



**Hasso  
Plattner  
Institut**

IT Systems Engineering | Universität Potsdam

Informationsintegration  
**Schema Mapping**

2.6.2008

Felix Naumann

## Erste Hälfte

2

Einführung in die Informationsintegration

Szenarien der Informationsintegration

Verteilung und Autonomie

Heterogenität

# Problemstellung

Materialisierte und virtuelle Integration

5-Schichten-Architektur

Mediator/Wrapper-Architektur / PDMS

# Architekturen

Schema Mapping

Schema Matching

SchemaSQL

# Mapping

Global-as-View und Lokal-as-View Modellierung

Global-as-View-Anfragebearbeitung

# Modellierung

3

Containment & Local-as-View Anfragebearbeitung  
Bucket Algorithmus  
Verteilte Anfragebearbeitung

## Anfragen

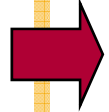
Duplikaterkennung  
Datensision - Union & Co.  
DWH, ETL & Data Lineage  
Informationsqualität

## Datenintegration

Hidden Web  
Semantic Web

## Anwendungen

4

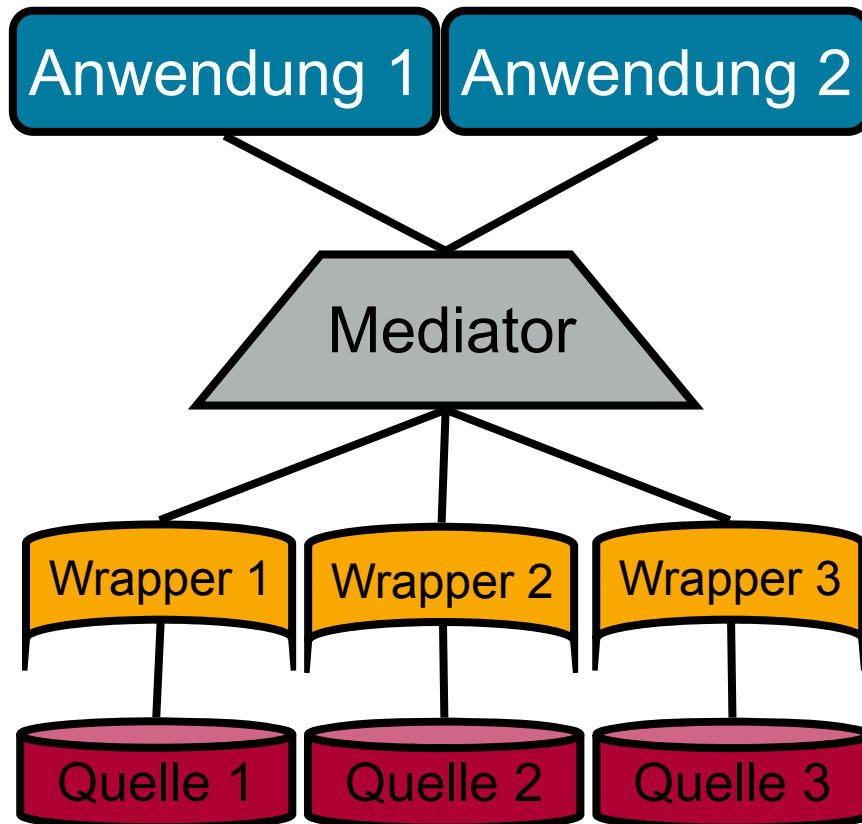


- Motivation
- Schema Mapping
- Schema Matching
- Mapping Interpretation
- Mapping Werkzeuge



# Wdh: Virtuelle Integration

5



- Datenfluss
- Anfragebearbeitung
- Entwicklung
  - Top-down
- Schema

# Wdh: Bottom-up oder Top-down Entwurf

6

- Beim Entwurf des integrierten Systems
- Bottom-up
  - Ausgelöst durch den Bedarf, mehrere Quellen integriert anzufragen
  - Schemaintegration ist wichtig.
  - Änderungen schwierig, da neu integriert werden muss.
  - Typisches Szenario: Data Warehouse
- Top-down
  - Ausgelöst durch globalen Informationsbedarf
  - Vorteilhaft bei labilen Quellen
  - Schemaintegration nicht nötig, bzw. leichter

# Wdh: Bottom-up Entwicklung

7

Wesentliche Merkmale

- Vorhandene lokale Schemata
- ALLE Daten sollen integriert werden

Deshalb

- Schema-integrativ

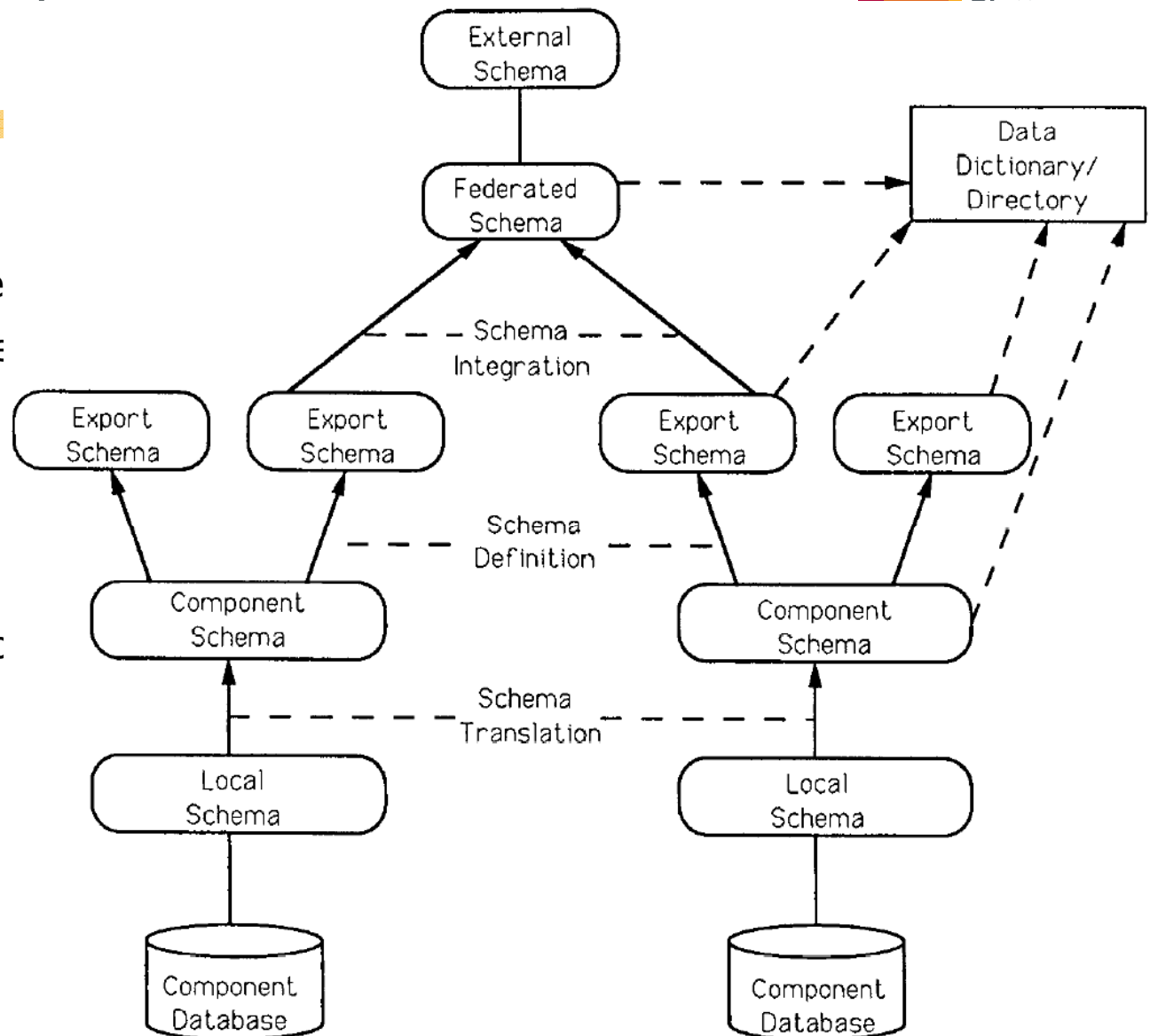


Figure 21. Bottom-up FDBS developing process.

# Wdh: Top-down Entwicklung

8

## Wesentliche Merkmale

- Vorhandenes globales
- Passende lokale Systeme werden gesucht

## Deshalb

- LaV
- und Schema Mapping

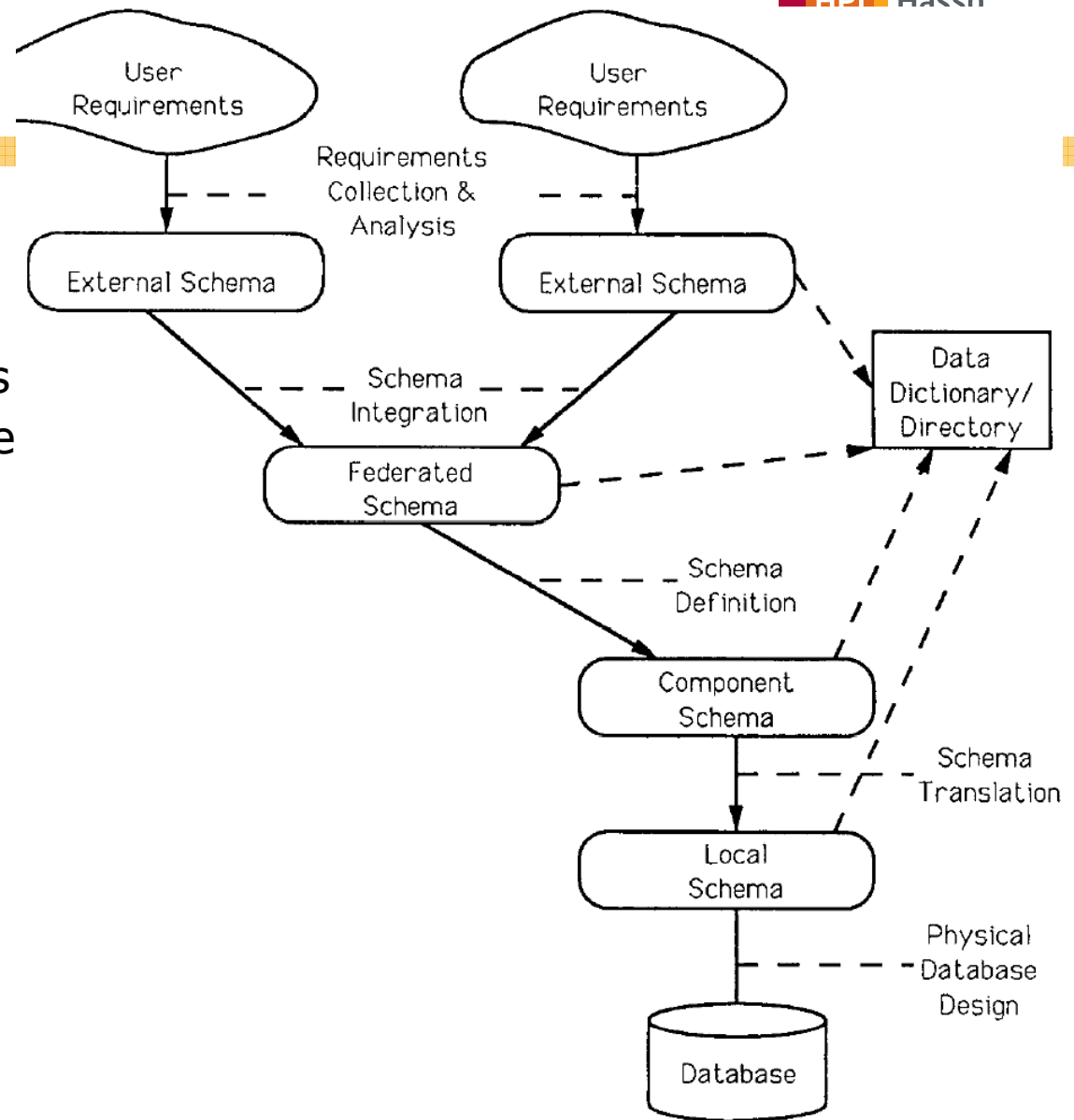


Figure 22. Top-down FDBS developing process.



# Schemaintegration vs. Schema Mapping

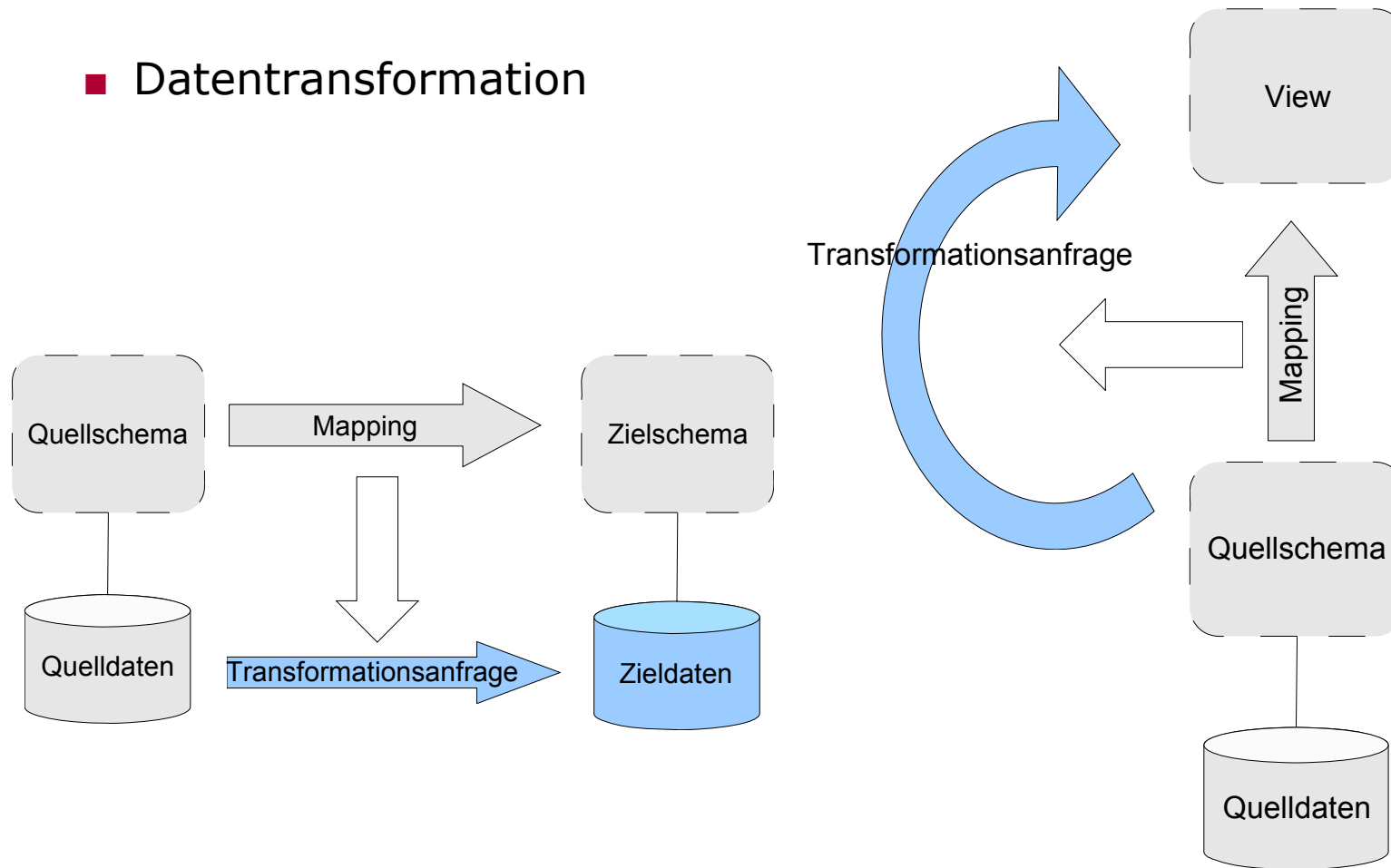
9

- Beide Probleme müssen strukturelle und semantische Heterogenität überwinden.
- Aber:
  - Schemaintegration liefert Schema Mapping „frei Haus“.
  - Zielschema hat keine eigene Semantik.
  - Schemaintegration ist unflexibel.
- Deshalb nun: Schema Mapping

# Verwendung von Schema Mappings

10

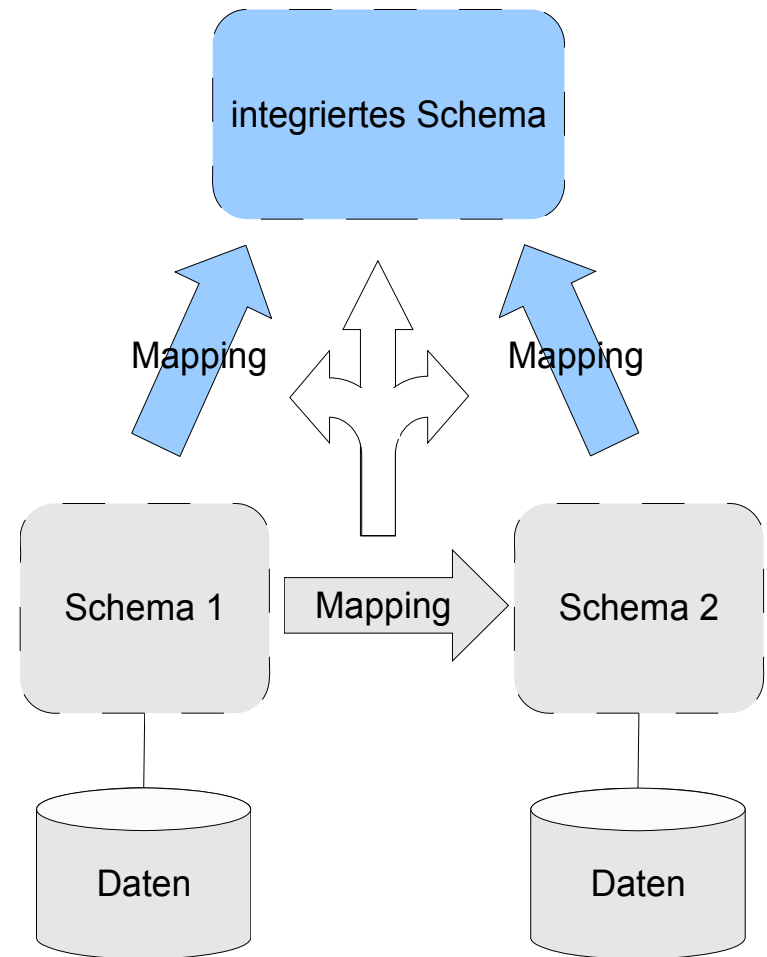
## ■ Datentransformation



# Verwendung von Schema Mappings

11

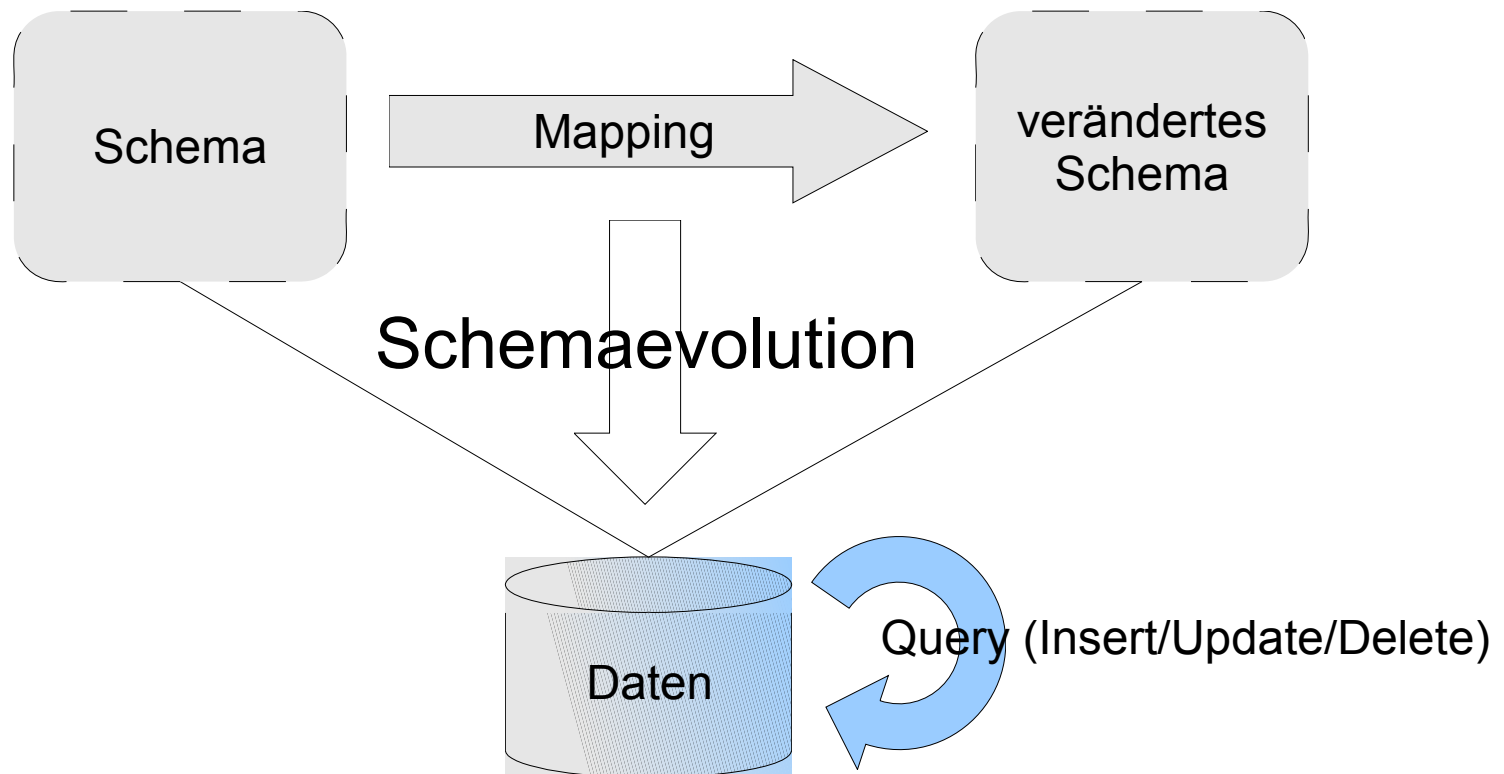
- Schemaintegration



# Verwendung von Schema Mappings

12

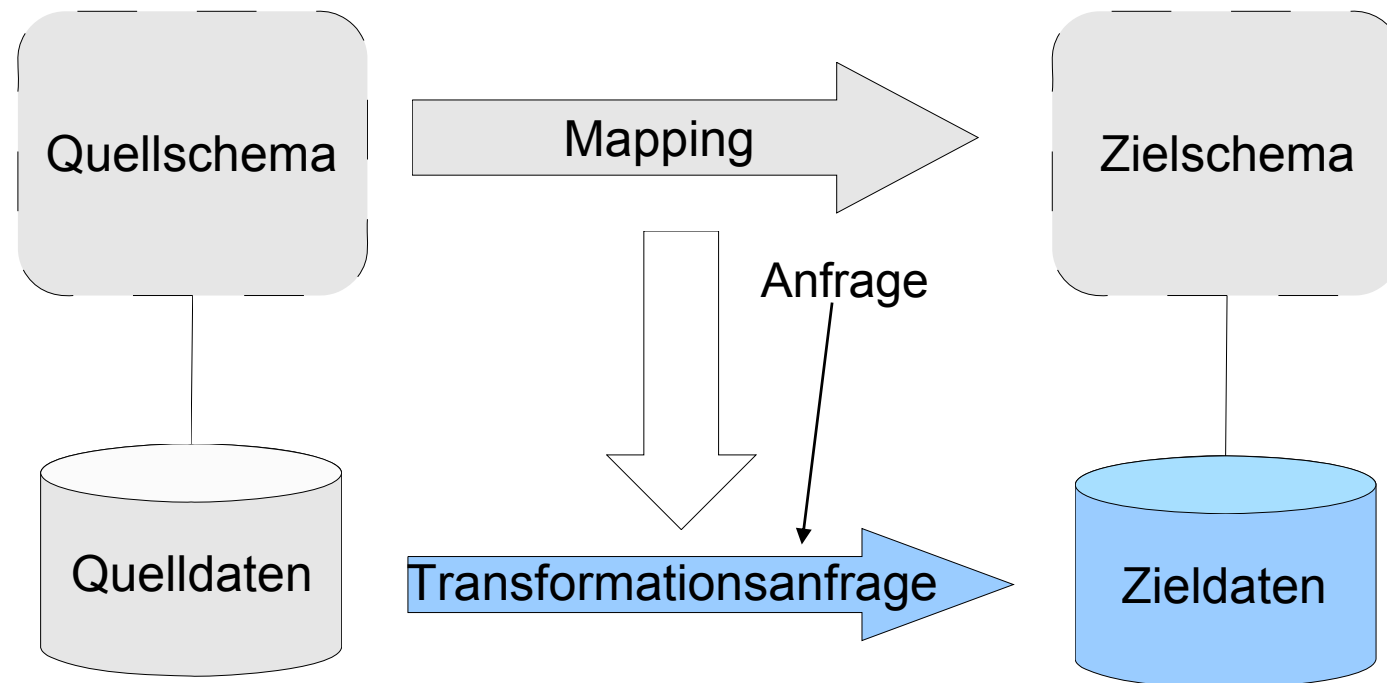
## ■ Schemaevolution



# Verwendung von Schema Mappings

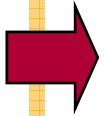
13

- Im weiteren: Datentransformation
  - Materialisierte Integration
  - Virtuelle Integration



14

- Motivation
- Schema Mapping
- Schema Matching
- Mapping Interpretation
- Mapping Werkzeuge



# Schema Mapping - Definitionen

15

- (Inter-Schema) **Korrespondenz**
  - Eine Zuordnung eines oder mehrerer Elemente eines (Quell-) Schemas zu einem oder mehreren Elementen eines anderen (Ziel-) Schemas
  - Auch: Value-correspondence
- (High-level) **Mapping**
  - Eine Menge von Korrespondenzen zwischen zwei Schemas.
- (Low-Level) **Logisches Mapping**
  - Logische Übersetzung eines oder mehrerer Mappings, die
    - ◇ den Integritätsbedingungen beider Schemas gehorcht und
    - ◇ die Intention des Nutzers widerspiegelt.
- **Interpretation**
  - Übersetzung eines Mappings in ein oder mehrere logische Mappings
  - Übersetzung eines logischen Mappings in eine Transformationsanfrage
- **Transformationsanfrage**
  - Anfrage in einer Anfragesprache (z.B. SQL), die Daten des Quellschemas in die Struktur des Zielschemas überführt

# Wdh: Schematische Heterogenität

16

## ■ Struktur

### □ Modellierung

- ◇ Relation vs. Attribut
- ◇ Attribut vs. Wert
- ◇ Relation vs. Wert

### □ Benennung

- ◇ Relationen
- ◇ Attribute

### □ Normalisiert vs. Denormalisiert

### □ Geschachtelt vs. Fremdschlüssel

High-order Mappings

hier



# Motivation

17

- Datentransformation zwischen heterogenen Schemas
  - Altes aber immer wiederkehrendes Problem
  - Üblicherweise schreiben Experten komplexe Anfragen oder Programme
    - ◇ Zeitintensiv
    - ◇ Experte für die Domäne, für Schemata und für Anfrage
    - ◇ XML macht alles noch schwieriger
      - XML Schema, XQuery
- Idee: Automatisierung
  - Gegeben: Zwei Schemata und ein high-level Mapping dazwischen.
  - Gesucht: Anfrage zur Datentransformation

