 - ◇ Konventions (country/region-specific spelling)
 - ◇ Ontologies
 - ◇ Thesauri, dictionaries for homonyms, synonyms, ...

- And data fusion...

Completeness and Conciseness

47



Overview

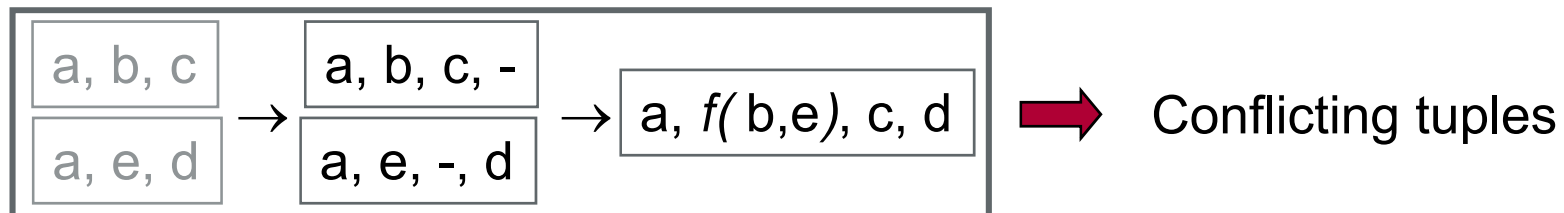
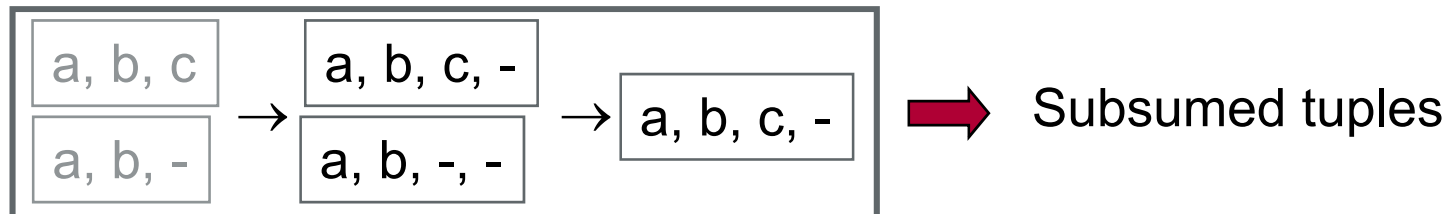
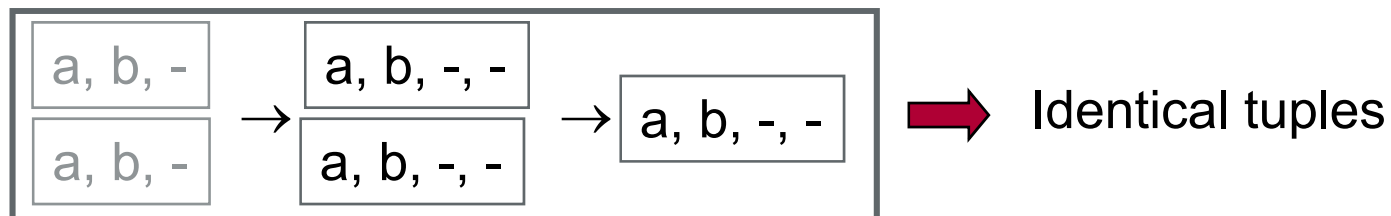
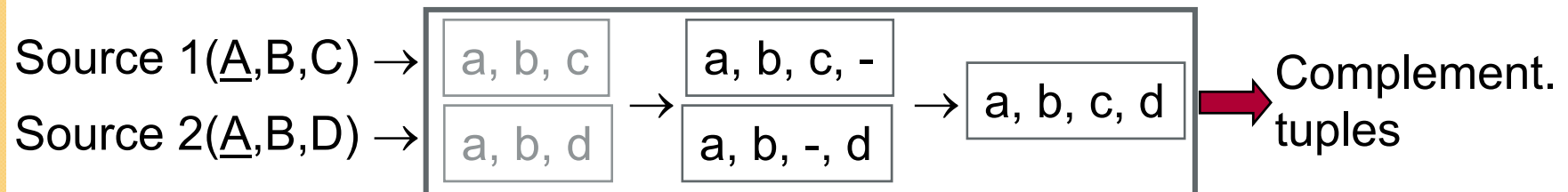
48

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
 - Data Conflicts
 - Relational Operators
 - Conflict Resolution
 - Tools
- Summary



"Proper" Data Fusion

49



Relational object integration

50

Union

+ Elimination of exact duplicates

Minimum Union, [GL94]

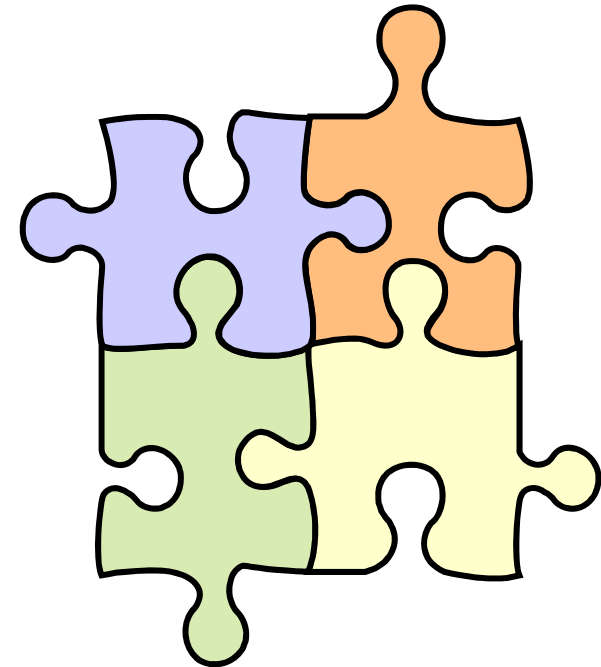
+ Elimination of subsumed tuples

But

✗ No duplicate integration

✗ Conflict resolution

Later: Join, Merge, Group, ...



Minimum Union – Outer Union

51

- Outer Union pads relations with NULL-values, to unify schemata.
- Then normal UNION.
- Usually not implemented in SQL

R			S	
A	B	C	B	D
P	1	2	2	U
P	2	1	3	V
Q	1	2		

$R \uplus S$			
A	B	C	D
P	1	2	⊥
P	2	1	⊥
Q	1	2	⊥
⊥	2	⊥	U
⊥	3	⊥	V

Minimum Union – Subsumption

52

- A tuple t_1 subsumes a tuple t_2 , if
 - Has same schema,
 - t_2 has more NULL-values than t_1 ,
 - $t_1 = t_2$ for all non-NULL-values of t_2 .

- Notation:

- $R\downarrow$ returns those tuples of R , that are not subsumed by any other tuple in R .

R			
p_id	fname	lname	age
1	Peter	Müller	32
1	Peter	Müller	⊥
1	Peter	⊥	⊥
1	Peter	⊥	32
1	Peter	⊥	42
2	Wiebke	⊥	2
2	⊥	Meyer	2

How many tuples does $R\downarrow$ have?

