



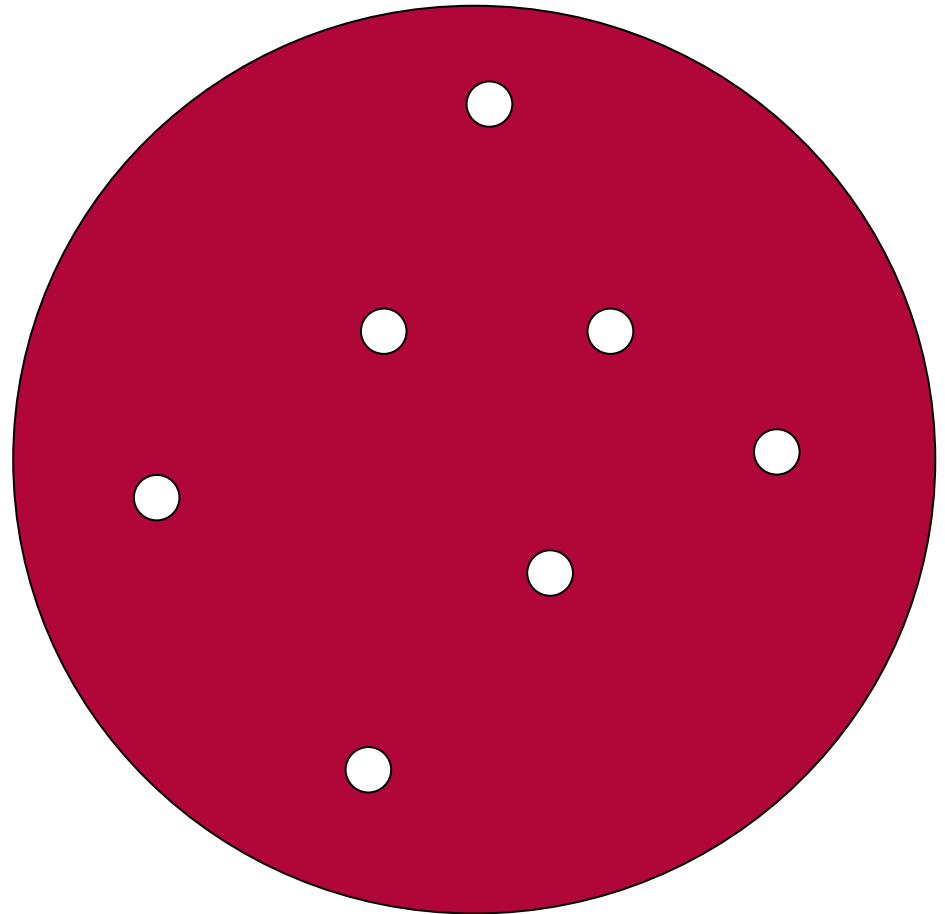
Schema Matching Seminar
WS 2007/08
Themen & Organisation

Prof. Felix Naumann,
Alexander Albrecht, Frank Kaufer, Melanie Weis

Hasso-Plattner-Institut, Potsdam
24. Oktober 2007

Seminarmodus - klassisch

- Großes Themengebiet
- 1-2 Artikel
- Seminararbeit
(20 Seiten in TeX)
- Vortrag



Seminarmodus - klassisch



Strg + C,V
modulo Thesaurus



Seminarmodus – FUBU Style



- Wissenschaftliche Kollaboration
 - Arbeitsgruppe, Forschungsprojekt
 - Forschungskolleg
- Gemeinsames Erarbeiten eines neuen, kohärenten Themengebiets
 - Theorie verstehen
 - Techniken herausarbeiten
 - Systeme evaluieren
- Ausarbeitung: Wiki (Wikipedia als Bonus)
- Talk im Tutorialstil
- Hands-on (6 CP)
 - Kleine Demoimplementationen
 - Testkollektionen erstellen (positive & negative Beispiele)
 - Evaluieren



Werde ich hier zum Fachidioten?

5

Mensch-Maschine-Schnittstelle

Profiling

Ontologien

Software Component Retrieval

P2P-Systeme

Maschinelles Schließen

Deep Web
Schemaintegration

Künstliche Intelligenz

Textklassifikation

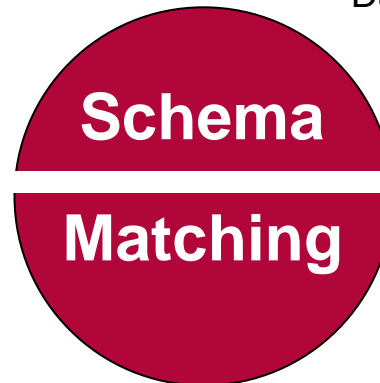
Data Translation

Graphentheorie

Information Retrieval

Katalogintegration

Datenreinigung



Service Matching

Service Komposition

Anwendungsintegration

Semantic Web

Query Mediation

Maschinelles Lernen

Computerlinguistik

Data Warehouse Loading

Web Community Systeme

Informationsintegration

Spam-Filter

Mustererkennung

XML Message Mapping

Web Query Interfaces

OLAP

Wissensrepräsentation

Business Intelligence



Themenkomplexe

6

A) Modellierung und Heterogenität

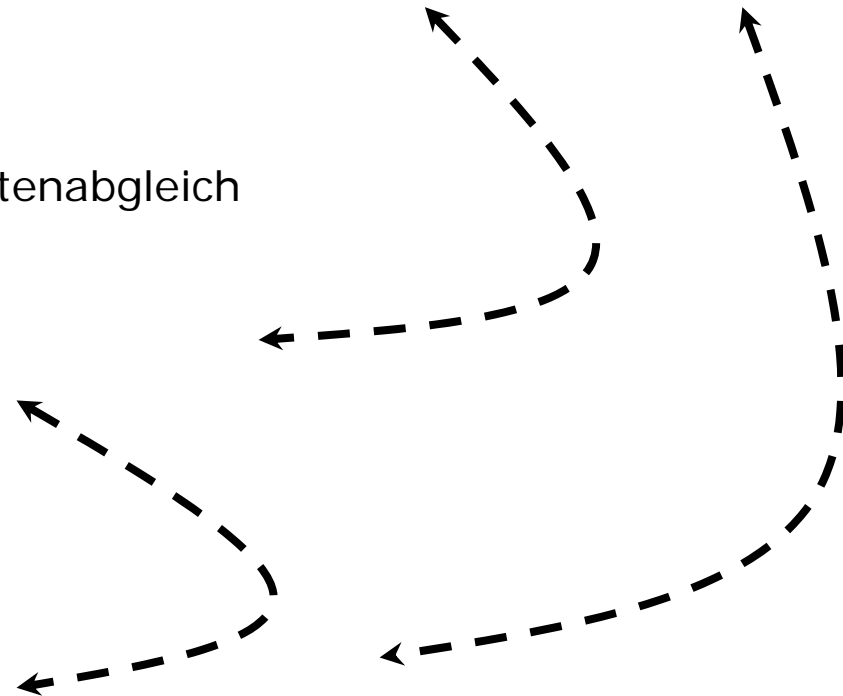
- Datenmodelle: Relational, XML-Schema, OWL, OO, ER, RDFS, F-Logic
- Heterogenitäten

