



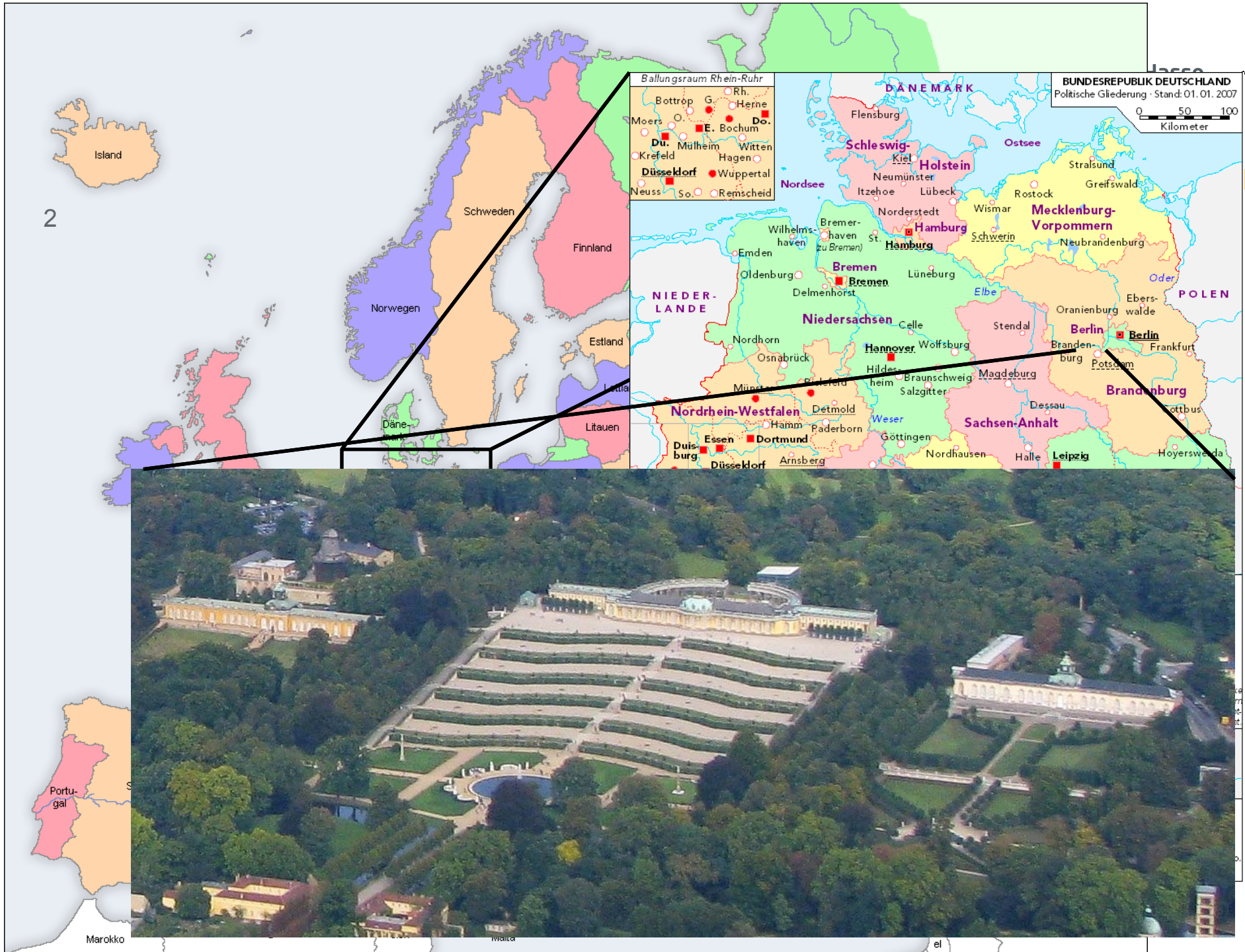
**Hasso
Plattner
Institut**

IT Systems Engineering | Universität Potsdam

Data Integration in Two Steps and From Above

Palo Alto, February 3, 2009

Felix Naumann



The HPI – Hasso Plattner Institut

3

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



Information systems team

4

project **ViQTOR**



Paul Führung



Pat Hobro

DQ Assessment



Prof. Felix Naumann

Information Integration

Information Quality



Jens Bleiholder

Data Fusion

project **HumMer**



Karsten Draba

Data Profiling & Cleaning



Christoph Böhm

Peer Data Management Systems



Armin Roth

project **System P**

Service-Oriented Systems

Matching

Data Integration for Life Science Data Sources

project **Aladin**



Alexander Albrecht

Personal Information Management



Mohammed AbuJarour

Ontologies



Frank Kaufer

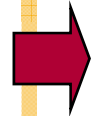


Jana Bauckmann

Data Profiling for Schema Management

Overview

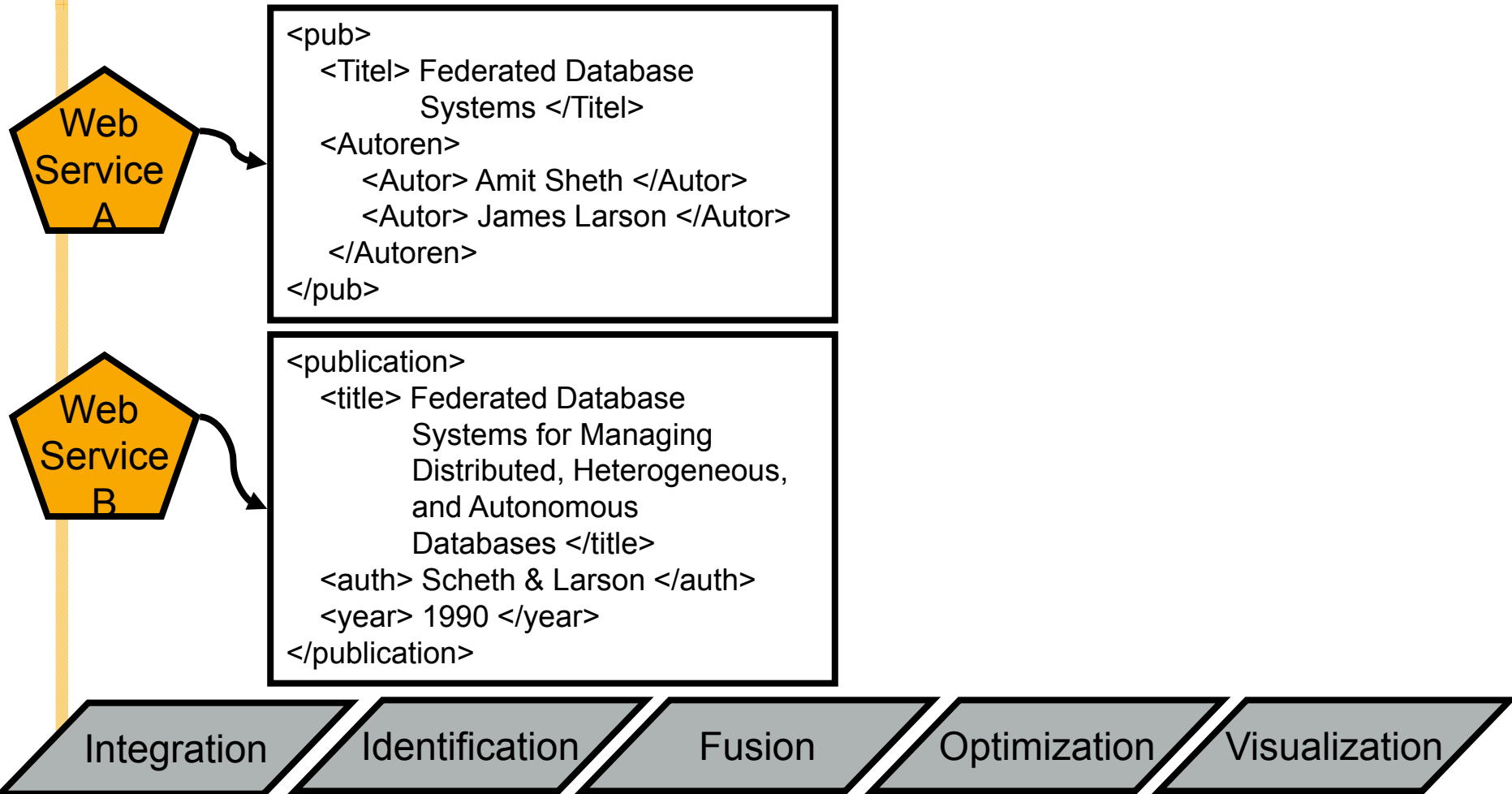
7



- Introductory example
- Step 1: Schema level integration
 - Schema Mapping
 - Schema Matching
- Step 2: Data level integration
 - Duplicate detection
 - Data fusion
- From above: ETL management