$R\downarrow$			
p_id	fname	lname	age
1	Peter	Müller	32
1	Peter	⊥	42
2	Wiebke	⊥	2
2	⊥	Meyer	2

Minimum Union – NULL-values

53

Semantics of NULL? [GUW02]

- *"unknown"*
 - There is a value, but I do not know it.
 - E.g.: Unknown birthday
- *"inapplicable"*
 - There is no meaningful value
 - E.g.: partner for singles
- *"withheld"*
 - There is a value, but we are not authorized to see it.
 - E.g.: Private phone

Minimum Union – NULL-values

54

"Value not supplied"

"Value does not exist"

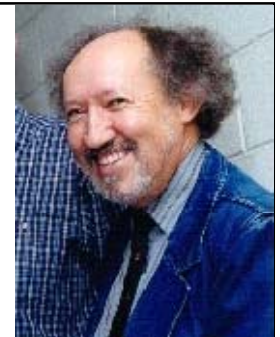
"Value undefined"

„Distinguished" NULL

"Total ignorance" NULL

C.J. Date:

- "Into the Unknown"
- "Much Ado About Nothing"
- "NOT Is Not Not!"
- "Oh No Not Nulls Again"
- ...



From now on: "Unknown"

Merge and Prioritized Merge

55

Merge (\boxtimes), [GPZ01]

- Mixes Join and Union to a new operator
 - Fuses complementary tuples (only from different sources)
- COALESCE removes NULL values
- Prioritization possible (\triangleright)
- Can be expressed with standard SQL

```
( SELECT K.p_id, K.fname, Coalesce(K.lname, C.lname), Coalesce(K.age, C.age)
FROM K LEFT OUTER JOIN C ON K.p_id = C.p_id )
```

UNION

```
( SELECT C.p_id, K.fname, Coalesce(C.lname, K.lname), Coalesce(C.age, K.age)
FROM K RIGHT OUTER JOIN C ON K.p_id = C.p_id )
```

Merge – example

56

Kunde K

p_id	fname	lname	age
1	Peter	Müller	32
2	Franz	Schmidt	55
3	Wiebke	Meyer	⊥
4	Klaus	Lehmann	28

Customer C

p_id	lname	age
1	⊥	32
2	Schmidt	⊥
3	Meier	56
5	Weger	47

C ⊗ K

p_id	fname	lname	age
1	Peter	Müller	32
2	Franz	Schmidt	55
3	Wiebke	Meier	56
3	Wiebke	Meyer	56
4	Klaus	Lehmann	28
5	⊥	Weger	47

What else is there?

57

Match Join [YaÖz99]

- Complex operator
- *HighConfidence, RandomEvidence, and PossibleAtAll*

ConQuer [FuFM05]

- „Consistent Query Answering“
- Rewriting of SQL queries

Burdick et. al. [BDJR05]

- Uncertainty in Data Warehouses
- „Possible Worlds“

Probabilistic Models [Mich89]

- Extending schema by probabilities

Grouping for Integration

58

$K \uplus C$

p_id	fname	lname	age
1	Peter	Müller	32
2	Franz	Schmidt	55
3	Wiebke	Meyer	55
1	⊥	Müller	32
2	⊥	Schmidt	⊥
3	⊥	Meier	56

p_id	fname	lname	age
1	Peter	Müller	32
2	Franz	Schmidt	55
3	Wiebke	Meier	56

```
SELECT      p_id, MAXLEN(fname), CHOOSE(lname,C), MAX(age)
FROM
GROUP BY    $K \uplus C$ 
           p_id
```

Longest
String

C is favored
source

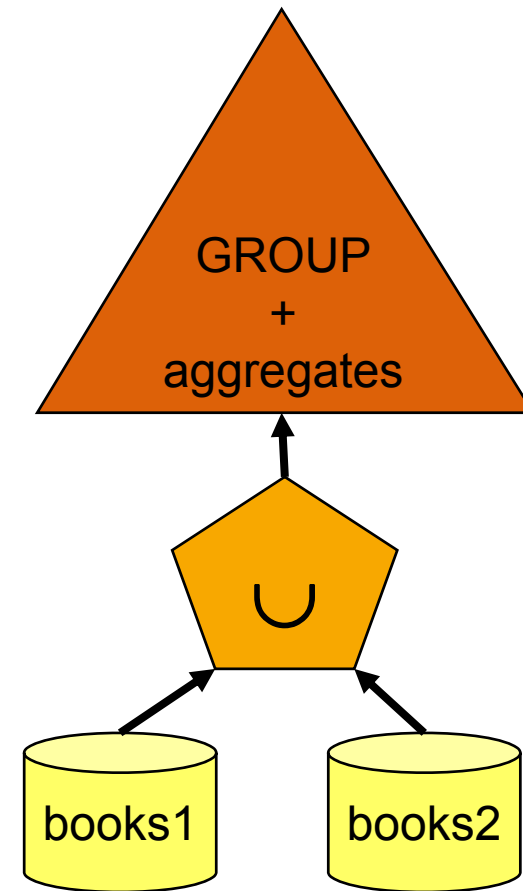
Highest
value

Grouping – example

59

```

SELECT books.isbn,
       MAXLEN(books.title),
       MIN(books.price)
FROM   (
        SELECT * FROM books1
        UNION
        SELECT * FROM books2
      )
AS books
GROUP BY
       books.isbn
  
```



Grouping – Pros and Cons

60

+ Pros

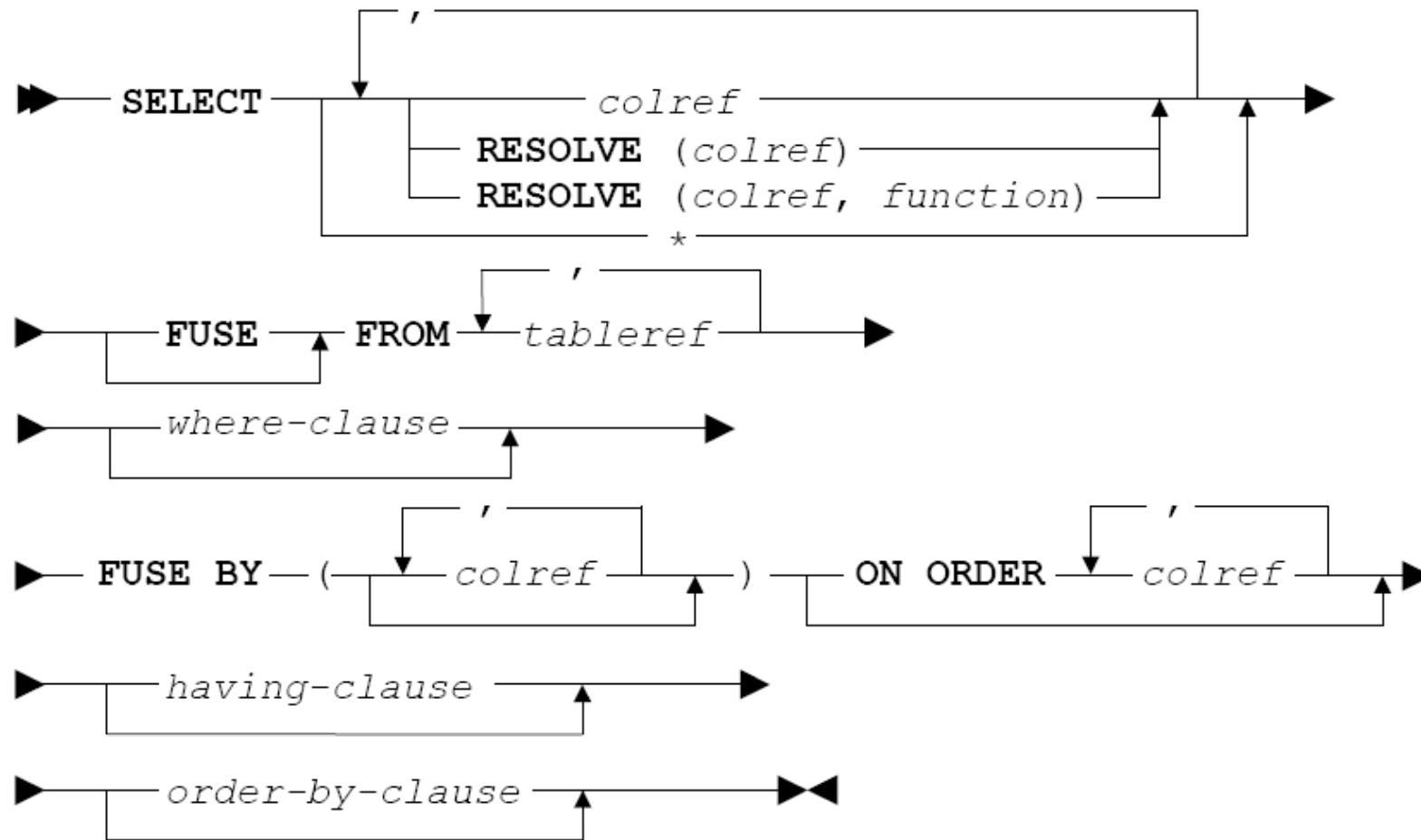
- Efficient
 - Implemented by sorting
- Catches duplicates within source and across sources
- Simple / short