## Motivation – Probleme

18

- Generierung der „richtigen“ Anfrage unter Berücksichtigung der Schemata und des Mappings
- Garantie, dass die transformierten Daten dem Zielschema entsprechen
- Effiziente Datentransformation
  - Für Materialisierung (Ausführung, inkrementell)
  - Für virtuelle Integration (query-unfolding)

**Hier: Nur Effektivität, nicht Effizienz**

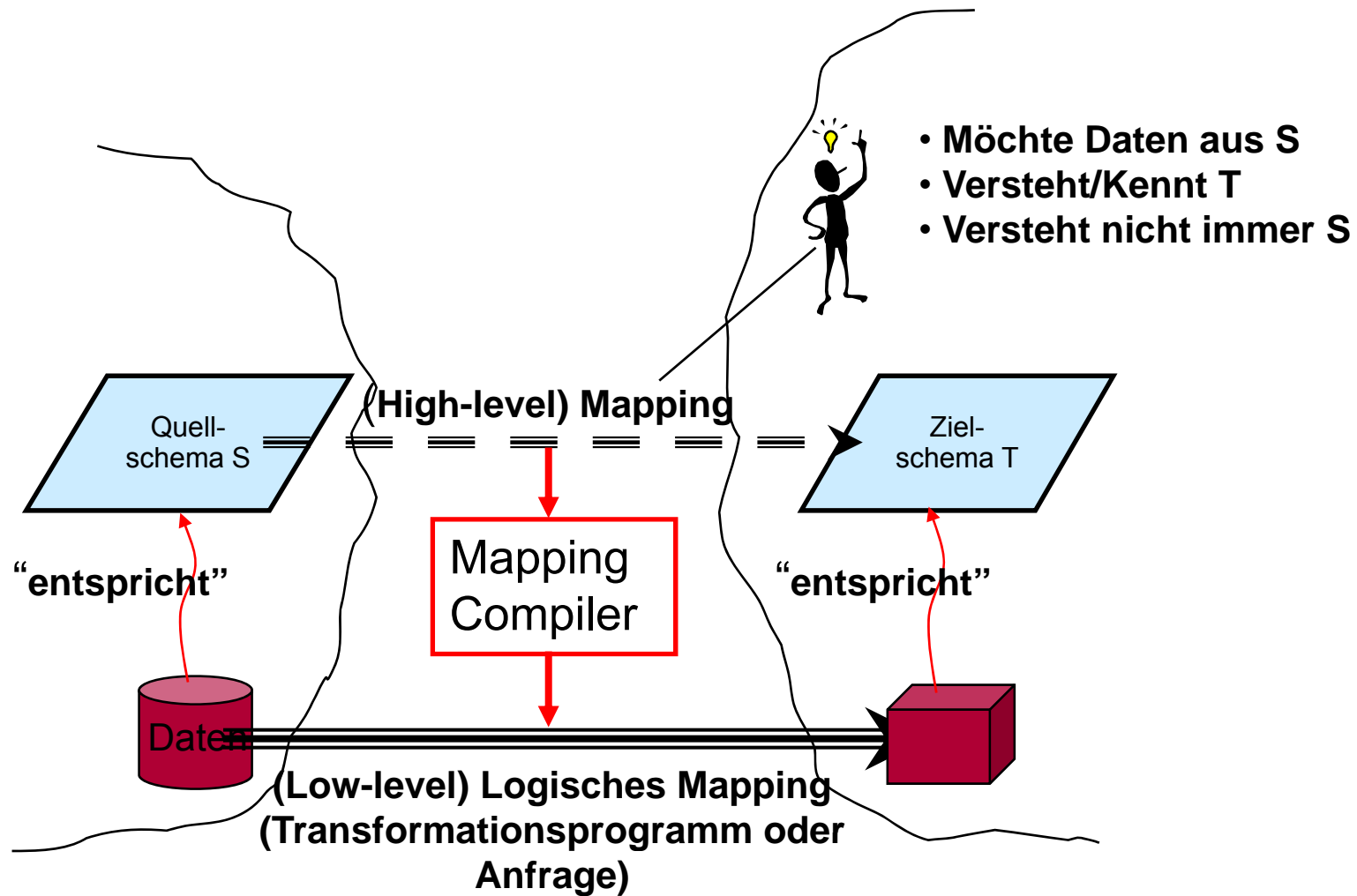
# Motivation – Weitere Probleme

19

- Geschachtelte Strukturen unterstützen
  - Geschachteltes, relationales Modell
  - XML
  - Geschachtelte Integritätsbedingungen
- Korrespondenzen
  - Nutzerfreundlich
  - Automatische Entdeckung (Schema Matching)
- Intention des Nutzers erkennen und repräsentieren
- Semantik der Daten erhalten
  - Assoziationen entdecken & erhalten
  - Schemata und deren Integritätsbedingungen nutzen
- Neue Datenwerte erzeugen
- Korrekte Gruppierungen erzeugen
- ...

# Schema Mapping im Kontext

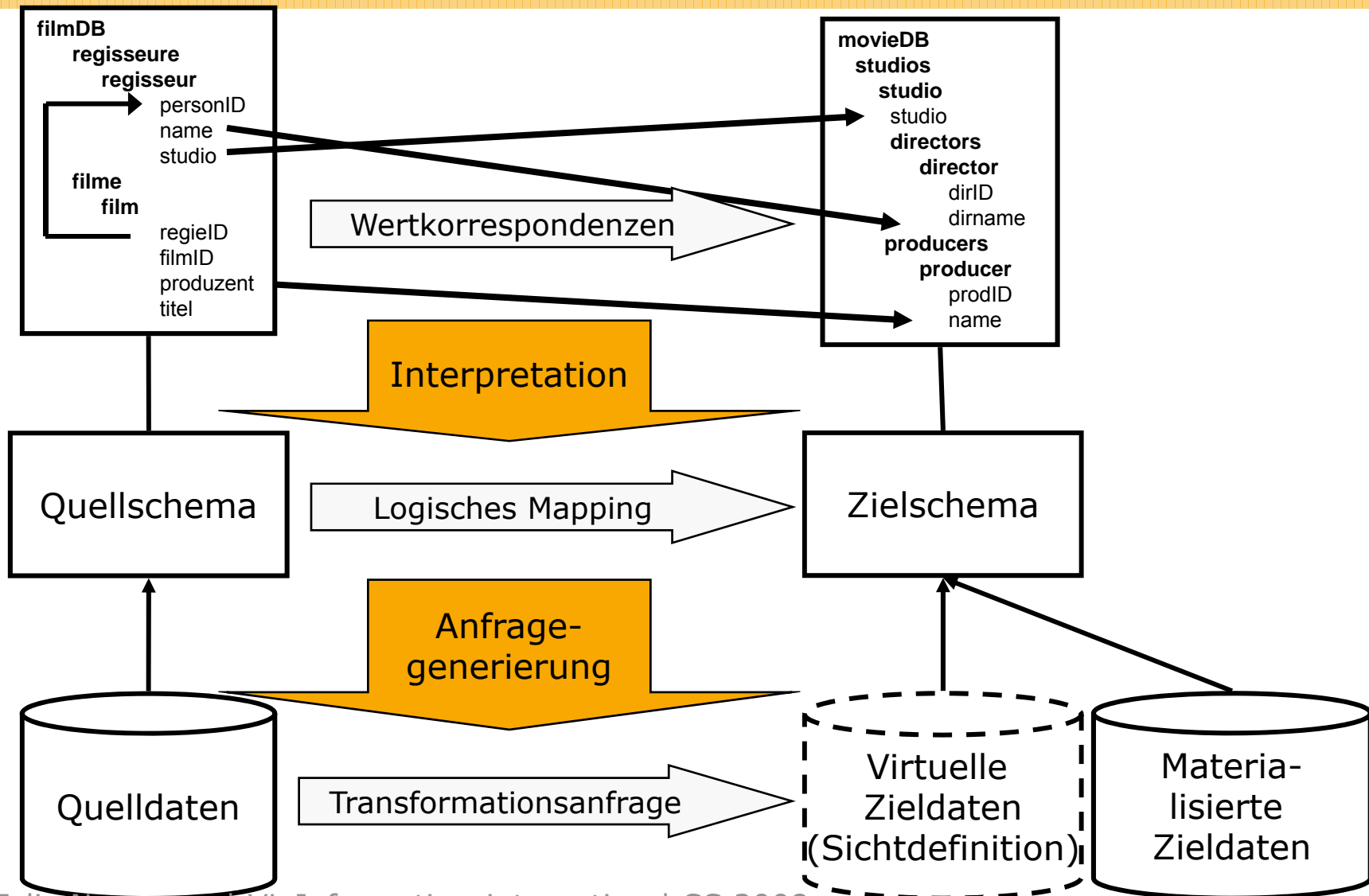
20



Quelle: [FHP+02]

# Schema Mapping im Kontext

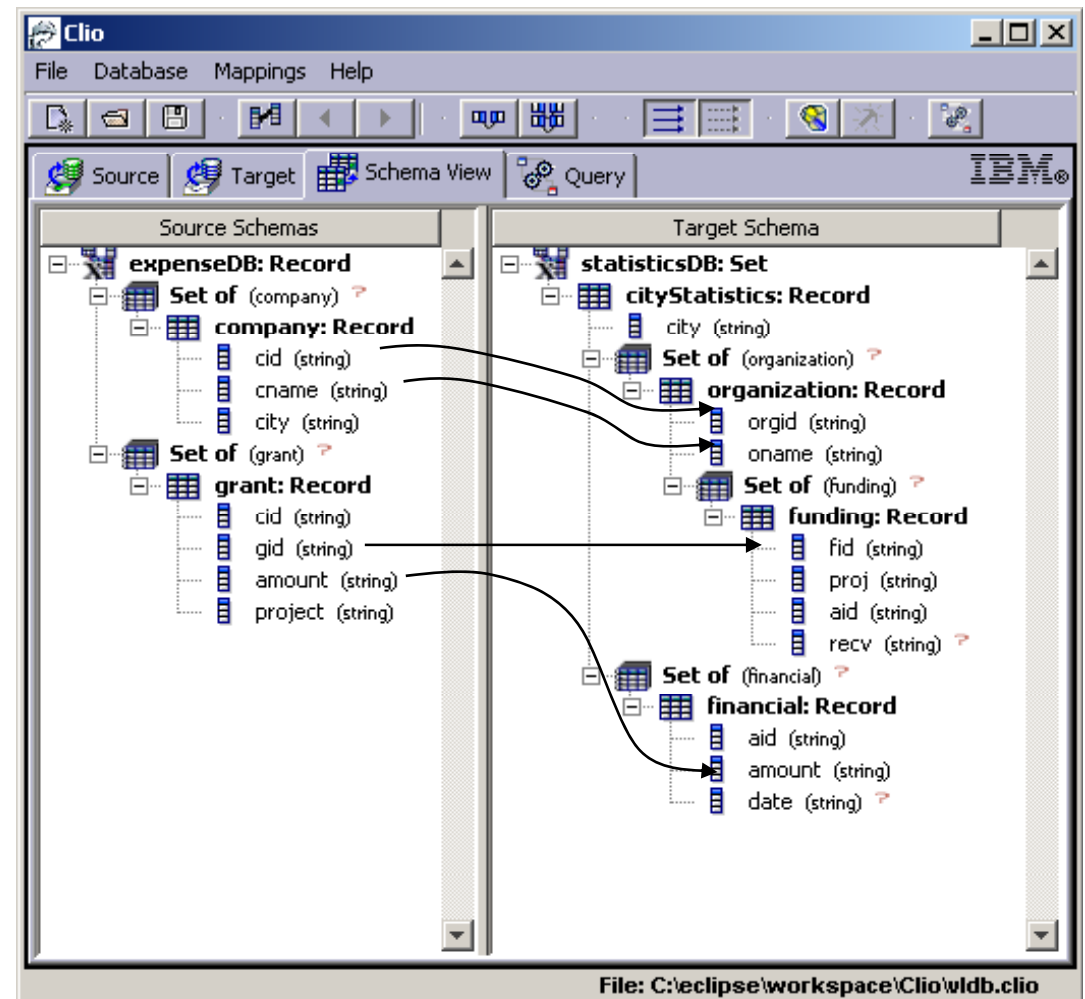
21



# Schema Mapping im Kontext

22

- Schema Matching & Korrespondenzen
- Schema Mapping
- Mapping Interpretation
- Daten-transformation



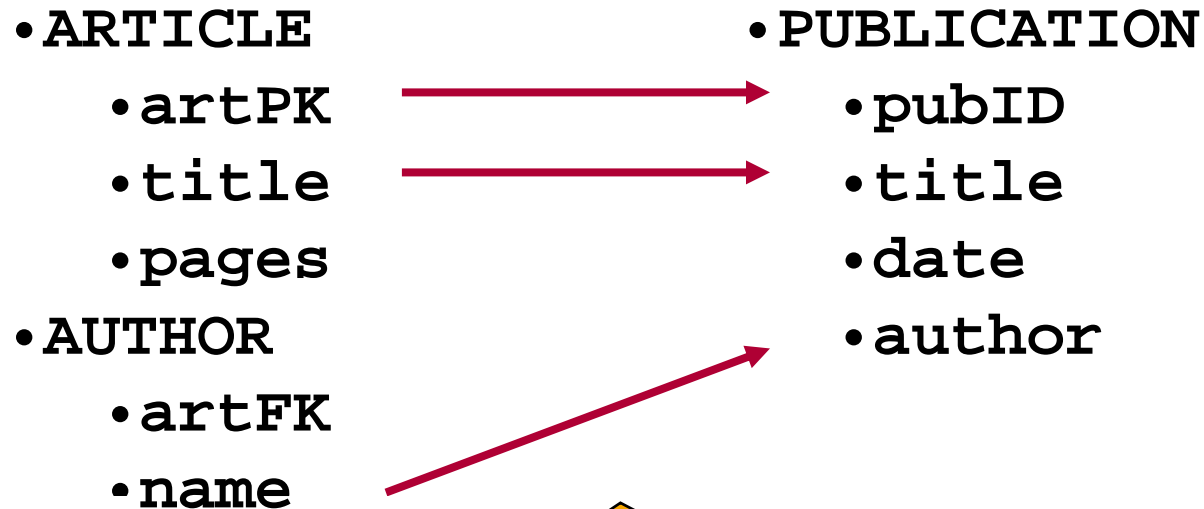
# Wdh: Schematische Heterogenität – Beispiel

23

- Normalisiert vs. Denormalisiert
  - 1:n Assoziationen zwischen Werten wird unterschiedlich dargestellt
    - ◇ Durch Vorkommen im gleichen Tupel
    - ◇ Durch Schlüssel-Fremdschlüssel Beziehung
  
- Lösung: Schema Mapping

# Schema Mapping Beispiel

24



```

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
  
```



# Schematische Heterogenität – Lösungen

25

• **ARTICLE**

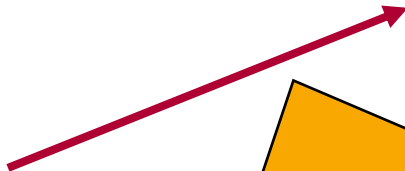
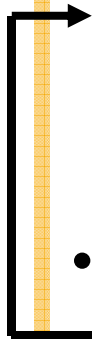
- artPK
- title
- pages

• **AUTHOR**

- artFK
- name

• **PUBLICATION**

- pubID
- title
- date
- author



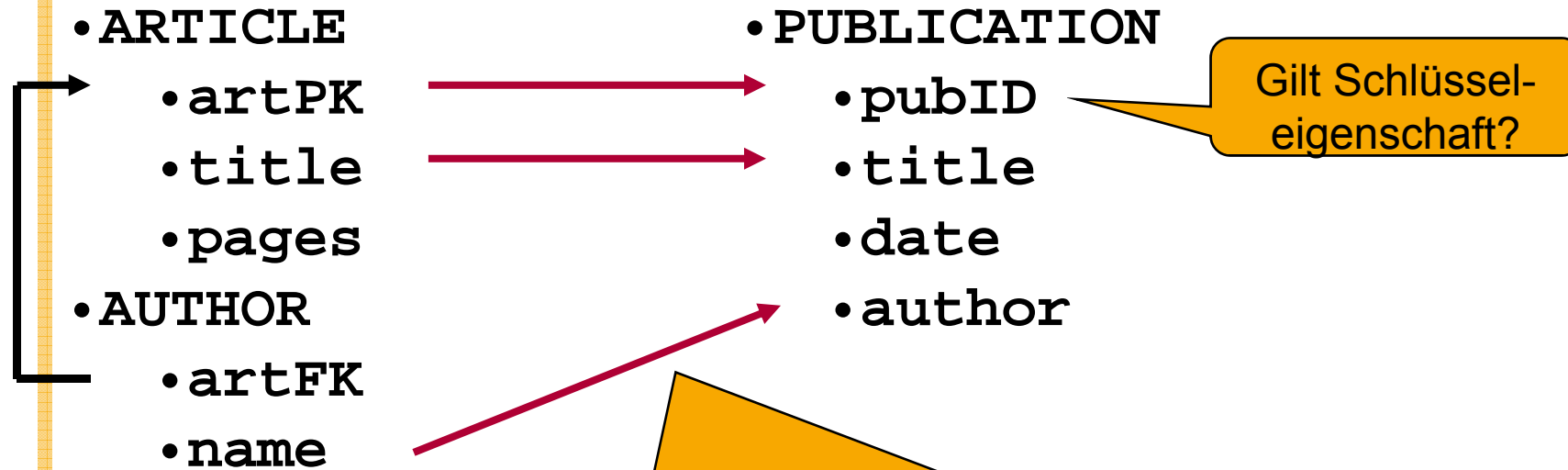
Weitere Interpretationen?

```

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

# Schematische Heterogenität – Lösungen

26



```

SELECT      artPK AS pubID
            title AS title
            null AS date
            name AS author
FROM        ARTICLE LEFT OUTER JOIN AUTHOR
ON          ARTICLE.artPK = AUTHOR.artFK
    
```

# Wdh: Schematische Heterogenität – Lösungen

27

## • PUBLICATION

- title
- date
- author

## • ARTICLE

- artPK
- title
- pages
- AUTHOR
- artFK
- name

```
SELECT SK(title) AS artPK
       title AS title
       null AS pages
FROM   PUBLICATION
```

```
SELECT SK(title) AS artFK
       author AS name
FROM   PUBLICATION
```

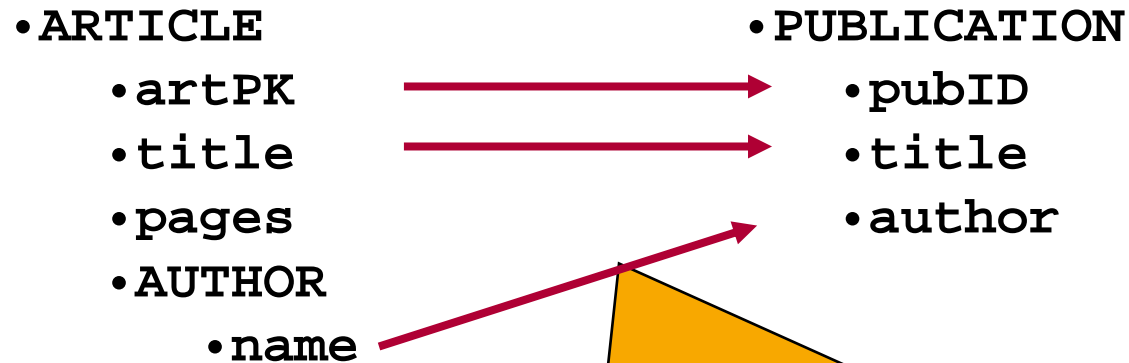
# Wdh: Schematische Heterogenität – Beispiel

28

- Geschachtelt vs. Flach
  - 1:n Assoziationen werden unterschiedlich dargestellt
    - ◇ Als geschachtelte Elemente
    - ◇ Als Schlüssel-Fremdschlüssel Beziehung
  
- Lösung: Schema Mapping

# Wdh: Schematische Heterogenität – Lösungen

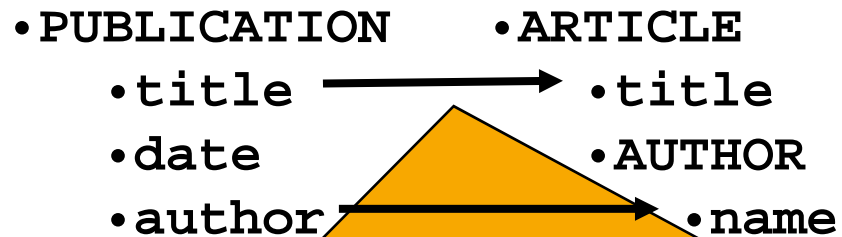
29



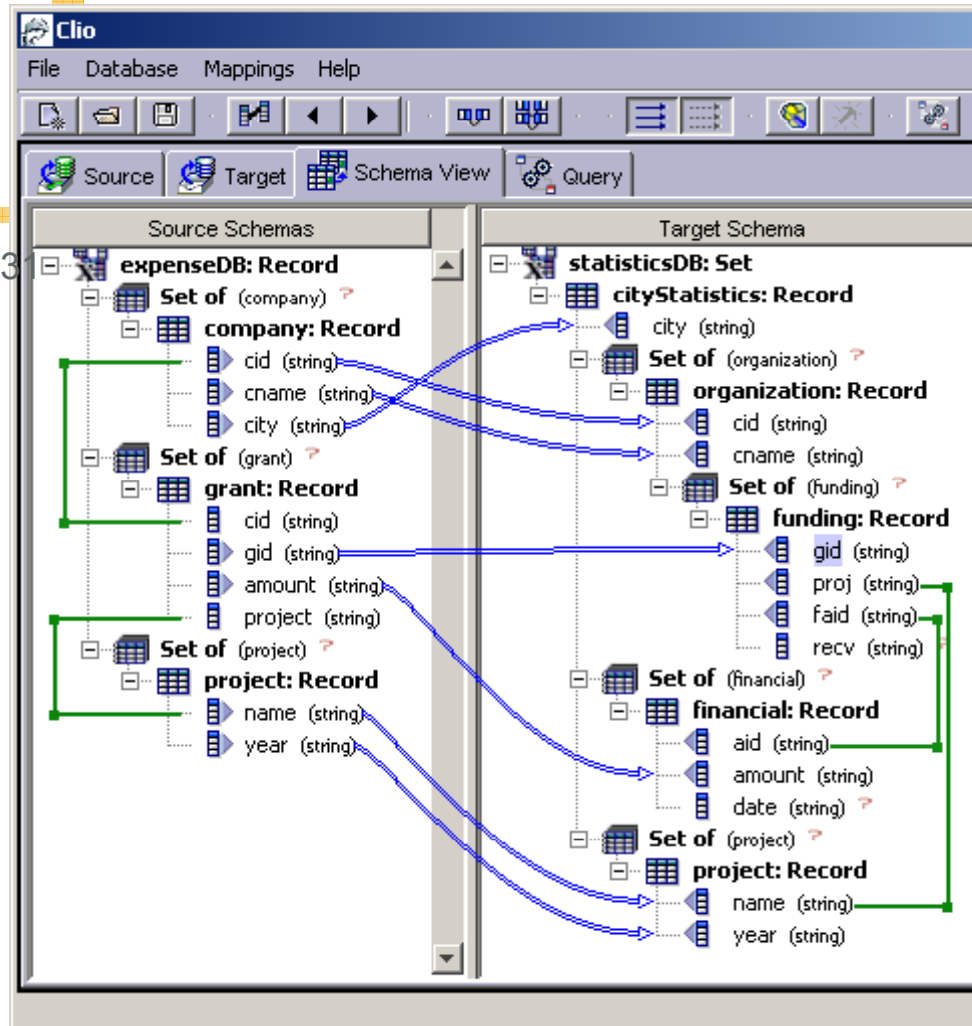
```
LET $doc0 := document("articles.xml") RETURN
<root> { distinct-values (
  FOR $x0 IN $doc0/authorDB/ARTICLE, $x1 IN
  $x0/AUTHOR
  RETURN
    <publication>
    <pubID> { $x0/artPK/text() } </pubID>
    <title> { $x0/title/text() } </title>
    <author> { $x1/name/text() } </author>
    </publication>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

# Wdh: Schematische Heterogenität – Lösungen

30



```
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/publication
      WHERE $x0/title/text() = $x0L1/title/text()
      RETURN
        <AUTHOR>
          <name> { $x0L1/author/text() } </name>
        </AUTHOR> ) }
  </ARTICLE> ) } </articles>
```



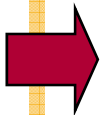
```

XQuery
RETURN
  $x0L1/project/text() = $x1L1/name/text() AND
  $x2L1/cid/text() = $x0L1/cid/text() AND
  $x2/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>,
      <faid> "SK267(", $x0L2/project/text(), " ,
    </funding> )
  </organization> ,
  distinct (
  FOR
    $x0L1 IN $doc/expenseDB/grant ,
    $x1L1 IN $doc/expenseDB/project ,
    $x2L1 IN $doc/expenseDB/company
  WHERE

```

Jetzt: Schema Matching

- Motivation
- Schema Mapping
- Schema Matching
  - Klassifikation von Schema Matching Methoden
  - Erweiterungen
  - Globales Matching
- Mapping Interpretation
- Mapping Werkzeuge





# Schema Matching – Motivation

33

- Große Schemas
  - > 100 Tabellen, viele Attribute
  - Bildschirm nicht lang genug
- Unübersichtliche Schemas
  - Tiefe Schachtelungen
  - Fremdschlüssel
  - Bildschirm nicht breit genug
  - XML Schema
- Fremde Schemas
  - Unbekannte Synonyme
- Irreführende Schemas
  - Unbekannte Homonyme
- Fremdsprachliche Schemas
- Kryptische Schemas
  - |Attributnamen| ≤ 8 Zeichen
  - |Tabellennamen| ≤ 8 Zeichen

Clio

File Schemas Mappings Help

Source Target Schema View Query

Source schemas

- author: Record (personType)
  - name (string)
- pubMed: Record
  - value (unsignedInt) ?
  - status (NMTOKEN) ?
  - medlineID (unsignedInt) ?
- Set [0..\*] ?
  - ev: Record (evidenceRefType)
    - type (string) ?
    - href (anyURI)
    - role (anyURI) ?
    - title (string) ?
    - label (string) ?
    - value() (string) ?
- citedInBook: Record (bookType)
  - linktype (string) ?
  - role (anyURI) ?
  - title (string) ?
  - volume (unsignedInt) ?
  - year (gYear)
  - first (unsignedInt)
  - last (unsignedInt)
  - ISDN (string) ?
  - publisher (string)
  - city (string)
  - country (string) ?
- Record
  - title (string)
  - bookTitle (string)
  - editors: Record (editorsType)
    - Set [1..\*]
      - editor: Record (personType)
        - name (string)
  - Set [1..\*]
    - author: Record (personType)
      - name (string)
  - Set [0..\*] ?
    - ev: Record (evidenceRefType)
      - type (string) ?
      - href (anyURI)
      - role (anyURI) ?
      - title (string) ?
      - label (string) ?
      - value() (string) ?
- Set [0..\*] ?
  - ev: Record (evidenceRefType)
    - type (string) ?
    - href (anyURI)
    - role (anyURI) ?
    - title (string) ?
    - label (string) ?
    - value() (string) ?
- observations: Record (observationType)
  - linktype (string) ?

Target Schema

34

Man beachte die Scrollbar!

Man beachte die Schachtelungstiefe!