B) Matchingtechniken

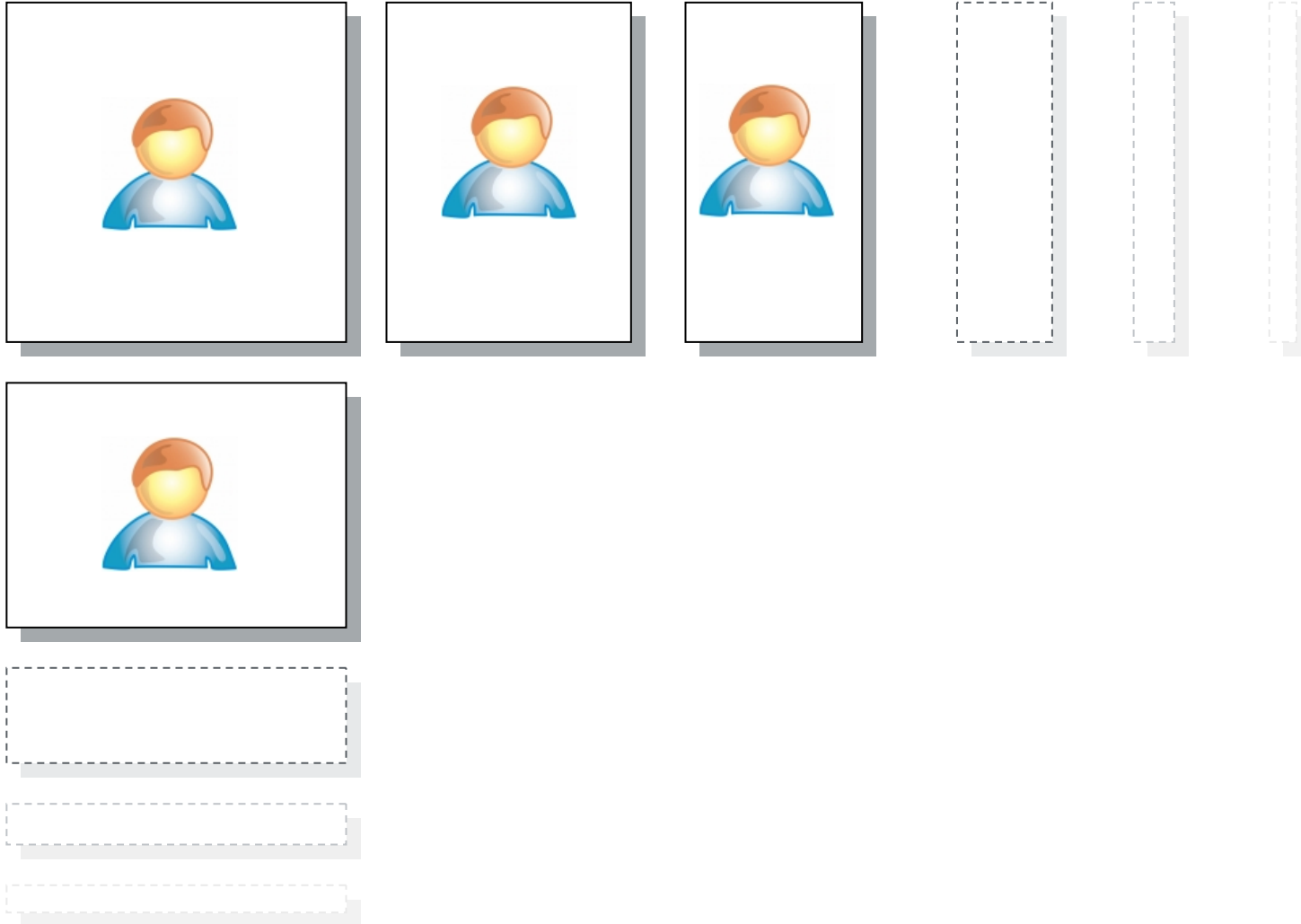
- Elementabgleich: Linguistisch, Metadatenabgleich
- Instanzabgleich
- Graphabgleich
- Lerntechniken
- Logisches Schlussfolgern

C) Matchingsysteme & Evaluation

- kommerzielle Systeme
- wissenschaftliche Prototypen
- Evaluation



Themenkomplexität



Die Themen



Datenmodelle

- Relational
- Bäume, XML-Schema
- OO (Objektrelationales Mapping, Java), Entity-Relationship
- Ontologien: OWL, RDF/S, F-Logic

Heterogenitäten ...

- in allen Modellen
- in bestimmten Modellen (Modelleigenheiten)
- zwischen den Modellen

A. Datenmodelle

Datenmodellheterogenität: XML-RDB

```
<Products>  
  <Product id="355156666" productgroup="DVD">  
    <title>Pulp Fiction</title>  
    <details>http://www.amazon.de/dp/B00004TH46</details>  
    <actor>John Travolta</actor>  
    <author>Quentin Tarantino</author>  
  </Product>  
  <Product id="355156666" productgroup="Book"> ... <Product>  
  <Product id="B00004TH46" productgroup="Book"> ... <Product>  
</Products>
```

Product		
id	title	details
355156666	Harry Potter ...	http://...
3899407288	Harry Potter ...	http://...
B00004TH46	Pulp Fiction	http://...

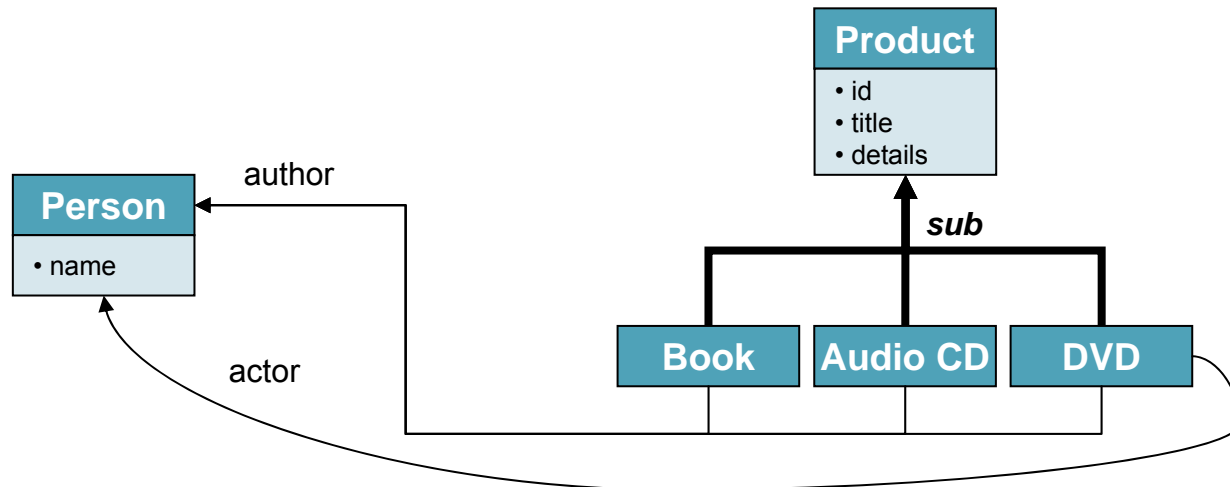
DVD		
id	actor	author
B00004TH46	John Travolta	Quentin Tarantino

Book	
id	author
355156666	Joanne K. Rowling

Audio CD	
id	author
3899407288	Rufus Beck

A. Datenmodelle

Datenmodellheterogenität: OO-RDB



Product		
id	title	details
3551566666	Harry Potter ...	http://...
3899407288	Harry Potter ...	http://...
B00004TH46	Pulp Fiction	http://...