✗ Cons

- Restricted to built-in standard aggregate-functions:
 - ◇ MAX, MIN, AVG, VAR, STDDEV, SUM, COUNT
- Grouping only by equality
 - ◇ ID attribute is necessary
- Outer Union usually not implemented

FUSE BY Syntax Diagram

61



Fuse By – Queries

62

```
SELECT *
FUSE FROM Q1
FUSE BY (Name)
```

```
SELECT *
FUSE FROM Q1
FUSE BY ()
```

```
SELECT *
FUSE FROM Q1, Q2
FUSE BY ()
```

Grouping with
coalesce
aggregation

Subsumption

Minimum
Union

```
SELECT Name, RESOLVE(Age, max), RESOLVE(Student,
    vote), RESOLVE(Place), RESOLVE(Phone)
FUSE FROM Q1, Q2
FUSE BY (Name) ON ORDER Q2.Age DESC
```

FUSE BY – Example

63

Name	Age	Student	Place
------	-----	---------	-------

<i>Felix</i>	⊥	No	Hamburg
<i>Melanie</i>	22	Yes	⊥
<i>Jens</i>	⊥	Yes	Karlsruhe
<i>Christoph</i>	25	Yes	Berlin
<i>Sven</i>	26	Yes	Berlin
<i>Sven</i>	⊥	Yes	Berlin

Name	Age	Student	Phone
------	-----	---------	-------

<i>Melanie</i>	⊥	Yes	030/12345
<i>Jens</i>	27	⊥	030/54321
<i>Christoph</i>	24	Yes	⊥
<i>Melanie</i>	21	No	030/98765
<i>Karsten</i>	24	Yes	⊥
<i>Karsten</i>	24	Yes	⊥

```

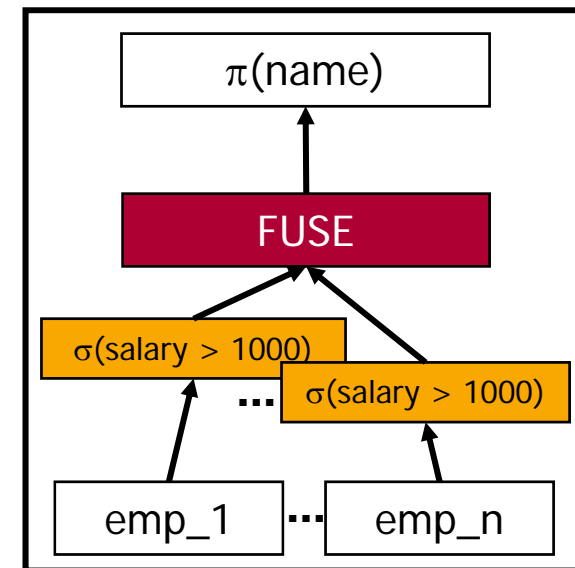
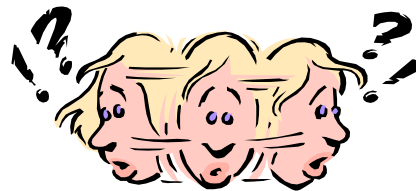
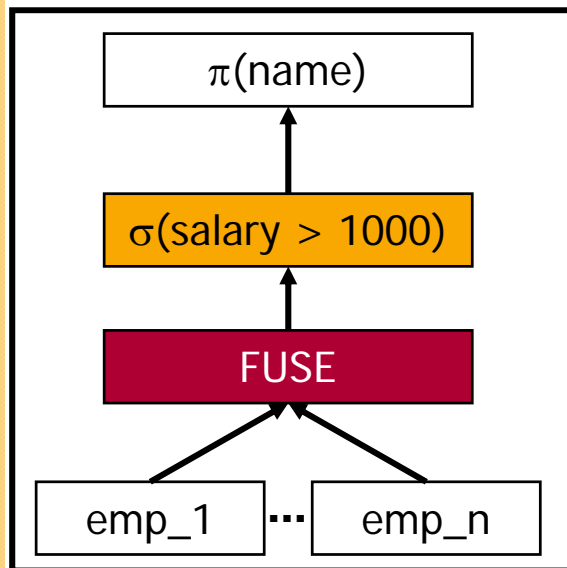
SELECT Name,
  RESOLVE(Age, max),
  RESOLVE(Student, vote),
  RESOLVE(Place),
  RESOLVE(Phone)
FUSE FROM Q1, Q2
FUSE BY (Name)
ON ORDER Q2.Alter DESC
  
```

Name	Age	Student	Place	Phone
------	-----	---------	-------	-------

<i>Felix</i>	⊥	No	Hamburg	⊥
<i>Melanie</i>	22	Yes	⊥	030/98765
<i>Jens</i>	27	Yes	Karlsruhe	030/54321
<i>Christoph</i>	25	Yes	Berlin	⊥
<i>Sven</i>	26	Yes	Berlin	⊥
<i>Karsten</i>	24	Yes	⊥	⊥

Query Optimization with Fusion

64



Correctness?
Completeness?
Efficiency?