Felix Naumann | VI Informationsintegration | SS 2008

er ut

# Schema Matching – Motivation

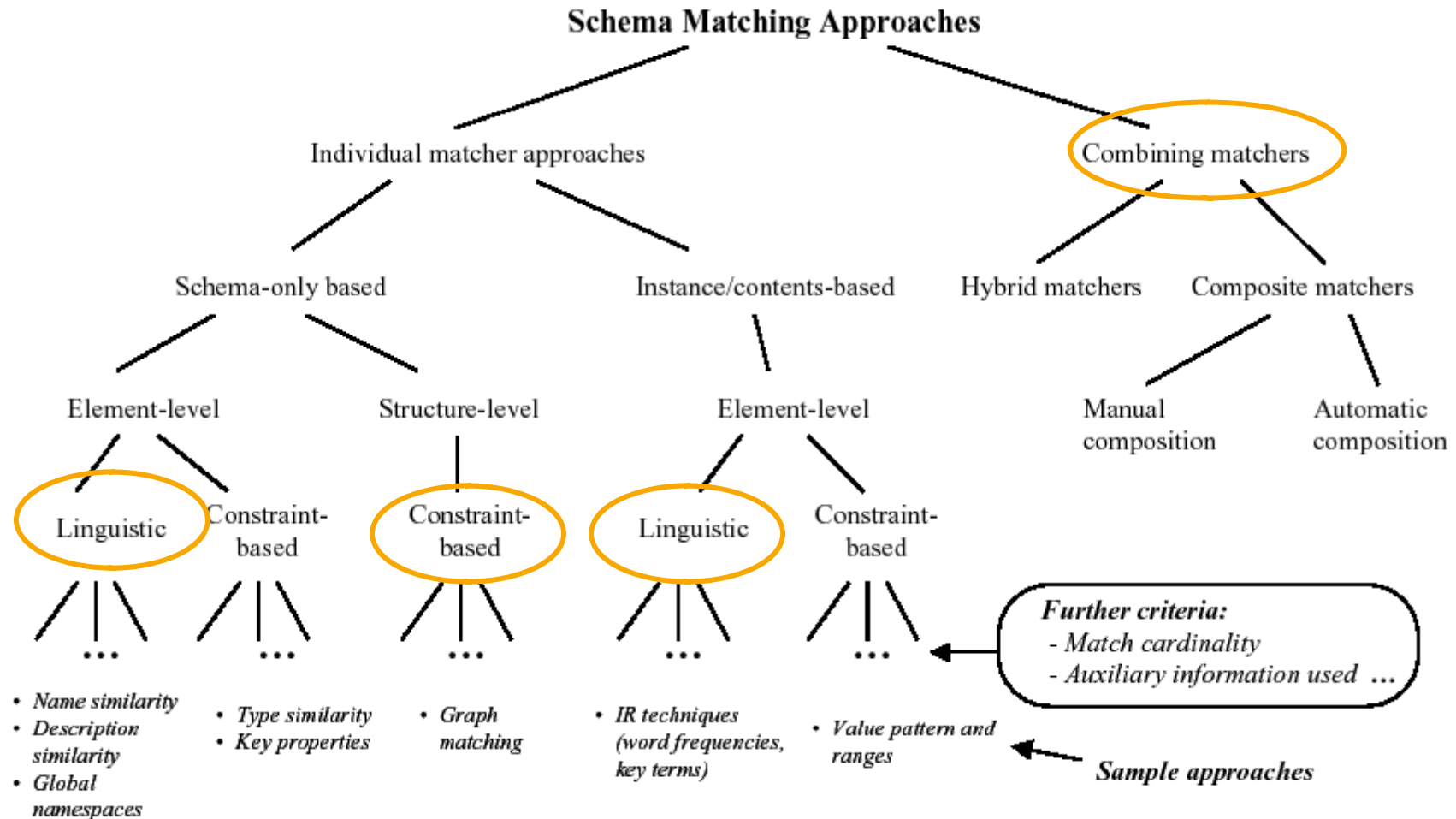
35

## Die Folgen

- Falsche Korrespondenzen (false positives)
- Fehlende Korrespondenzen (false negatives)
- Frustration
  - User verlieren sich im Schema
  - User verstehen Semantik der Schemas nicht

# Schema Matching Klassifikation nach [RB01]

36



# Schema Matching Klassifikation

37

Schema Matching basierend auf

- Namen der Schemaelemente (*label-based*)
- Darunterliegende Daten (*instance-based*)
- Struktur des Schemas (*structure-based*)
- Mischformen

# Schema Matching – Label-based

38

Gegeben zwei Schemata mit Attributmengen A und B

Kernidee:

- Bilde Kreuzprodukt aller Attribute aus A und B.
- Für jedes Paar vergleiche Ähnlichkeit bezgl. Attributnamen (Label).
  - Z.B. Edit-distance
- Ähnlichste Paare sind Matches

Probleme:

- Effizienz
- Auswahl der besten Matches (globales Matching)
  - Iterativ?
  - Stable Marriage?
- Synonyme und Homonyme werden nicht erkannt

# Schema Matching – Label-based

39

## Stand der Technik in kommerziellen Produkten

- Label-based
- Namensgleichheit
- Kein globales Matching
- Keine Ähnlichkeitsmaße
- Kein Instanz-basiertes Matching

# Altova MapForce 2007

40

**Connection settings:**  
Opens the Connections Settings dialog, in which you can define the specific mixed content settings as well as the connector annotation settings, please see the [Connection](#) section in the Reference section.

**Connect matching Children dialog box**  
This command allows you to create multiple connectors between items of the same name in both the source and target components.

1. Connect two (parent) items that share identically named child items in both components.
2. Right click the connector and select the **Connect matching child elements** option.

**Connect Matching Children**

Ignore Case  
 Ignore Namespaces  
 Recursive  
 Mix Attributes and Elements

Existing Connections  
 Ignore existing output connections  
 Retain  
 Overwrite  
 Delete all existing

OK Cancel

3. Select the required options discussed in the text below, and click OK to create the mappings.  
Mappings are created for all the child items that have identical names and adhere to the settings defined in the dialog box.



# Schema Matching – Instance-based

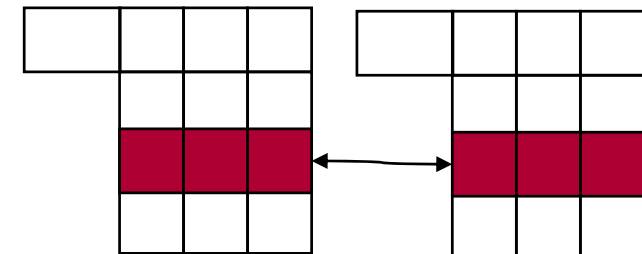
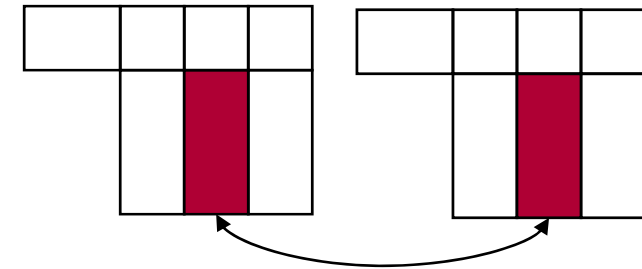
41

- Gegeben zwei Schemata mit Attributmengen A und B, jeweils mit darunterliegenden Daten.
- Kernidee
  - Für jedes Attribute extrahiere interessante Eigenschaften der Daten
    - ◇ Buchstabenverteilung, Länge, etc.
  - Bilde Kreuzprodukt aller Attribute aus A und B.
  - Für jedes Paar vergleiche Ähnlichkeit bzgl. der Eigenschaften
- Probleme
  - Auswahl der Eigenschaften
  - Datenmenge: Sampling
  - Vergleichsmethode, z.B. Naive Bayes
  - Gewichtung (Maschinelles Lernen)

# Instanz-basiertes Schema Matching

42

- **Herkömmliche Lösung: Vertikal**
  - Vergleich von Spalten
  - = **Attribut-Klassifikation**
  - [ICDE'02] u.v.a.m.
- **Neue Lösung: Horizontal**
  - Vergleich von Zeilen
  - = **Duplikaterkennung**
    - ◇ trotz fehlender Attribut-korrespondenzen
  - [ICDE'05]

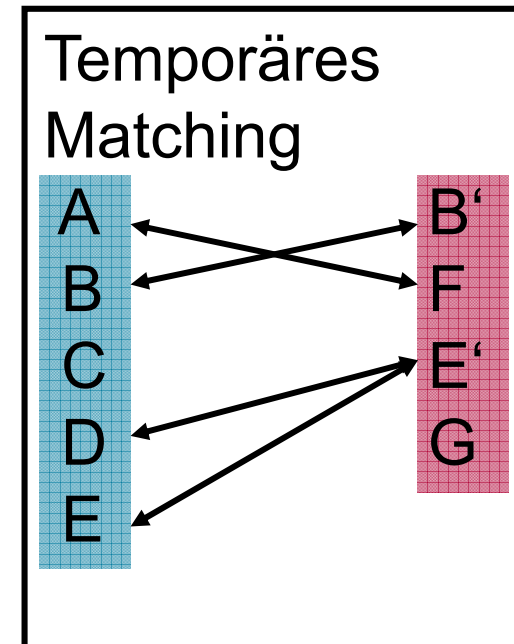
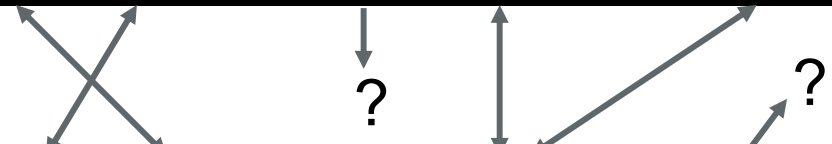


# Duplicate-driven Schema Matching

43

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

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

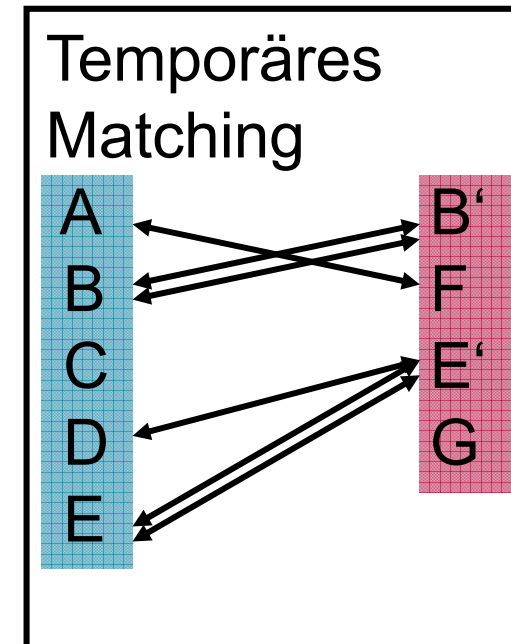
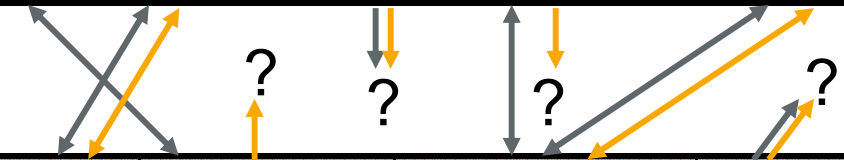


# Duplicate-driven Schema Matching

44

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



## Annahmen

- Es existieren Daten in beiden DBs.
- Es existieren (wenigstens ein paar) Duplikate in beiden DBs.

# Schema Matching – Structure-based

45

- Gegeben zwei Schemata mit Elementmengen A und B.
- Kernidee
  - Nutze (komplexe) Struktur des Schemas aus.
  - Hierarchieebene
  - Elementtyp (Attribut, Relation, ...)
  - Nachbarschaftsbeziehungen

# Schema Matching – Structure-based

46

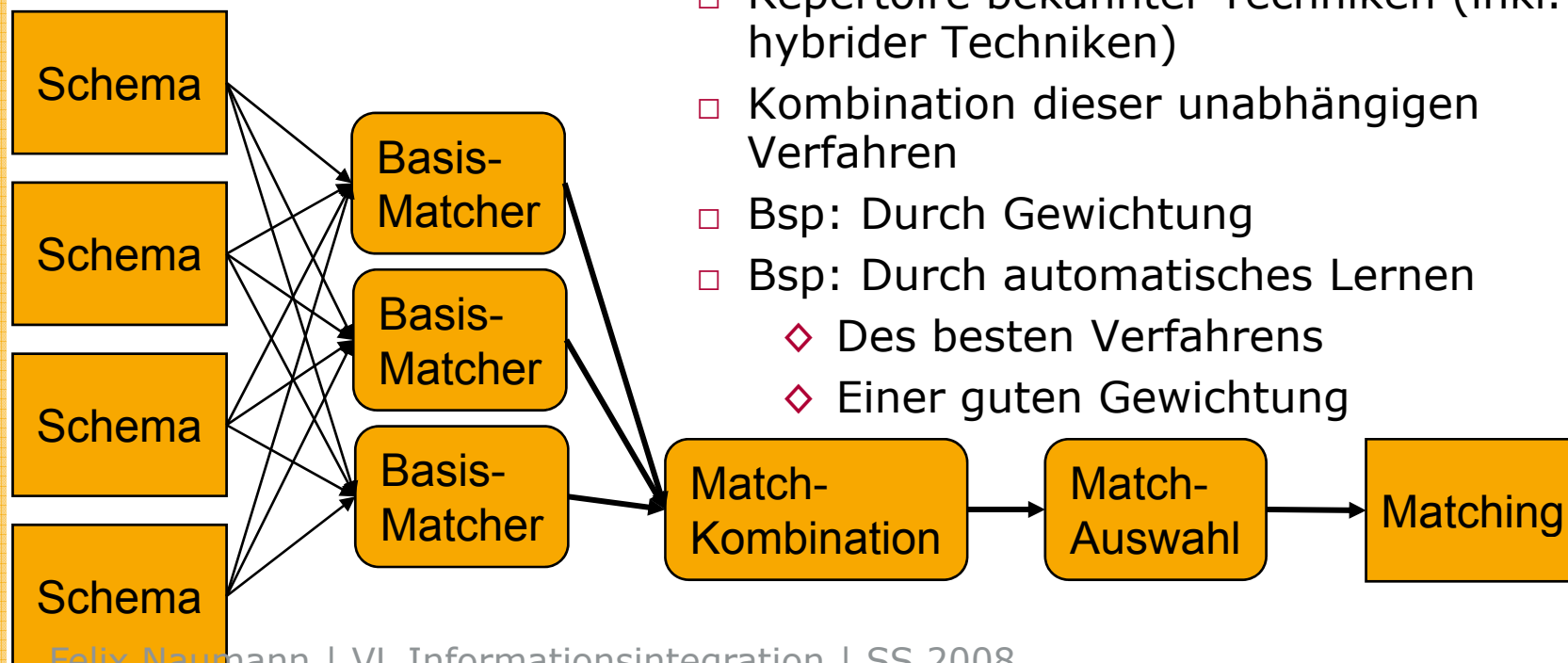
Beispiel: Similarity Flooding nach [MGMR02]

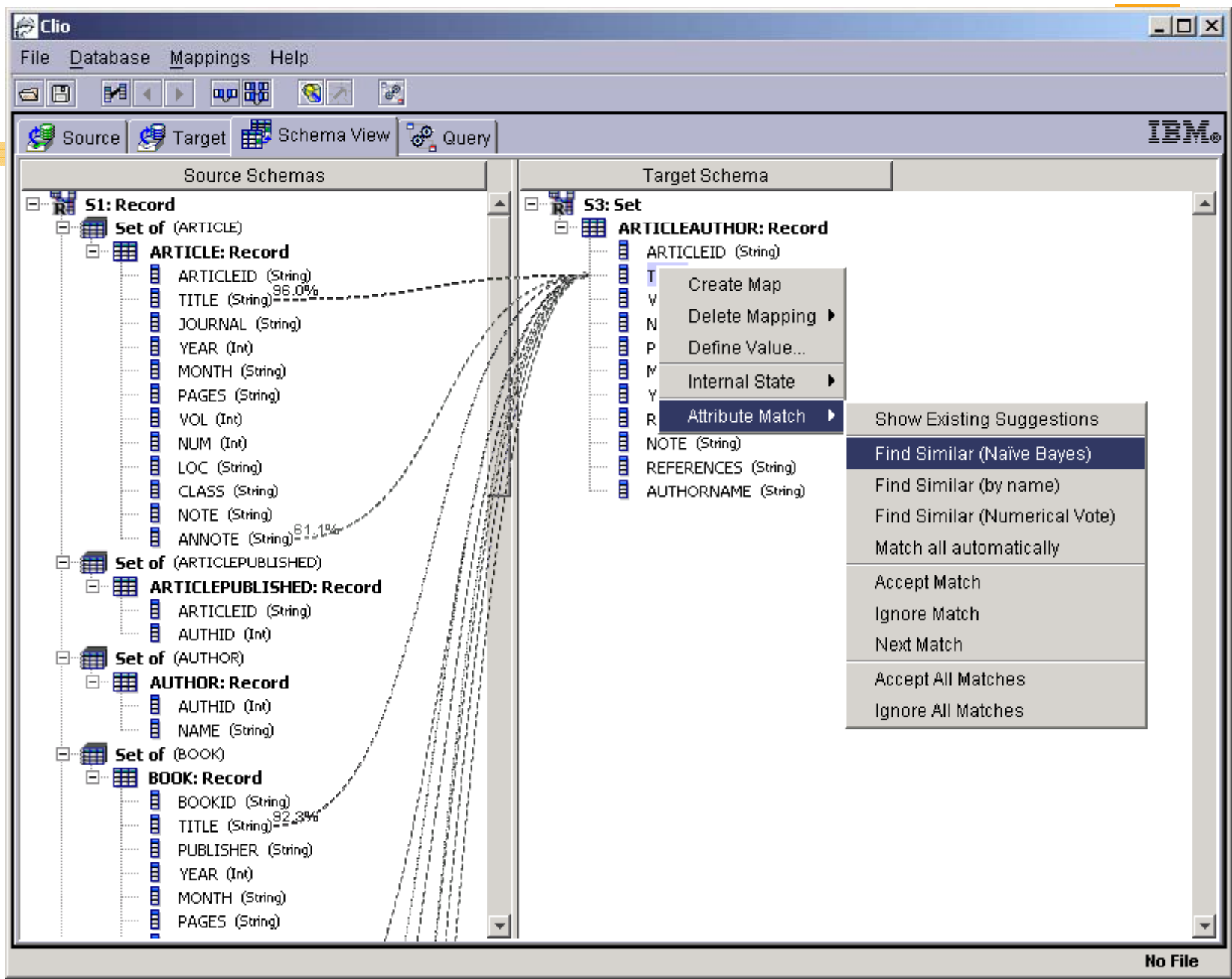
- Gegeben initiale Ähnlichkeit zwischen Schemaelementen (z.B. durch edit-distance oder durch Analyse der darunterliegenden Daten)
- Lasse Ähnlichkeiten „abfärben“ auf die Nachbarn
  - Nachbarn sind durch Struktur definiert
  - Sind alle Nachbarn von x und y ähnlich zueinander, sind (vielleicht) auch x und y ein match.
- Analogie: Man „flutet“ das Netzwerk der Ähnlichkeiten bis ein Gleichgewicht erreicht ist.

# Schema Matching – Mischformen

47

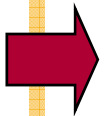
- Hybrid
  - Gleichzeitige Anwendung mehrerer Techniken
  - Bsp: Instance-based + Datentypvergleich
- Composite
  - Repertoire bekannter Techniken (inkl. hybrider Techniken)
  - Kombination dieser unabhängigen Verfahren
  - Bsp: Durch Gewichtung
  - Bsp: Durch automatisches Lernen
    - ◇ Des besten Verfahrens
    - ◇ Einer guten Gewichtung







- Motivation
- Schema Mapping
- Schema Matching
  - Klassifikation von Schema Matching Methoden
  - Erweiterungen
  - Globales Matching
- Mapping Interpretation
- Mapping Werkzeuge



# Schema Matching – Erweiterungen

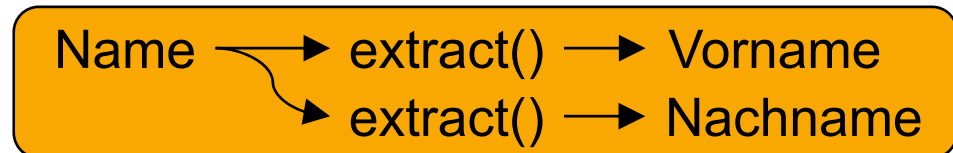
50

- n:1 und 1:n Matches
  - Viele Kombinationsmöglichkeiten
  - Viele Funktionen denkbar
    - ◇ Mathematische Operatoren, Konkatination, etc.
  - Parsingregeln
- n:m Matching?
- Matching in komplexen Schemata
  - Ziel: Finde Mapping, nicht Korrespondenzen

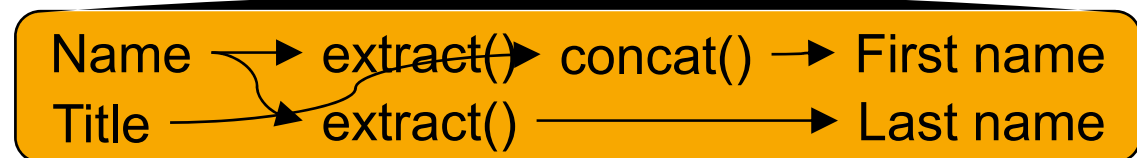
## n:1 Matching



## 1:n Matching



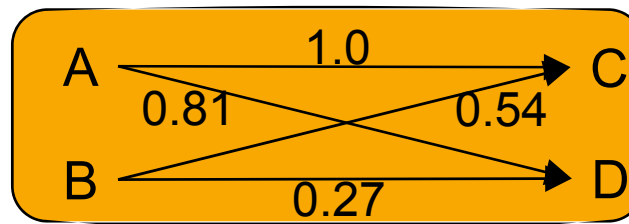
## m:n matching



# Schema Matching – Erweiterungen

51

- Global matching (gleich)
  - Matche nicht nur einzelne Attribute (oder Attributmengen)
  - Sondern komplette Tabellen oder komplette Schemata
  - Stable Marriage Problem

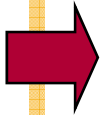


# Schema Matching – Weitere Anwendungen

52

- Herkömmlich: Korrespondenzen finden
- Schlüssel – Fremdschlüssel finden
  - Ähnliche Attribute innerhalb eines Schemas sind gute Kandidaten
- Höher-stufige Korrespondenzen finden
  - Ähnlichkeiten von Tabellen durch Aggregation der Matches ihrer Attribute

- Motivation
- Schema Mapping
- Schema Matching
  - Klassifikation von Schema Matching Methoden
  - Erweiterungen
  - Globales Matching
- Mapping Interpretation
- Mapping Werkzeuge



# Schema Matching – Stable Marriage

54

- Gegeben
  - $n$  Frauen (Attribute in Schema A) und  $m$  Männer (Attribute in Schema B)
- Monogamie
  - Je eine Frau kann nur mit je einem Mann verheiratet sein (nur 1:1 matches)
- Jede Frau hat eine Rangliste der Männer und umgekehrt
  - Bei Schema Matching
    - ◇ Attribut-Ähnlichkeit gemäß eines der vorigen Verfahren
    - ◇ Rangliste ist (normalerweise) symmetrisch
- Gesucht: Paarung (globales Matching), so dass niemals gilt
  - $f_1$  heiratet  $m_1$ ,  $f_2$  heiratet  $m_2$ ,
  - aber  $f_1$  bevorzugt  $m_2$  und  $m_2$  bevorzugt  $f_1$  (Instabil!)

# Stable Marriage – Beispiel

55

Männer (1-4)

1: B, D, A, C

2: C, A, D, B

3: B, C, A, D

4: D, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

Beispiel aus: David Toth,

["The Stable Marriage Problem: More Marital Happiness than Reality TV"](#)

April 25, 2003, Connecticut College, New London, CT, USA,

# Stable Marriage – Beispiel

56

Männer (1-4)

1: B, D, A, C

2: C, A, D, B

3: B, C, A, D

4: D, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein: (1, B)



# Stable Marriage – Beispiel

57

Männer (1-4)

1: B, D, A, C

2: C, A, D, B

3: B, C, A, D

4: D, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein: (1, B)

2 stellt Antrag an C, sie willigt ein: (1, B) (2, C)

# Stable Marriage – Beispiel

58

Männer (1-4)

1: ~~B~~, D, A, C

2: C, A, D, B

3: B, C, A, D

4: D, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein: (1, B)

2 stellt Antrag an C, sie willigt ein: (1, B) (2, C)

3 stellt Antrag an B, sie willigt ein & verlässt 1: (2, C) (3, B)

# Stable Marriage – Beispiel

59

Männer (1-4)

1: ~~B~~, D, A, C

2: C, A, D, B

3: B, C, A, D

4: D, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein : (1, B)

2 stellt Antrag an C, sie willigt ein : (1, B) (2, C)

3 stellt Antrag an B, sie willigt ein & verlässt 1: (2, C) (3, B)

1 stellt Antrag an D, sie willigt ein : (1, D) (2, C) (3, B)

# Stable Marriage – Beispiel

60

Männer (1-4)

1: ~~B~~, D, A, C

2: C, A, D, B

3: B, C, A, D

4: ~~D~~, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein : (1, B)

2 stellt Antrag an C, sie willigt ein : (1, B) (2, C)

3 stellt Antrag an B, sie willigt ein & verlässt 1: (2, C) (3, B)

1 stellt Antrag an D, sie willigt ein : (1, D) (2, C) (3, B)

4 stellt Antrag an D, sie lehnt ab : (1, D) (2, C) (3, B)

# Stable Marriage – Beispiel

61

Männer (1-4)

1: ~~B~~, D, A, C

2: C, A, D, B

3: B, C, A, D

4: ~~D~~, A, C, B

Frauen (A-D)

A: 2, 1, 4, 3

B: 4, 3, 1, 2

C: 1, 4, 3, 2

D: 2, 1, 4, 3

1 stellt Antrag an B, sie willigt ein : (1, B)

2 stellt Antrag an C, sie willigt ein : (1, B) (2, C)

3 stellt Antrag an B, sie willigt ein & verlässt 1: (2, C) (3, B)

1 stellt Antrag an D, sie willigt ein : (1, D) (2, C) (3, B)

4 stellt Antrag an D, sie lehnt ab: (1, D) (2, C) (3, B)

4 stellt Antrag an A, sie willigt ein : (1, D) (2, C) (3, B) (4, A)

## Vier stabile Paare!

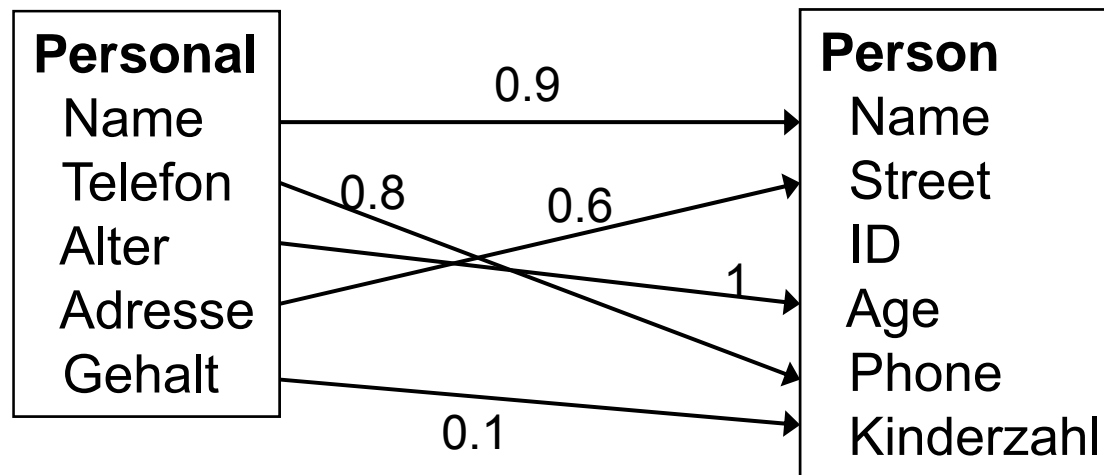
# Maximum Weighted Matching

62

- Alternative zu Stable Marriage
- Suche Matching mit maximalem Gewicht in bipartiten Graphen
  - Bipartit:
    - ◇ Knoten in zwei Klassen (Quelle & Ziel)
    - ◇ Kanten nur zwischen Knoten verschiedener Klassen (Korrespondenzen)
  - Maximiere Summe der einzelnen Gewichte/Ähnlichkeiten
- $O(n^3)$  („Ungarische Methode“)

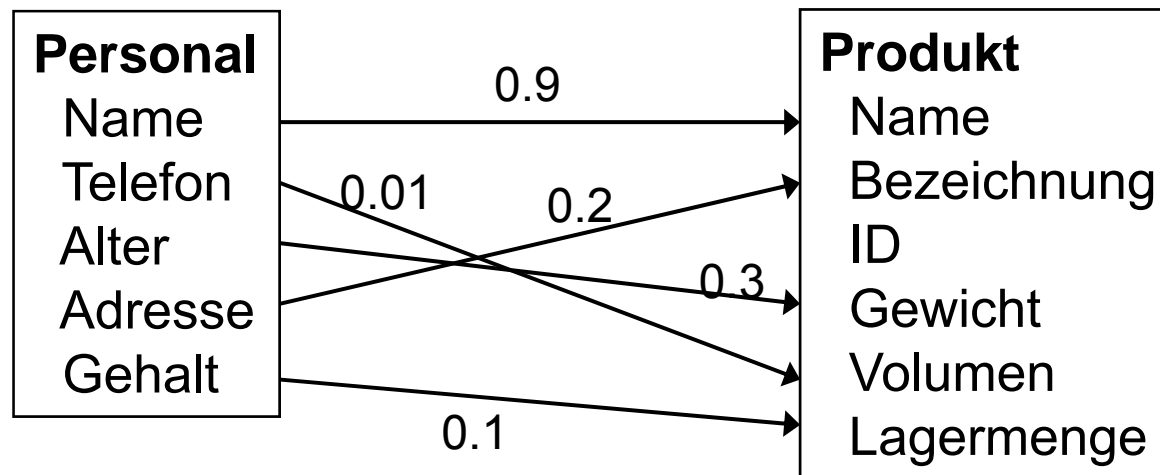
# Diskussion: Globales Matching

63



# Diskussion: Globales Matching

64



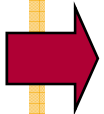


# Zusammenfassung – Schema Matching

65

- Schema Matching basierend auf
  - Namen der Schemaelemente (label-based)
  - Darunterliegende Daten (instance-based)
  - Struktur des Schemas (structure-based)
  - Mischformen, Meta-Matcher
- High-order Matching
- n:m Matching
- Globales Matching

- Motivation
- Schema Mapping
- Schema Matching
- Mapping Interpretation
  - Aus [FHP+02] und VLDB 2002 Vortragsfolien
- Mapping Werkzeuge



# Mapping – Das Problem

67

- Gegeben: Zwei Schemata
  - Unabhängig voneinander erzeugt
  - Relational
  - Geschachtelt
  - Mit Integritätsbedingungen (Schlüssel/Fremdschlüssel)
  - Stellen teilweise unterschiedliche Daten dar
- Gegeben: Eine Menge von Korrespondenzen zwischen den Schemata
- Gesucht: Anfrage, die Daten des einen in Daten des anderen Schemas transformiert, wobei
  - Semantik des Quellschemas erhalten bleibt,
  - Integritätsbedingungen des Zielschemas berücksichtigt werden,
  - und möglichst alle Korrespondenzen berücksichtigt werden.