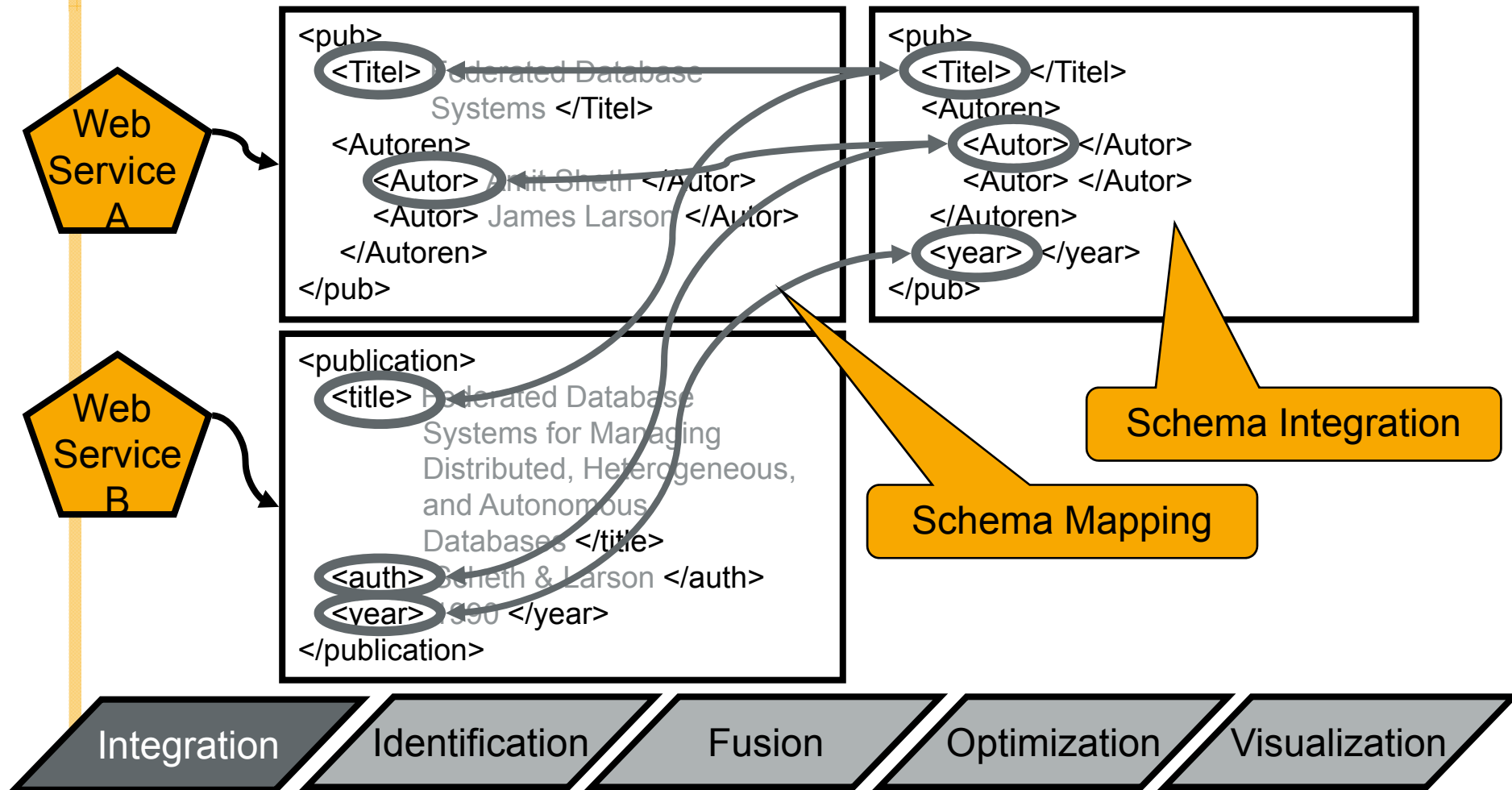
Information Integration

8



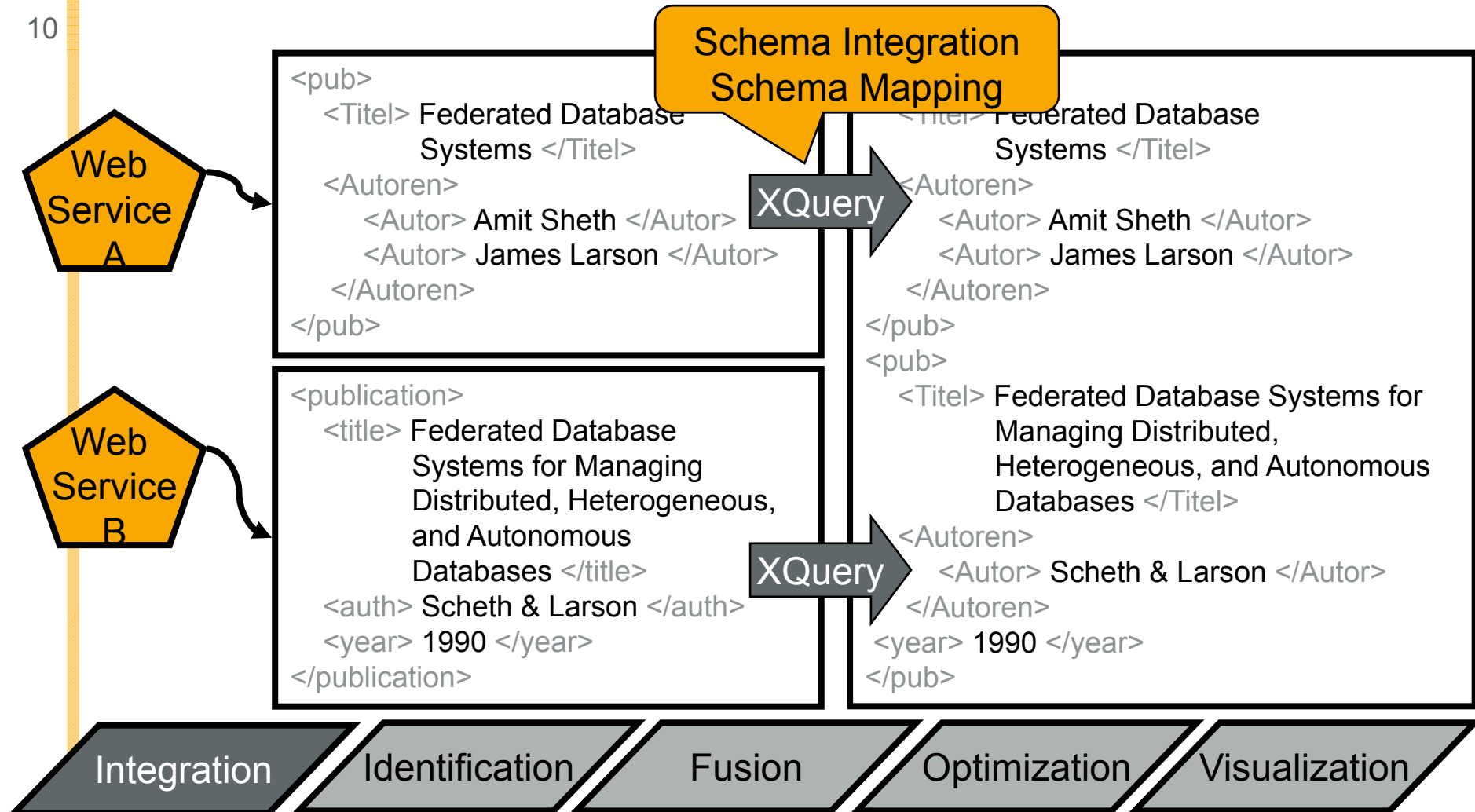
Information Integration

9



Information Integration

10



Information Integration

11

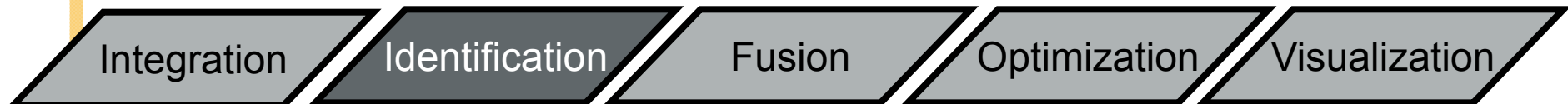


```
<pub>
  <Titel> Federated Database
    Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```



```
<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

12

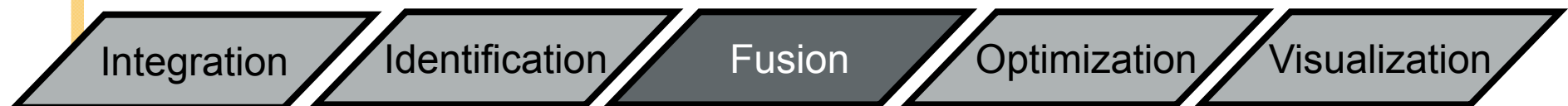


```
<pub>
  <Titel> Federated Database
    Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```



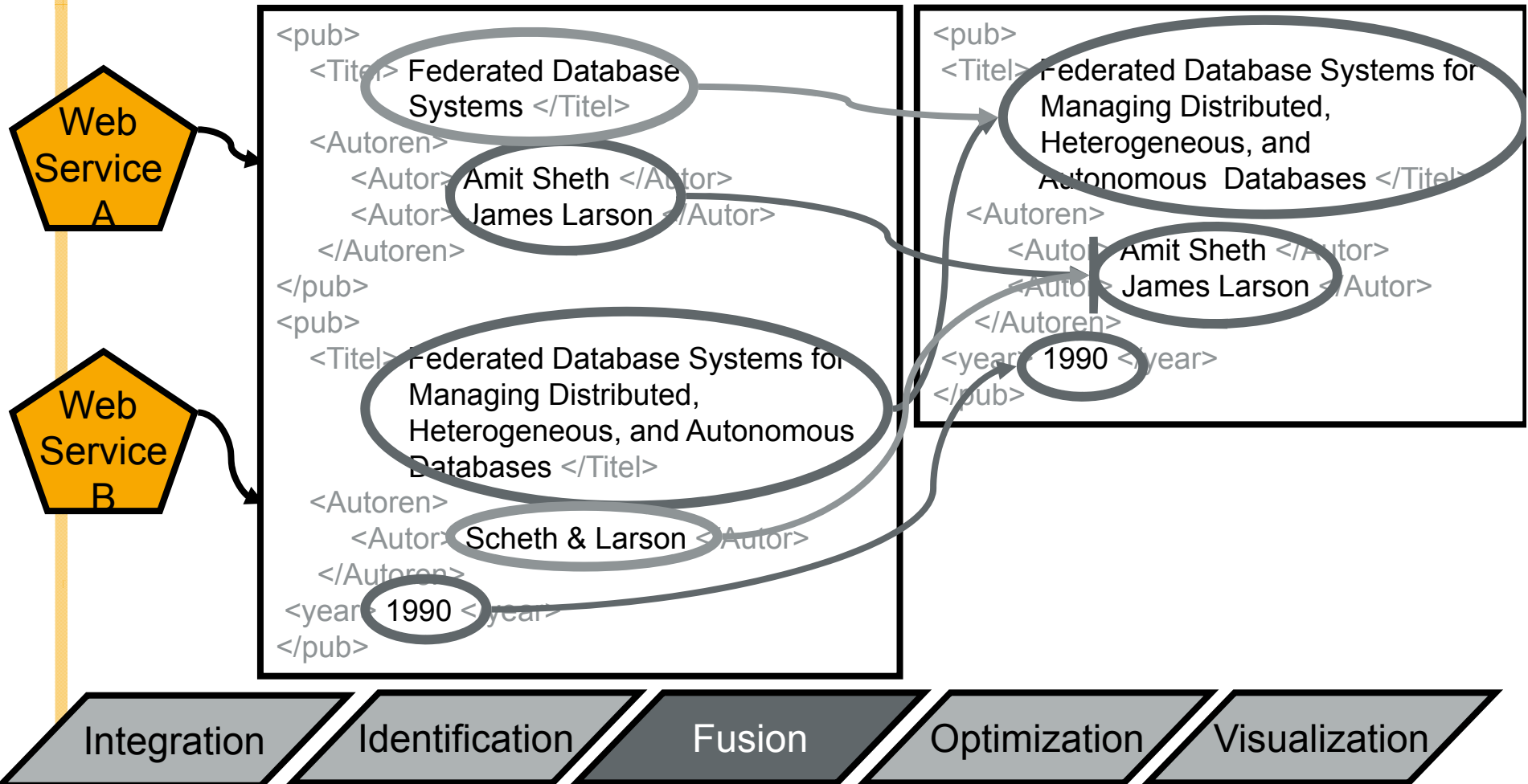
```
<publication>
  <title> Federated Database
    Systems for Managing
    Distributed, Heterogeneous,
    and Autonomous
    Databases </title>
  <auth> Scheth & Larson </auth>
  <year> 1990 </year>
</publication>
```