Overview

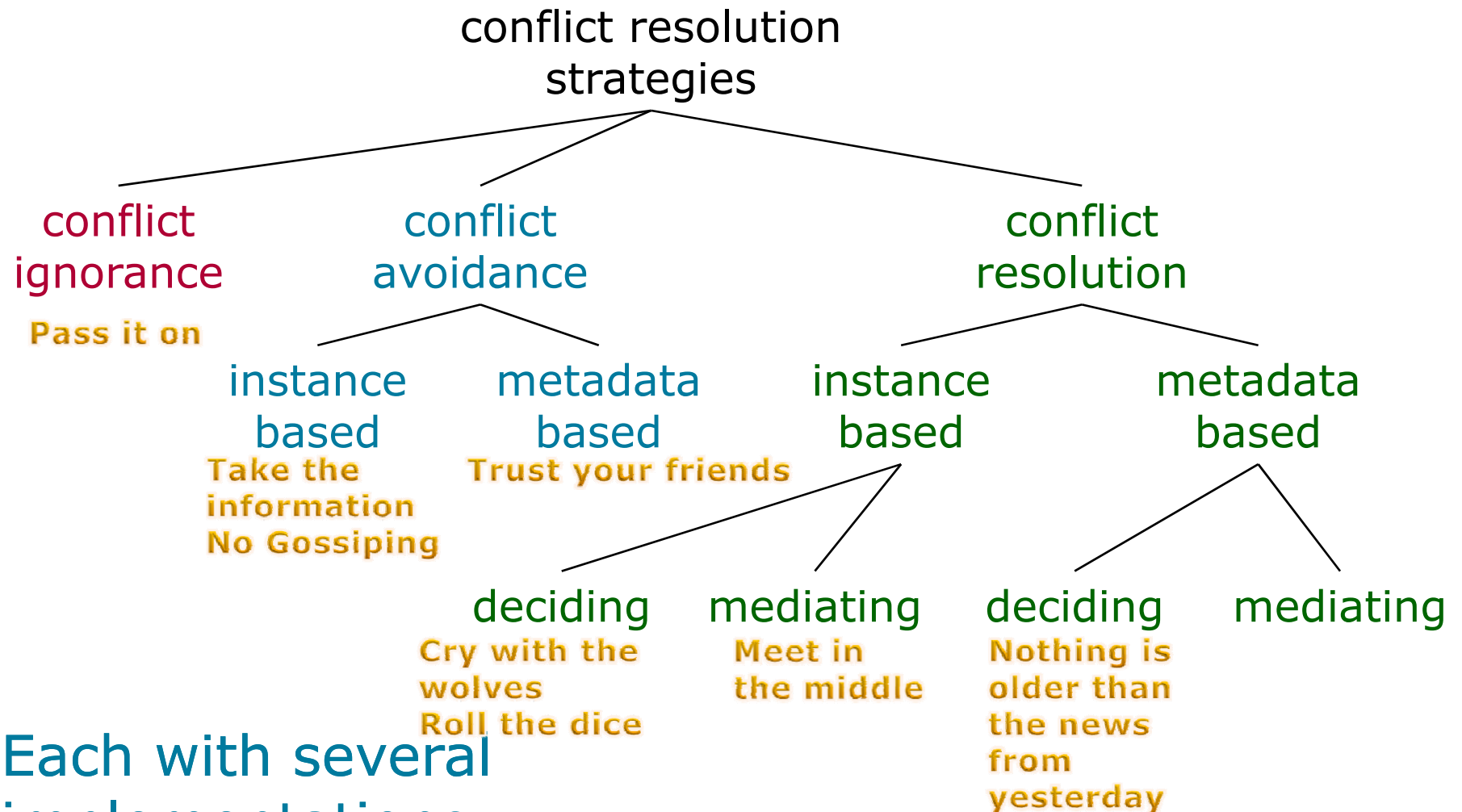
65

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
 - Data Conflicts
 - Relational Operators
 - Conflict Resolution
 - Tools
- Summary



Classification of strategies

66



Each with several implementations

Conflict Resolution Functions

67

Min, Max, Sum, Count, Avg, StdDev	Standard aggregation
Random	Random choice
First, Last	Choose first/last value; depends on order
Longest, Shortest	Choose longest/shortest value
Choose(<i>source</i>)	Choose value from a particular source
ChooseDepending(<i>col</i> , <i>val</i>)	Choose depending on <i>val</i> in other column <i>col</i>
Vote	Majority decision
Coalesce	Choose first non-null value
Group, Concat	Group or concatenate all values
MostRecent	Choose most recent (up-to-date) value
MostAbstract, MostSpecific	Use a taxonomy / ontology
....

Overview

68

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
 - Data Conflicts
 - Relational Operators
 - Conflict Resolution
 - Tools
- Summary



Visualization of Integrated Data

69

HumMer-Demo File Extra Help

0. Sources
 1. Matching
 2. Duplicate Definition
 3. Duplicate Detection
 4. Conflict
 5. Result

Result

Choose the fusion implementation to use **default**

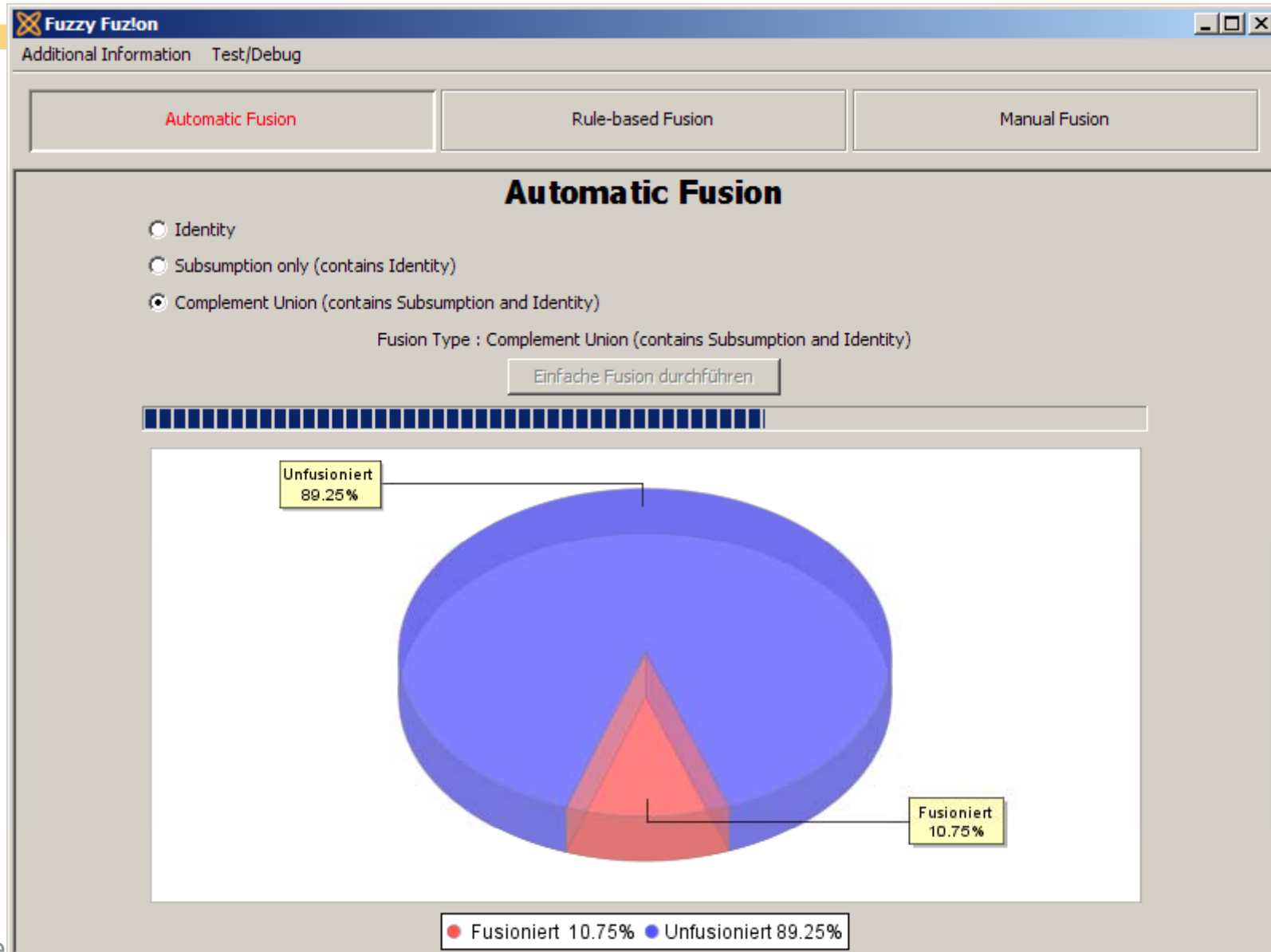
#	CLU...	TITLE	VERSI...	COUN...	YEAR	ORIGI...	GENRE	DIREC.
	VOTE	COALESCE	COALES...	MAX	COALES...	LAST	COALES.	
13	87	HOPE FLOATS	engl...	USA	1998	Hop...	Unterhaltu...	Fore...
14	84	GOOD WILL H...	engl...	USA	1998	Goo...	Drama	Gus .
15	83	GODZILLA	engl...	USA	1998	God...	Fantasy, S...	Rola.
16	80	Gadjo Dilo Gadjo Dilo GADJO DILO	franz... franz.&r...	F/Rum	1998 1998 1997	Gadj... Gadjo ...	Unterhaltu... Unterhaltung Drama	Ton.
17	77	Deconstructin...	engl...	USA	1998	Dec...	Komödie/...	Woo.
18	74	City Of Angels	engl...	USA	1998	City ...	Drama	Brad.
19	69	BOOGIE NIGH...	engl...	USA	1998	Boo...	biografisc...	Paul .
20	65	Antz	engl...	USA	1998	Antz	Animation...	Darn.
21	57	SPIDER			2002		Drama	
22	51	SECRETARY			2002		Komödie	
23	49	S.F.W.			1994		Komödie	
24	31	Intolerable Cr...			2003		Komödie	
25	25	GANGSTER N...			2000		Gangsterfi...	
26	24	From Hell			2001			
27	17	DEATHWATCH			2002		Kriegsfilm	
28	15	CHARLOTTE ...			2001		Melodram	
29	11	Big Fish			2003		Drama	