DVD		
id	actor	author
B00004TH46	John Travolta	Quentin Tarantino

Book	
id	author
3551566666	Joanne K. Rowling

Audio CD	
id	author
3899407288	Rufus Beck

A. Datenmodelle

Schematische Heterogenität in XML

```
<Products>
  <Product id="3551566666" productgroup="DVD">
    <title>Pulp Fiction</title>
    <details>http://www.amazon.de/dp/B00004TH46/</details>
    <actor>John Travolta</actor>
    <author>Quentin Tarantino</author>
  </Product>
  <Product id="3551566666" productgroup="Book"> ... <Product>
  <Product id="B00004TH46" productgroup="Book"> ... <Product>
</Products>
```

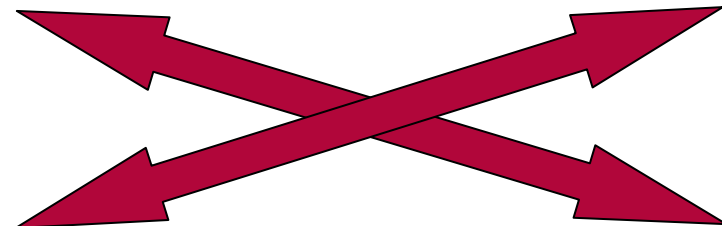
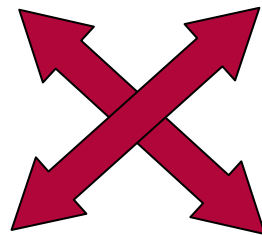
```
<Products>
  <Product title="Pulp Fiction">
    <id>3551566666</id>
    <attribute>
      <name>productGroup</name>
      <value>DVD</value>
    </attribute>
    <details>http://www.amazon.de/dp/B00004TH46/</details>
    <title>Pulp Fiction</title>
    <actor>John Travolta</actor>
    <author>Quentin Tarantino</author>
  </Product>
  <Product id="3551566666" productgroup="Book"> ... <Product>
  <Product id="B00004TH46" productgroup="Book"> ... <Product>
</Products>
```

B.1. Intrinsisches Elementmatching

13

- Ähnlichkeit der Attributnamen wird gemessen
- Ähnliche Attributnamen sind potentiell bessere Matches als unähnliche Attributnamen.

ID	NName	VName	Tel	Adresse
1	Müller	Michael	123456	Hauptstr. 3, Berlin
2	Schmidt	Maria	987654	Dorfstr. 45, Golzow



PNR	Vorname	Nachname	Adresse	Telefon
1001	Michel	Müller	Hauptstrasse 3, Berlin	123456
1002	Peter	Schmidt	Forststr. 74, Hintertupfingen	918736

B.1. Intrinsisches Elementmatching

String-Ähnlichkeit

14

Berechnung der Ähnlichkeit von Attributnamen

- Durch ein Ähnlichkeitsmaß
 - Je größer das Ergebnis des Ähnlichkeitsmaßes, desto ähnlicher sind die Strings.
 - Jaccard similarity
 - Cosine similarity
 - ...
- Durch ein Distanzmaß
 - Je kleiner das Ergebnis des Distanzmaßes, desto ähnlicher sind die Strings.
 - Edit distance
 - Q-Gram distance
 - ...
- Bei normalisierten Maßen: $\text{Ähnlichkeit} = 1 - \text{Distanz}$

B.1. Intrinsisches Elementmatching

Edit-Distance

15

Beispiel „HASE“ \Rightarrow „RASEN“

- triviale Umformung durch Einfügung und Löschung mittels Leerzeichen # am Anfang und Ende
- HASE##### #####RASEN

Transkript: DDDDIIIII

- Kosten 9
- Nicht minimal!

B.1. Intrinsisches Elementmatching

Edit-Distance contd.

16

Beispiel „HASE“ \Rightarrow „RASEN“

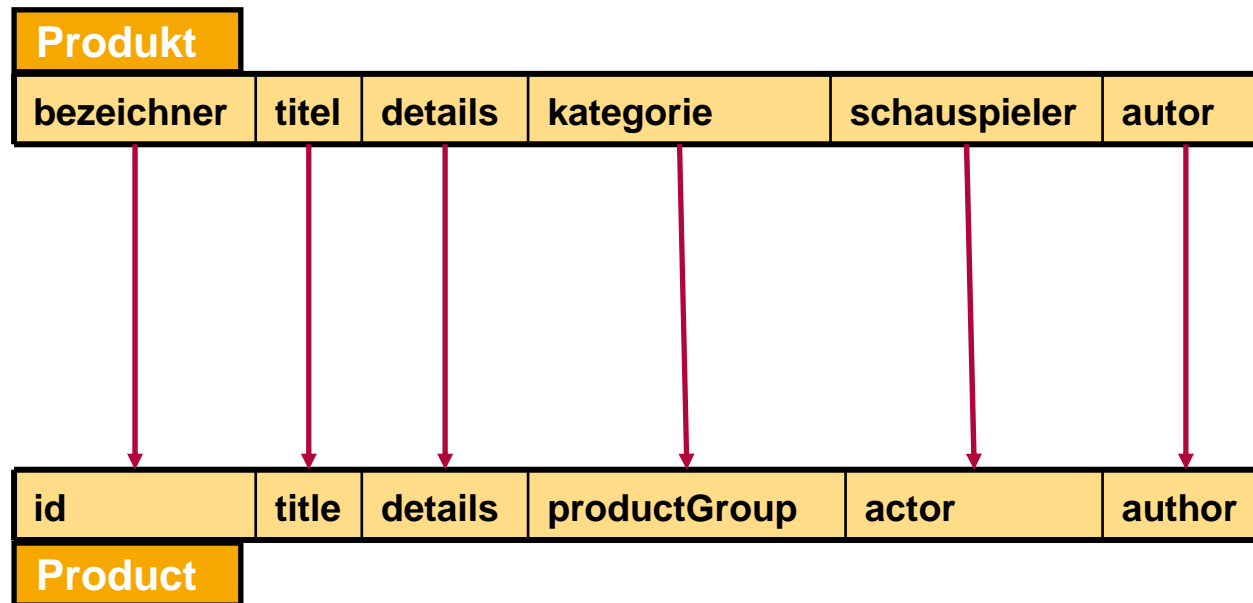
- H \rightarrow R durch Ersetzung (R)
- A \rightarrow A durch Übereinstimmung (M)
- S \rightarrow S durch Übereinstimmung (M)
- E \rightarrow E durch Übereinstimmung (M)
- „ “ \rightarrow N durch Einfügung (I)

Transkript: RMMMI

- $\text{EditDistance}(\text{HASE}, \text{RASEN}) = 2$
- Minimal! Woher weiß man das?