# Relationale vs. XML Schemata

68

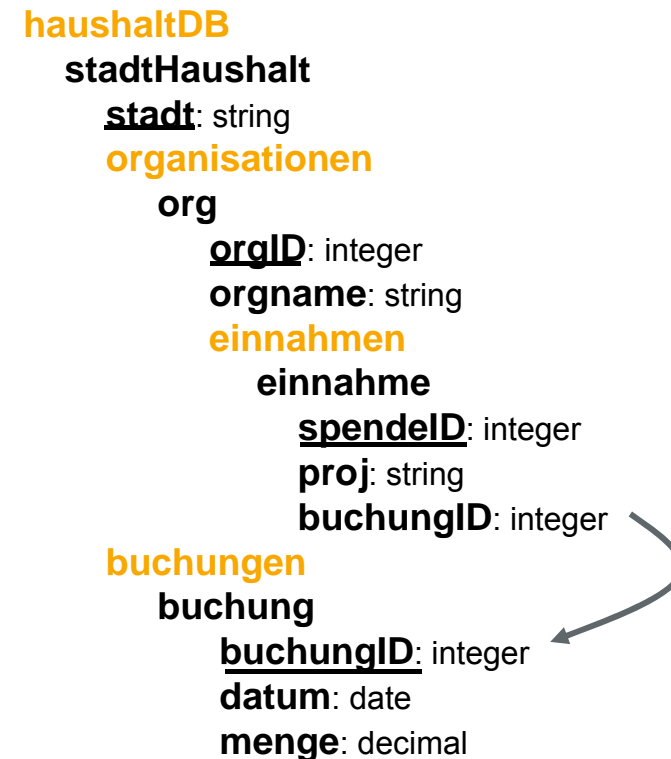
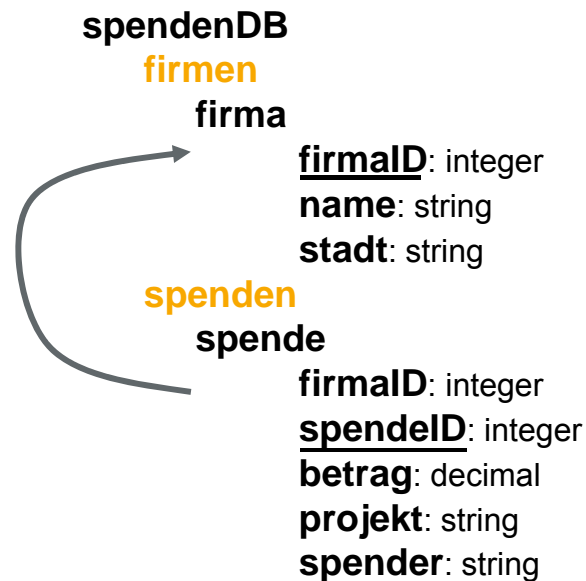
- Relationale Schemata
  - Flach
- XML Schemata
  - Flach oder geschachtelt
- NF2

Employee

ename	Children		Skills		
	name	dob	type	Exams	
				year	city
Smith	Sam	2/10/84	typing	1984	Atlanta
	Sue	1/20/85		1985	Dallas
			dictation	1984	Atlanta
Watson	Sam	3/12/78	filing	1984	Atlanta
				1975	Austin
				1971	Austin
			typing	1962	Waco

# Mapping – Beispiel

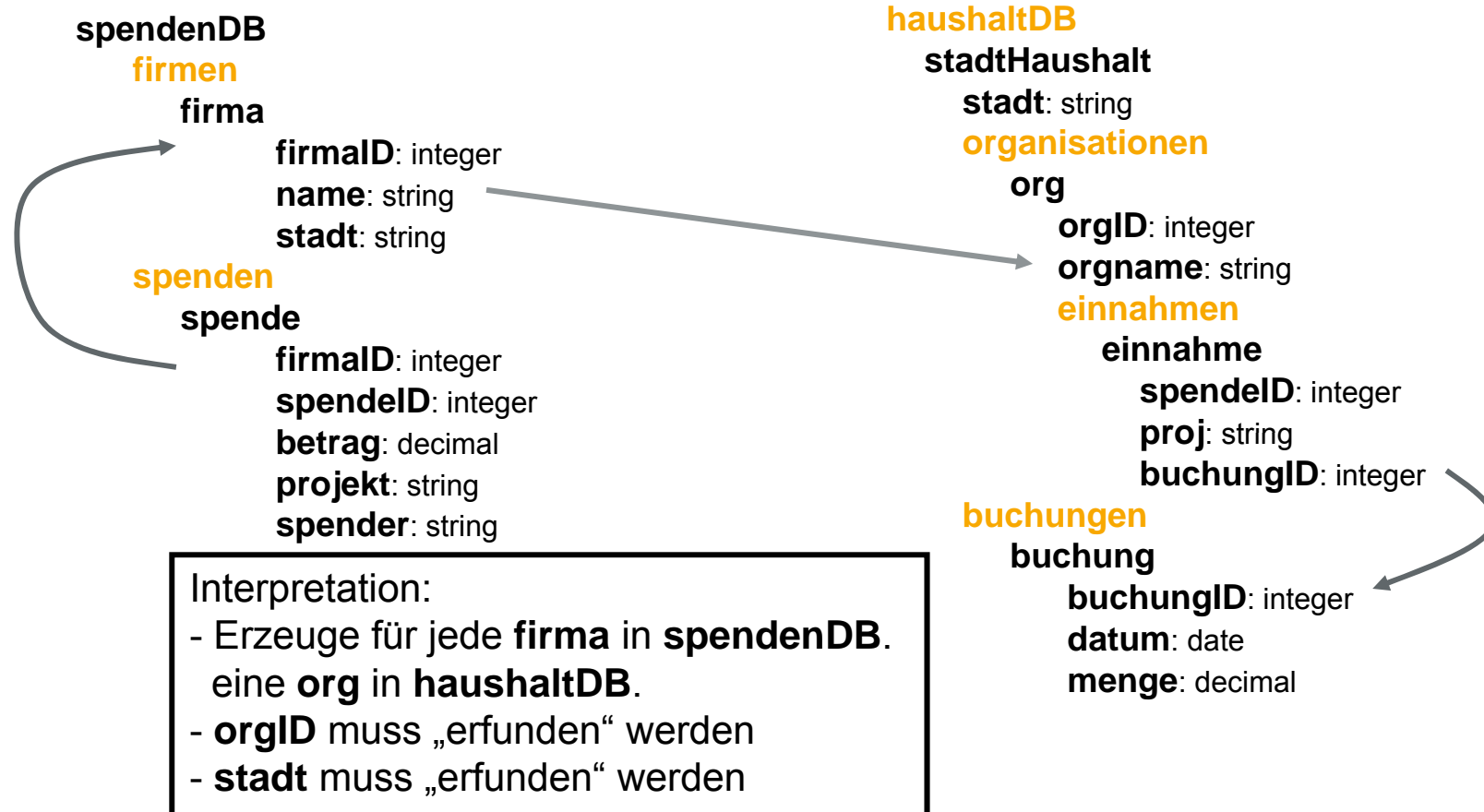
69



Quelle für Beispiel: [FHP+02]

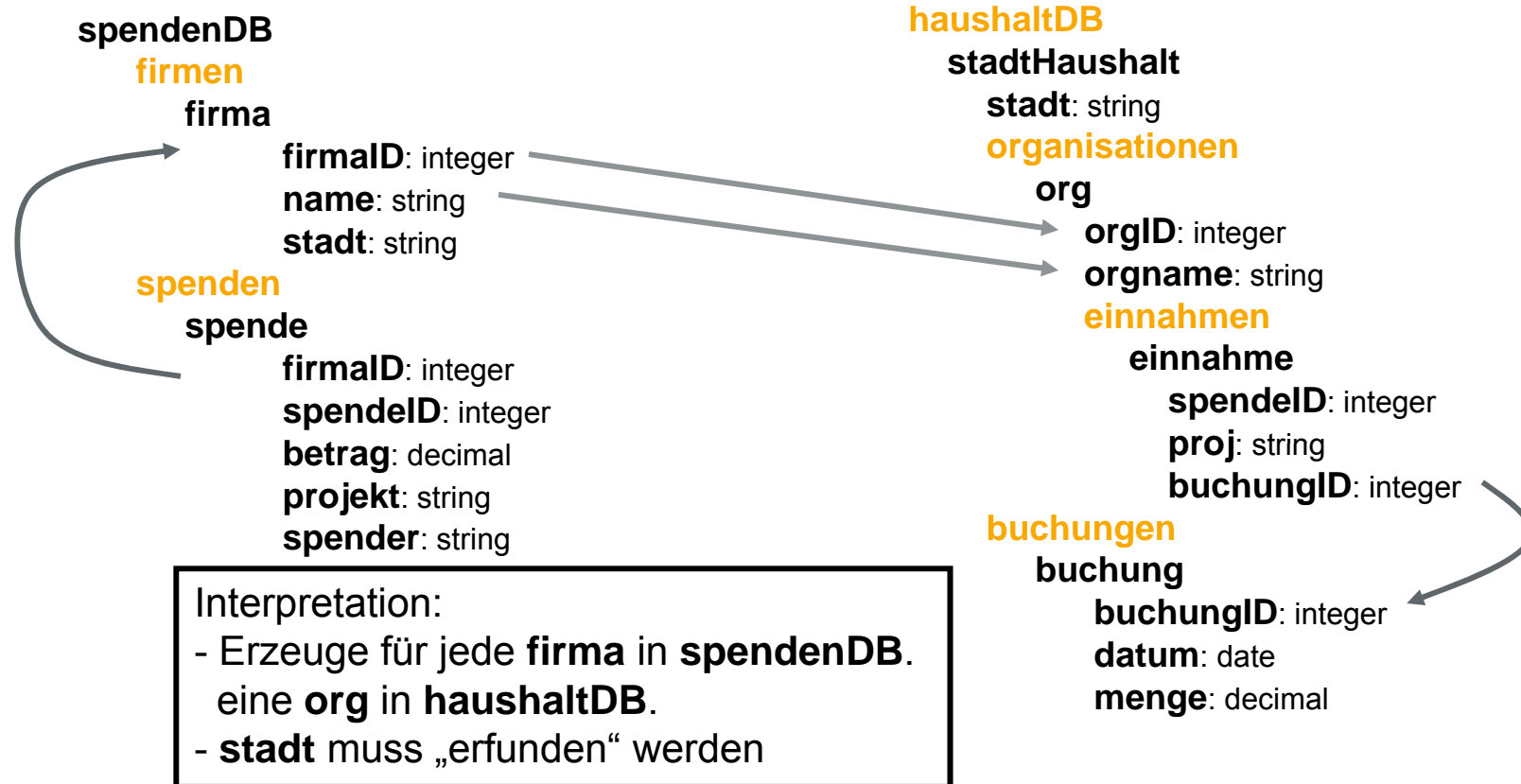
# Mapping – Beispiel

70



# Mapping – Beispiel

71



# „Erfinden“ von Werten

72

Zwei Gründe zum Erfinden: „non-null“ und Identität  
„non-null“ Werte

- Erfundener Wert egal
- Z.B. „unbekannt“ oder „null“ (oder „Berlin“)

ID Werte

- Skolemfunktion:
  - Input: n Werte (beliebige Domäne)
  - Output: bezgl. Input eindeutiger Wert (beliebiger Domäne)
  - Beispiel: Konkatenation aller Inputwerte als String

**firma**

**firmaID**: integer  
**name**: string  
**stadt**: string

**stadt**: string

**organisationen**

**org**

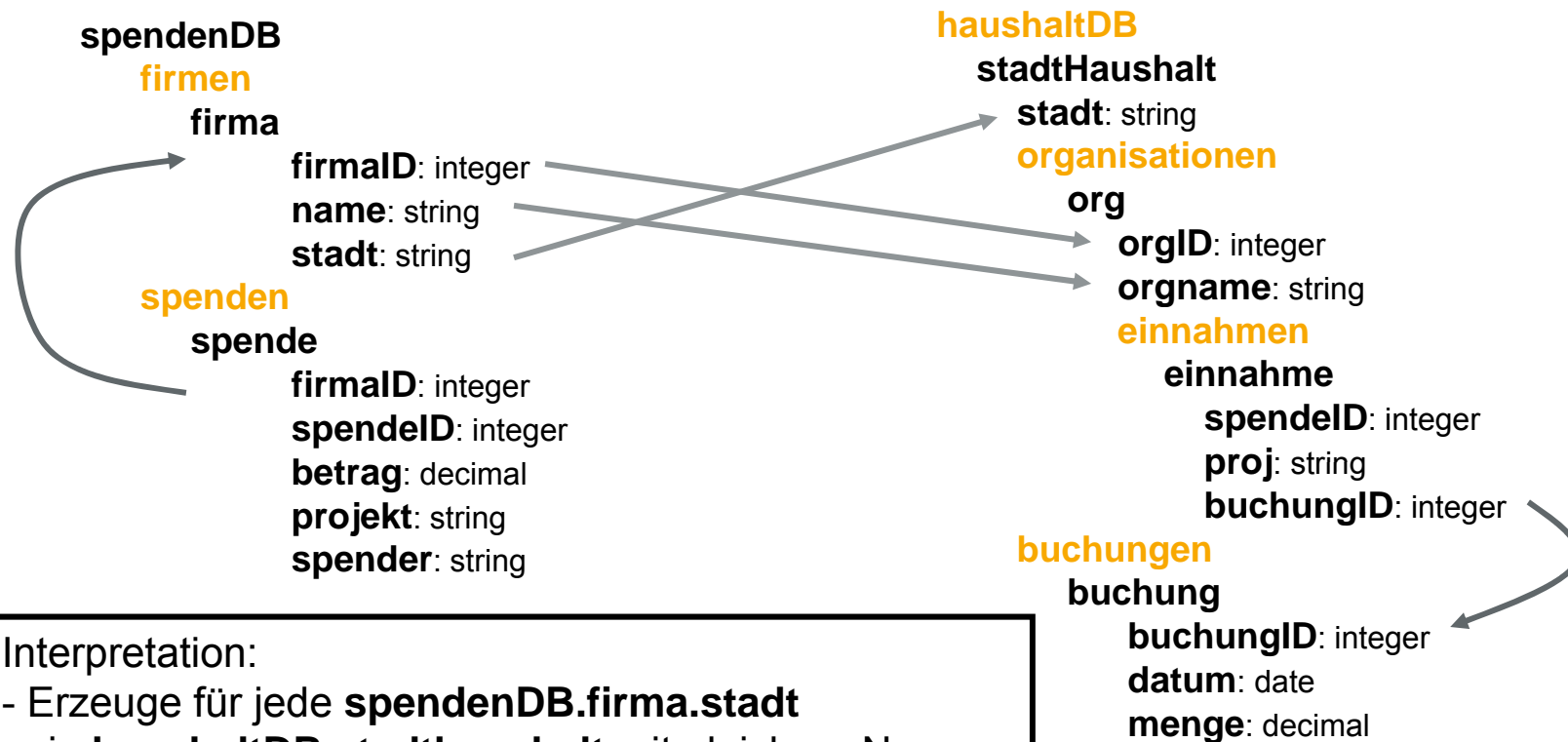
**orgID**: integer  
**orgname**: string

Wert für org.orgID nicht egal, sondern je nach firma.name eindeutig!



# Mapping – Beispiel

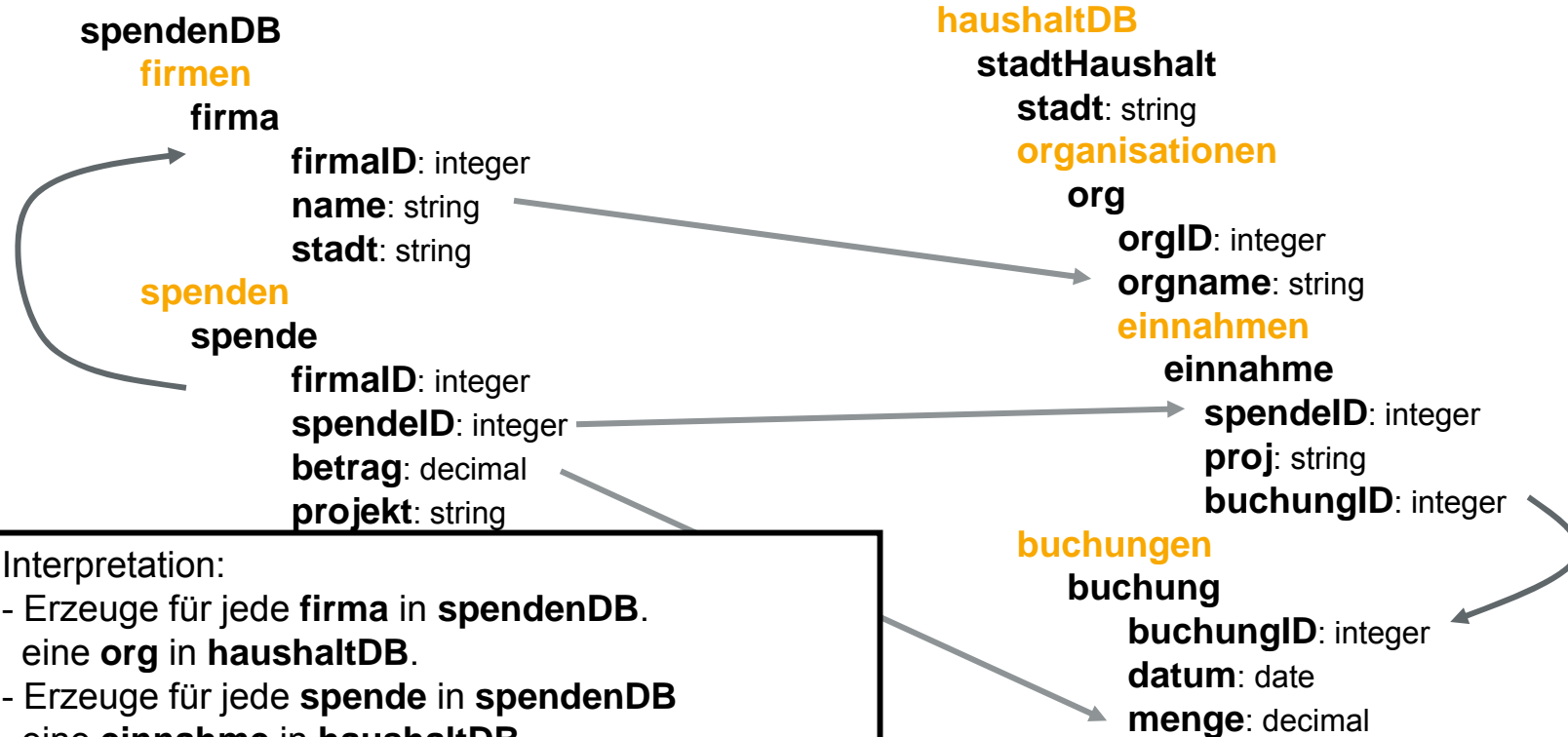
73



Interpretation:  
 - Erzeuge für jede **spendenDB.firma.stadt** ein **haushaltDB.stadthaushalt** mit gleichem Namen.  
 - Gruppriere jede **firma** unter den entsprechenden **stadtHaushalt**.

# Mapping – Beispiel

74

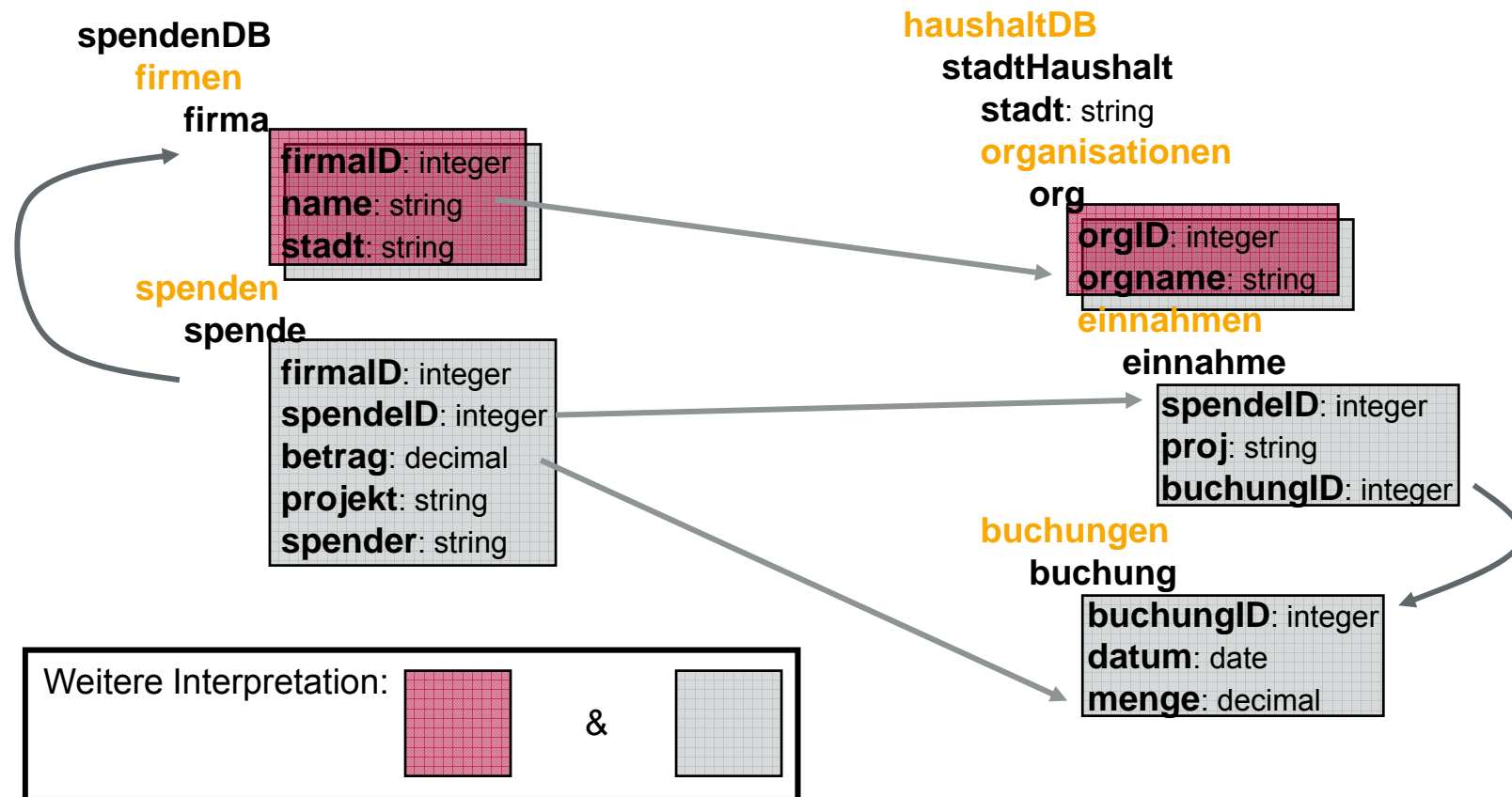


Interpretation:

- Erzeuge für jede **firma** in **spendenDB**. eine **org** in **haushaltDB**.
- Erzeuge für jede **spende** in **spendenDB** eine **einnahme** in **haushaltDB**
- Erzeuge für jede **spende** in **spendenDB** eine **buchung** in **haushaltDB**
- Gruppiere korrekt: Schachtelung & Fremdschlüssel!