Rows: 0:99

Duplicate **Contradiction** **Uncertainty** **Unique**

Bachelorprojekt „Fuz!on“

70



Bachelorprojekt „Fuz!on“

71

Fuzzy Fuzion Additional Information Test/Debug

Automatic Fusion **Rule-based Fusion** Manual Fusion

Rule Matrix

	Firstname	Lastname	Street	houenumber	postcode	city	ignore	phone
None	66105	68111	58872	66404	63121	71285	100000	73936
Null values	5671	6402	6116	16746	12208	5643	0	26064
Case Variance	10835	12745	14563	0	0	11330	0	0
Abbreviation	7095	1170	8256	16850	12364	942	0	0
Tokenization	0	0	0	0	0	0	0	0
Substrings	2122	2091	1088	0	12307	1701	0	0
Dominance	2170	2424	2883	0	0	2434	0	0
Low edit distance	5913	7057	7101	0	0	6664	0	0
Global dominance	88	0	762	0	0	1	0	0
Undefined	1	0	359	0	0	0	0	0

Actions

Fusionsregel(n) anzeigen/erzeugen Nur aktuelle Markierung anzeigen **WEITER -->**

Selected Rules

Regeldefinition (Status: neu)

Spalten: Firstname, Lastname

Konflikttypen: Low edit distance

Primäre Konfliktauflösung: Vote
 Minimum fraction of solution (in %) : 50

Ignore case
 Ignore null-values

Sekundäre Konfliktauflösung: First

Aktionen: Übernehmen, Ausblenden, Spalte hinzufügen, Konflikttyp hinzufügen

Bachelorprojekt „Fuz!on“

72

Fuzzy Fuzlon Additional Information Test/Debug

Automatic Fusion Rule-based Fusion Manual Fusion

Groups 0 to 50 of 100000 All Groups Filter Mode

fdb.group	Firstname	Lastname	Street	houenumber	postcode	city	ignore	phone
31750025-01	Werner	Trimpert	Thomas-Man...	89	24943	Kiel	19470524	0461
31758055-01	Artur	Heiser	Kalkgrund	4	24939	Kiel	19360106	
31765505-01	Siegfried	Aswegen	Mürwiker Str.	6	4943	Flensburg	19250404	0461
31772625-01	M.	Blankenburg	Harmsstr.	48	24116	Kiel	19610727	0461
31780965-01	K	Degen	Peter-Chr.-H...	5	24114	Flensburg	19630331	0461
31789325-01	Manh The	Knaut	Wiedeberger ...	37	24943	Flensburg	19280312	0461
31798345-01	horst	Booitsmann		6	24937	Flensburg	19281225	0461

Back Next

21. Group :

Firstname	Lastname	Street	houenumber	postcode	city	ignore	phone
Manh The	Knaut	Wiedeberger Weg	37	24943	Flensburg	19280312	0461
Manh The	KNAUT	Wiedeberger Weg		24943	Flensburg	19280312	0461
Manh	Knaut	WIEDEBERGER WEG	37	24943	Flensburg	19280312	0461
First	Mixed ...	Vote	First non-null value	First	First	First	First
Manh The	Knaut	Wiedeberger Weg	37	24943	Flensburg	19280312	0461

Merge Save Configurations

Tool: fusem

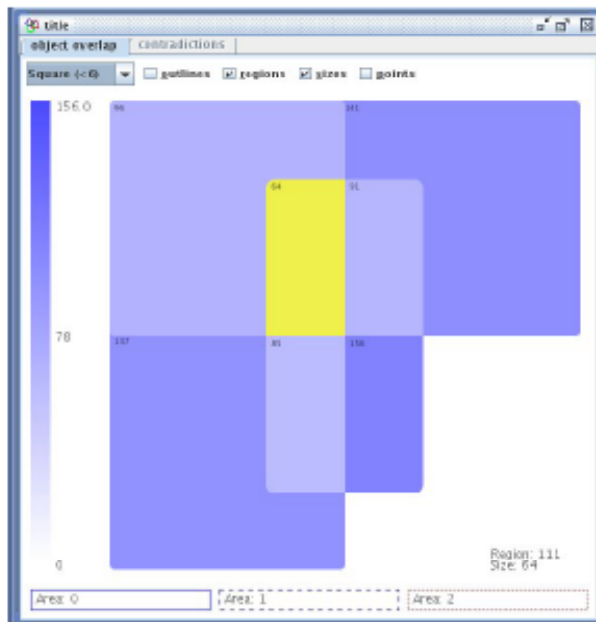
73

The screenshot shows the 'fusem' application window. The main window has a menu bar (File, Extra, Window, Test, Help) and a toolbar with buttons for SQL, FB, CQ, ME, and MJ. A 'Match Join' window is open, displaying a SQL query: `SELECT [ANY] OBJECTID, TUPLEID, NAME, AGE, PHONE, CAI FROM CSSTU WHERE Age> WITH HIGHC`. A 'Fuseby' window is also open, showing a query: `SELECT * FUSE FROM CSSTUDENTS, EESTUDENTS FUSE BY (Name)`. The 'Fuseby' window has a 'Run Query' button and 'Fusion' options: Rewrite, normal opt., and ext. opt. Below the query window is a 'Result Table' tab showing a table with columns TUPLEID, OBJECTID, NAME, and AGI. The table contains 6 rows of data.