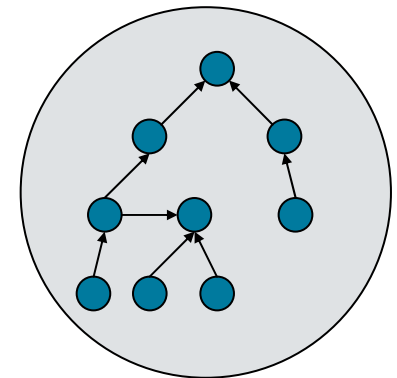
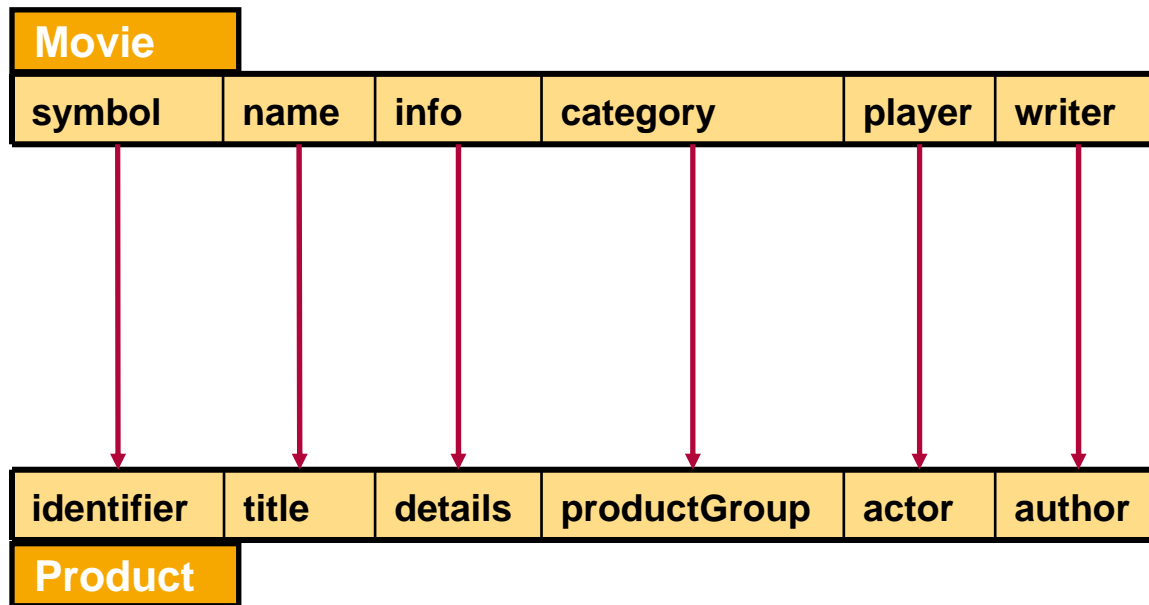
B.2. Extrinsisches Elementmatching Wörterbücher

17



B.2. Extrinsic Elementmatching Lexika

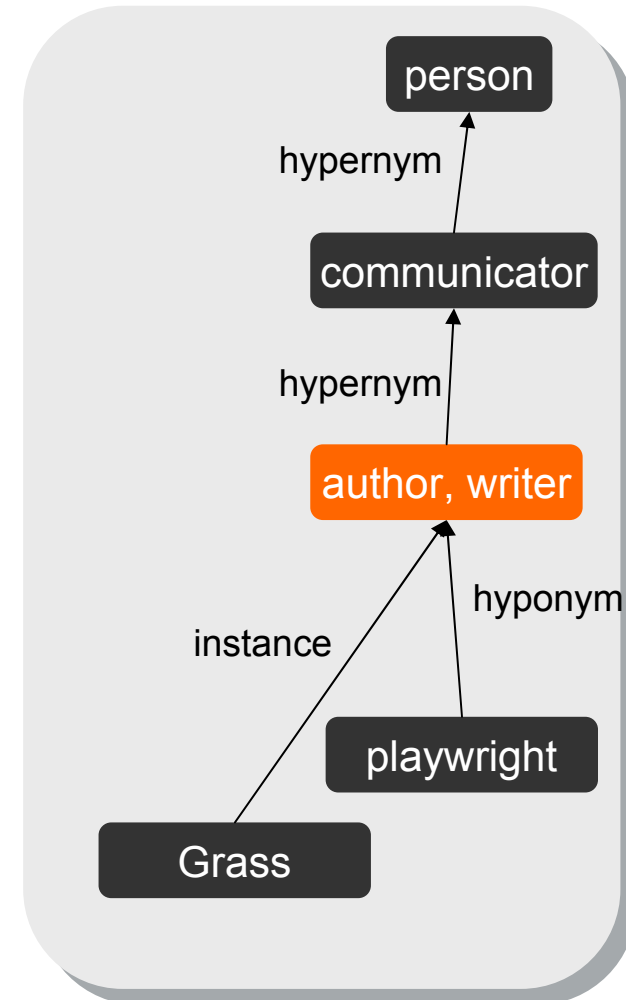
18



B.2. Extrinsisches Elementmatching WordNet

19

- Linguistische Ontologie
- Konzepte: Synonymmengen vs. Wörter
- **author**
 - **{author¹, writer}**
Schriftsteller, Drehbuchautor (noun)
 - **{author², generator, source}**
Verursacher, Initiator (noun)
 - **{author³}**
verfassen (verb)
- Semantische Relationen
 - hypernym, hyponym, holonym, ...

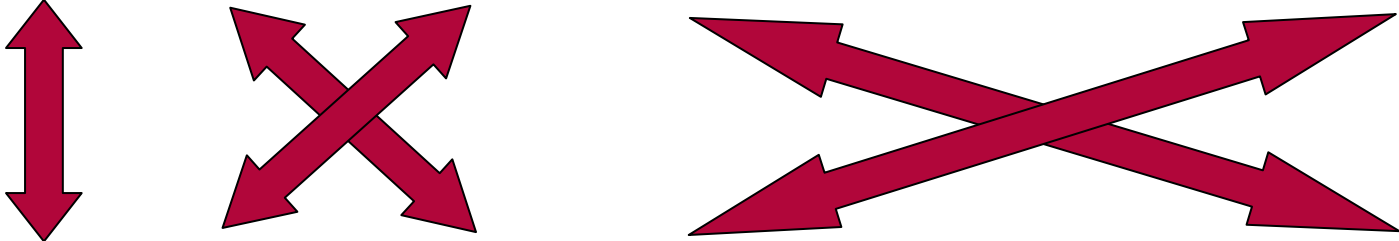


B.3. Extensionales Matching

20

- Es werden ähnliche Daten oder Datencharakteristika ermittelt.
- Attribute / Attributkombinationen mit ähnlichen Daten sind potentiell bessere Matches als Attribute mit unähnlichen Daten.

ID	NName	VName	Tel	Adresse
1	Müller	Michael	123456	Hauptstr. 3, Berlin
2	Schmidt	Maria	987654	Dorfstr. 45, Golzow



PNR	Vorname	Nachname	Adresse	Telefon
1001	Michel	Müller	Hauptstrasse 3, Berlin	123456
1002	Peter	Schmidt	Forststr. 74, Hintertupfingen	918736

Betrachtung der Instanzdaten

- Betrachtung von Statistiken (Verteilung, Muster, ...)
- Erkennung von Duplikaten zwischen den zu matchenden Tabellen
 - auf Attributebene
 - auf Tupelebene
- Erkennen von Transformationsregeln zwischen den Instanzen

Ähnlichkeitsmaße für

- lange Strings
- Strings, die aus mehreren Worten (Token) bestehen
- komplexe Strukturen (Tupel statt Strings, XML Daten)

Beispiele

- Jaccard Similarity
- Cosine Similarity
- SoftTFIDF
- ...