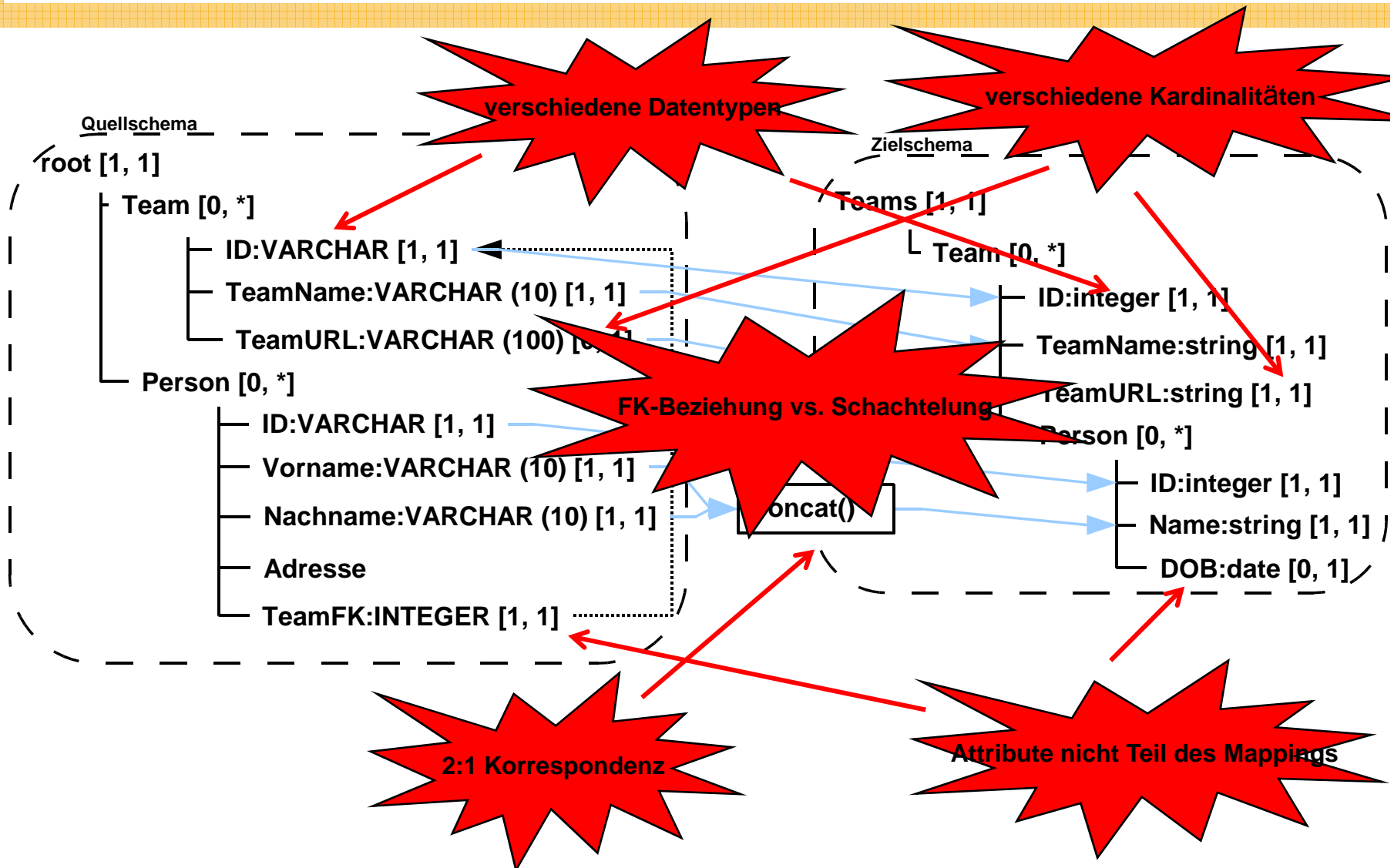
# Mapping – Beispiel

75



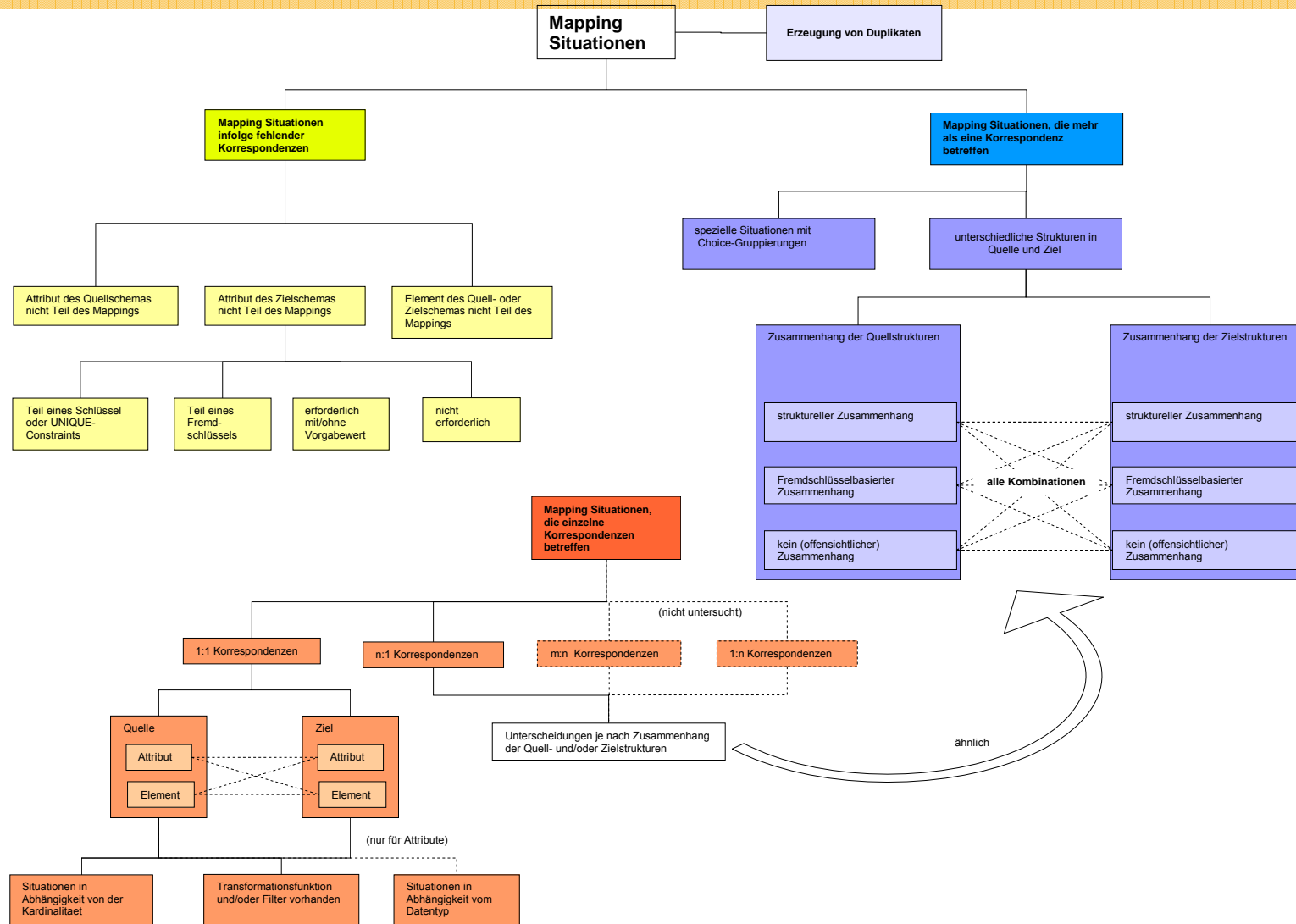
# Verschiedene Mapping Situationen

76



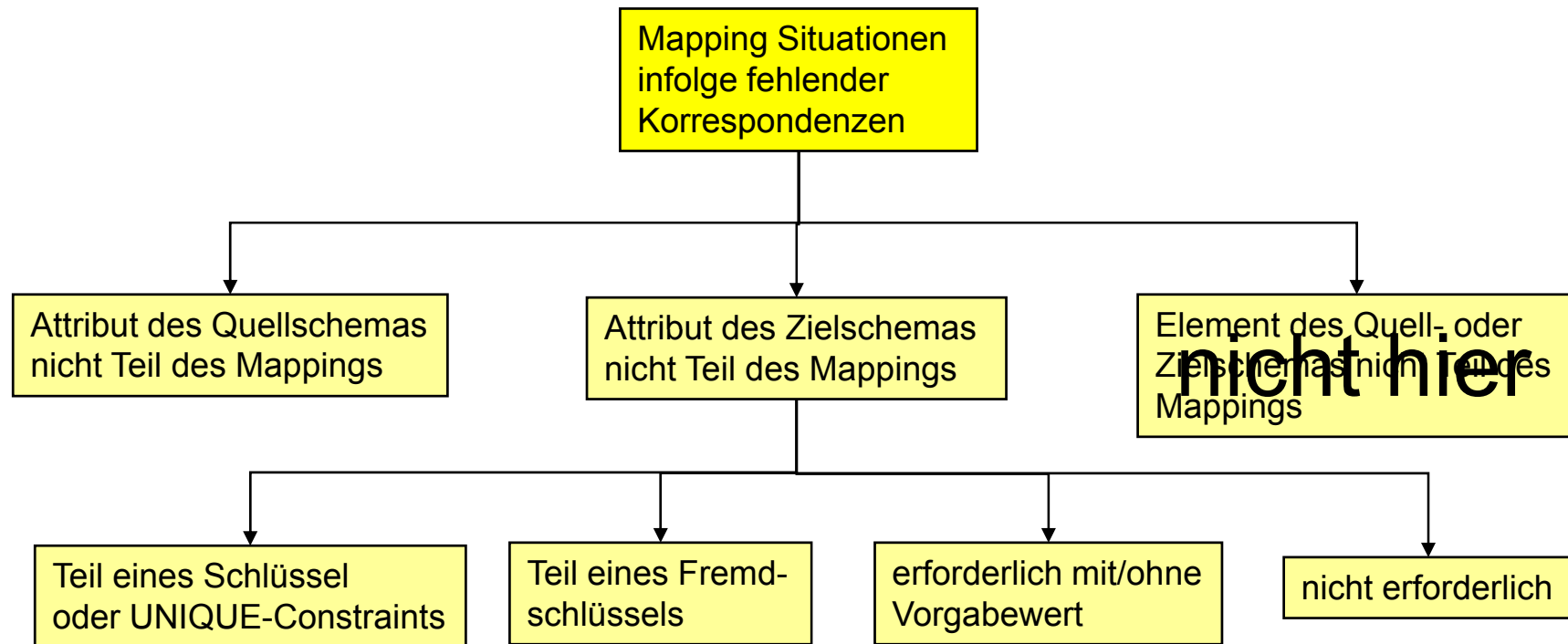
# Klassifikation von Mapping Situationen

77



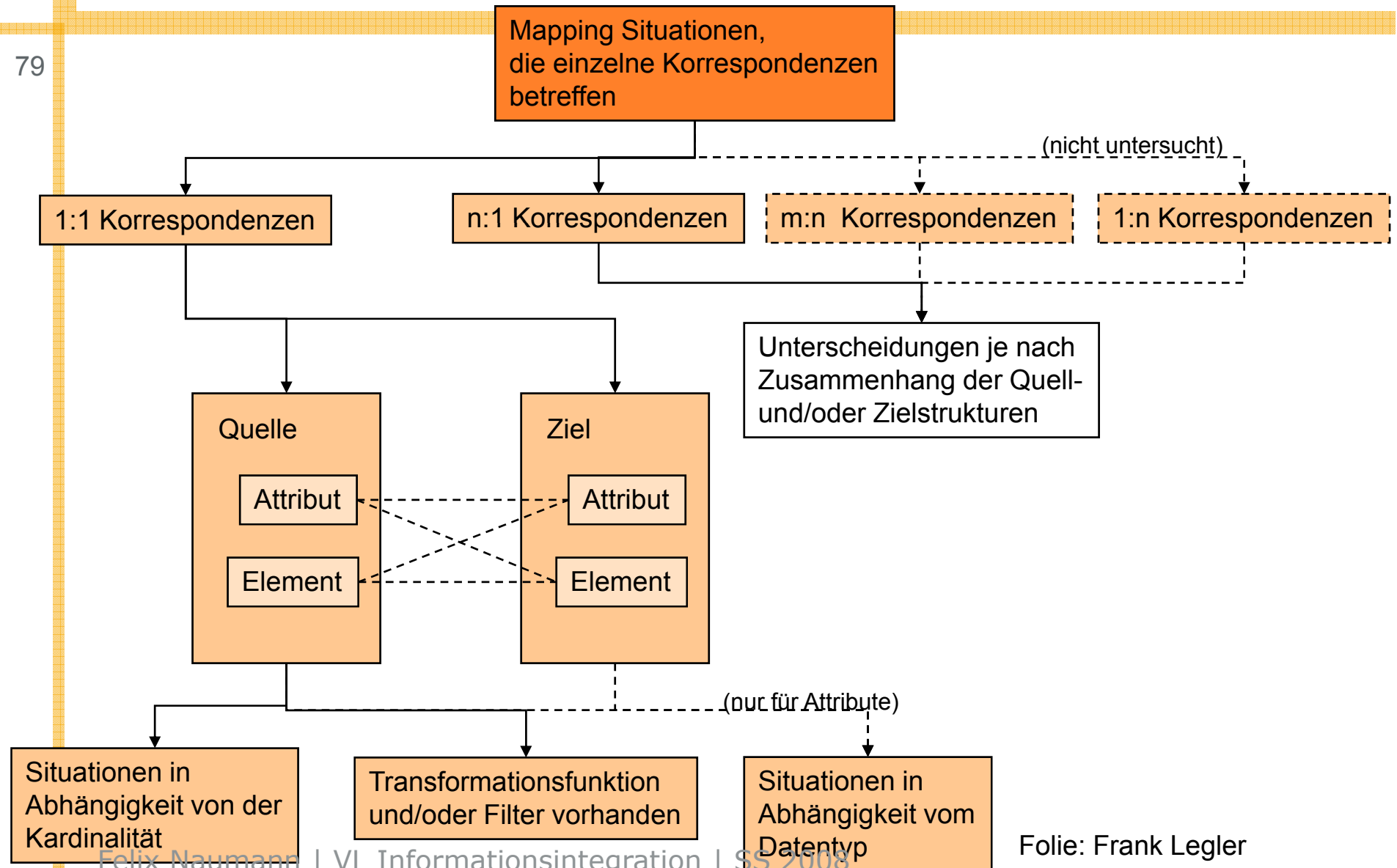
# Klassifikation von Mapping Situationen

78



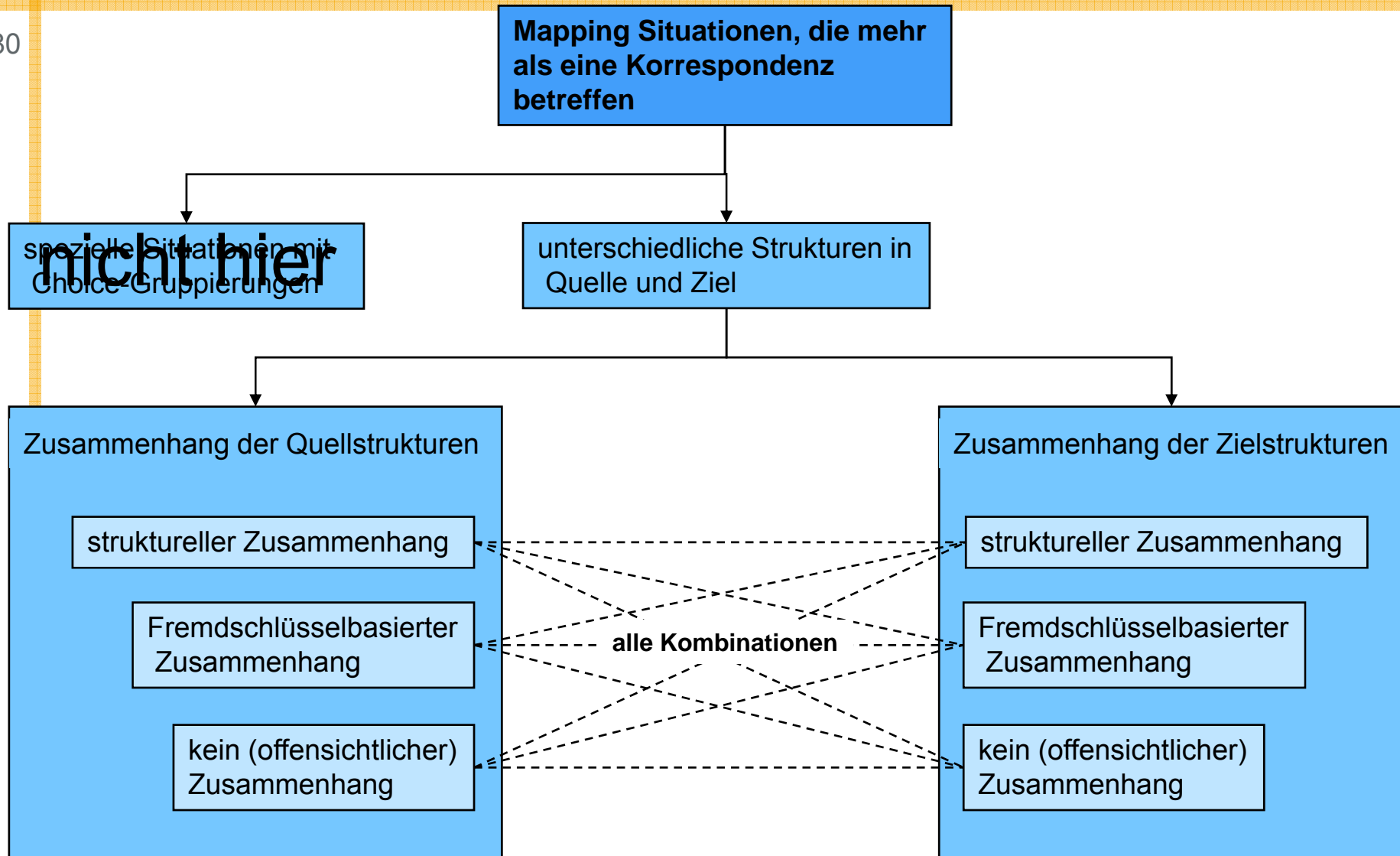
# Klassifikation von Mapping Situationen

79



# Klassifikation von Mapping Situationen

80



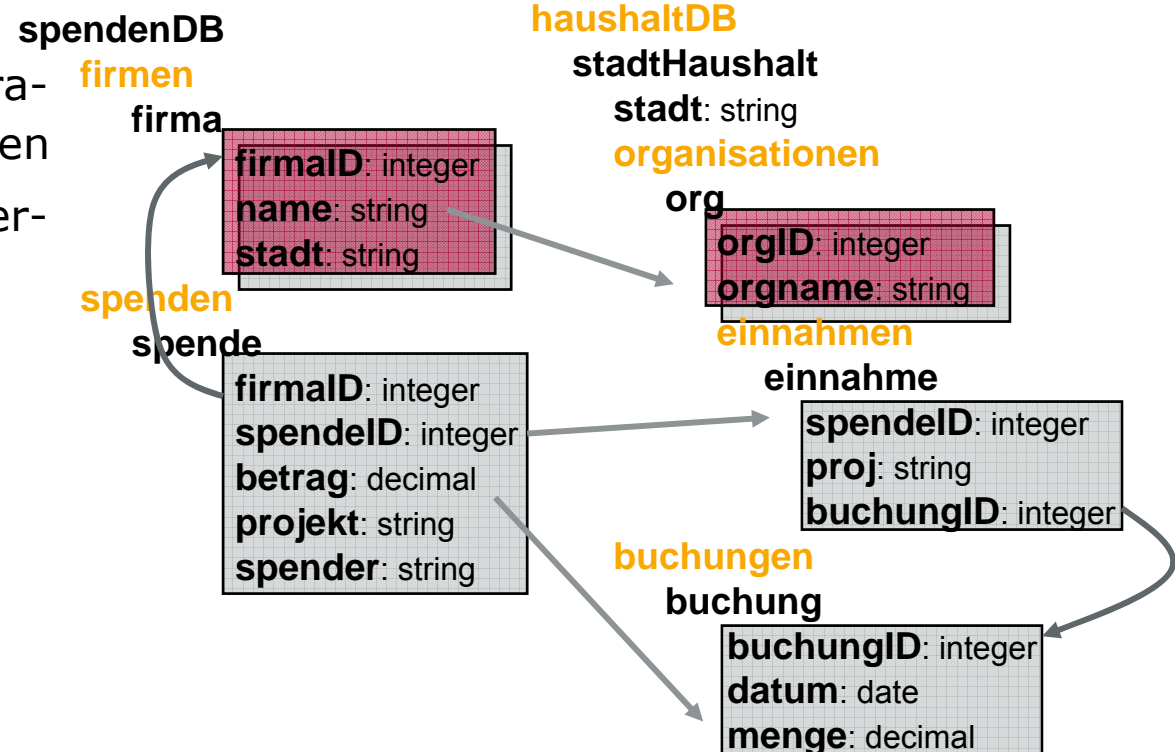


# Mapping – Algorithmus

81

## Drei Schritte

1. Entdeckung von intra-Schema Assoziationen
2. Entdeckung von inter-Schema logischen Mappings
3. Anfrage-erzeugung



# Entdeckung von Assoziationen

82

## Schritt 1

- Intra-schema Assoziationen zwischen Schemaelementen
- Relationale Sichten enthalten maximale Gruppen assoziierter Elemente
- Jede Sicht repräsentiert eine eigene „Kategorie“ an Daten der Datenquelle
- Unabhängig vom Mapping (aber beschränkt auf „gemappte“ Elemente)



Quelle: [FHP+02]

# Entdeckung von Assoziationen

83

Start: Alle „primären“ Pfade (*primary paths*)

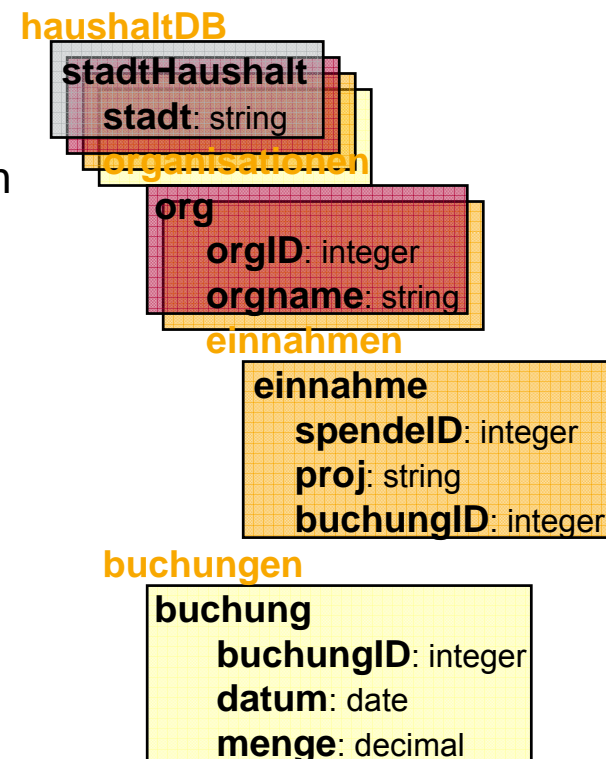
- Assoziationen im Schema ohne Integritätsbedingungen

Relationale Schemas

- Jede Relation entspricht einem primären Pfad

Geschachtelte Schemas

- Attribute einer Ebene
- Attribute geschachtelter Ebenen



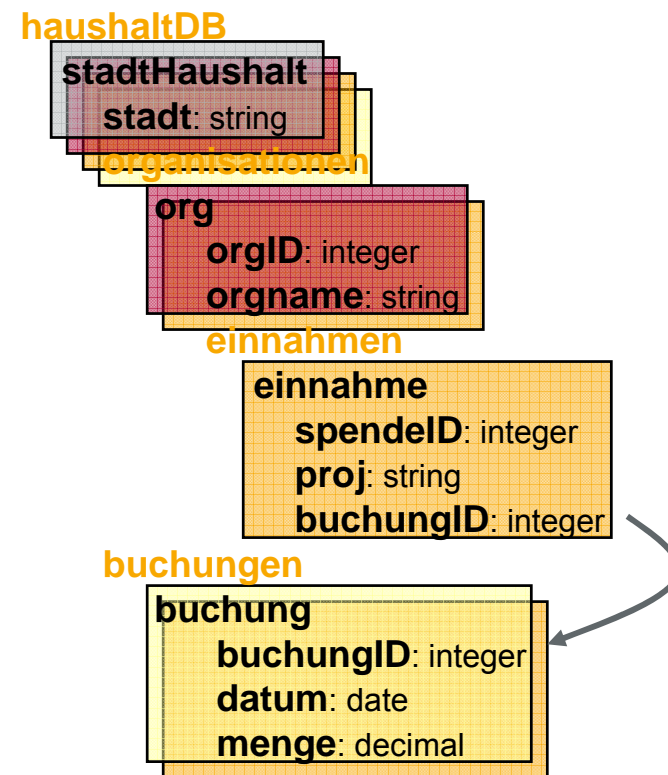
# Entdeckung von Assoziationen

84

Betrachte nun Schlüssel / Fremdschlüssel (ICs)

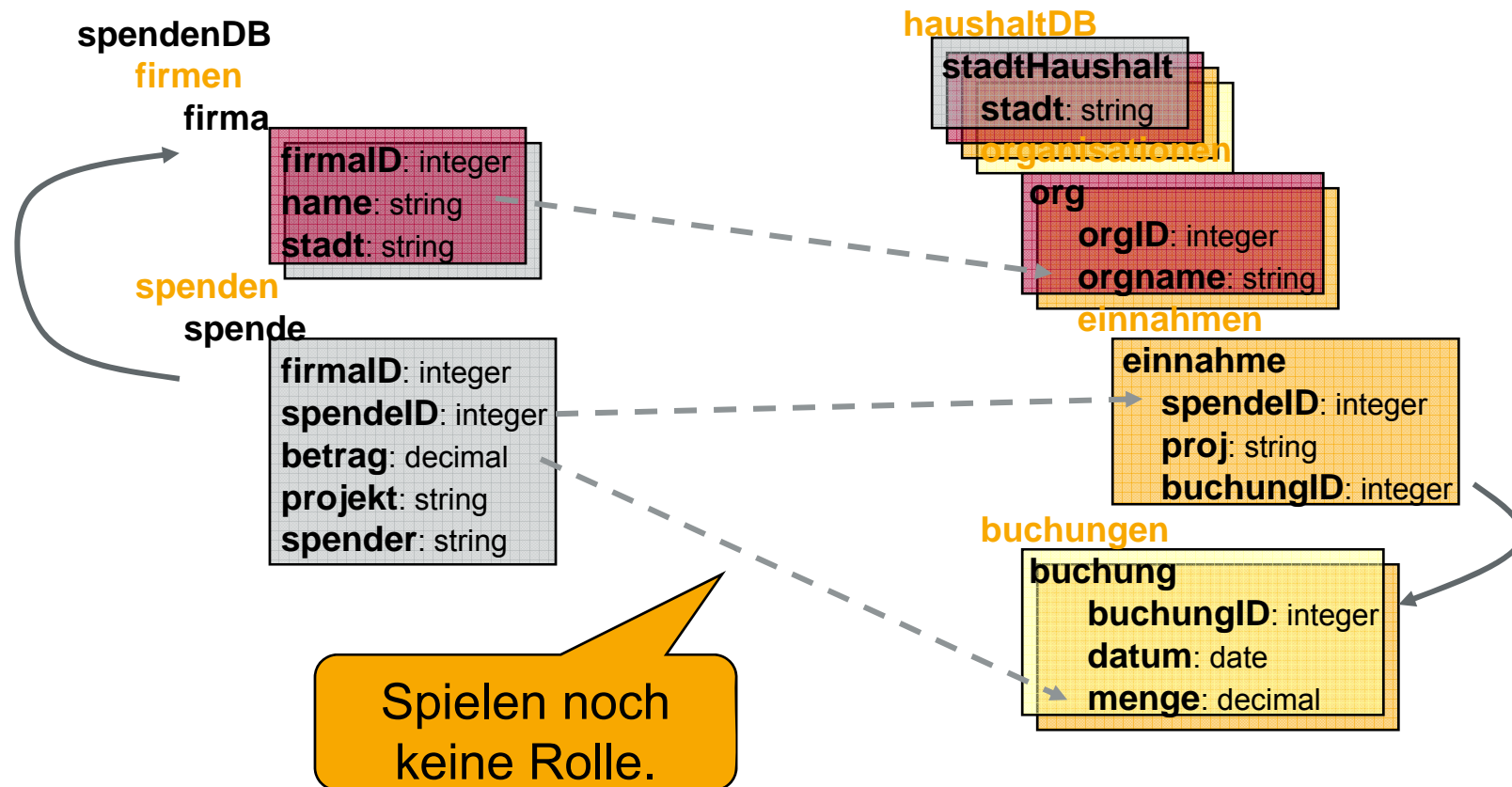
Logische Relation

- Erweitere jeden primären Pfad durch „Verfolgen“ der ICs (*chase*)