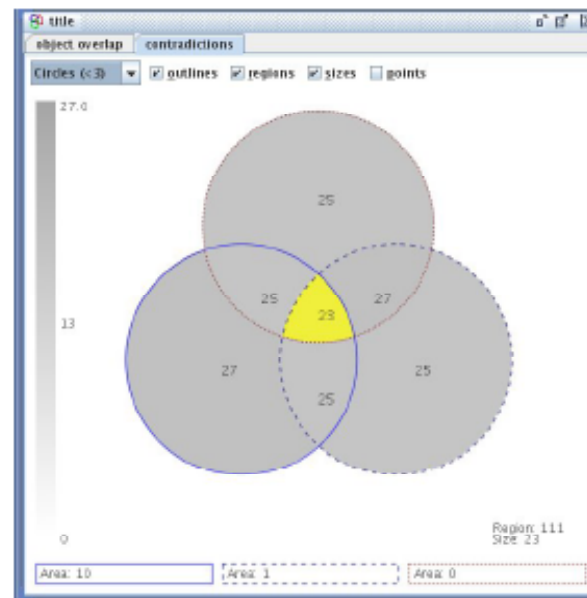
TUPLEID	OBJECTID	NAME	AGI
1	1	Peter	
5	5	Paul	
11	6	Mary	
13	7	Frank	
4	4	Charly	
3	3	Bob	