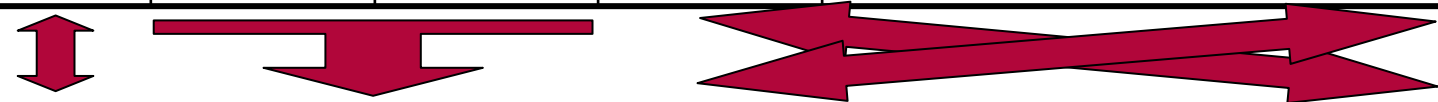
B.3. Extensionales Matching

Jaccard Similarity

23

ID	NName	VName	Tel	Adresse
1	Müller	Michael	123456	Hauptstr. 3, Berlin

PNR	Name	Anschrift	Telefon
1001	Michel Müller	Hauptstrasse 3, Berlin	123456



1. Tokenmenge (TM) erzeugen

$$TM(T1) = \{1, Müller, Michael, 123456, Hauptstr., 3, Berlin\}$$

$$TM(T2) = \{1001, Michael, Müller, Hauptstrasse, 3, Berlin, 123456\}$$

2. Ähnliche Token identifizieren

Michael = Michel

Hauptstr. = Hauptstrasse

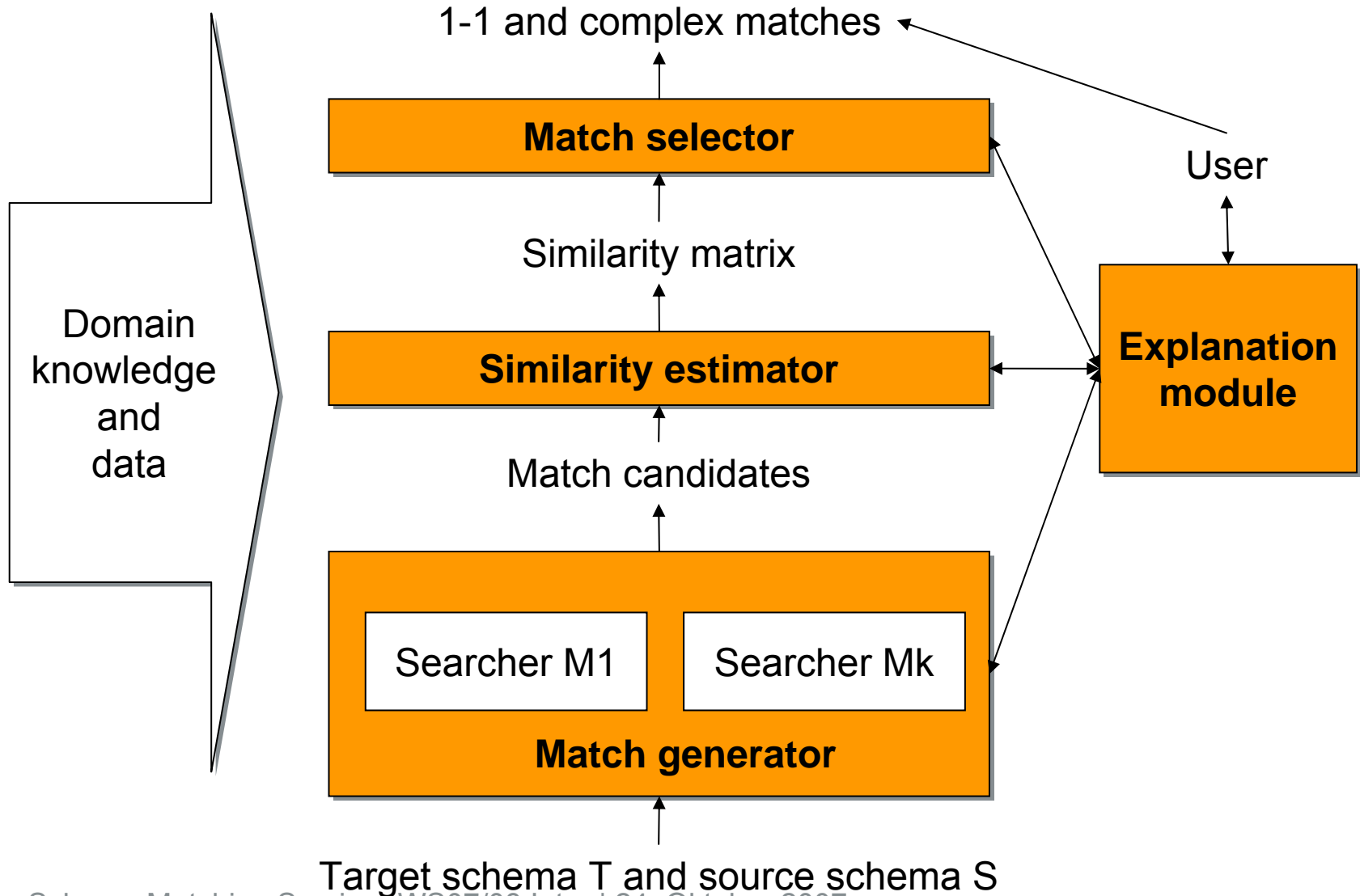
3. Jaccard Ähnlichkeit berechnen

$$Jaccard(TM(T1), TM(T2)) = \frac{|TM(T1) \cap TM(T2)|}{|TM(T1) \cup TM(T2)|} = \frac{6}{7}$$

- Finden von 1:1 und komplexen Matchings
- Durchsucht den Raum aller möglichen Matches anhand von spezialisierten Searcher Modulen
 - text searcher: `concat(Name, Vorname)`
 - numeric searcher: `preis = preis + Mwst`
 - date searcher: `datum = concat(Monat, „.“, Jahr)`
- Setzt frühzeitig Domänenwissen ein
 - Integritätsbedingungen
 - Betrachtet überlappende Tupel (Duplikate)
 - Erlaubt, Daten außerhalb der Datenbank zu verwenden