# Entdeckung von Assoziationen

85

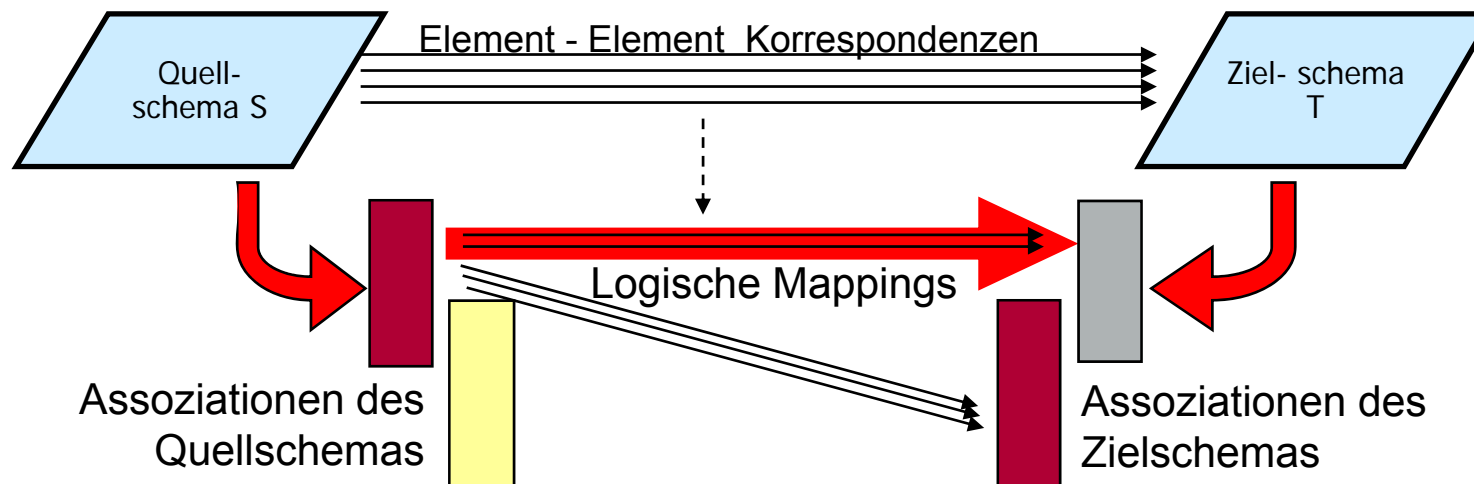


# Entdeckung von logischen Mappings

86

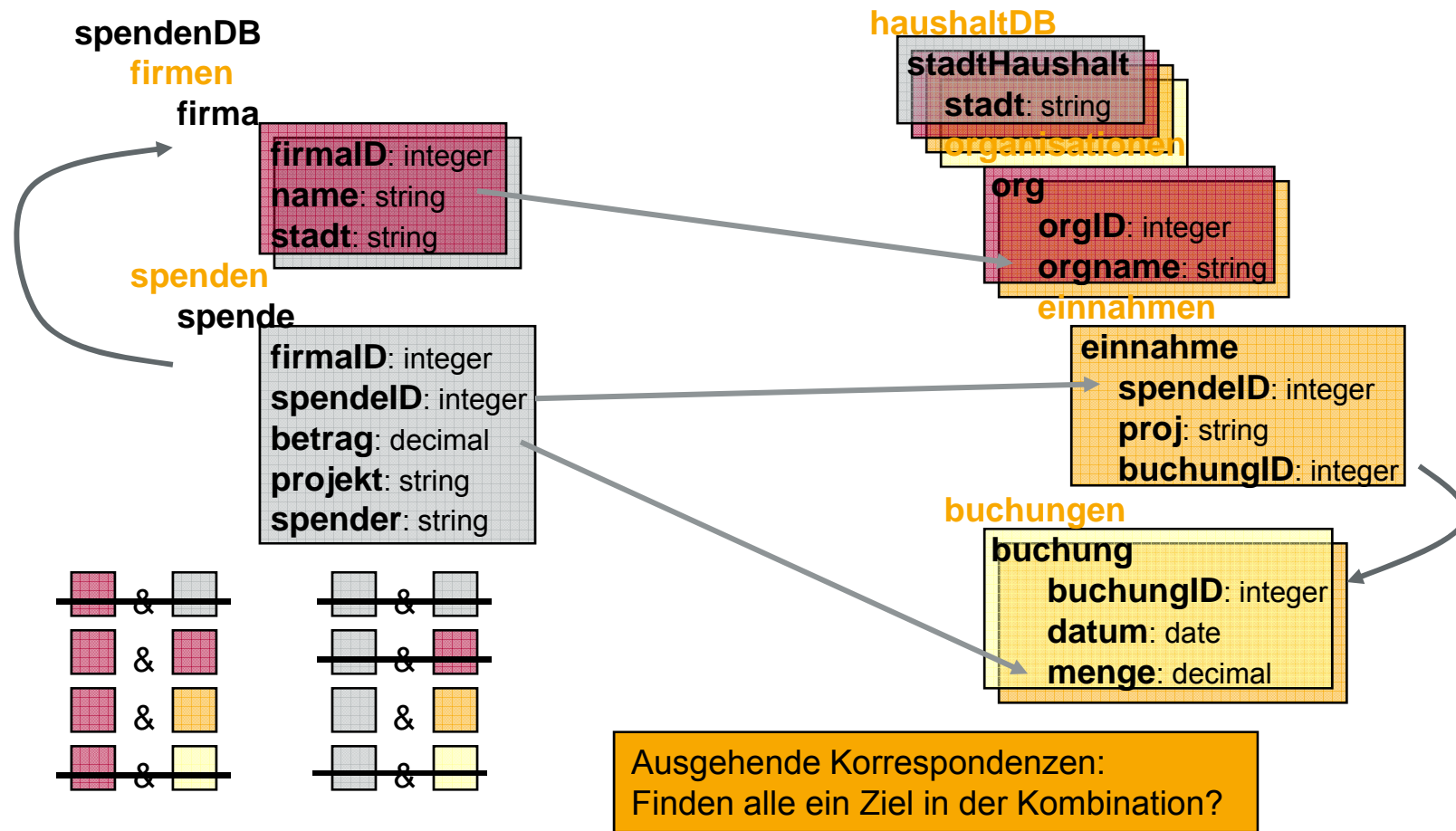
## Schritt 2

- Entdecke logische Mappings zwischen Quell- und Zielschema
- Betrachte alle Kombinationen aus Assoziationen des Quellschemas und Assoziationen des Zielschemas
  - Unter Berücksichtigung aller Korrespondenzen (sofern Korrespondenzen zwischen ihnen überhaupt existieren)



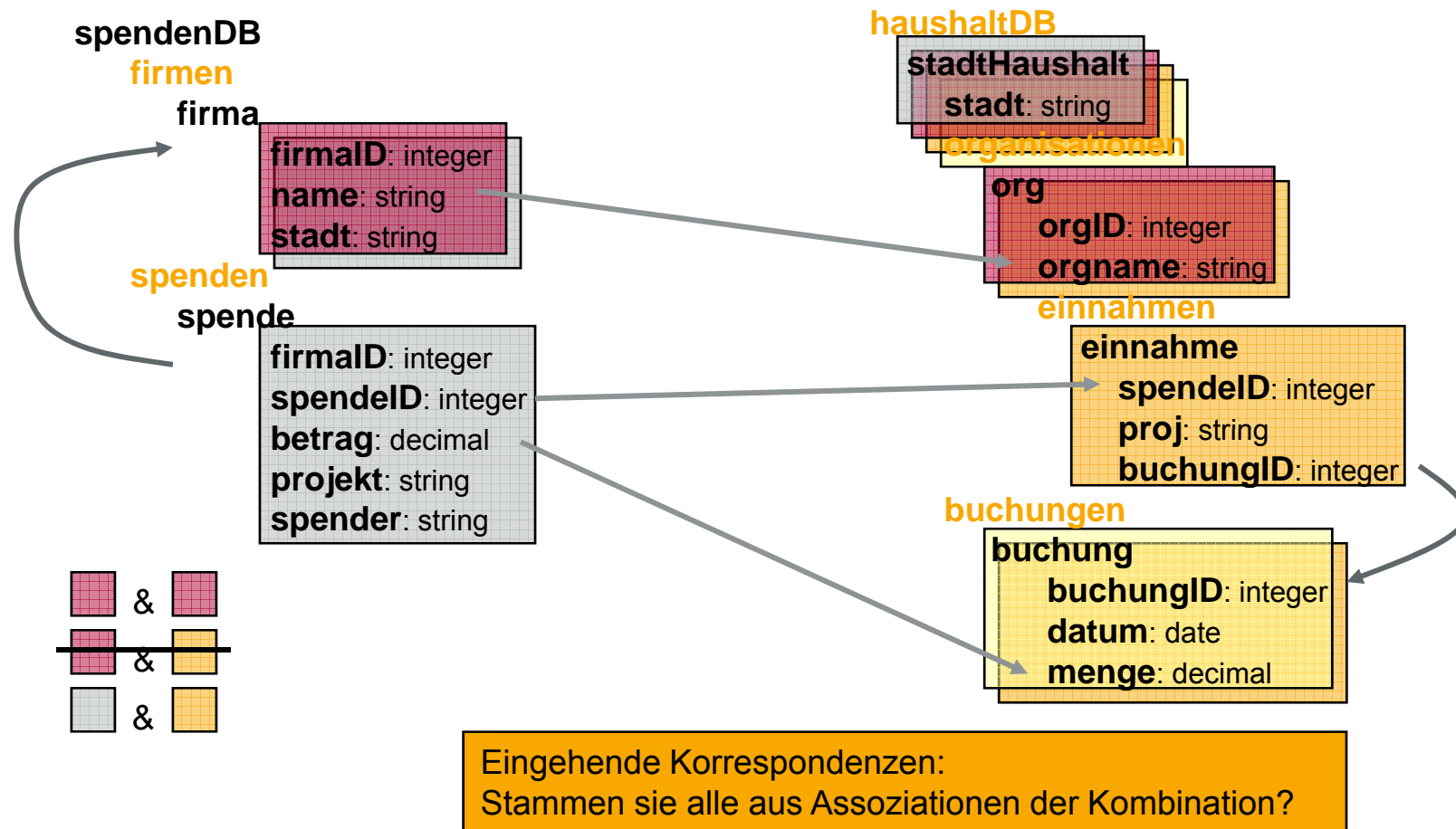
# Entdeckung von logischen Mappings

87



# Entdeckung von logischen Mappings

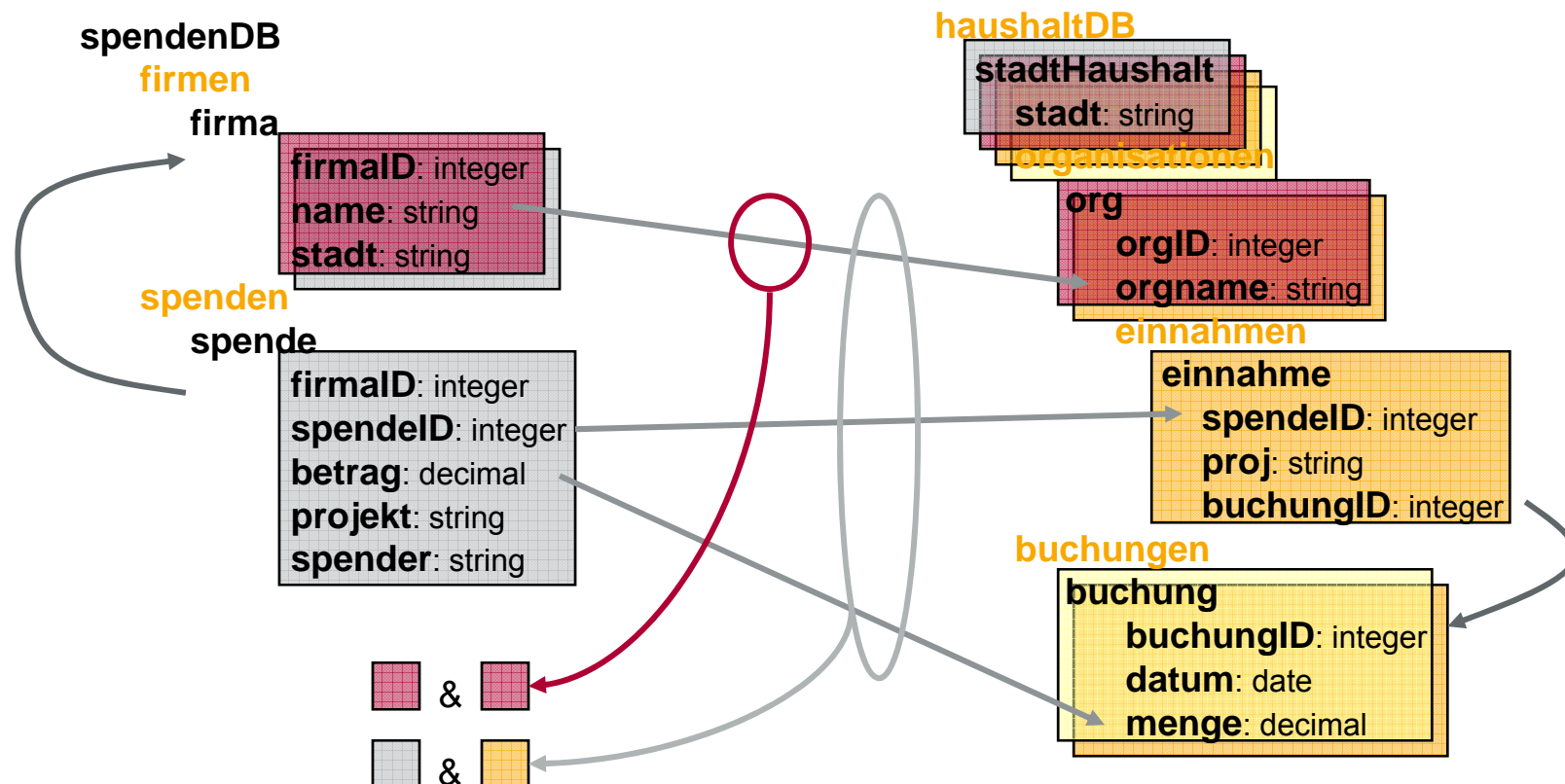
88





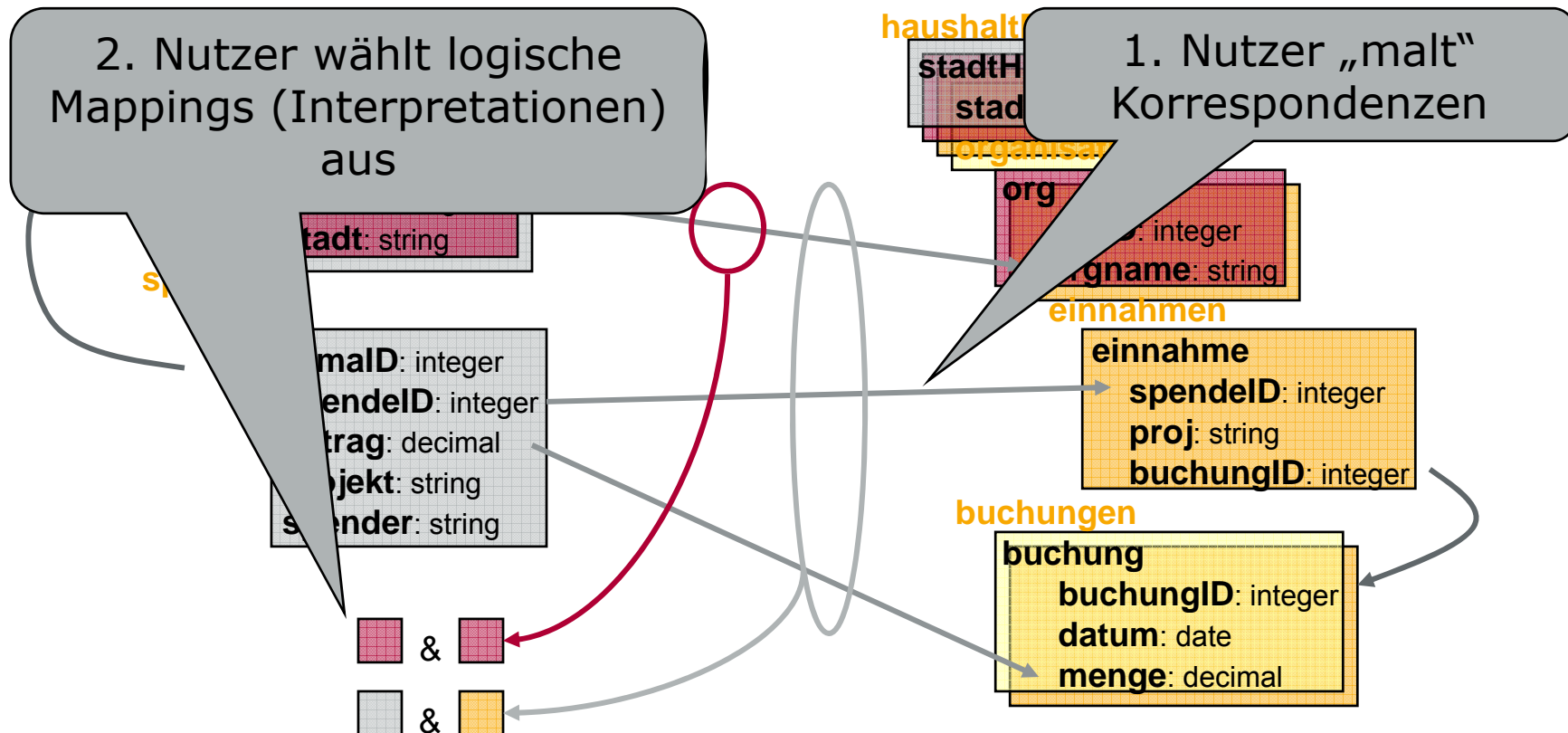
# Entdeckung von logischen Mappings

89



# Input des Nutzers bzw. Domänenexperten

90

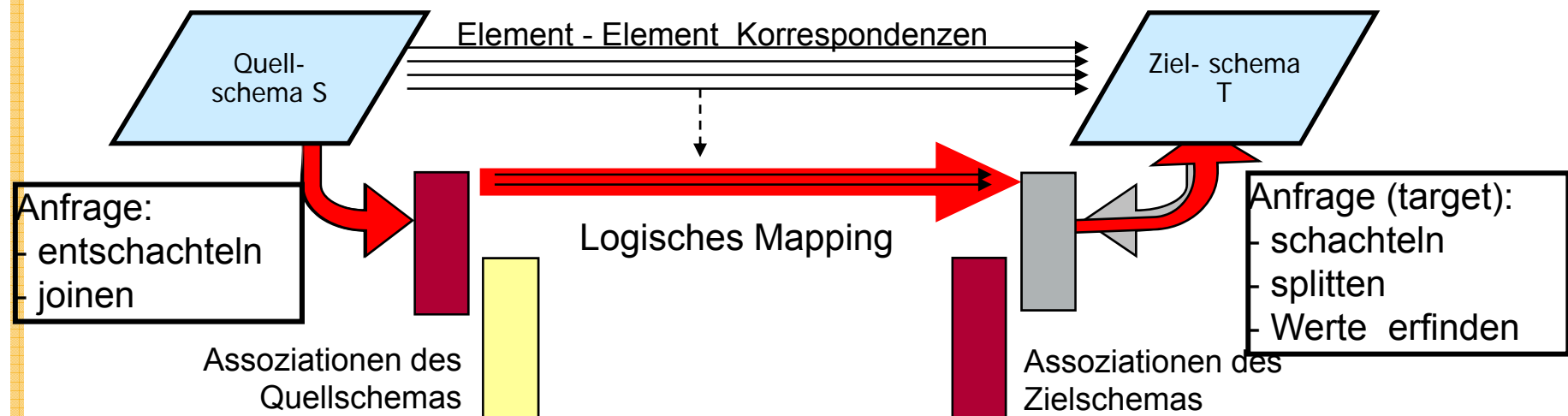


# Erzeugung der Anfragen

91

## Schritt 3

- Erzeuge für jedes (ausgewählte) logische Mapping eine Anfrage
  - Auswahl und verknüpfen der entsprechenden Quelldaten
  - Generierung der entsprechenden Zieldaten



# Erzeugung der Anfragen

92

## 2 Probleme

- Erfinden von Werten
  - NULL-Werte nicht immer ausreichend.
  - Schlüssel und passende Fremdschlüssel müssen erzeugt werden.
- Schachtelung
  - Es soll nicht eine logische Relation (flach) materialisiert werden, sondern geschachtelte Strukturen.
  - Es muss entsprechend gruppiert werden.

# Erfinden neuer Werte (Wdh.)

93

- Logische Relation: buchungen
- Wert für buchungID wird nicht gefüllt
  - Entweder: Egal
  - Oder: Not-null: Dann Default-Wert, z.B. „null“
  - Oder ID: Dann eindeutigen Werte erzeugen
  - Skolemfunktion, basierend auf allen Werten des Mappings dieser logischen Relation

$$\text{buchungID} = \text{Sk}(\text{betrag})$$

spendenDB

**spenden**

spende

<b>firmaID</b> : integer
<b>spendeID</b> : integer
<b>projekt</b> : string
<b>betrag</b> : decimal
<b>spender</b> : string

haushaltDB

**einnahmen**

einnahme

**spendeID**: integer

**proj**: string

**buchungID**: integer

**buchungen**

buchung

**buchungID**: integer

**datum**: date

**menge**: decimal

# Erfinden neuer Werte

94

- Logische Relation: einnahmen, buchungen
- buchungID-Attribute haben keine Korrespondenz
  - Assoziationen gingen verloren
  - Also neue ID erfinden
  - Wieder: Skolemfunktion basierend auf allen Werten des Mappings dieser logischen Relation
  - Trick wie gehabt: Gleiche Funktion für Schlüssel und Fremdschlüssel

spendenDB

spenden

spende

```

firmaID: integer
spendeID: integer
projekt: string
betrag: decimal
spender: string
    
```

haushaltDB

einnahmen

einnahme

```

spendeID: integer
proj: string
buchungID: integer
    
```

buchungen

buchung

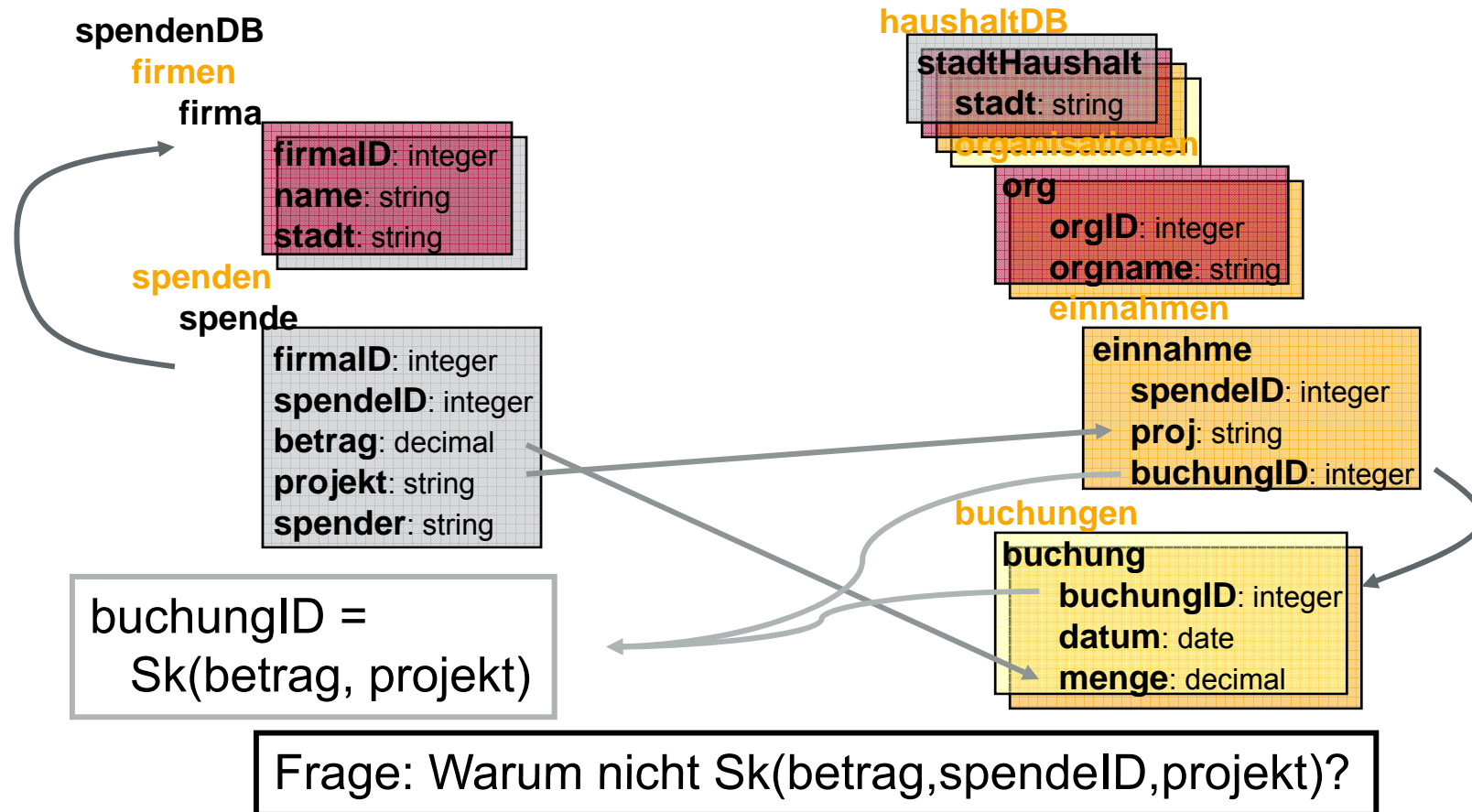
```

buchungID: integer
datum: date
menge: decimal
    
```

buchungID =  
Sk(spendeID, projekt, betrag)

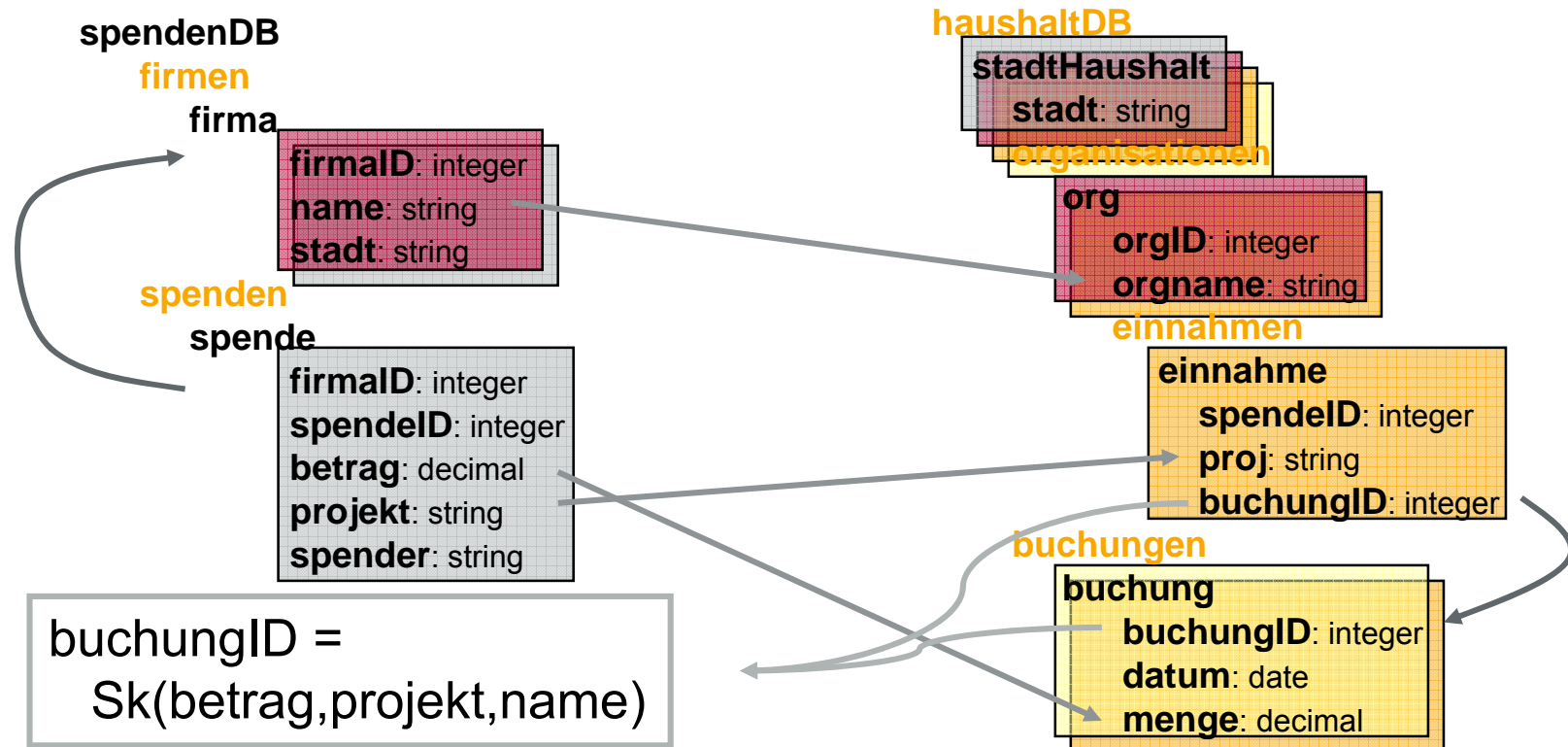
# Erfinden neuer Werte

95



# Erfinden neuer Werte

96



Jetzt erst recht: Warum nicht Sk(betrag, projekt, name, firmID)?