Visualizing semantics and conflicts

74



(a) Object overlap

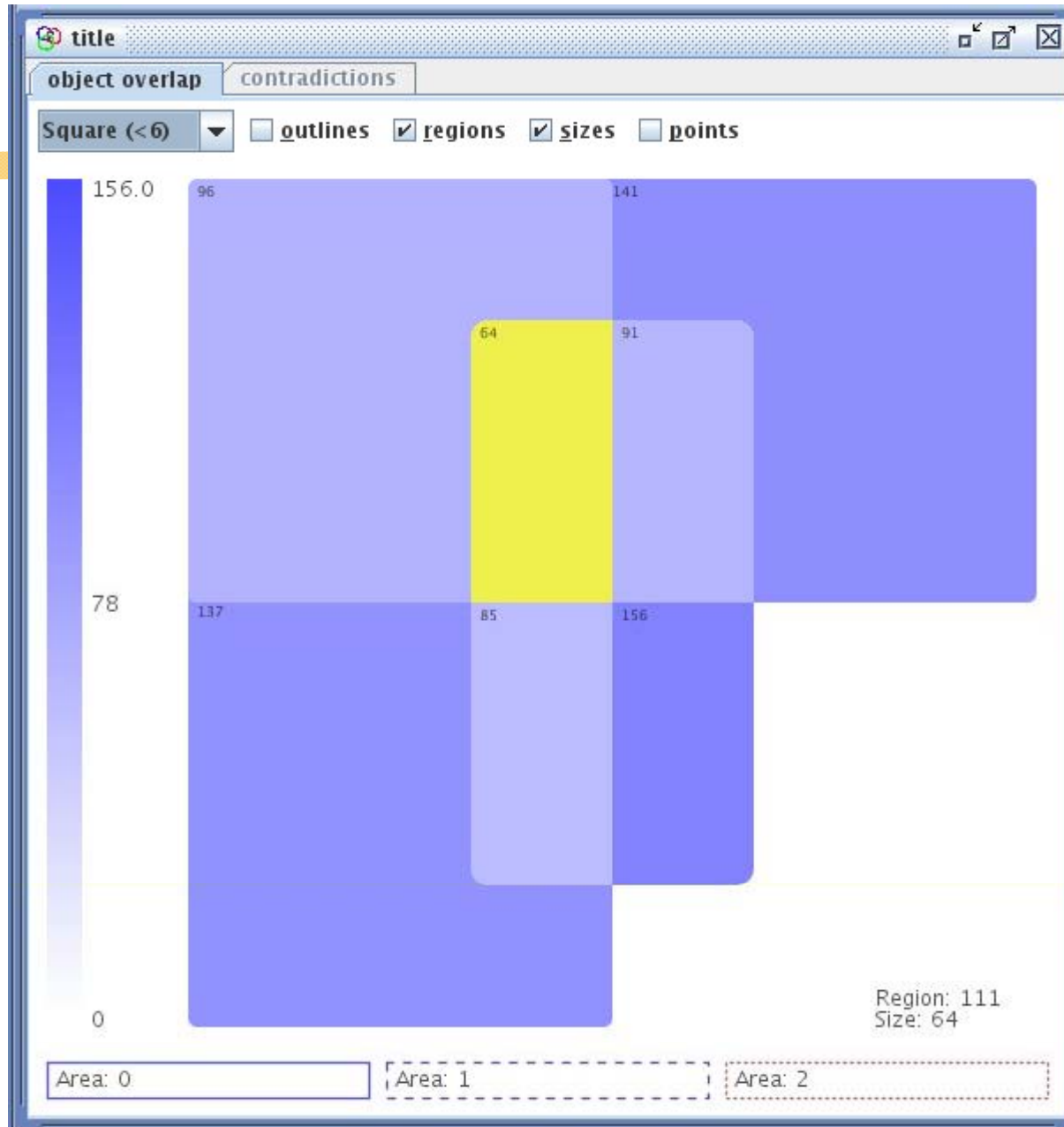


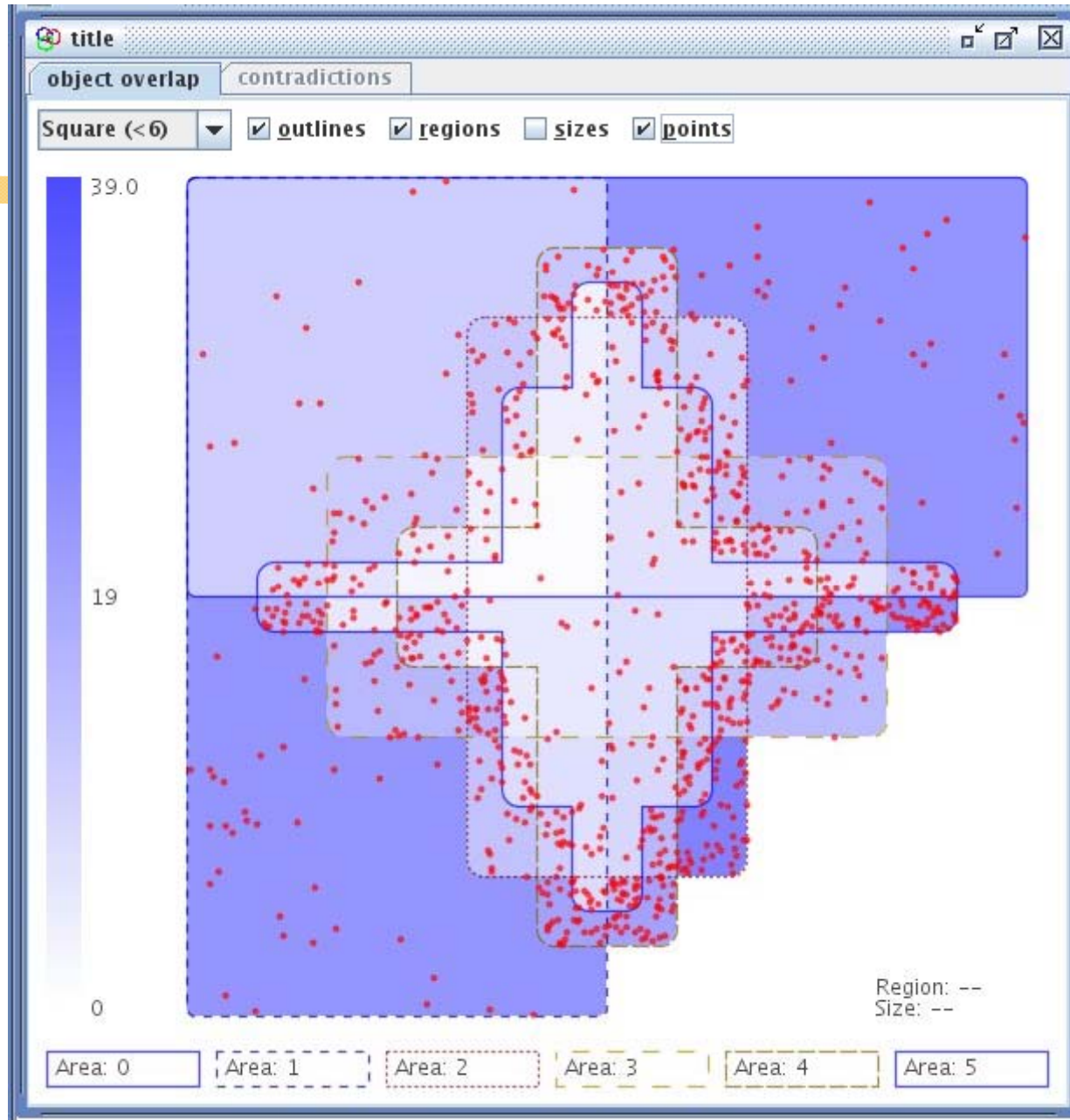
(b) Contradiction Count

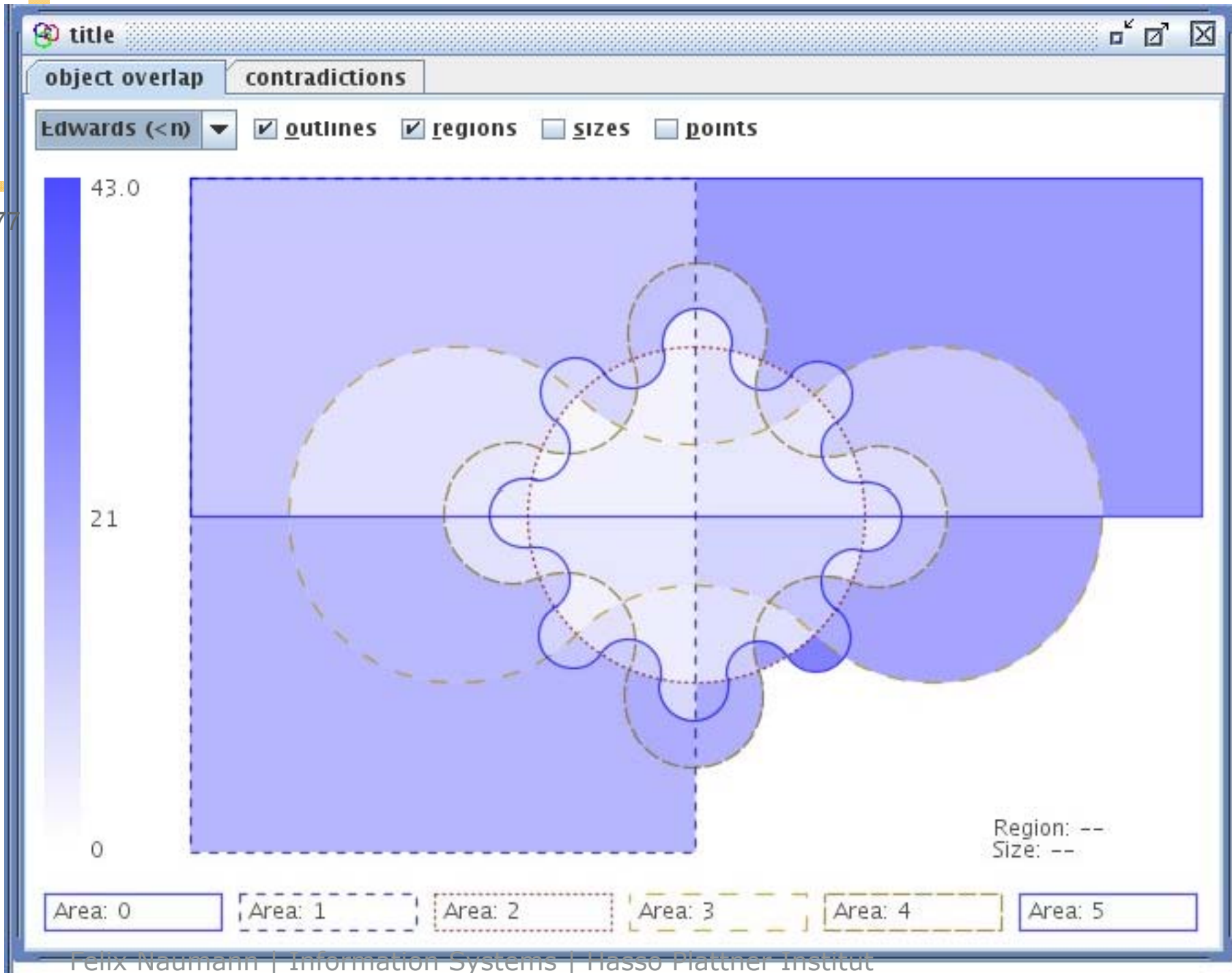
Object	Semantics	X	Y
O ₁	Merge	1	2
	FuseBy	1	2
	MatchJoin	1	3
O ₂	Merge	1	1
	FuseBy	1	1
	MatchJoin	1	1
O ₃	Merge	1	2
	FuseBy	3	4
	MatchJoin	5	6
O ₄
...