```
<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>
```



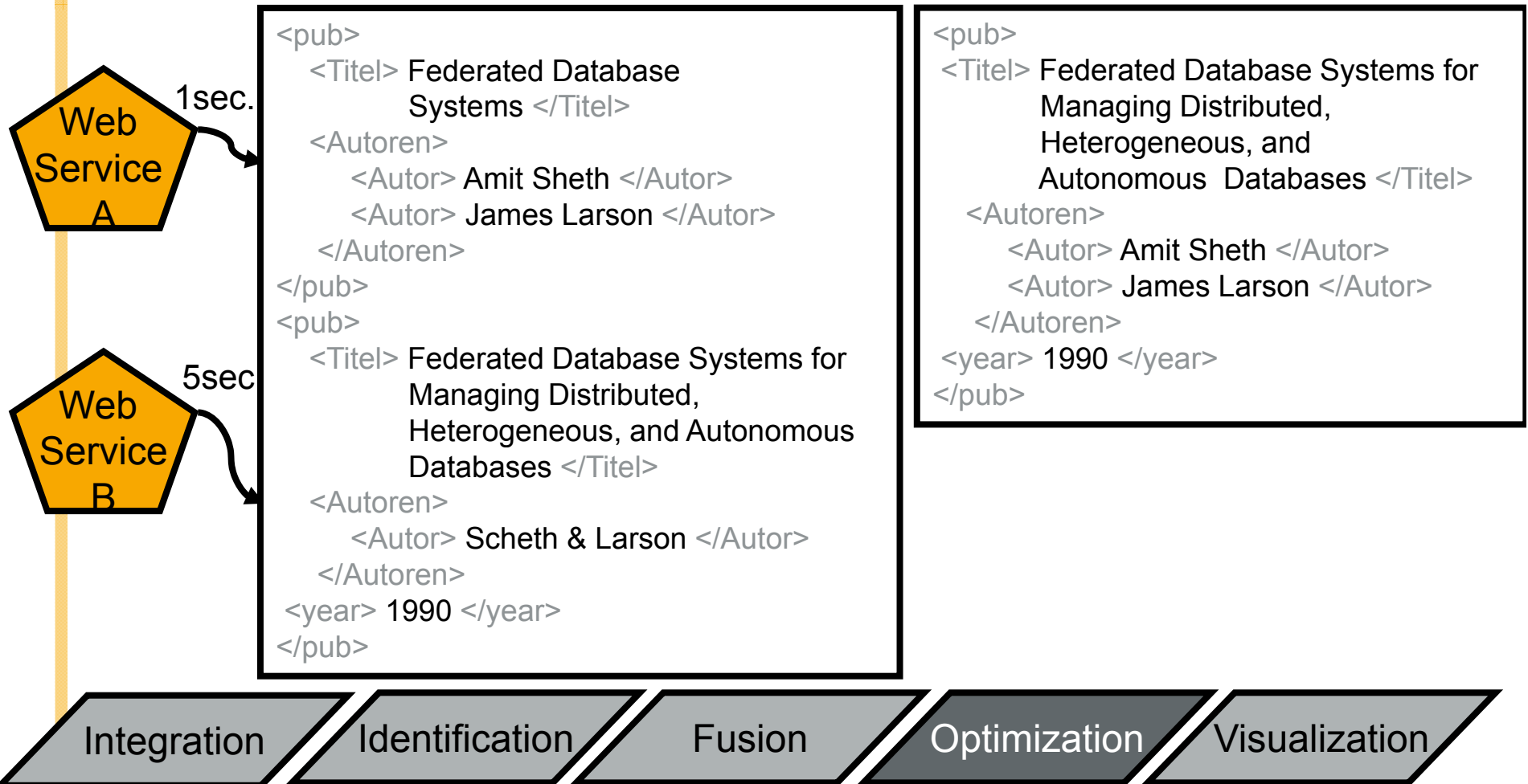
Information Integration

13



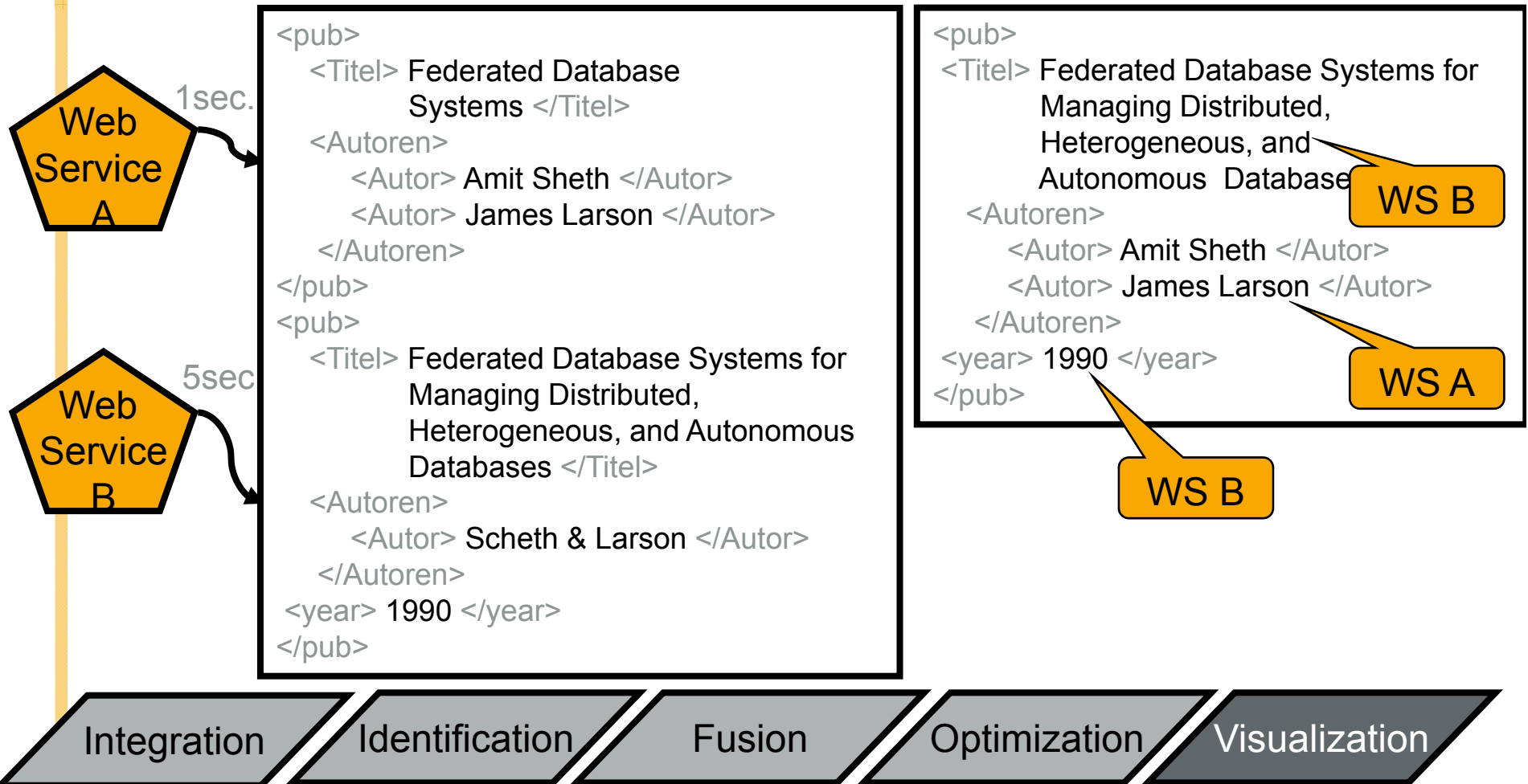
Information Integration

14



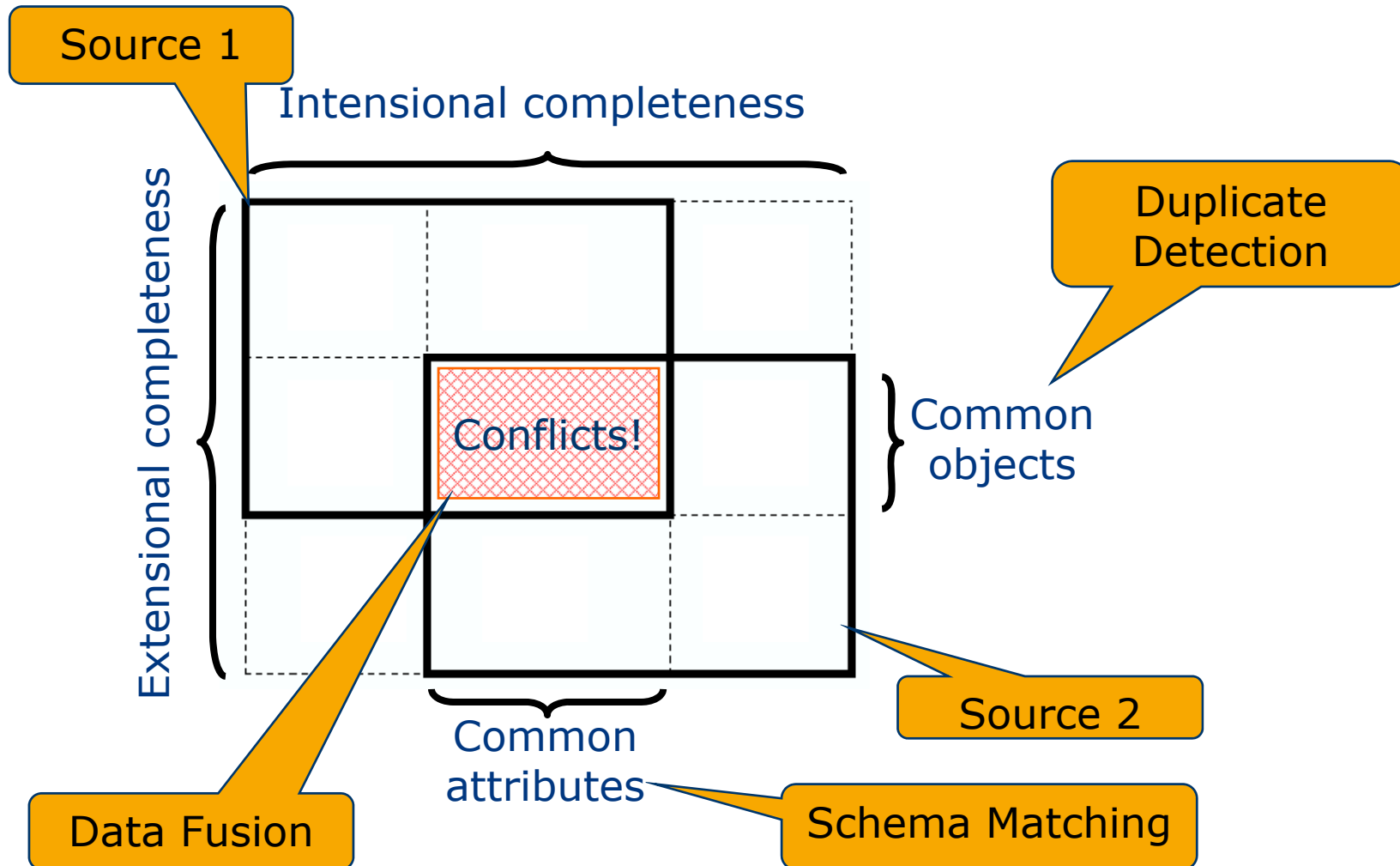
Information Integration

15



Completeness and Conciseness

16



Overview

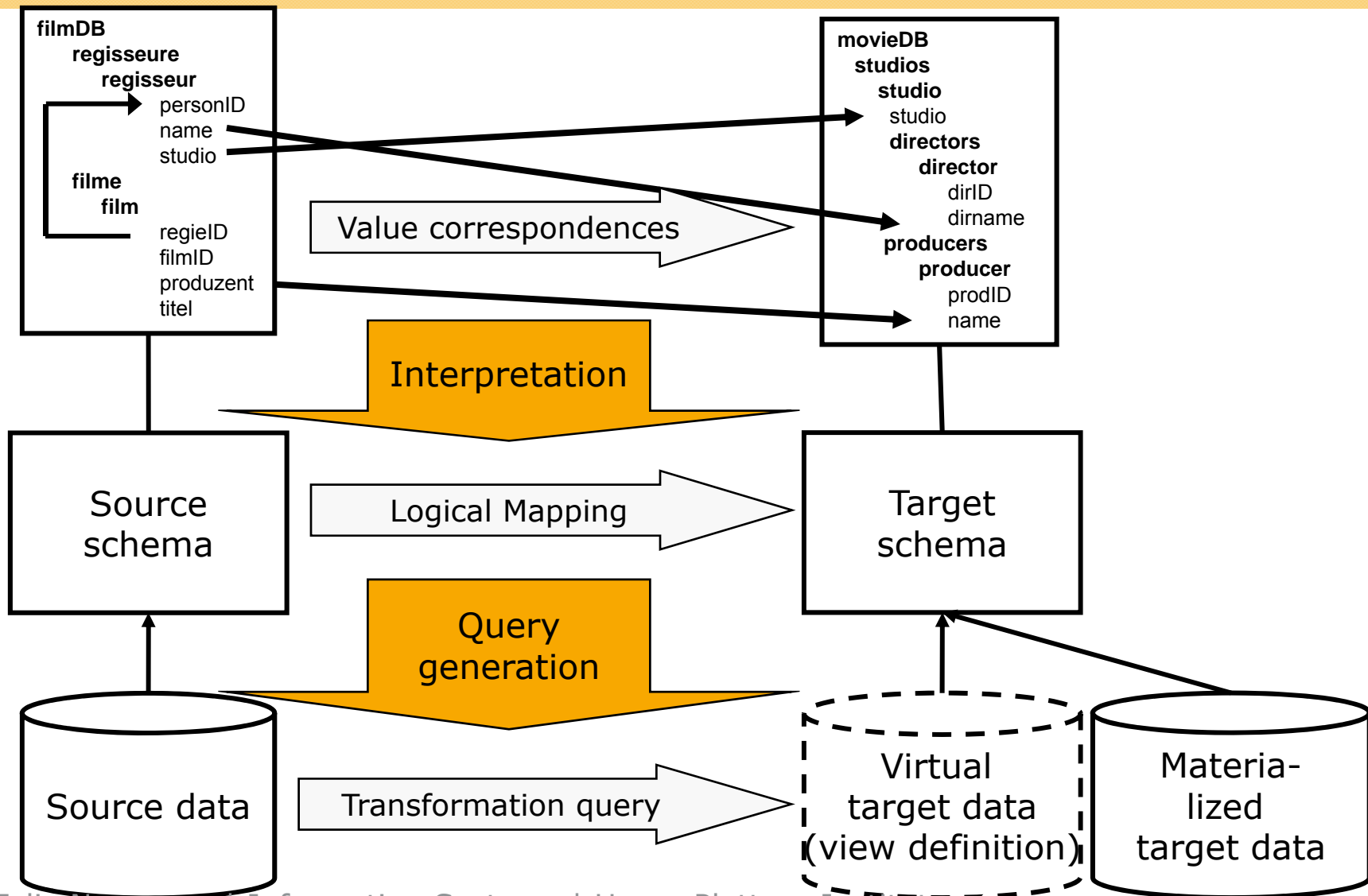
17

- Introductory example
- Schema level integration
 - Schema Mapping
 - Schema Matching
- Data level integration
 - Duplicate detection
 - Data fusion
- ETL management



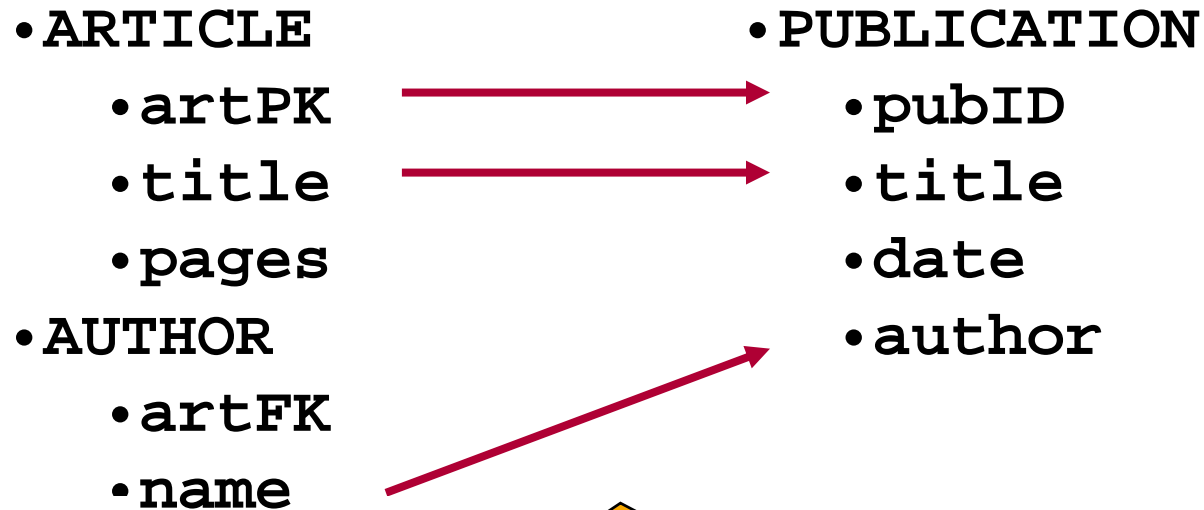
Schema Mapping

18



Schema Mapping Example

19

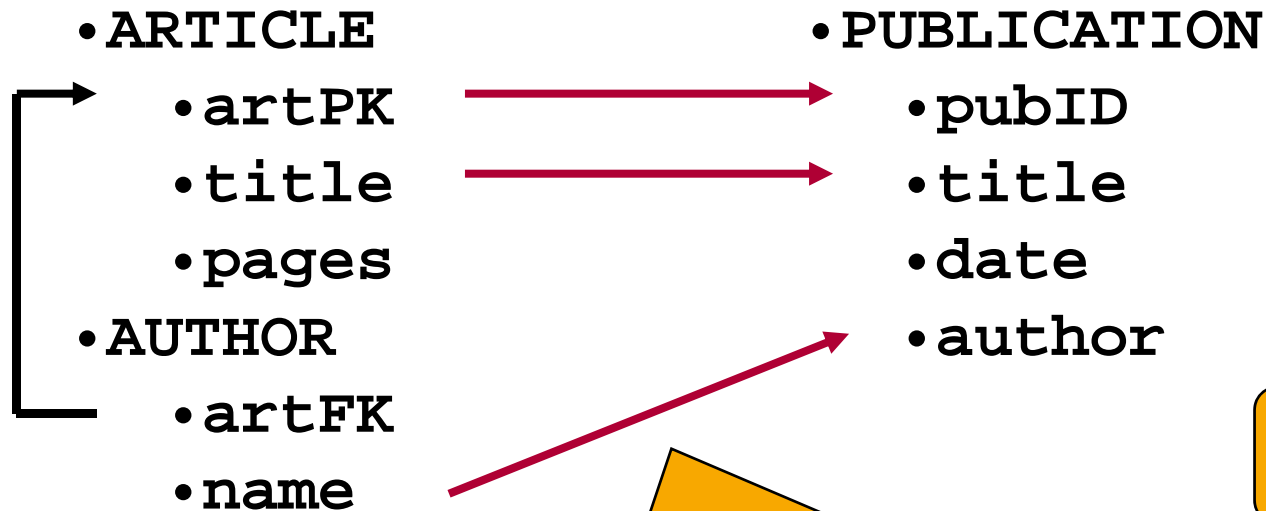