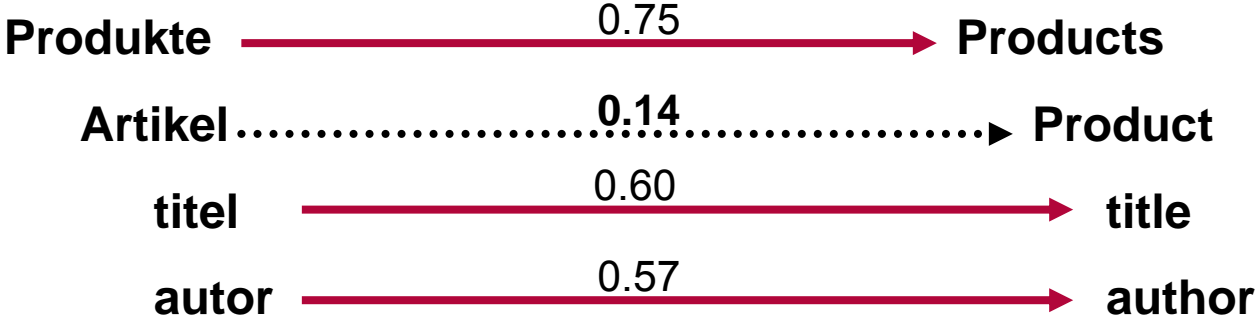
B.3. Extensionales Matching System iMap

25



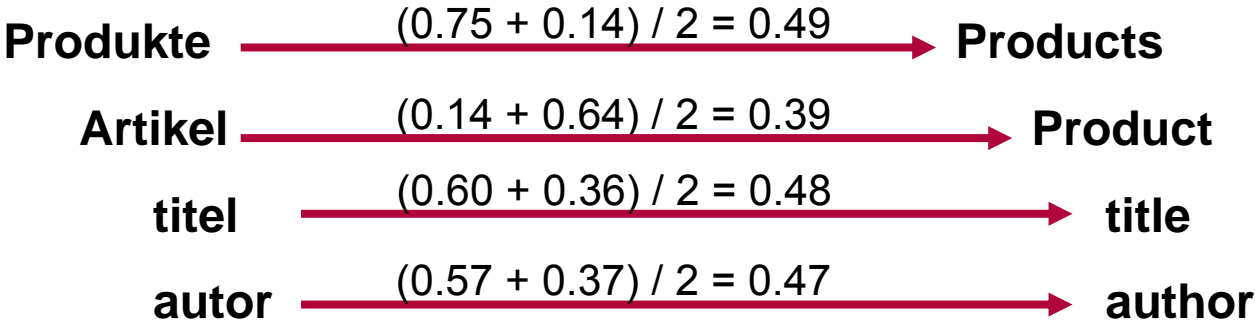
B.4. Strukturbasiertes Matching

26

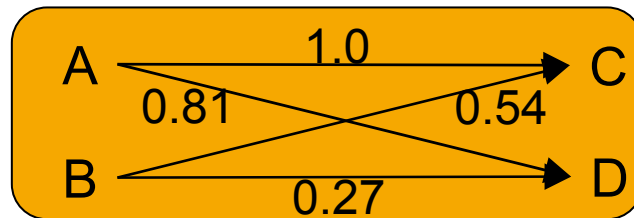


Min: 0.30

Idee: Ähnlichkeit = (Elementähnlichkeit + durchschnittliche Elementähnlichkeit der direkten Nachbarn) / 2



- Ergebnis der vorherigen Verfahren



- Algorithmen zur Extraktion des eigentlichen (besten) Mappings
 - Stable Marriage
 - Maximum Weighted Bipartite Graph Matching
- Vergleich der beiden Verfahren
- Praktische Aufgabe:
 - Enummerierung mehrerer Stable Marriages
 - Finden der besten Stable Marriage

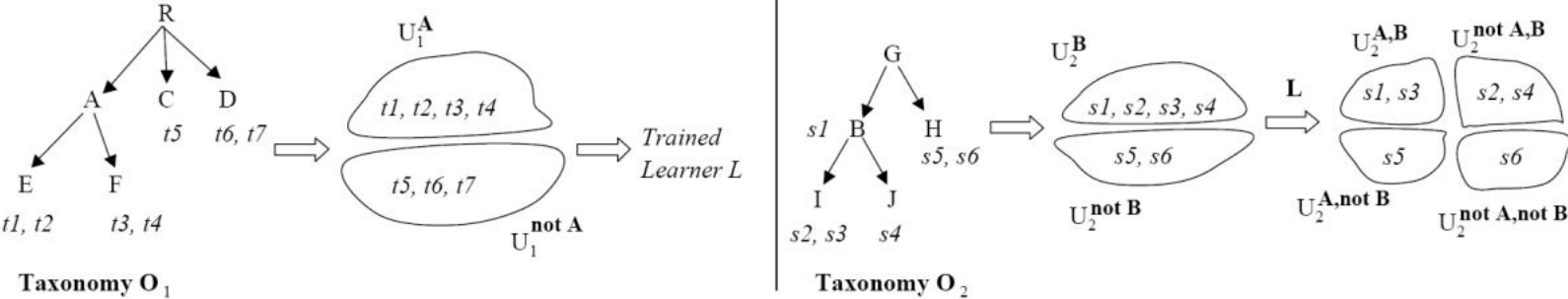
B.6. Maschinelles Lernen

28

- Stochastik
- Künstliche Intelligenz
- Unüberwachtes vs. Überwachtes Lernen

- Ansätze
 - Naive Bayes - **GLUE (LSD)**
 - Neuronale Netzwerke (Semint)
 - Entscheidungsbäume
 - ...

- Survey



Ontologien: formale Spezifikation einer gemeinsamen Konzeptualisierung

- Schemas (Konzepte, Relationen)
- Instanzen
- Logische Axiome

Unterschiede

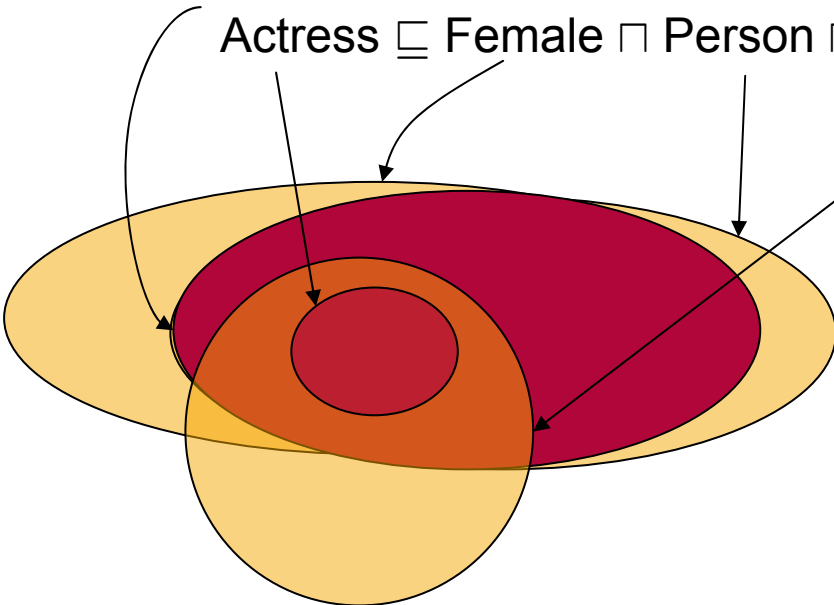
- Objekt-Orientierung
- Betonung auf konzeptionelle Hierarchien (subclass, hypernym)
- Metamodellierung (Instanzen \leftrightarrow Konzepte)
- Anwendung: Wissensmodell vs. Datenmodell (expliziter)
- **Formale Logik und Semantik**