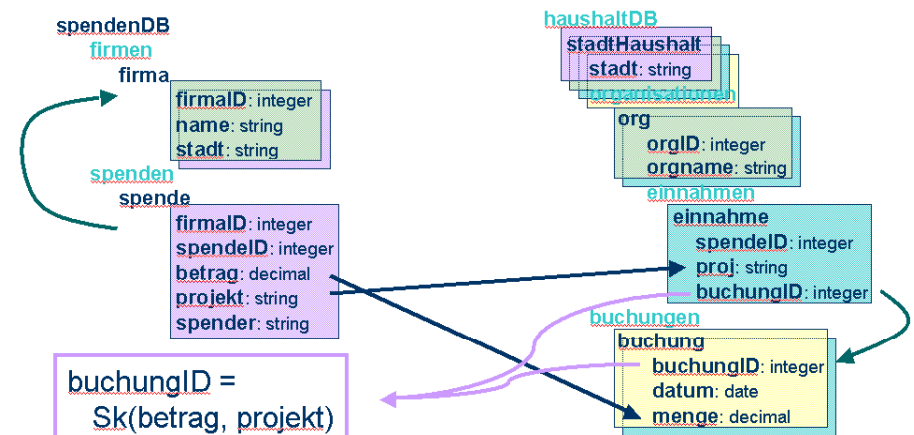
# Erfinden neuer Werte

97

Gegenfrage: Warum nicht gleiche ALLE ungemappten Werte?

Vier Antworten (von Lucian Popa)

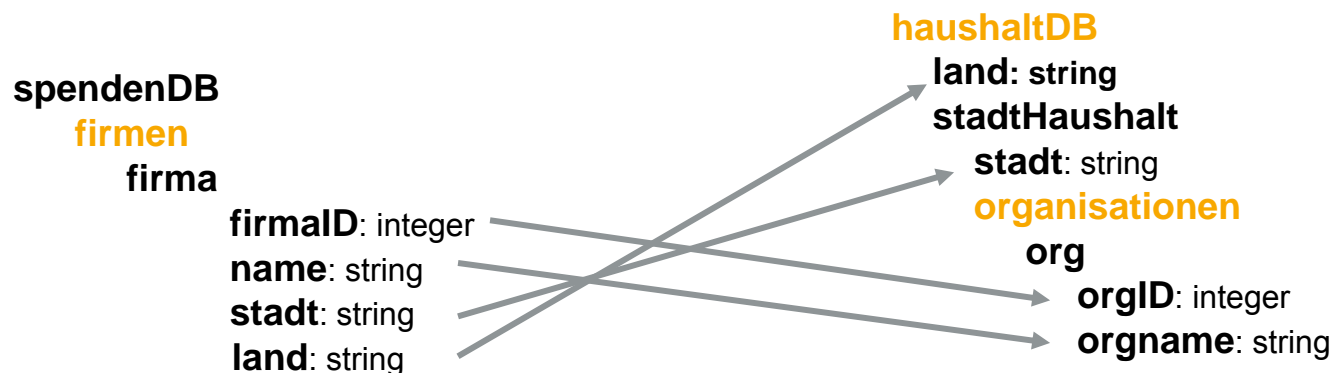
- Prinzipiell ginge das, aber
- „Duplikate“ könnten entstehen
  - 2 Spenden, mit gleichem Betrag und gleichem Projekt aber von unterschiedlichen Firmen
- firmenID könnte für Nutzer völlig uninteressant sein.
  - Mapping hätte dann die Rolle einer Projektion, in der Duplikate fehl am Platz wären
  - Sie wären sogar völlig unerklärlich für den Nutzer.
- Minimalität
  - Beide Varianten sind in Ordnung, aber eine ist kleiner.



# Gruppierung

98

- Alle Attribute erhalten Werte
- Aber:
  - Assoziationen könnten verloren gehen.
  - Neue (falsche) Assoziationen könnten erzeugt werden.
- Deshalb: Gruppierung notwendig
- Trick: Virtuelle ID Generierung mittels Skolemfunktion basierend auf allen (gemapten) Werten hierarchisch über der aktuellen Relation.
  - Im Beispiel: Jedes Tupel `haushaltDB` erhält ID `Sk(land, stadt)`
  - Jedes weitere Tupel aus `firmen` mit gleichen Werten für `stadt` und `land` errechnet gleiche ID und wird unter gleichem Element geschachtelt.



# Erzeugung der Anfragen

99

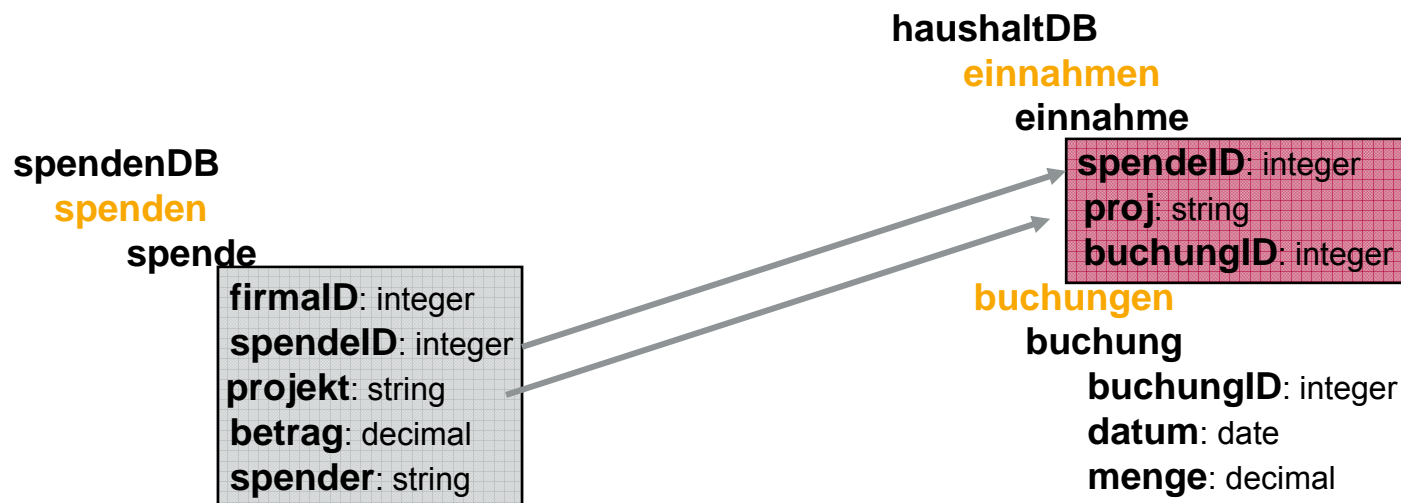
- Implementationsspezifisch
- Abhängig von Datenmodell
  - Relational → Relational: SQL
  - Relational → XML: SQL/XML
    - ◇ Schachtelung und tagging des Ergebnisses
  - XML → Relational: XQuery oder XSLT
    - ◇ Tags weglassen
  - XML → XML: XQuery oder XSLT
- In Clio
  - Erzeugung proprietärer Regeln
  - Dann: Übersetzung der Regeln in jeweilige Anfragesprache

# Erzeugung der Anfragen

100

Erzeugung proprietärer Regeln in Clio

- for x in spenden
- let a = x.spende.spendeID, b = x.spende.projekt
- return <einnahme = <spendeID = a, proj = b, buchungID = null>>
- in einnahmen

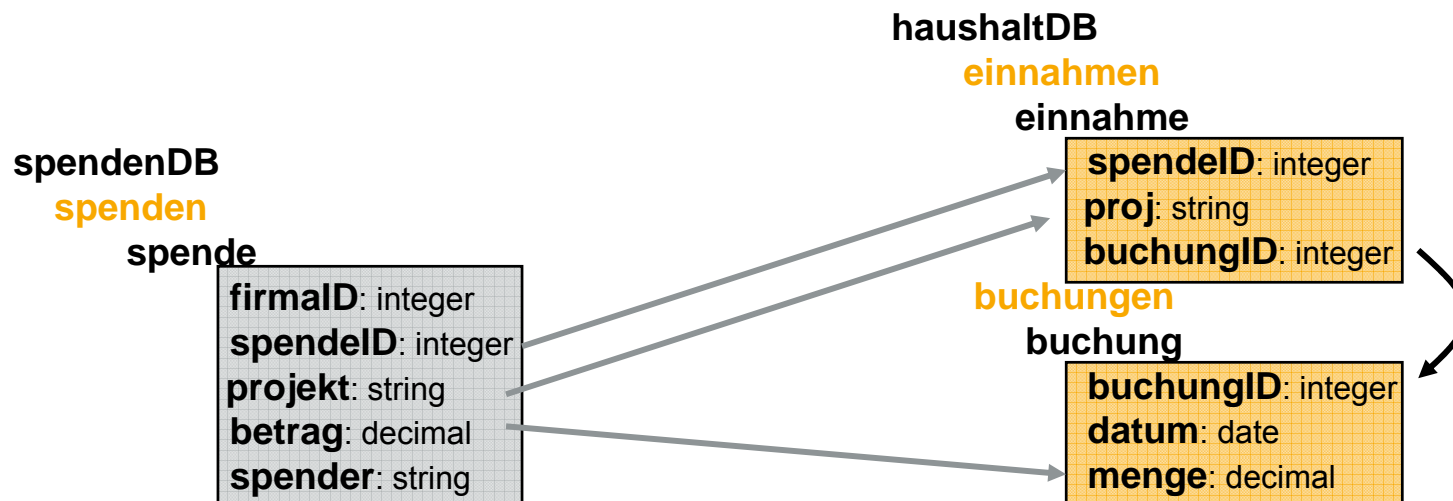


# Erzeugung der Anfragen

101

Erzeugung proprietärer Regeln in Clio

- for x in spenden
- let a = x.spende.spendeID, b = x.spende.projekt, c = x.spende.betrag
- return <einnahme = <spendeID = a, proj = b, buchungID = Sk(a,b,c)>>
- in einnahmen,
- <buchung = <buchungID = Sk(a,b,c), datum = null, menge = c>>
- in buchungen



# Erzeugung der Anfragen – XQuery

102

```

for x in spenden
let a = x.spende.spendeID, b = x.spende.projekt, c = x.spende.betrag
return <einnahme = <spendeID = a, proj = b, buchungID = Sk(a,b,c)>> in einnahmen,
       <buchung = <buchungID = Sk(a,b,c), datum = null, menge = c>> in buchungen
  
```

```

LET $doc0 := document("input XML file goes here")
RETURN <haushaltDB> {
  distinct-values (
    FOR $x0 IN $doc0/spendenDB/spende
    RETURN <einnahme>
      <spendeID> { $x0/spendeID/text() } </spendeID>
      <proj> { $x0/projekt/text() } </proj>
      <buchungID> { "Sk32(", $x0/betrag/text(), ", ", $x0/spendeID/text(), ", ", $x0/projekt/text(), ")" } </buchungID>
      </einnahme> ) }
  { distinct-values (
    FOR $x0 IN $doc0/spendenDB/spende
    RETURN <buchung>
      <buchungID> { "Sk32(", $x0/betrag/text(), ", ", $x0/spendeID/text(), ", ", $x0/projekt/text(), ")" } </buchungID>
      <menge> { $x0/betrag/text() } </menge>
      </buchung> ) }
</haushaltDB>
  
```

# Erzeugung der Anfragen – SQL

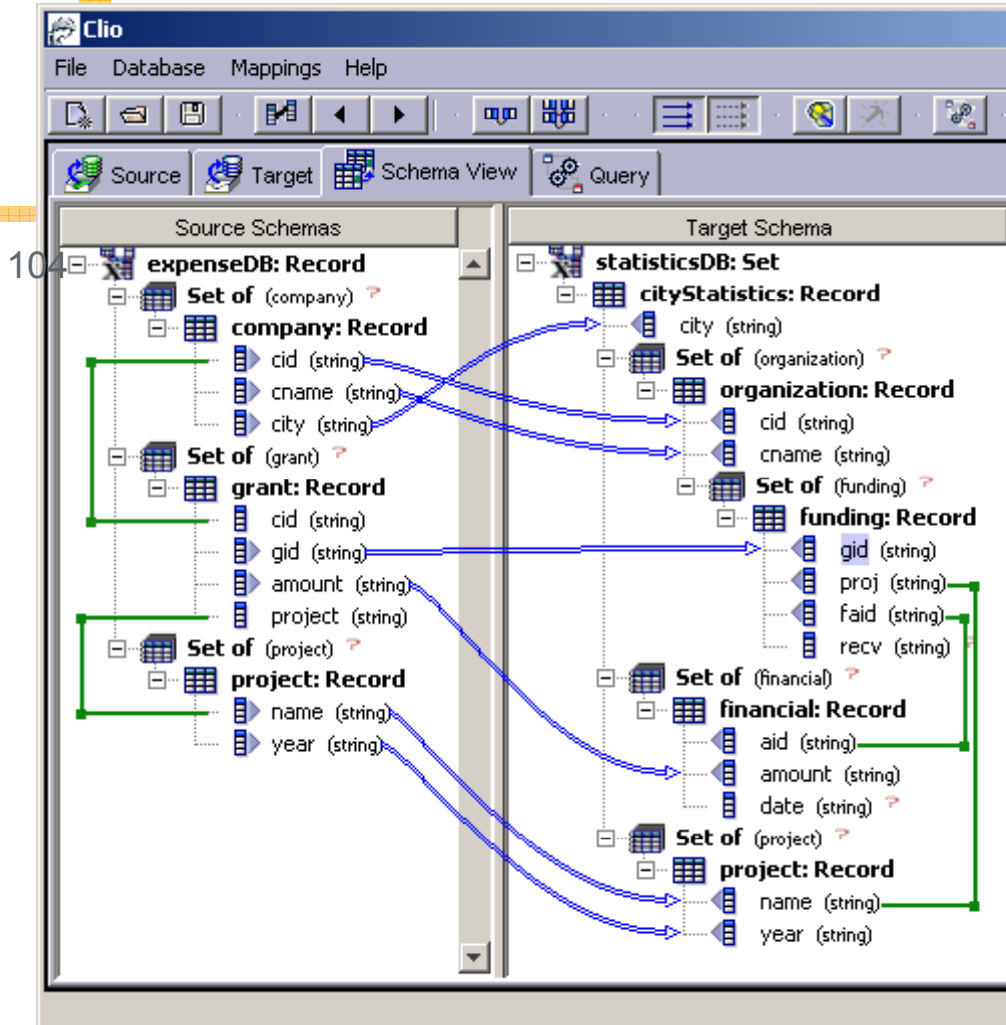
103

```

for x in spenden
let a = x.spende.spendeID, b = x.spende.projekt, c = x.spende.betrag
return <einnahme = <spendeID = a, proj = b, buchungID = Sk(a,b,c)>> in einnahmen,
      <buchung = <buchungID = Sk(a,b,c), datum = null, menge = c>> in buchungen
  
```

```

-- CREATE VIEW einnahme AS
SELECT
  x0.spendeID AS spendeID,
  x0.projekt AS proj,
  RTRIM('Sk32(' || CHAR(x0.betrag) || ',' || (CHAR(x0.spendeID) || ',' || CHAR(x0.projekt) || '')) AS buchungID
FROM
  spendenDB.spende x0
-----
-- CREATE VIEW buchung AS
SELECT
  RTRIM('Sk32(' || CHAR(x0.betrag) || ',' || CHAR(x0.spendeID) || ',' || CHAR(x0.projekt) || '')) AS buchungID,
  x0.betrag AS menge
FROM
  spendenDB.spende x0
  
```



Clio

File Database Mappings Help

Source Target Schema View Query

XQuery XSL SQL Query

```

RETURN
  $x0L1/project/text() = $x1L1/name/text() AND
  $x2L1/cid/text() = $x0L1/cid/text() AND
  $x2/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>,
      <faid> "SK267(", $x0L2/project/text(), " ,
    </funding> )
  </organization> ,
  distinct (
  FOR
    $x0L1 IN $doc/expenseDB/grant ,
    $x1L1 IN $doc/expenseDB/project ,
    $x2L1 IN $doc/expenseDB/company
  WHERE

```

Preview Execute Query Copy to Clipboard

No File



# Vorschau: Global as View / Local as View

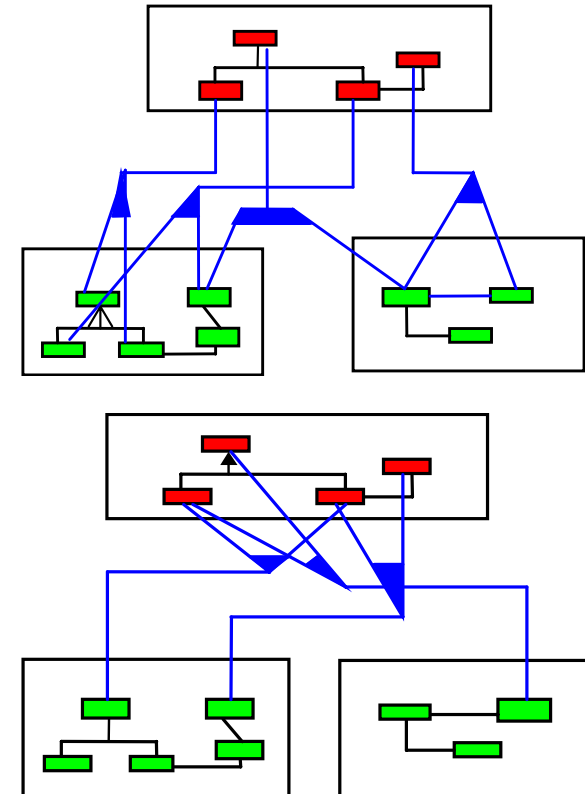
105

## Global as View (GaV)

- Relationen des globalen Schemas werden als Sichten auf die lokalen Schemas der Quellen ausgedrückt.

## Local as View (LaV)

- Relationen der Schemas der Quellen werden als Sichten auf das globale Schema ausgedrückt.



# Einbettung in GaV/LaV

106

Ergebnis des Schema Mapping Prozess sind Anfragen

- Eine oder mehrere Anfragen pro Assoziation

Relationales Modell:

- Jede Ziel-Relation entspricht einem primären Pfad
- Verwerfe Fremdschlüssel
- Jede Ziel-Relation entspricht einer Assoziation
- Jede Anfrage (oder Vereinigung von Anfrage) liefert Daten für eine Relation, also wie GaV.

Schema Mapping leistet aber mehr!

- Erzeugung von Daten für mehrere Zielrelationen zugleich
- Unter Berücksichtigung von Integritätsbedingungen in Quell- und Zielschema
- „Kombination von LaV und GaV“: GLaV

# Schematische Heterogenität

107

## Struktur

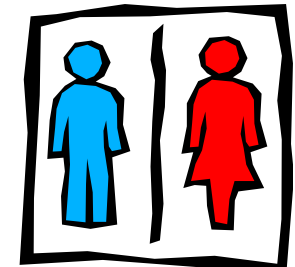
- Modellierung
  - Relation vs. Attribut
  - Attribut vs. Wert
  - Relation vs. Wert
- Benennung
  - Relationen
  - Attribute
- Normalisiert vs. Denormalisiert
- Geschachtelt vs. Fremdschlüssel
- Geschachtelt vs. Flach

} Higher-order  
Mappings

} bisher

# High-order Mappings

108



```
Männer( Id, Vorname, Nachname)
Frauen( Id, Vorname, Nachname)
```

Relation vs. Attribut

```
Person( Id, Vorname,
        Nachname, männlich,
        weiblich)
```

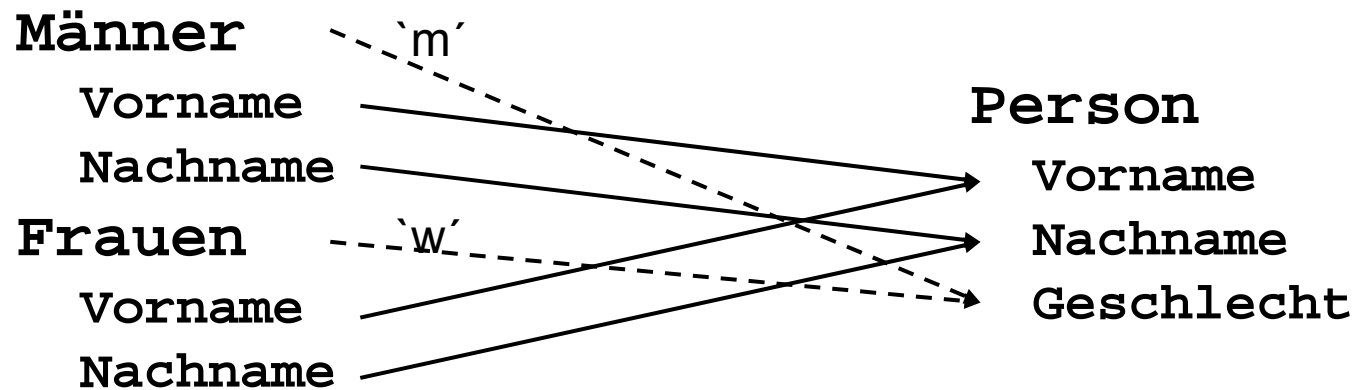
Relation vs. Wert

```
Person( Id, Vorname,
        Nachname, Geschlecht)
```

Attribut vs. Wert

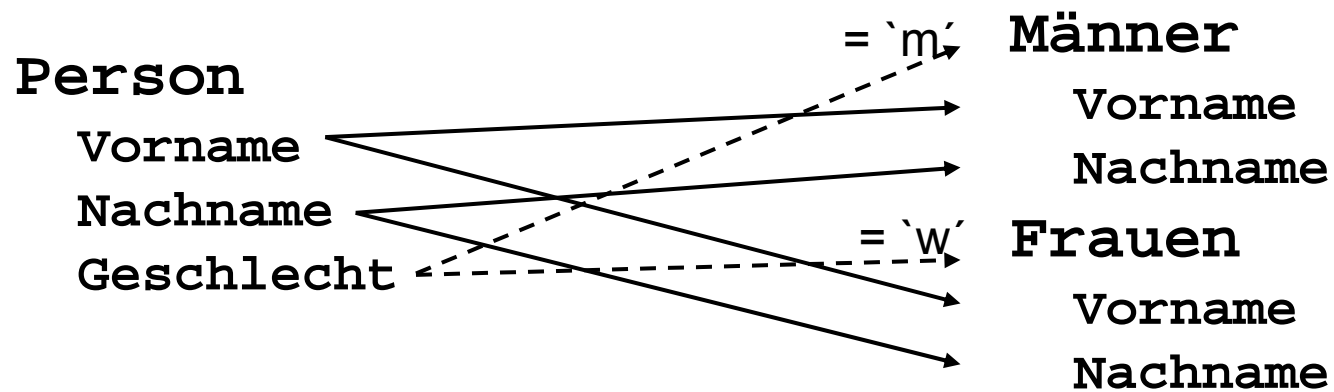
# High-order Mappings

109



# High-order Mappings

110

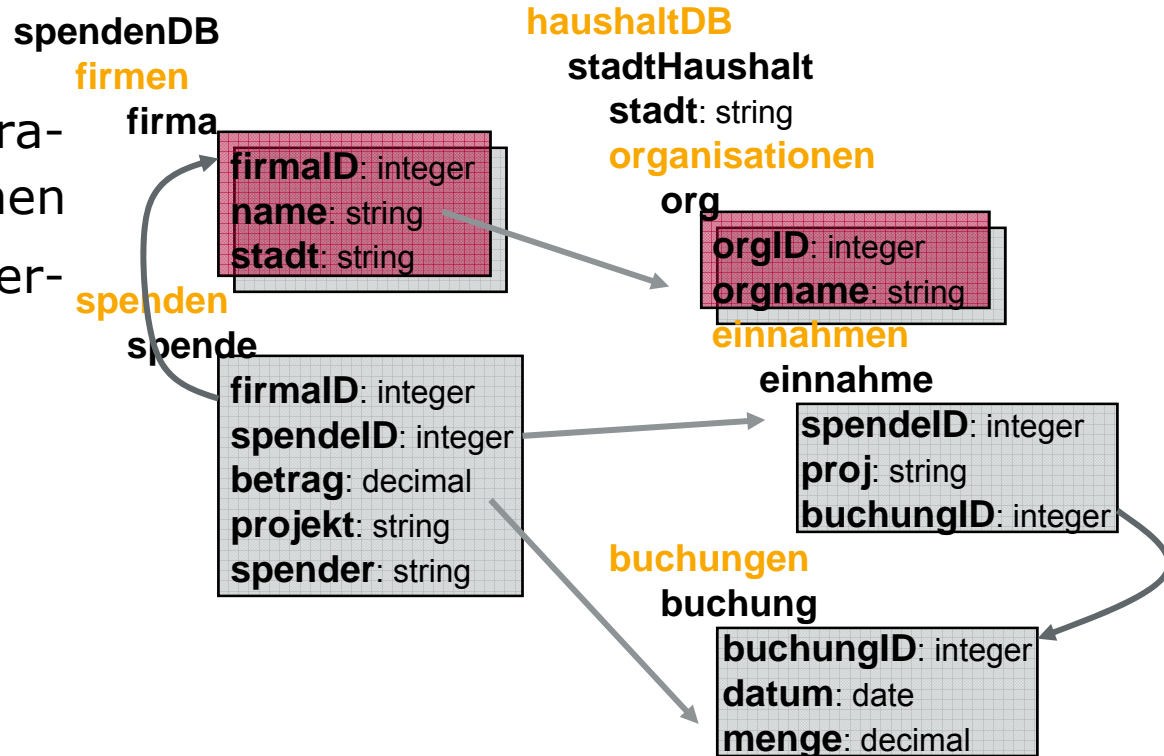


# Zusammenfassung

111

## Drei Schritte

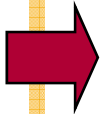
1. Entdeckung von intra-Schema Assoziationen
2. Entdeckung von inter-Schema logischen Mappings
3. Anfrage-erzeugung