```

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

20



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

21

Schemata are

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

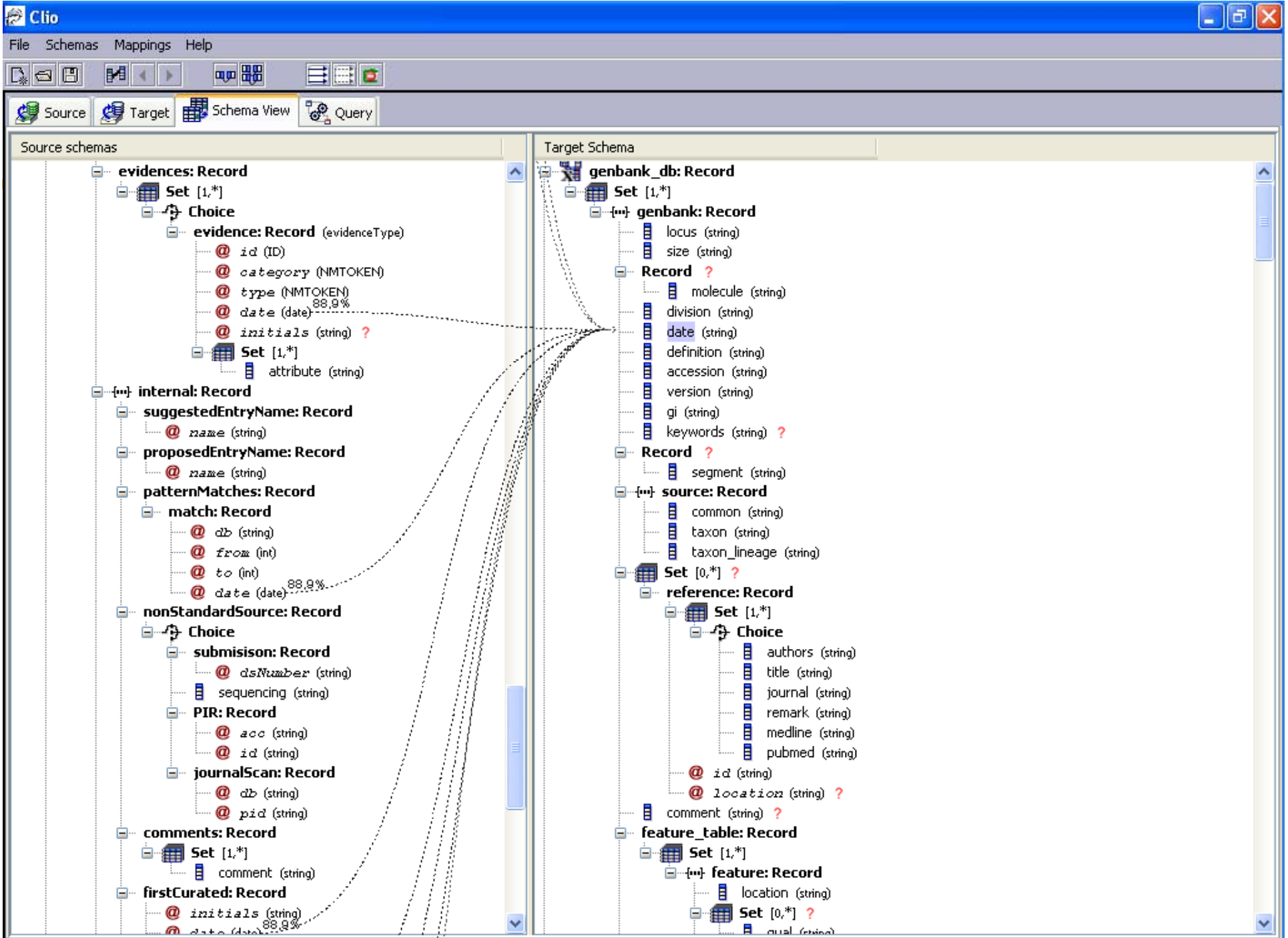
> 100 tables, many attributes

Deep nesting
Foreign keys
XML Schema

Unknown synonyms

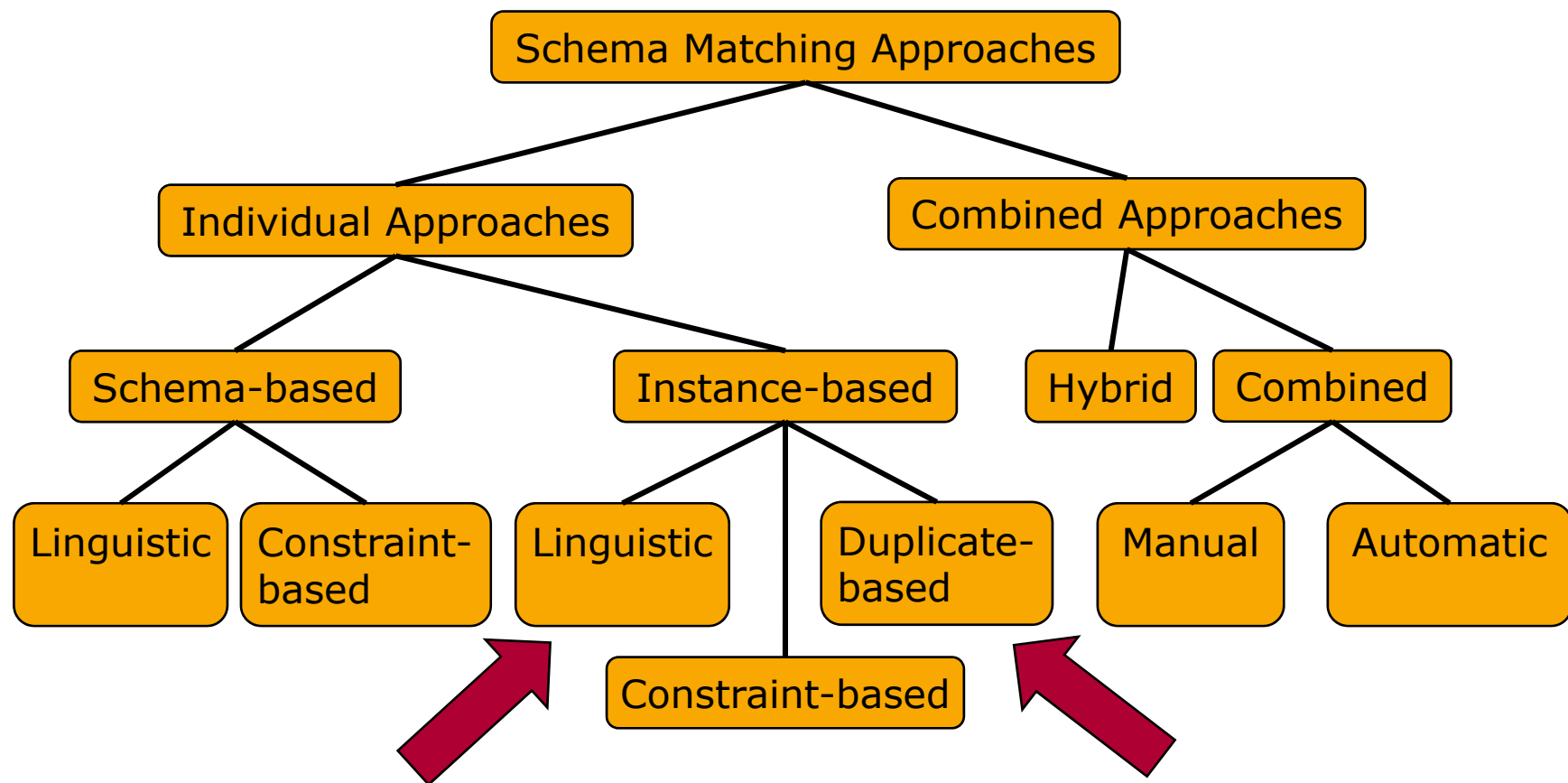
Unknown homonyms

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



Schema Matching Classification [RB01]

23



Instance-based Schema Matching

24

Instance-based Schema Matching:

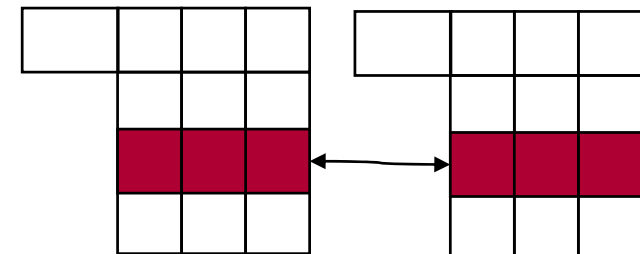
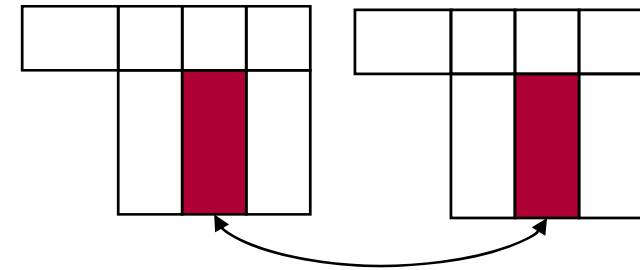
- Correspondences based on similar data values or their properties