B.7. Semantisches Matching

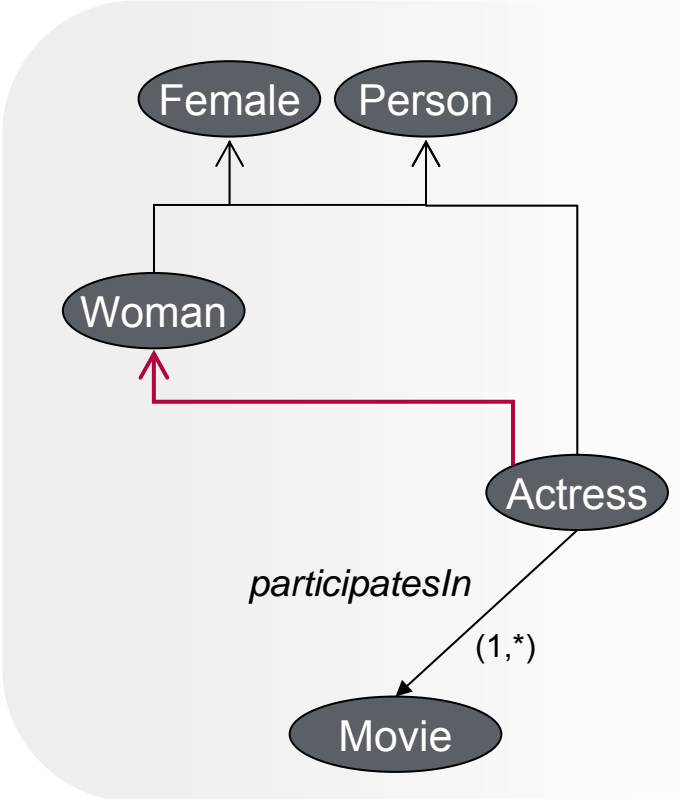
Logik, Semantik, Inferenz

30

$Woman \equiv Female \sqcap Person.$
 $Actress \sqsubseteq Female \sqcap Person \sqcap \exists \text{ participatesIn.Movie.}$



$\emptyset \models Actress \sqsubseteq Woman$



B.7. Semantisches Matching

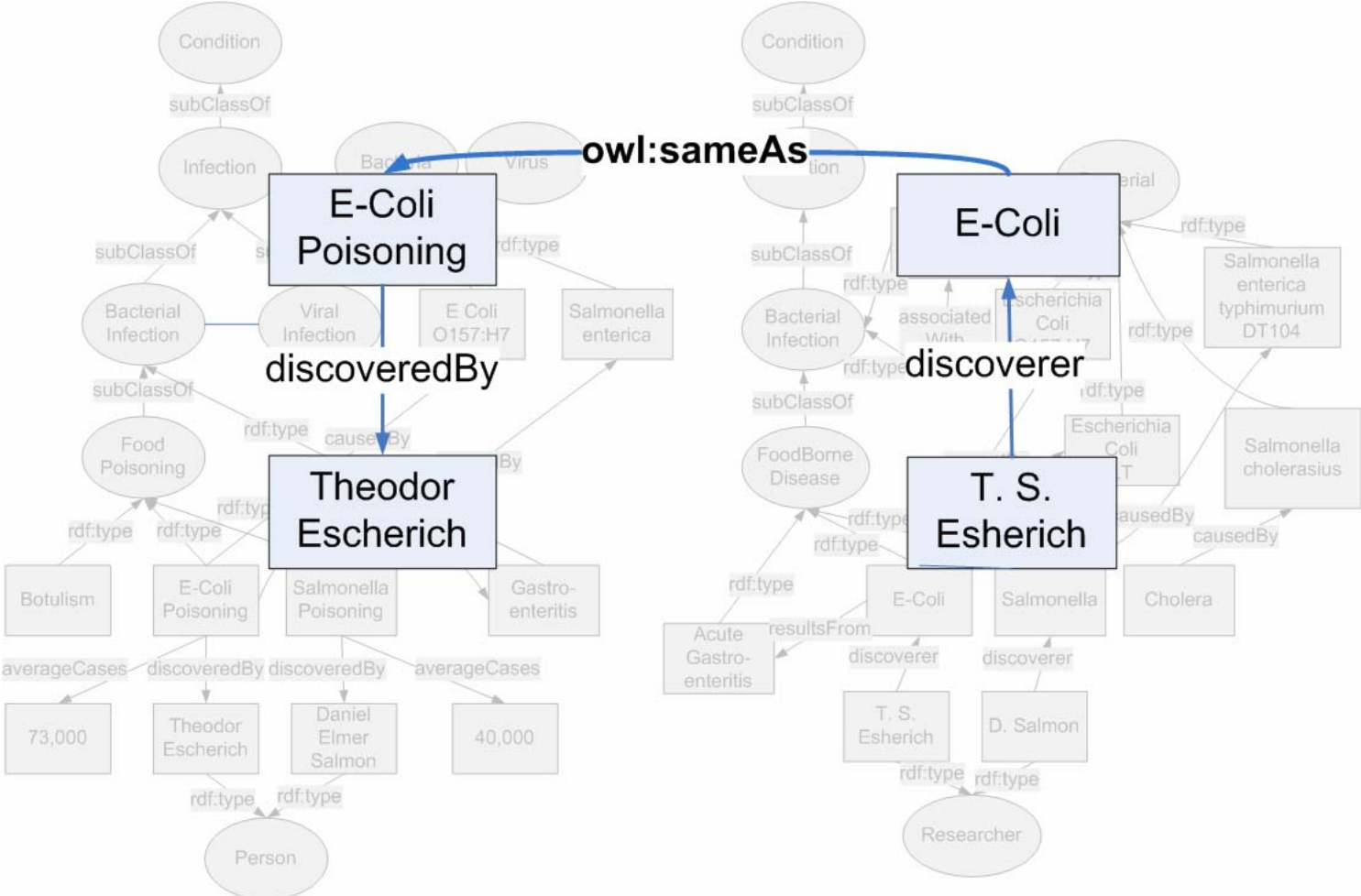
Ontology Matching

Schema Matching mit zusätzlichen Methoden & Herausforderungen

- Semantisches Matching ~ Logische Inferenz (Schlussfolgern)
- Semantische Relationen statt $[0,1]$ -Werte
 - Äquivalenz, Implikation, Disjunktion, Konjunktion
 - Gleichheit, Untermenge, Vereinigung, Schnittmenge, Disjunktheit
- Integration von syntaktischem und semantischem Matching
- Systeme
 - S-Match (Aussagenlogik)
 - ILIADS (OWL/Beschreibungslogik + syntaktisches Matching)