# Überblick

112

- Motivation
- Schema Mapping
- Schema Matching
- Mapping Interpretation
  - Aus [FHP+02] und VLDB 2002 Vortragsfolien
- Mapping Werkzeuge





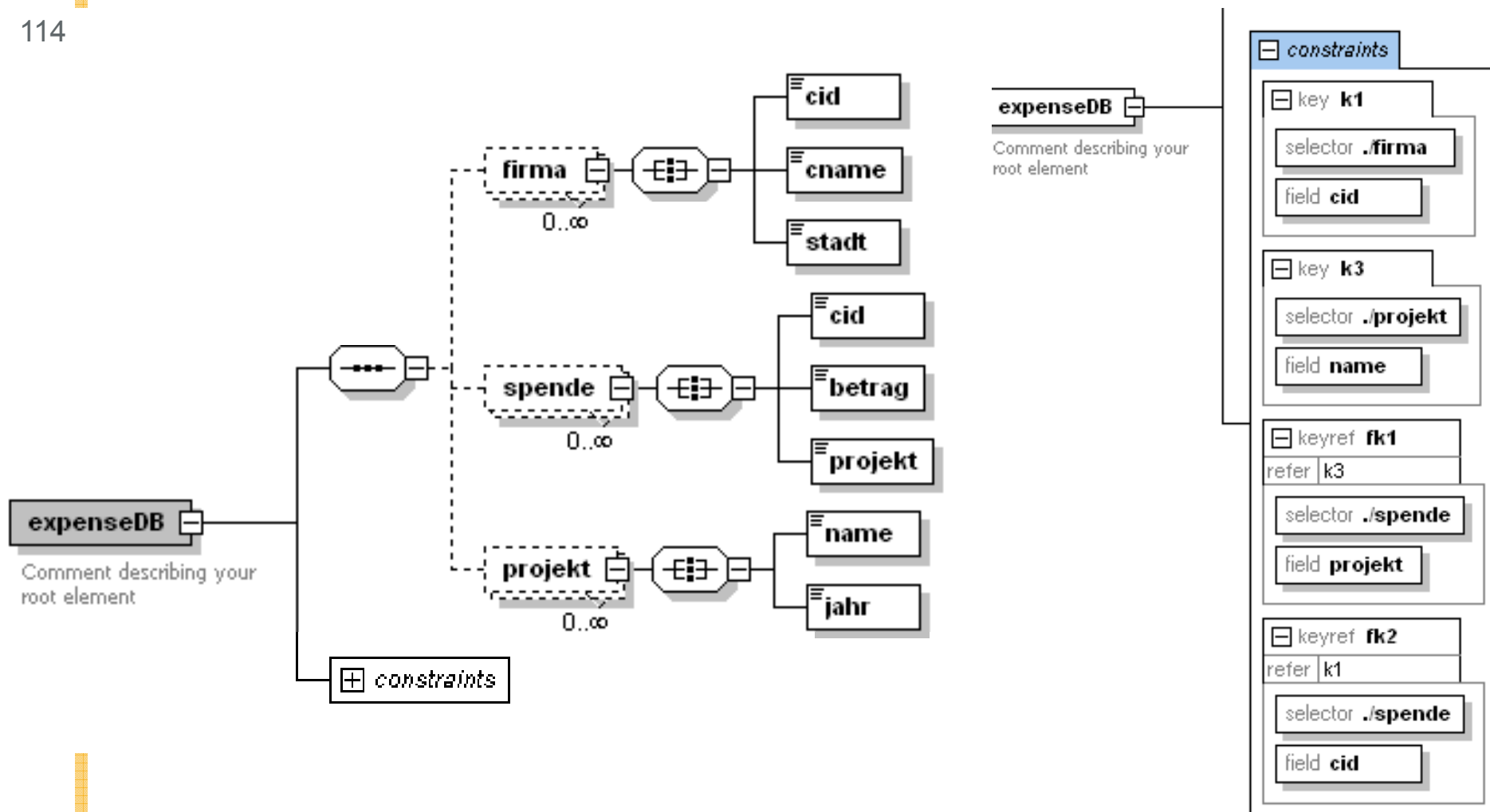
# Vergleich dreier Tools

113

- Altovas MapForce
- IBMs Clio
- [IBMs WSAD]
- Szenario
  - expenseDB -> statisticsDB
  - Definition des Mappings
  - Generierung der Transformation
  - Vergleich der XML Ergebnisse
- XML Screenshots aus XMLSpy
  - [www.xmlspy.com](http://www.xmlspy.com)

# Szenario - ExpenseDB

114



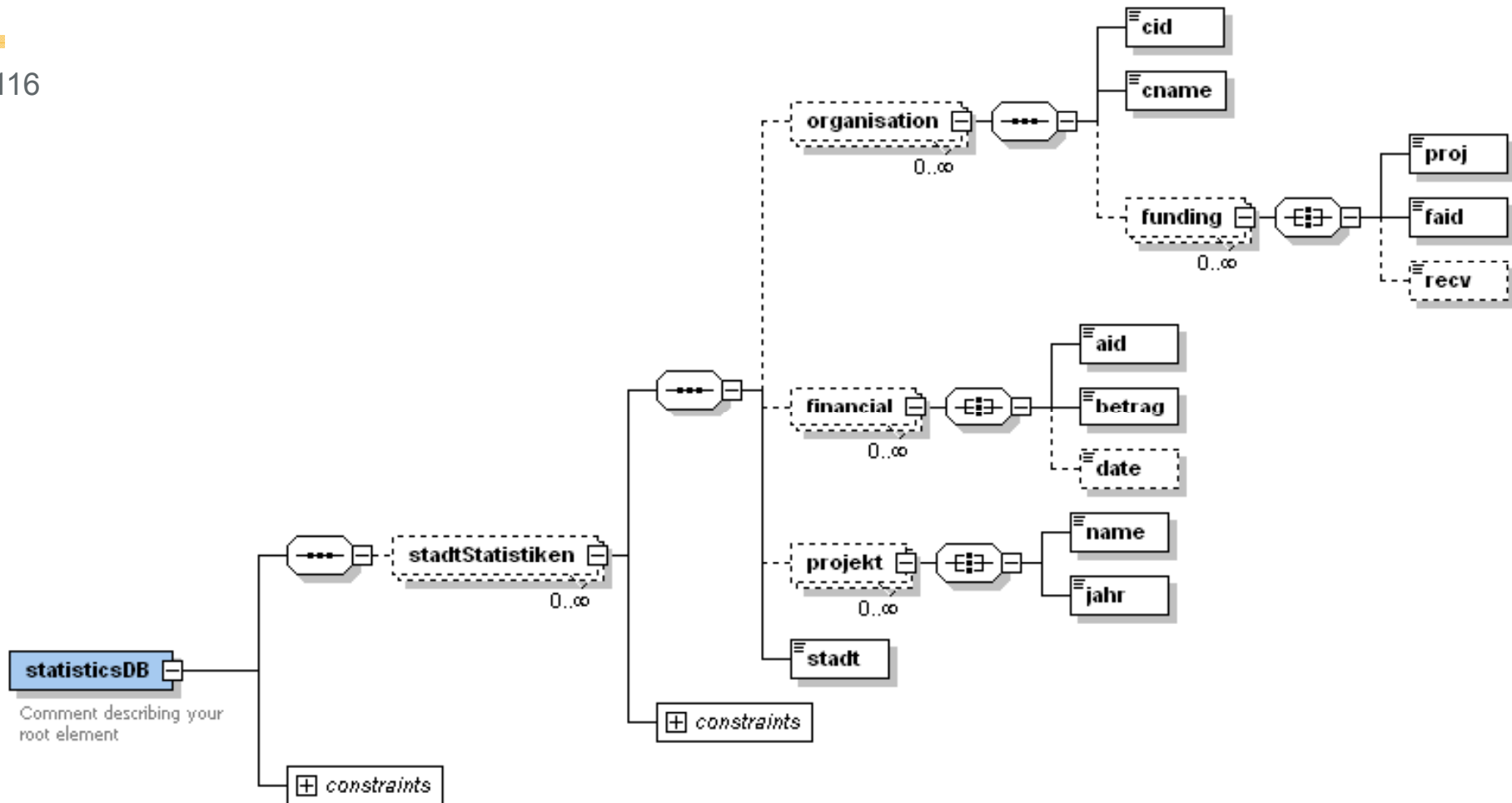
# Szenario - ExpenseDB

115

expenseDB			
<b>xmlns:xsi</b>	http://www.w3.org/2001/XMLSchema-instance		
<b>xsi:nillamespac...</b>	workspace\Clio\schemas\expenseDB.xsd		
▲ firma (5)			
	cid	cname	stadt
1	GG	Garlic	Gilroy
2	CMPQ	Compaq	San Jose
3	NA	Network Associates	San Jose
4	IBM	International Business Machines	San Jose
5	MSFT	Microsoft	Redmond
▲ spende (6)			
	cid	betrag	projekt
1	MSFT	12350	SQLServer
2	IBM	3500	Clio
3	IBM	25000	DB2
4	NA	5000	FastNetwork
5	MSFT	500	Leo
6	IBM	10000	Clio
▲ projekt (5)			
	name	jahr	
1	SQLServer	1997	
2	FastNetwork	1996	
3	DB2	1987	
4	Clio	1999	
5	Leo	2000	

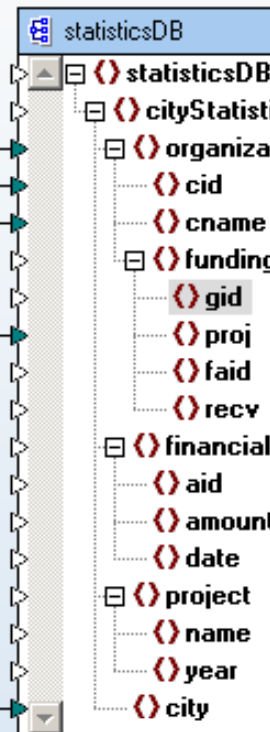
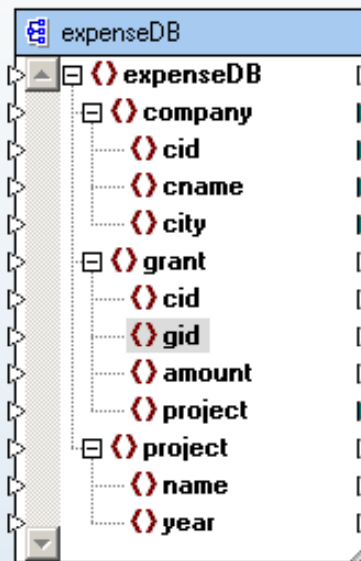
# Szenario StatisticsDB

116



# Altovas MapForce

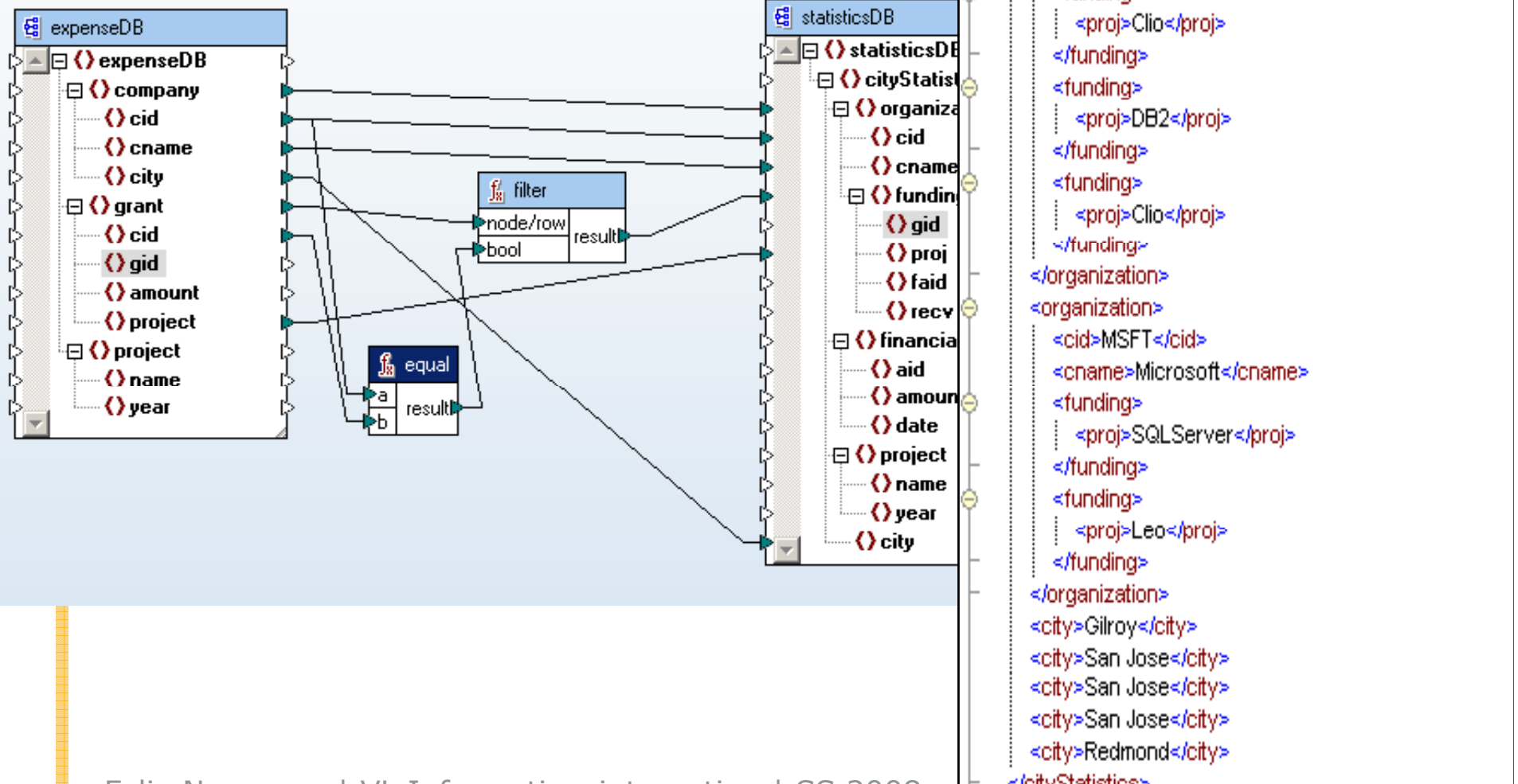
117



```
<?xml version="1.0" encoding="UTF-8"?>
<statisticsDB xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://www.w3.org/2001/XMLSchema-instance statisticsDB.xsd">
  <cityStatistics>
    <organization>
      <cid>GG</cid>
      <cname>Garlic</cname>
      <funding>
        <proj>SQLServer</proj>
        <proj>Clio</proj>
        <proj>DB2</proj>
        <proj>FastNetwork</proj>
        <proj>Leo</proj>
        <proj>Clio</proj>
      </funding>
    </organization>
    <organization>
      <cid>CMPQ</cid>
      <cname>Compaq</cname>
      <funding>
        <proj>SQLServer</proj>
        <proj>Clio</proj>
        <proj>DB2</proj>
        <proj>FastNetwork</proj>
        <proj>Leo</proj>
        <proj>Clio</proj>
      </funding>
    </organization>
    <organization>
    </organization>
    <organization>
    </organization>
    <city>Gilroy</city>
    <city>San Jose</city>
    <city>San Jose</city>
    <city>San Jose</city>
    <city>Redmond</city>
  </cityStatistics>
</statisticsDB>
```

# Altovas MapForce

118



# Altovas MapForce

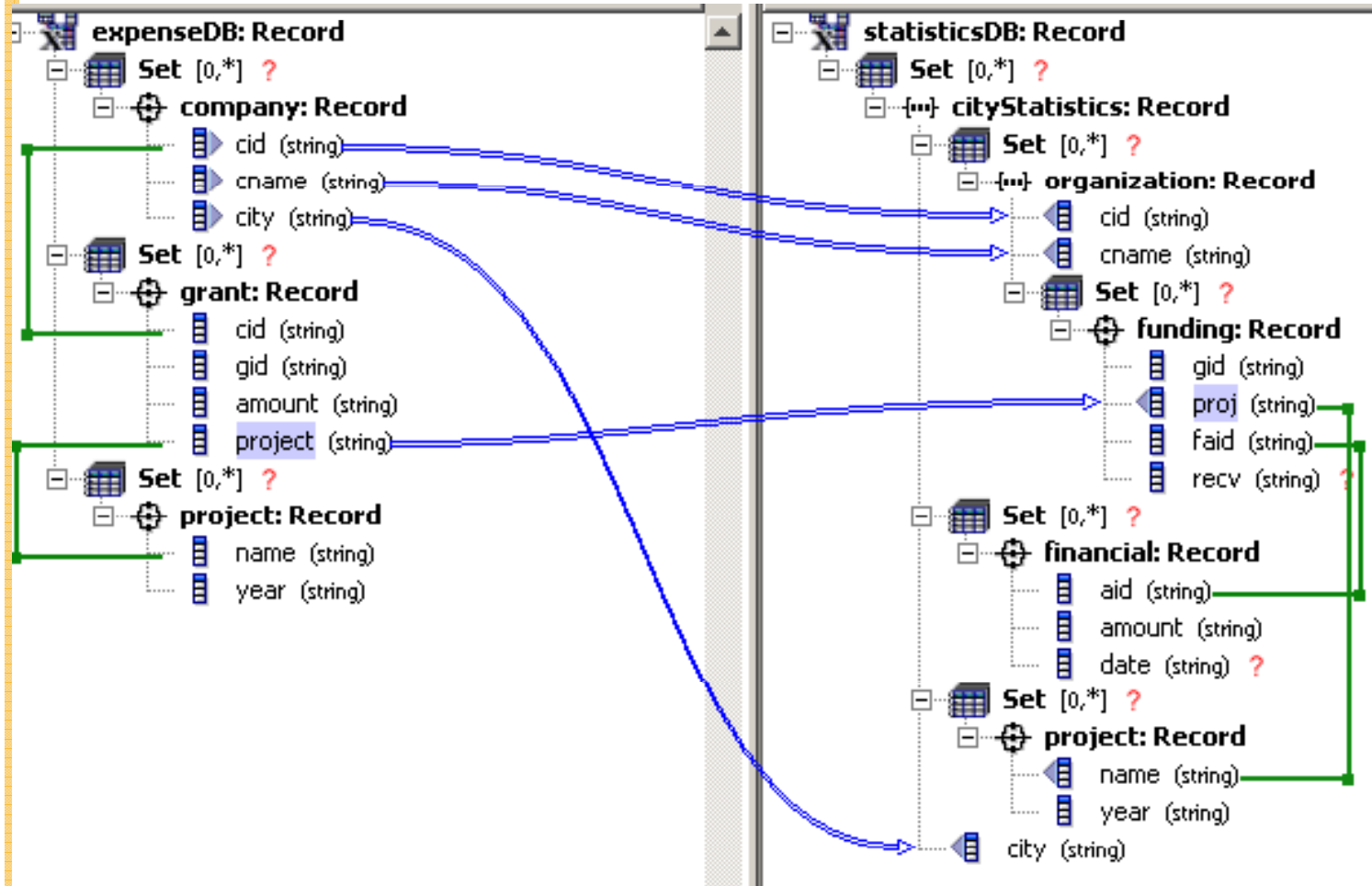
119

## Zusammenfassung

- Informationen gehen nicht verloren
- Joins
  - werden nicht erkannt
  - können manuell in GUI erstellt werden
- Gruppierung wird nicht erkannt

# IBMs Clio

120





# IBMs Clio

121

```

<?xml version="1.0" encoding="UTF-8" ?>
- <quip:result xmlns:quip="http://namespaces.softwareag.com/tamino/quip/">
- <statisticsDB>
+ <cityStatistics>
- <cityStatistics>
- <organization>
  <cid>IBM</cid>
  <cname>International Business Machines</cname>
- <funding>
  <gid>Sk72(Clio, International Business Machines, San Jose, IBM)</gid>
  <proj>Clio</proj>
  <faid>Sk67(Clio, International Business Machines, San Jose, IBM)</faid>
</funding>
- <funding>
  <gid>Sk72(DB2, International Business Machines, San Jose, IBM)</gid>
  <proj>DB2</proj>
  <faid>Sk67(DB2, International Business Machines, San Jose, IBM)</faid>
</funding>
</organization>
- <organization>
  <cid>NA</cid>
  <cname>Network Associates</cname>
- <funding>
  <gid>Sk72(FastNetwork, Network Associates, San Jose, NA)</gid>
  <proj>FastNetwork</proj>
  <faid>Sk67(FastNetwork, Network Associates, San Jose, NA)</faid>
</funding>
</organization>
+ <financial>
+ <financial>
+ <financial>
+ <project>
+ <project>
+ <project>
  <city>San Jose</city>
</cityStatistics>
</statisticsDB>
</quip:result>

```

# IBMs Clio

122

## Zusammenfassung

- Joins werden erkannt
  - Aber nicht erzwungen: Interpretation
- Gruppierung wird erkannt

# IBMs WSAD

Target	Source	Applied Function/Gr
statisticsDB	expenseDB	
statisticsDB		
cityStatistics		
organization	company	
cid	cid	
cname	cname	
funding	grant	
proj	name	
city	city	

```

<?xml version="1.0" encoding="UTF-8" ?>
- <statisticsDB>
- <cityStatistics>
  - <organization>
    <cid>GG</cid>
    <cname>Garlic</cname>
  + <funding>
  + <funding>
  + <funding>
  + <funding>
  + <funding>
  + <funding>
  - <funding>
    <gid />
    <proj>SQLServer</proj>
    <faid />
  </funding>
  - <funding>
    <gid />
    <proj>FastNetwork</proj>
    <faid />
  </funding>
  + <funding>
  + <funding>
  + <funding>
  </organization>
+ <organization>
+ <organization>
+ <organization>
+ <organization>
+ <financial>
- <project>
  <name />
  <year />
</project>
  <city>Gilroy</city>
</cityStatistics>
</statisticsDB>
  
```

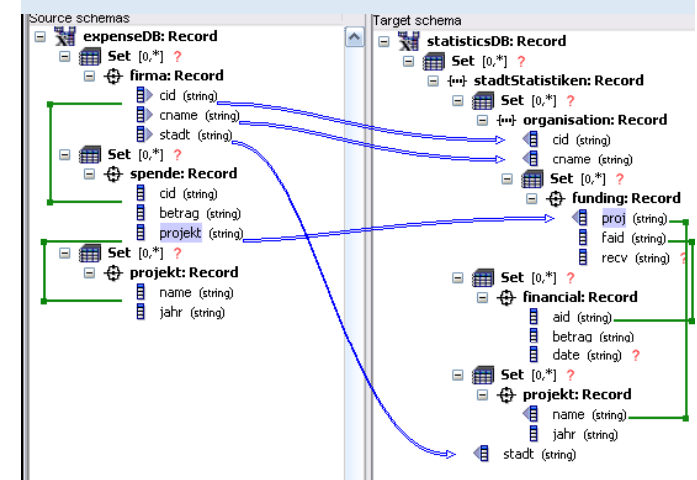
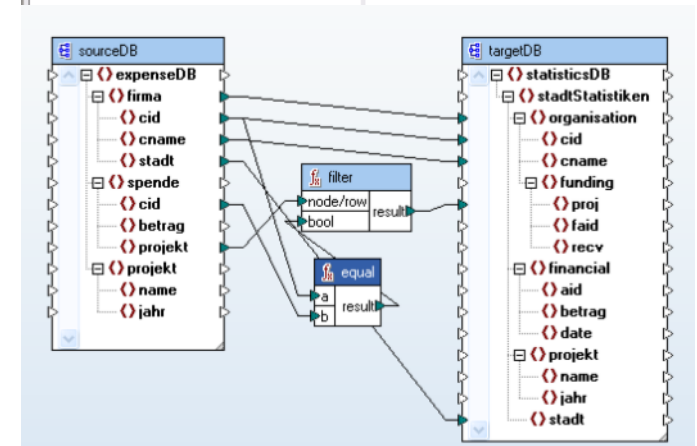
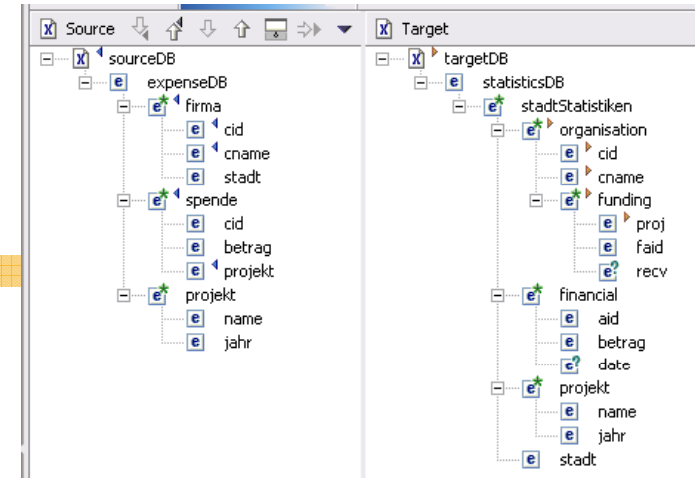
## Zusammenfassung

- Mögliche Joins werden nicht erkannt
- Gruppierung wird nicht erkannt
- Informationen werden verworfen
  - Nur erste Stadt
- Information wird „erfunden“
  - Jede Organization erhält jedes Projekt

# Zusammenfassung der Demo

125

- IBMs WSAD
  - Mögliche Joins werden nicht erkannt.
  - Gruppierung wird nicht erkannt.
  - Daten werden verworfen.
  - Assoziationen zw. Daten werden „erfunden“.
- Altovas MapForce
  - Daten gehen nicht verloren.
  - Joins werden nicht erkannt, aber können manuell in GUI erstellt werden.
  - Gruppierung wird nicht erkannt.
- IBMs Clio
  - Joins werden erkannt, aber nicht erzwungen: Interpretationen möglich
  - Gruppierung wird erkannt.



# Literatur – Schema Mapping

126

- [FHP+02] Ron Fagin, Mauricio Hernandez, Lucian Popa, Renee Miller, and Yannis Velegarakis, Translating Web Data, VLDB 2002, Hong Kong, China.
- [RB01] Erhard Rahm and Philip Bernstein, A survey of approaches to automatic schema matching, VLDB Journal 10(4), 2001.
- Zu der Problematik, ob Duplikate Teil des Outputs sein sollten:
  - [Ronald Fagin](#), [Phokion G. Kolaitis](#), [Renée J. Miller](#), Lucian Popa: Data Exchange: Semantics and Query Answering. [ICDT 2003](#): 207-224
- Clio:
  - <http://www.almaden.ibm.com/software/projects/criollo/>
  - oder <http://www.cs.toronto.edu/db/cliol/>

# Literatur – Schema Matching

127

Artikel mit der Klassifikation:

- [RB01] Erhard Rahm and Philip Bernstein, A survey of approaches to automatic schema matching, *VLDB Journal* 10(4), 2001.

Spezielle Algorithmen

- [MGMR02] Sergey Melnik, [Hector Garcia-Molina](#), [Erhard Rahm](#): Similarity Flooding: A Versatile Graph Matching Algorithm and Its Application to Schema Matching. [ICDE 2002](#): 117-128
- [MH03] Jayant Madhavan, [Alon Y. Halevy](#): Composing Mappings Among Data Sources. [VLDB 2003](#): 572-583.
- uvam.