Conventional solution: Vertical

- Comparison of columns
- = Attribute classification

Our solution: Horizontal

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

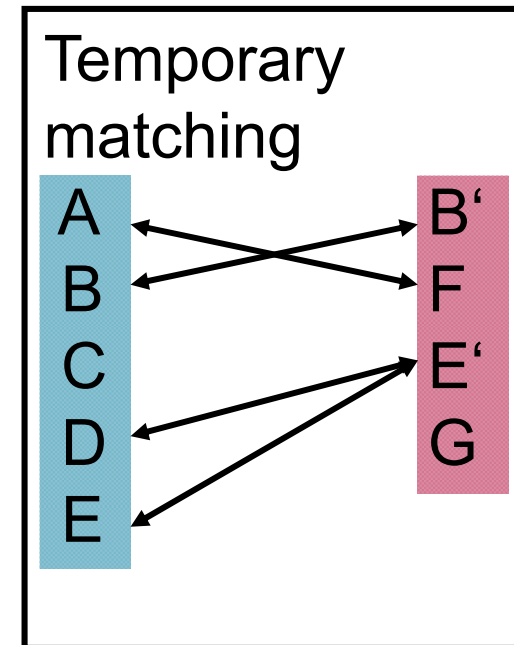
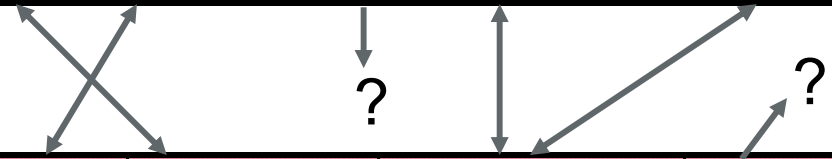


DUMAS Matcher

25

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

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

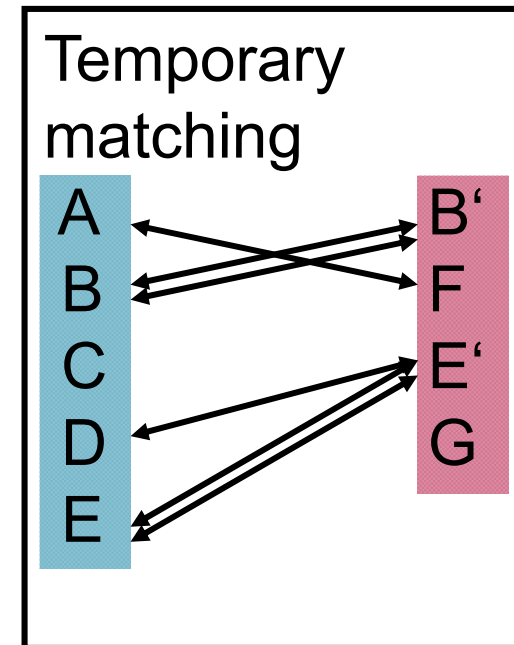
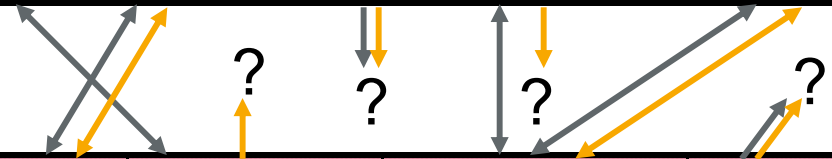


DUMAS Matcher

26

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 – open problems

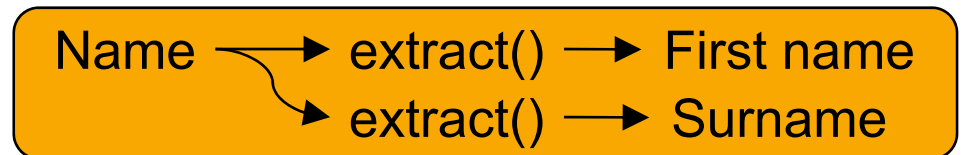
27

- n:1 und 1:n matches
 - Many combinations
 - Many functions
 - Parsing
- Matching in complex schemata
 - Find mapping, not only correspondences
 - Unions and joins

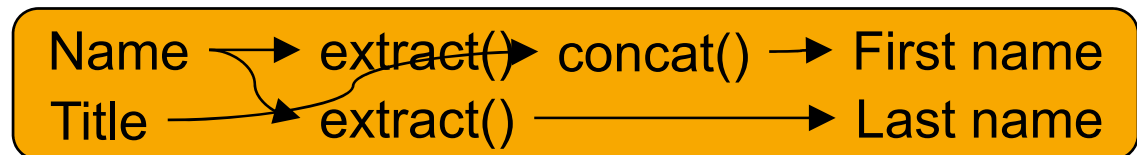
n:1 Matching



1:n Matching



m:n matching



- Tooling: User in the loop!

Overview

28

- Introductory example
- Schema level integration
 - Schema Mapping
 - Schema Matching
- Data level integration
 - Duplicate detection
 - Data fusion
- ETL management



Duplicate Detection

29

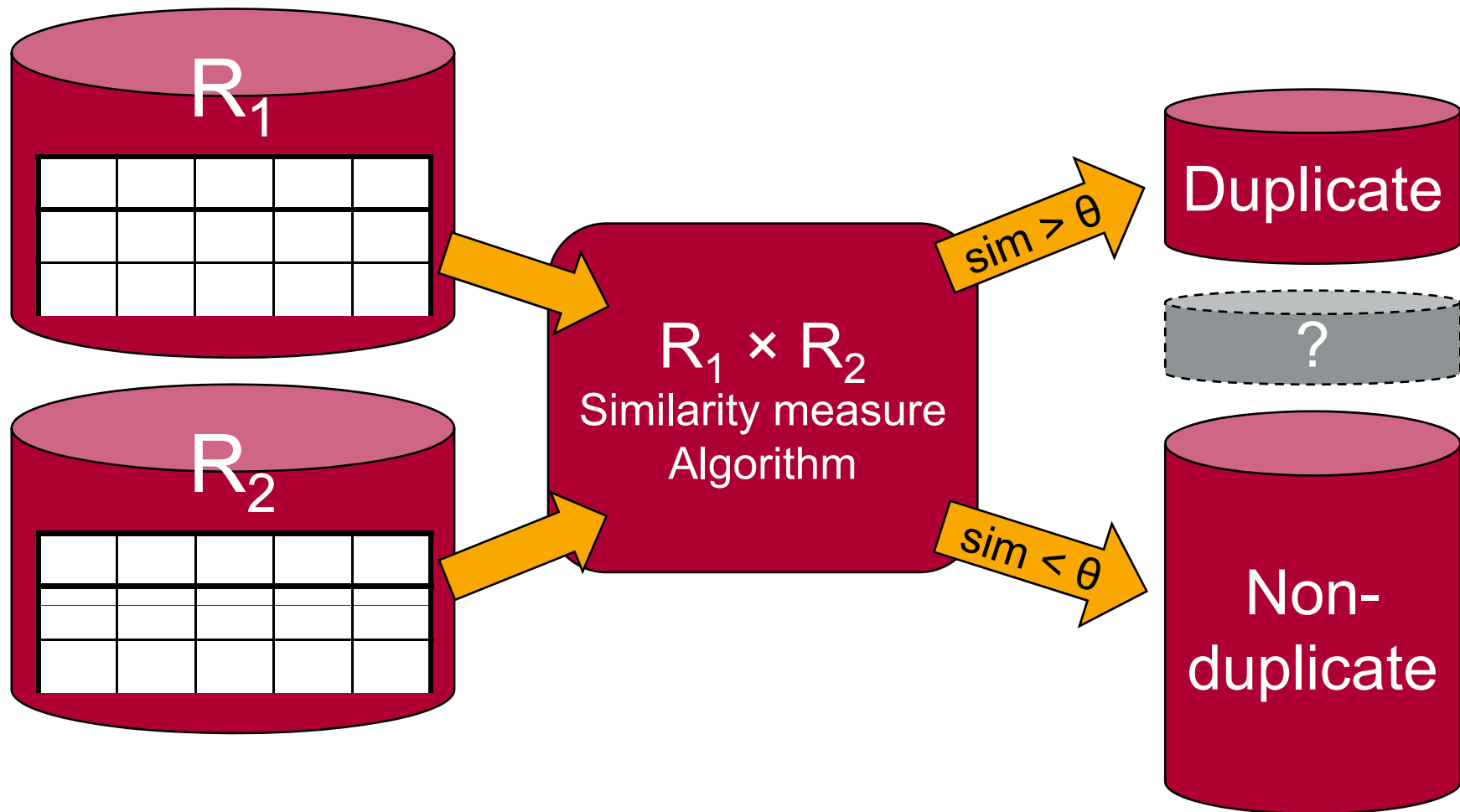
Duplicate detection is the discovery of multiple representations of the same real-world object.

- 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

30



Motivation

31

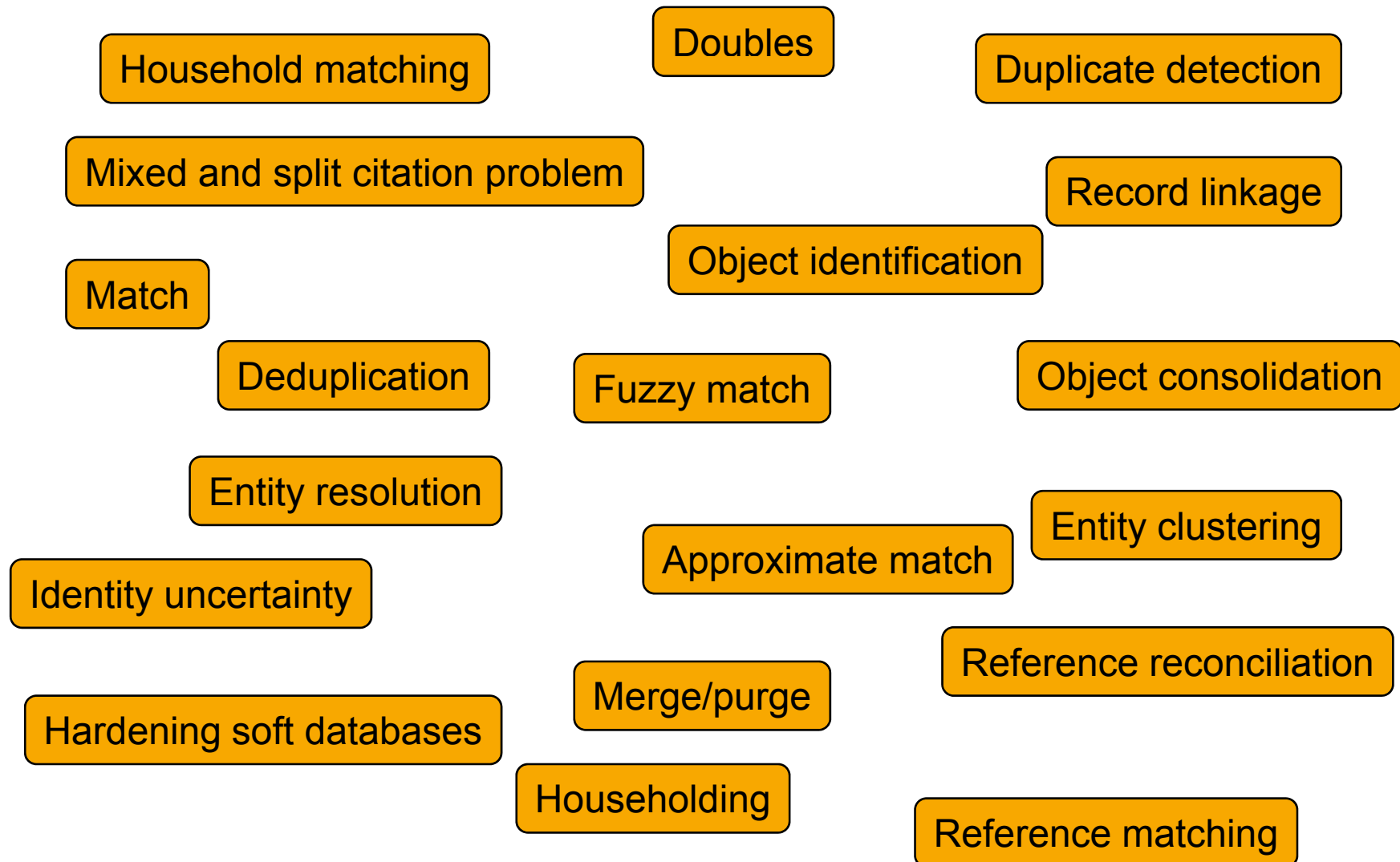
- Possible effects
 - Example: Portfolio Management Offers
 - Credit maximum not detected
 - Too low inventory levels
 - No quantity discount for multiple orders
 - Total revenue of preferred customers unknown
 - Multiple mailings of same catalog to same household
- General problems
 - Additional, unnecessary IT expenses
 - Low customer satisfaction
 - Potentials and dangers not detected
 - Poor quality financial data