(c) Contradicting tuples

75



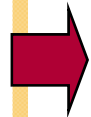




Overview

78

- Introductory example
- Step 1: Schema Matching
- Step 2: Duplicate detection
- Step 3: Data fusion
 - Data Conflicts
 - Relational Operators
 - Conflict Resolution
 - Tools
- Summary



Summary

79

- Step 1: Schema Matching
 - Similarity Measure
 - Combination of methods
- Step 2: Duplicate Detection
 - Similarity Measure
 - Algorithm
 - Data Model
- Step 3: Data Fusion
 - Relational Operators
 - Conflict Resolution
 - Visualization of Semantics and Overlap

VIQTOR: Quality Annotations

80

enter quality value for selected data: Reliability

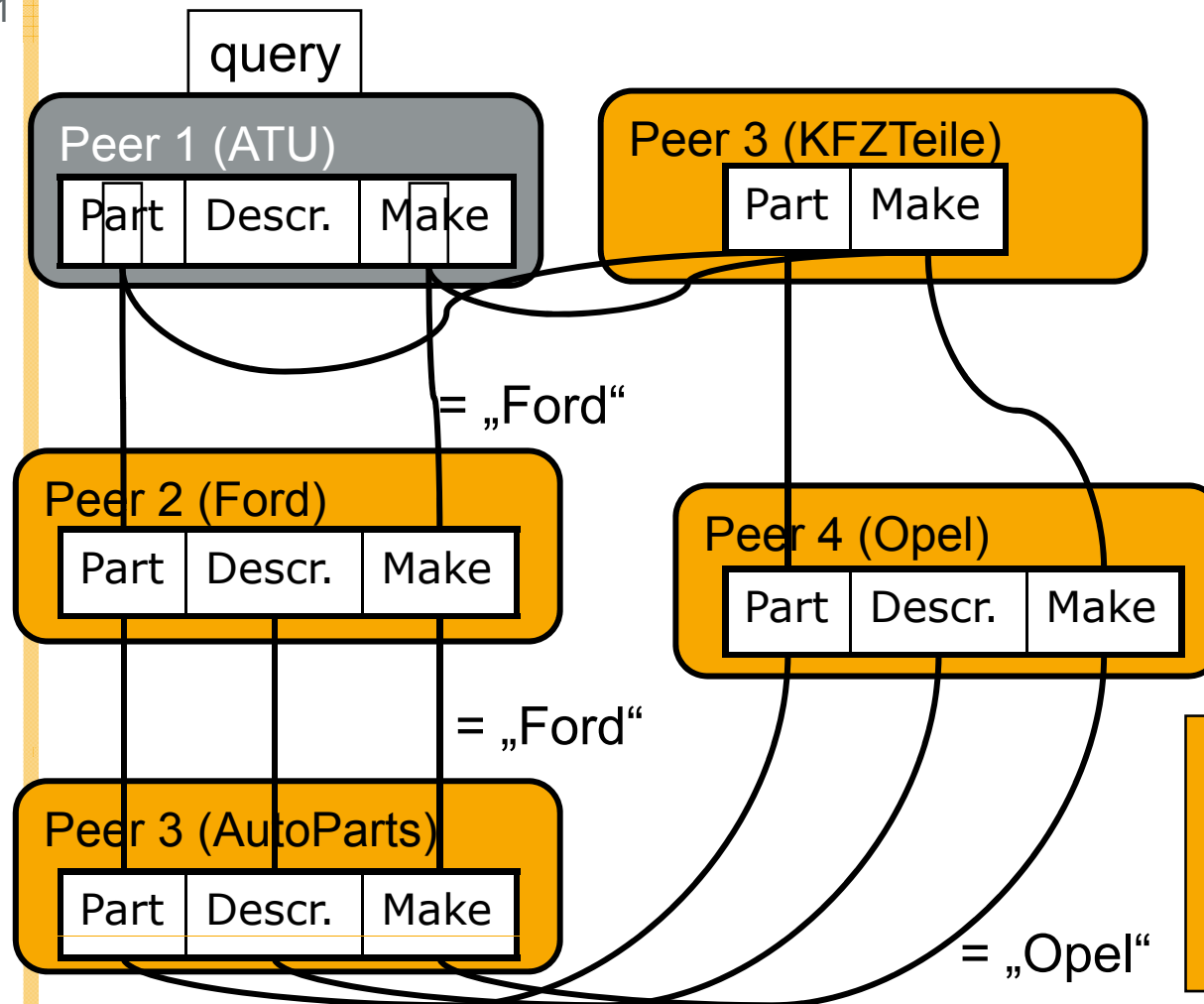
did	artist	title	category	genre	year	cdextra
<input type="checkbox"/> a810b60d	Various xes	Peppered Wit...	<input type="checkbox"/> reggae	<input type="checkbox"/> Electronic	<input type="checkbox"/> 2003	
<input checked="" type="checkbox"/> a810e20d	The Gladiators	Sold Out	<input checked="" type="checkbox"/> reggae			
<input type="checkbox"/> a90a270c	Jackie Brown	Look Pon You	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 1999	Joe Gibbs En...
<input checked="" type="checkbox"/> a90a6c0e	Byron Lee an...	Soca Engine	<input checked="" type="checkbox"/> reggae			
<input type="checkbox"/> a90b260c	Lee Perry	Dub Fire	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 2001	YEAR: 2001
<input checked="" type="checkbox"/> a90b340c	CHM	PNG XMAS 2...	<input checked="" type="checkbox"/> reggae			
<input checked="" type="checkbox"/> a90b5f0c	FEEL THE R...	FEEL THE R...	<input checked="" type="checkbox"/> reggae			
<input type="checkbox"/> a90cdb0c	Natural Zion ...	L'homme Co...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae		
<input checked="" type="checkbox"/> aa0a6c0d	Pressure Co...	I Want To Tell...	<input checked="" type="checkbox"/> reggae			
<input type="checkbox"/> aa0b590c	Best Of World...	Jamaïque	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 1994	
<input type="checkbox"/> aa0bfc0a	Lionel Richie	Live & Alive	<input type="checkbox"/> reggae	<input type="checkbox"/> Soul		
<input type="checkbox"/> aa0ceb0d	Latin Prince	Multiple Vibes	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 1993	
<input type="checkbox"/> aa0d200b	Various Artists	Rite Sound R...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae		
<input type="checkbox"/> aa0d5c0d	Charlie Chap...	Kings of Reg...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 2002	
<input type="checkbox"/> aa0de80c	The Dead Bill...	Heartfelt Ses...	<input type="checkbox"/> reggae	<input type="checkbox"/> Rock	<input type="checkbox"/> 1999	
<input checked="" type="checkbox"/> aa0f2a0c	Bob Marley	1976-04-23 -	<input checked="" type="checkbox"/> reggae			
<input type="checkbox"/> aa0fce0e	Baha Men	Who Let the ...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 2000	YEAR: 2000
<input type="checkbox"/> aa10550d	SES	Surprise	<input type="checkbox"/> reggae	<input type="checkbox"/> reggae	<input type="checkbox"/> 2001	
<input type="checkbox"/> aa11d00d	Ras	Rhythmic Alte...	<input type="checkbox"/> reggae	<input type="checkbox"/> Slow Jam		
<input type="checkbox"/> ab08280e	Siemens	Caribbean S...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 2001	YEAR: 2001
<input type="checkbox"/> ab09ca0d	La Factoria	Total Trance	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae		
<input checked="" type="checkbox"/> ab09da0e	King Tubby & ...	Foundation O...	<input checked="" type="checkbox"/> reggae			Produced by ...
<input checked="" type="checkbox"/> ab09e20c	Tappa Zukie	From the Arc...	<input checked="" type="checkbox"/> reggae			©1995 RAS ...
<input type="checkbox"/> ab0b490c	Jimmy Cliff	Higher & Hig...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 1996	YEAR: 1996 L...
<input type="checkbox"/> ab0d470c	Showtime	Snowtime - J...	<input type="checkbox"/> reggae	<input type="checkbox"/> Reggae	<input type="checkbox"/> 1998	Showtime Ju...

to filter data: category reggae

more Q-values

PDMS: Incomplete and Selective Mappings

81



Problem:
Cumulated selections

- implicit in schemata
- explicit in mappings
- Point selections and range selections

Problem:
Cumulated projections

- in schemata
- in mappings