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



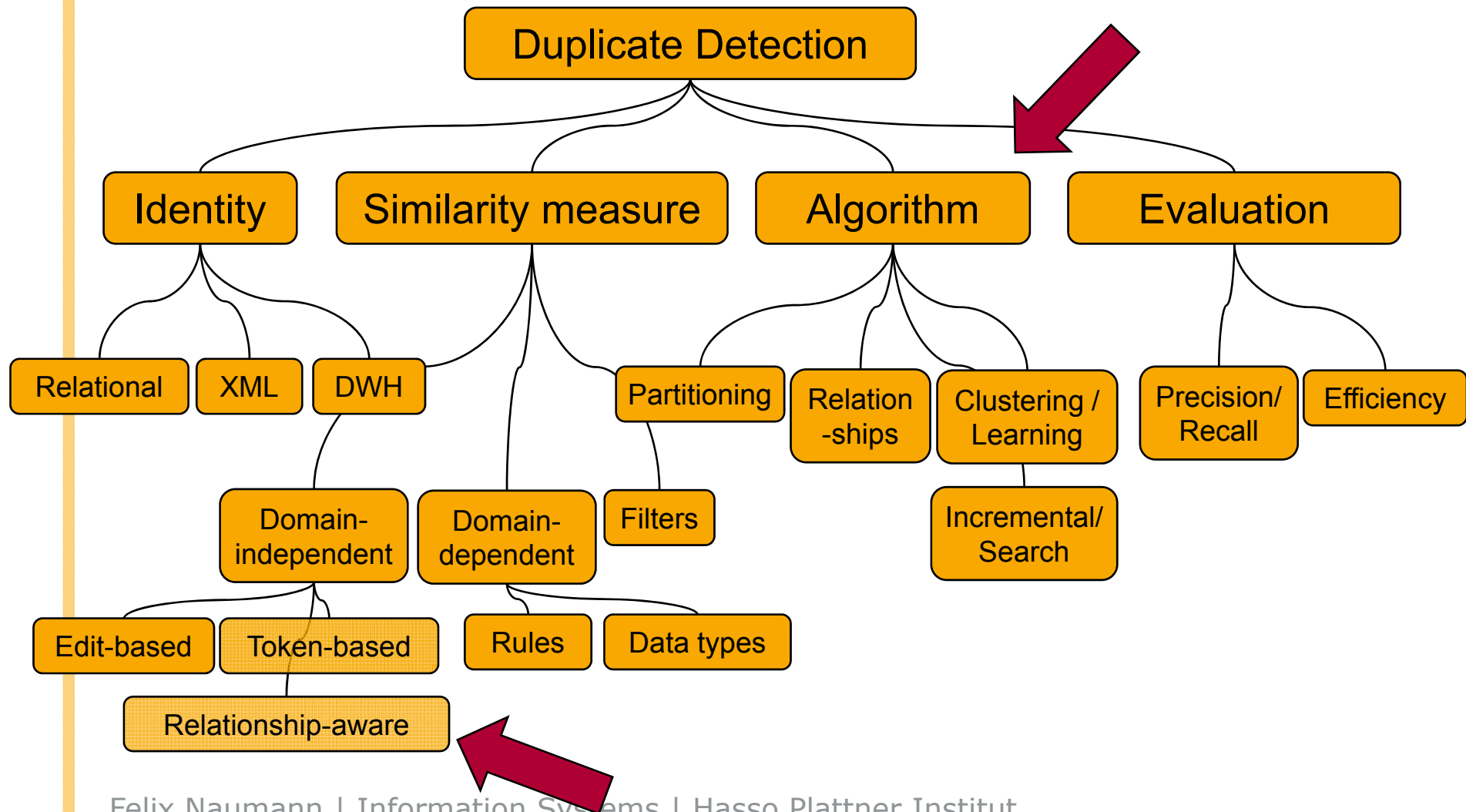
Ironically, “Duplicate Detection” has many Duplicates

32



Duplicate Detection – Research

33

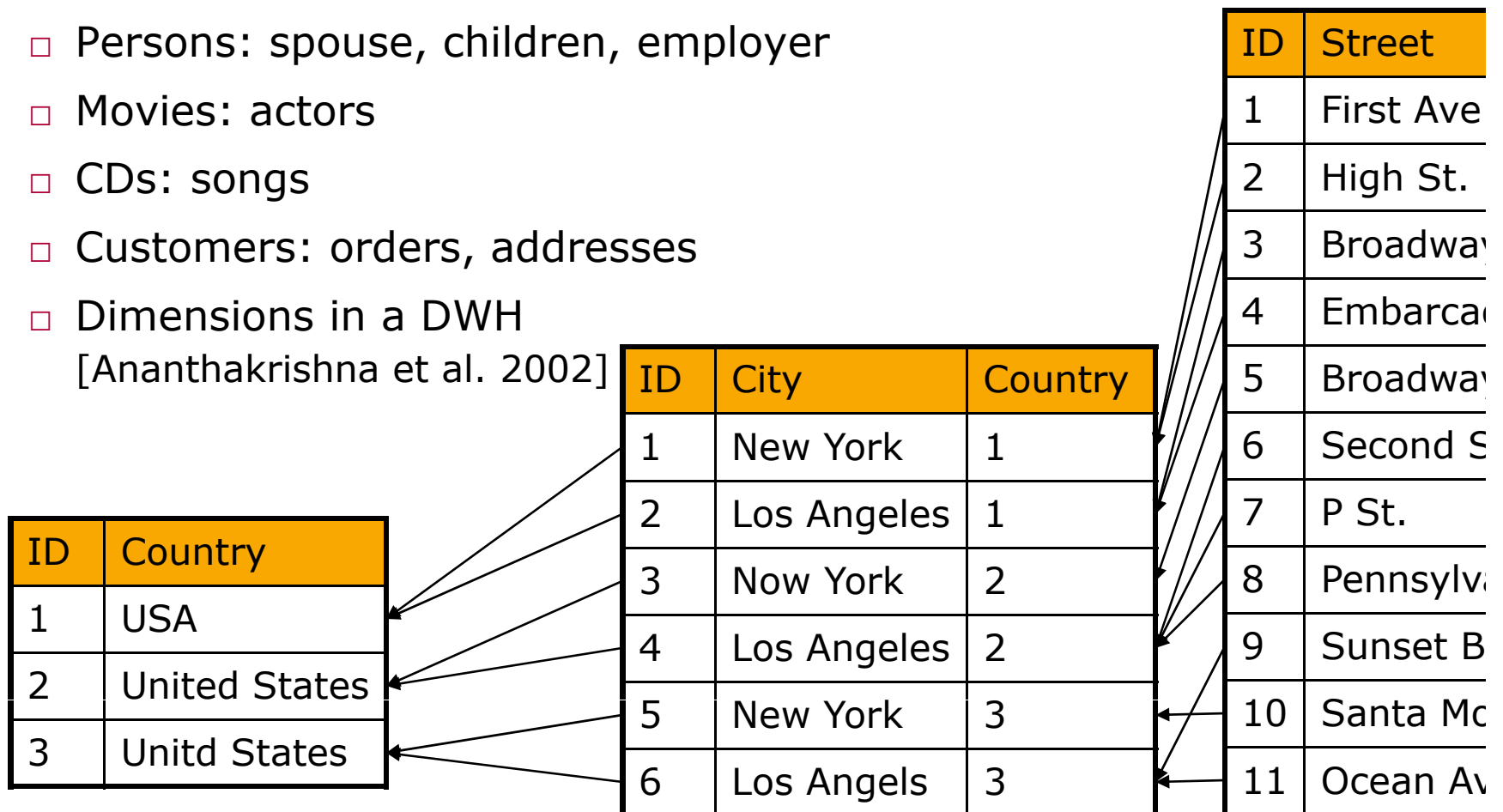


Relationship-aware Similarity Measures

34

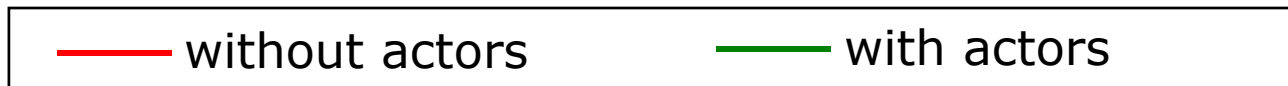
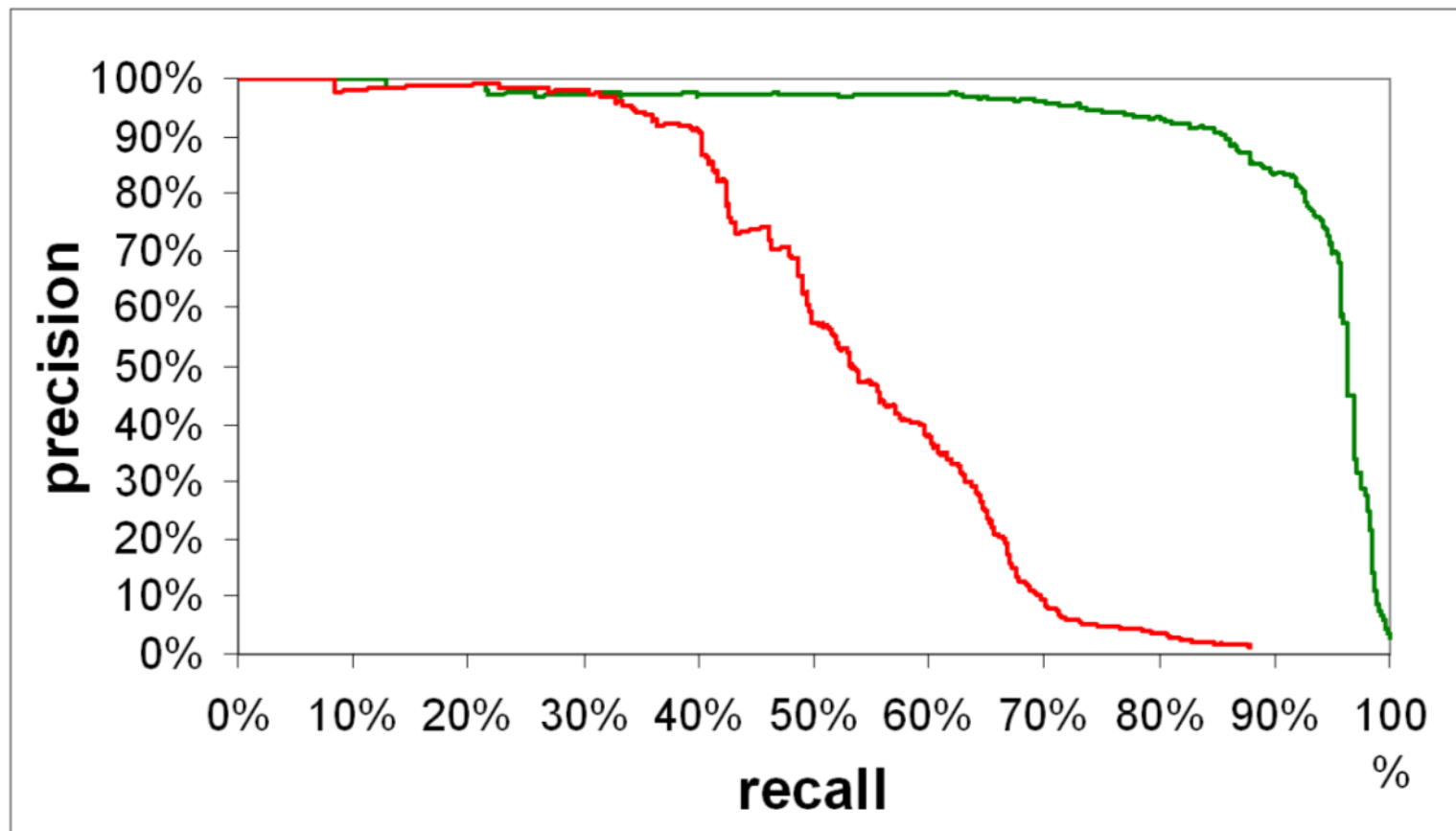
- Idea: Not only values of the records, but values of related records are relevant for similarity.

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



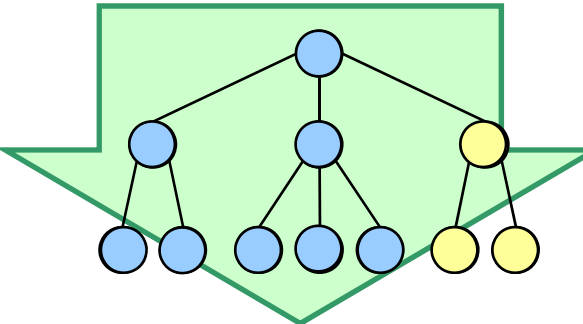
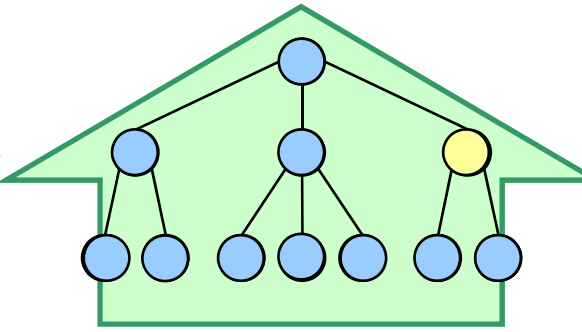
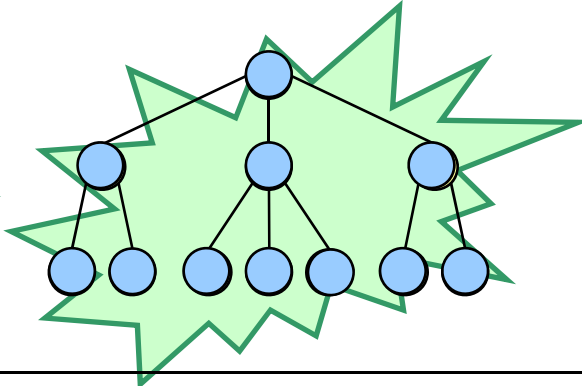
Relationship-aware Similarity Measures – Evaluation

35



Iterative RADD

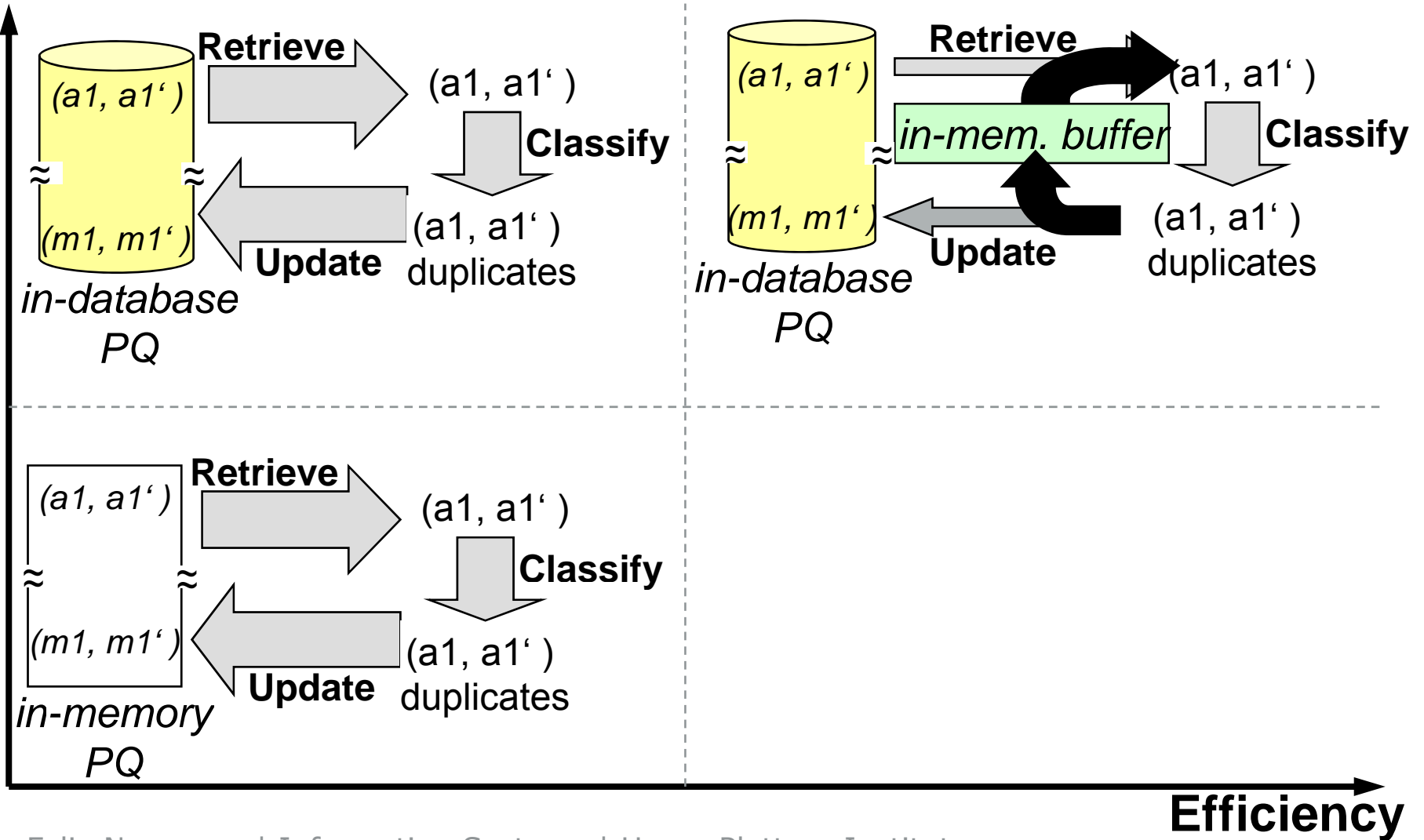
36

Top-down [SIGMOD'05]	Bottom-up [EDBT'06]	From-the-middle [ICDE'06]
		
Effectiveness ★ Efficiency ★★	Effectiveness ★★ Efficiency ★★	Effectiveness ★★★ Efficiency ★
Further techniques: Object filter; Edit distance filter; Transitivity		

Scalability of the generalized algorithm

37

Scalability

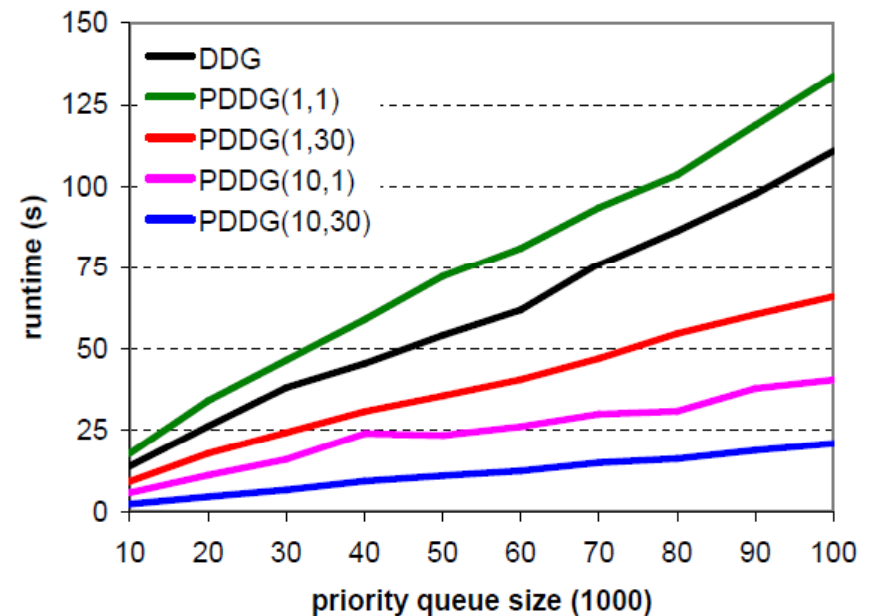


Scalability – Comparison

38

Approach	# candidates	Connectivity	Classif. Time (s)
Singla [27]	1,295	n.a.	n.a.
Dong05 [12]	40,516	13.4	n.a.
RC-ER [2]	45,000/65,000	1.9 / 5.3	100 / 890
RC-ER [4]	97,270	1.9	543 - 690
ReIDC [19]	75,000	low - high	180 - 13,000
LinkClus [31]	100,000	10	900
RECUS/BUFF	1,000,000	1.7	24,433 (7h)

Parallelization increases performance:

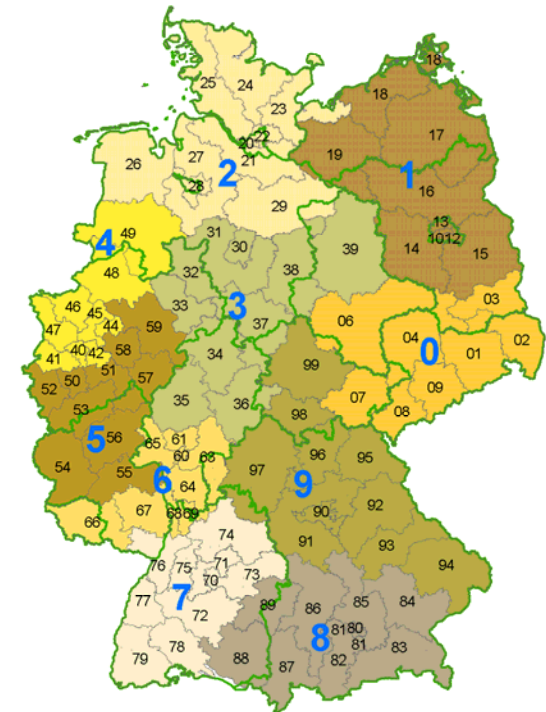


Partitioning / Blocking

39

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

- Partitioning by first two zip-digits
- or partition by first letter of surname
- or ...



Source: wikipedia.de

- Idea: Partition multiple times by different criteria.

- Then apply transitive closure on discovered duplicates.

Sorted Neighborhood

[Hernandez Stolfo 1998]

40

- 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

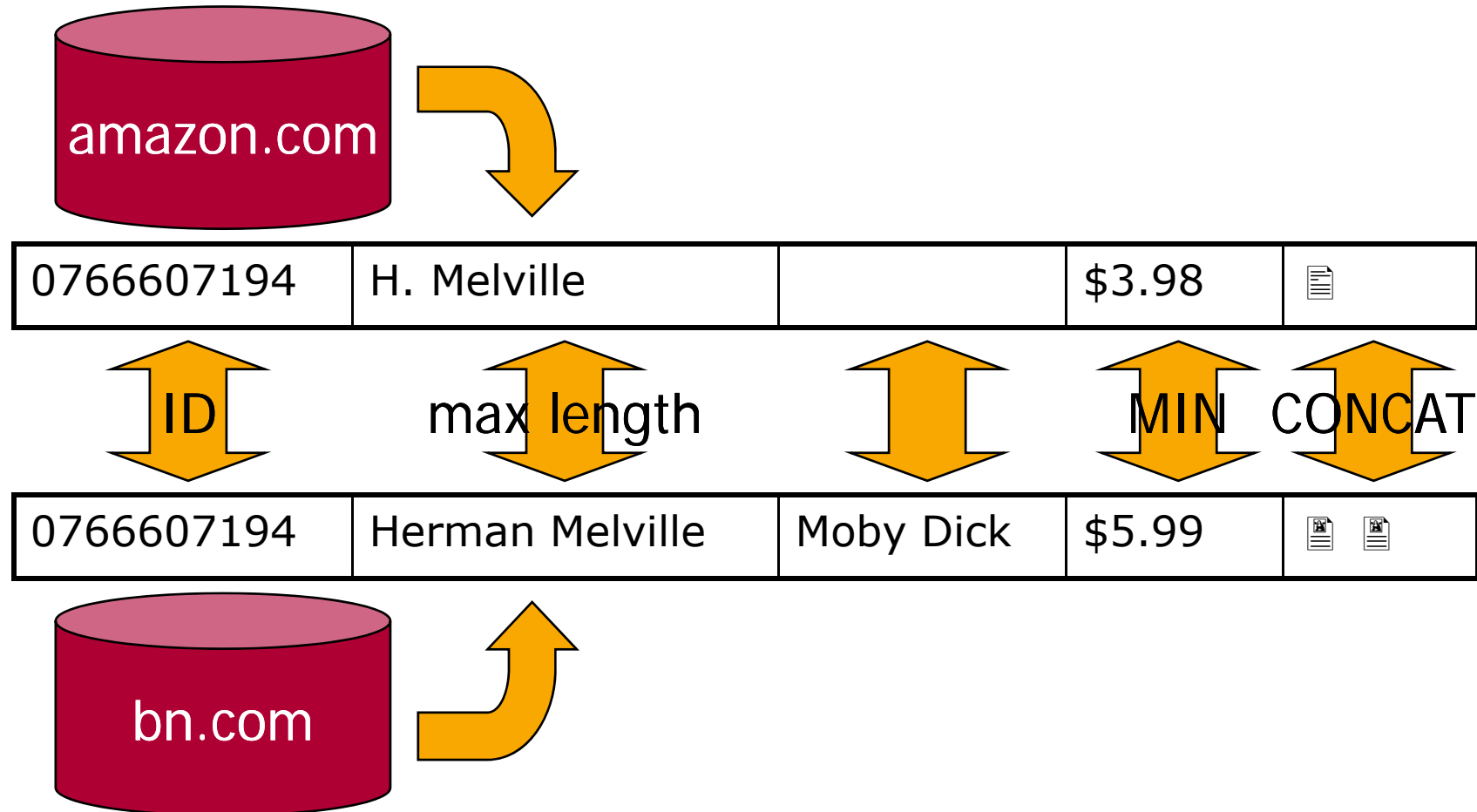
42

- Introductory example
- Schema level integration
 - Schema Mapping
 - Schema Matching
- Data level integration
 - Duplicate detection
 - Data fusion
- ETL management



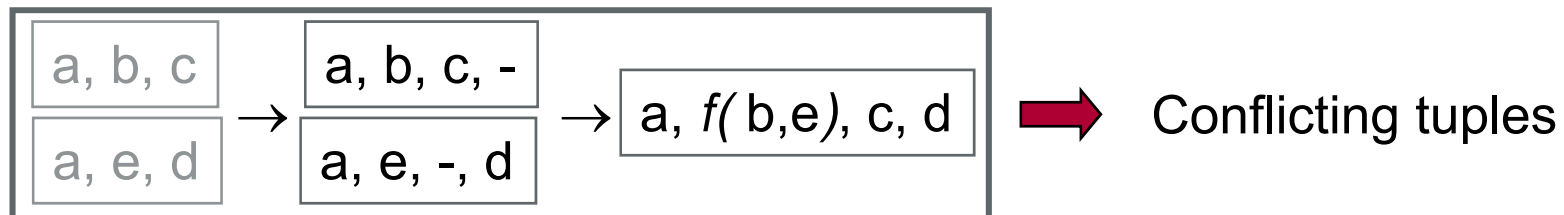
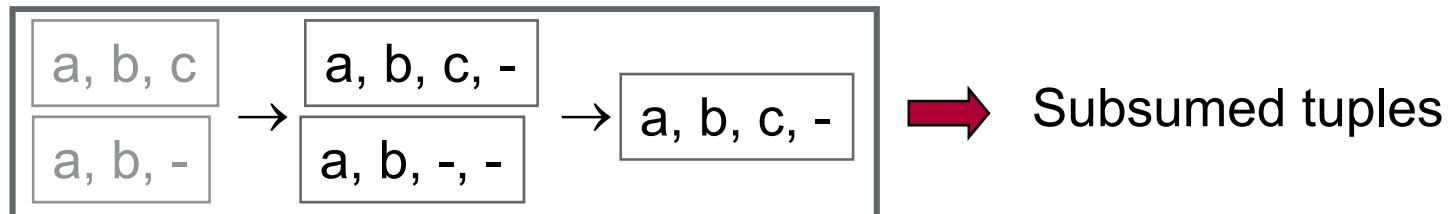
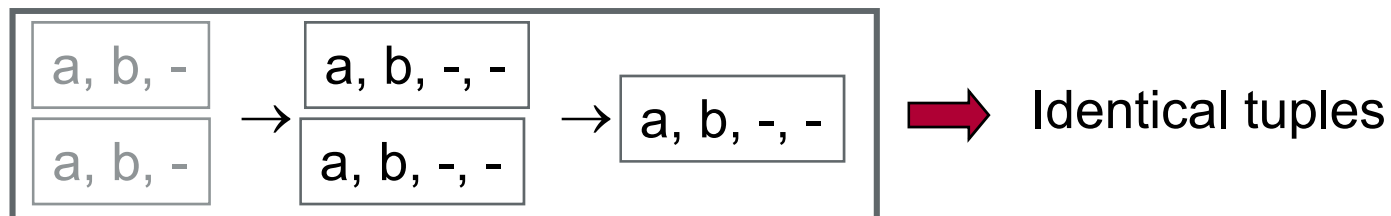
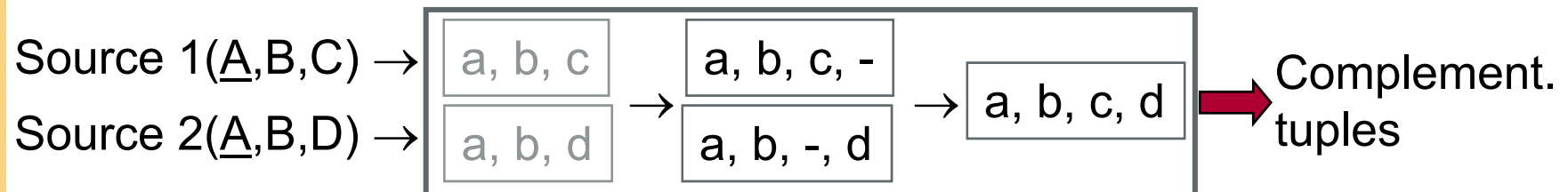
Data Fusion

43



“Proper” Data Fusion

44



Conflict Resolution Functions

45

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
....

Visualization of Integrated Data

46

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 VOTE	VERSI... COALESCE	COUN... COALES...	YEAR MAX	ORIGI... COALES...	GENRE LAST	DIREC. 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	Gadjo... 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

Felix

Tool-based Data Fusion

47

Fuzzy Fuzion
Additional Information Test/Debug

Gruppen 0 bis 50 von 39449 Filtermodus Wert:

fdb.gr...	TITLE	SALUT...	FIRST...	LASTN...	COMP...	COUN...	STREET	STREE...	ZIP	CITY	ADR1	ADR2	ADR3	ADR4	ADR5
1253		Frau u...		Koste...	Daimle...	D	Alt-Mo...	96 A	10559	Berlin	Frau u...	Kosten...	Daimle...	Alt-Mo...	
1333		Herr	Markus	Bauer		D	HPC V...		10878	Berlin	Herr	Marku...	HPC V...		D - 10...
1782		Herr	Frank	Leusc...		D	Arenh...		12103	Berlin	Herr	Frank ...	Arenh...		D - 12...
1874		Monsieur	Frank	Eichler		D	Falken...	78A	13589	Berlin	Monsieur	Frank ...	Falken...		D - 13...
2159		Herr	Horst	Fucks		D	Nassa...		10717	Berlin	Herr	Horst ...	Nassa...		D - 10...
2196	Dr.	Frau	Christa	Schün...		D	Uhlan...	121	10717	Berlin	Frau	Dr. Ch...	Uhlan...		D - 10...
2217	Dr.	Familie		Hofma...		D	Hede...	13	10969	Berlin	Familie	Dr. H...	Hede...		D - 10...
2498		Frau	Julia	Görsc...		D	Regin...	20	13409	Berlin	Frau	Julia G...	Regin...	13409 ...	
2552		Herr		Ehlers		D	Über F...		10587	Berlin	Herr	Ehlers	Über F...	10587 ...	

10. Gruppe :

fdb.group	TITLE	SALUTATION	FIRSTNAME	LASTNAME	COMPAN...	COUNTRY_CODE	STREET	STREET_NUMBER	ZIP	
1874		Monsieur	Frank	Eichler		D	Falkenseer Chaussee	78A	13589	Berlin
1874		Firma		Eichler		D	Fkenseer Chaussee 78A		13589	Berlin
1874		Herr	Frank	Eichler	Zres	D	Falkenseer Chaussee 78 A		13589	Berlin

..

1874		Monsieur	Frank	Eichler	Zres			78A	13589	Berlin
------	--	----------	-------	---------	------	--	--	-----	-------	--------

- Längster
- Längste gemei
- Supersequenz
- Verkettung
- Minimum
- Maximum
- Globale Mehrh
- GLOBALSIM

Overview

48

- Introductory example
- Schema level integration
 - Schema Mapping
 - Schema Matching
- Data level integration
 - Duplicate detection
 - Data fusion
- ETL management (with Alexander Albrecht)



ETL Process Management

49

- State of the art
 - Many data sources from different organizational areas are involved in a variety of data integration projects using ETL.
 - ETL processes may encompass shared data sources, same data targets, common sub-processes, and transformations

- Diagnosis
 - No common method, approach, or framework to uniformly manage entire ETL processes

- With METL (Managing ETL) we are implementing a next generation ETL tool that supports high-level ETL management.

Management Operators in METL

50

- **Search** – retrieves all ETL processes that satisfy the specified search query.
- **Match** – given an ETL (sub-)process it finds all corresponding ETL (sub-)processes that extract, transform, or load common data in a similar way.
- **Merge** – takes one or more ETL processes as input and returns a merged ETL process.
- **Invert** – feeds the output of one ETL process back to its sources
- **Create** – populates system from a variety of sources (mappings, SQL, scripts, mappings, ...)
- **Import/Deploy** – interface for existing ETL tools, based on a common representation of ETL processes

METL Prototype

51

Searching for ETL processes with following features:

Required	<input type="text" value="L"/>
Optional	<input checked="" type="checkbox"/> Authors
Unwanted	<input type="checkbox"/> Linda
	<input checked="" type="checkbox"/> Database Server

localhost
localhost:1234
localhost:8080

Files
longIDs_cust_AL.csv
longIDs_cust_MR.csv
longIDs_cust_SZ.csv

Transformations
LDAP Reader
LDAP Writer
Lookup

```

    graph LR
      Catalog[(Catalog)] --> Join[Join]
      Addresses[(Addresses)] --> Join
      Join --> Split[Split]
      LDAP[(LDAP)] --> Lookup[Lookup]
      Split --> Lookup
      Lookup --> Gather[Gather]
      Split --> Gather
      Gather --> CRM[(CRM)]
  
```

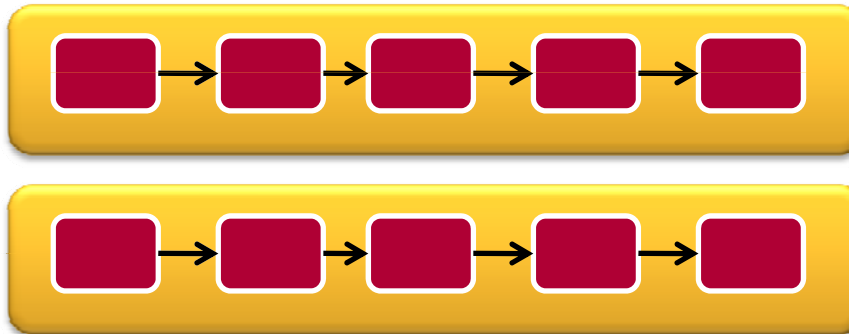
Search & Match

52

- **Search** – similar to search engines
 - Search predicates
 - Optionally prepended by a context-label
- **phone +table:addresses +stage:lookup - author:Alice**
- Ranking based on relevance to specified search terms using TF/IDF-like term weighting
- **Match** – easy-to-use access to ETL processes in repository
- Problem: Suitable similarity measure for ETL processes
 - Variety of ETL features
 - Semantic or syntactic heterogeneity
- Structure-aware Match operator
 - Position in ETL process
 - Data-Schema
 - Type of operator

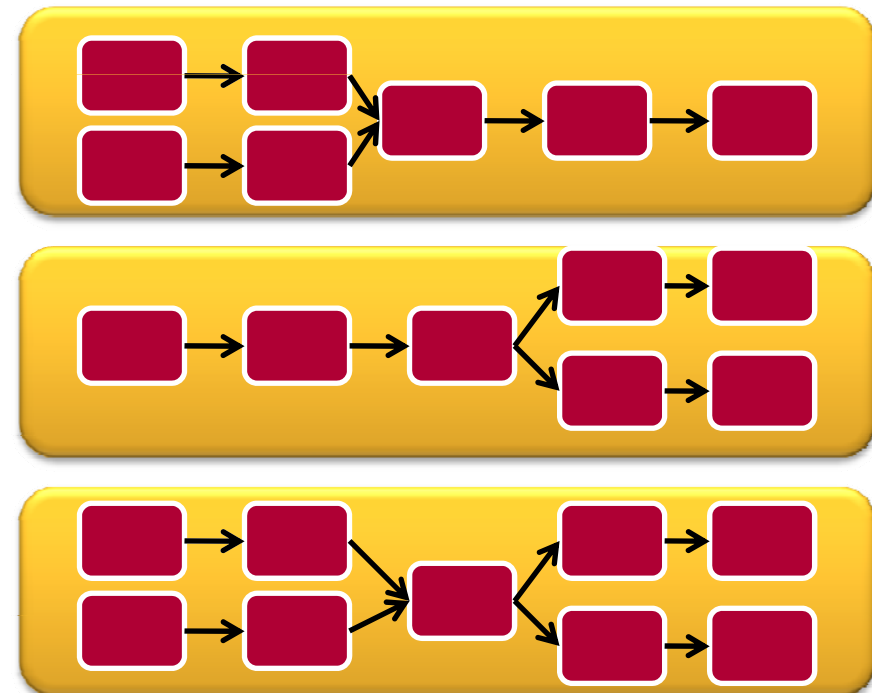
Merge

53



Benefits

- Better utilization of shared resources
- Latency improvement and reduced amount of data transmission
- Enhancement of performance compared to performing all processes in a separate run
- Single view of all information that was originally processed separately



Invert

54

- Invert output at a given transformation within the ETL process
- Generates for source table S a corresponding source table S' with cleansed data.
 - ETL-equivalence
- Benefits
 - Consolidated and cleaned data is fed back to the sources to ensure data quality for applications on top of original sources.
 - Propagation of corrections to sources improves future ETL projects
 - ◇ Avoid multiple corrections of same error
 - Application area: MDM