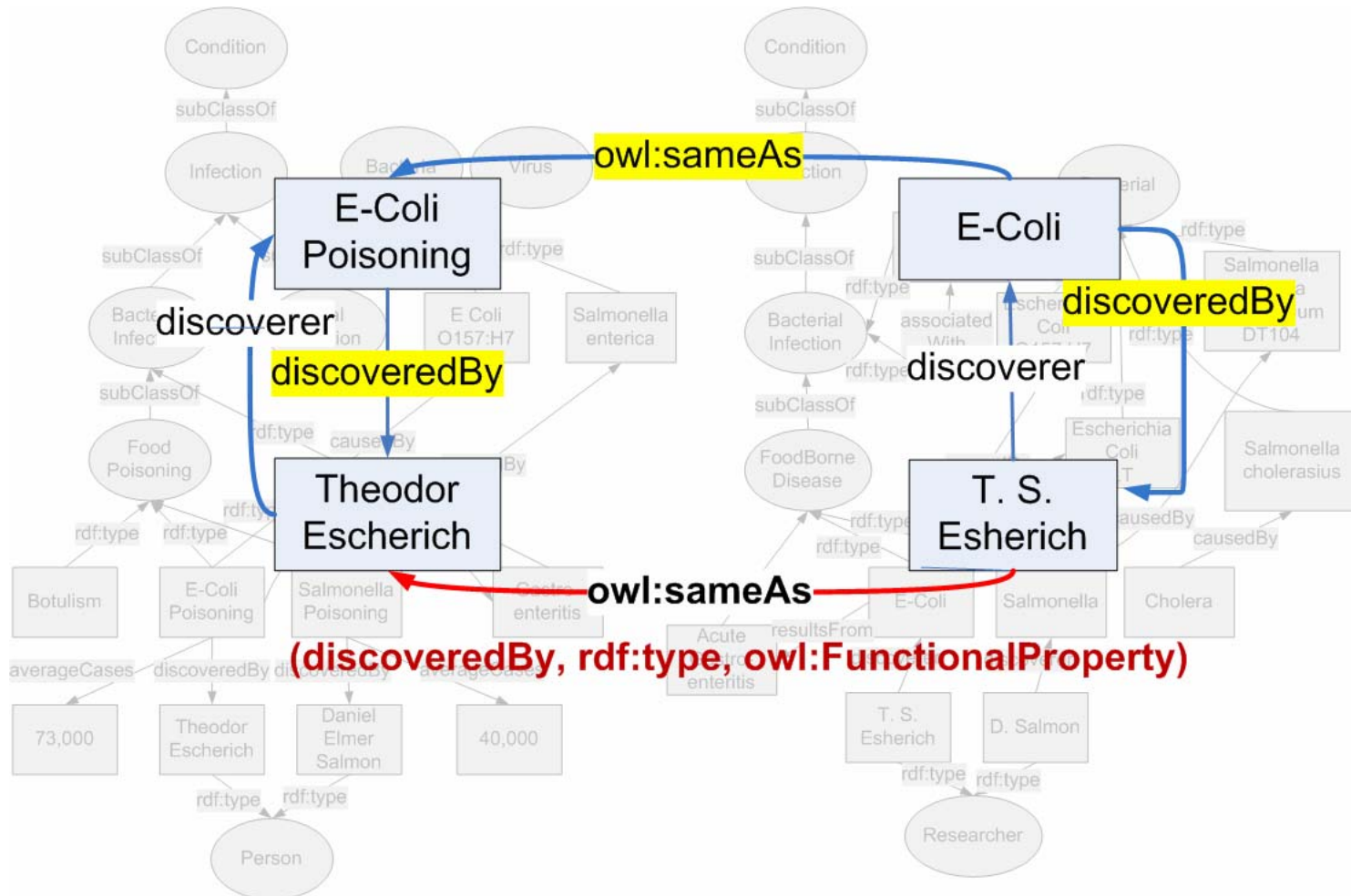
B.7. Semantisches Matching ILIADS

32



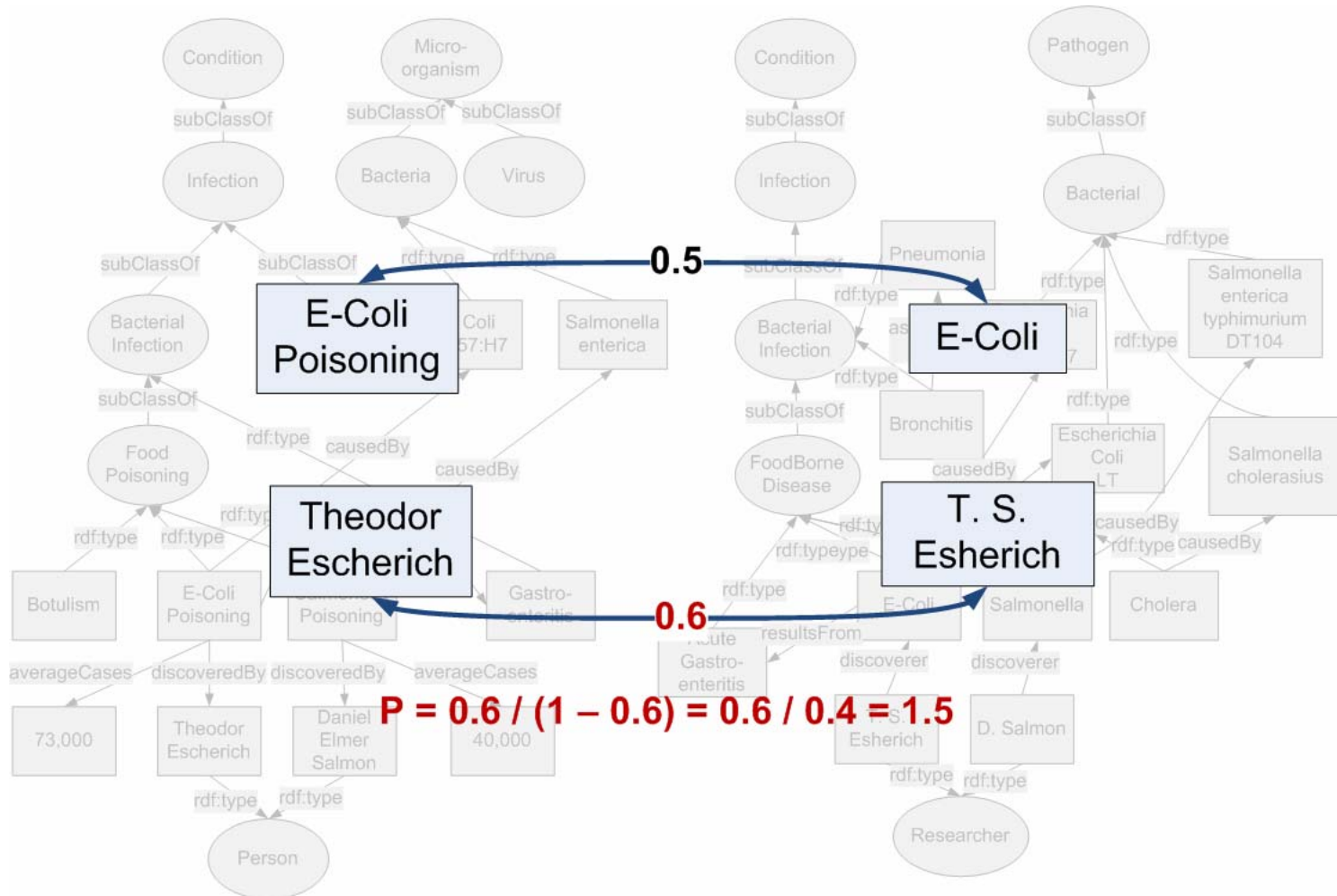
Source: <http://www.cs.umd.edu/projects/linqs/iliads>

B.7. Semantisches Matching ILIADS



Source: <http://www.cs.umd.edu/projects/linqs/iliads>

B.7. Semantisches Matching ILIADS



Source: <http://www.cs.umd.edu/projects/linqs/iliads>

C.1. Evaluationsverfahren

37

- Methoden

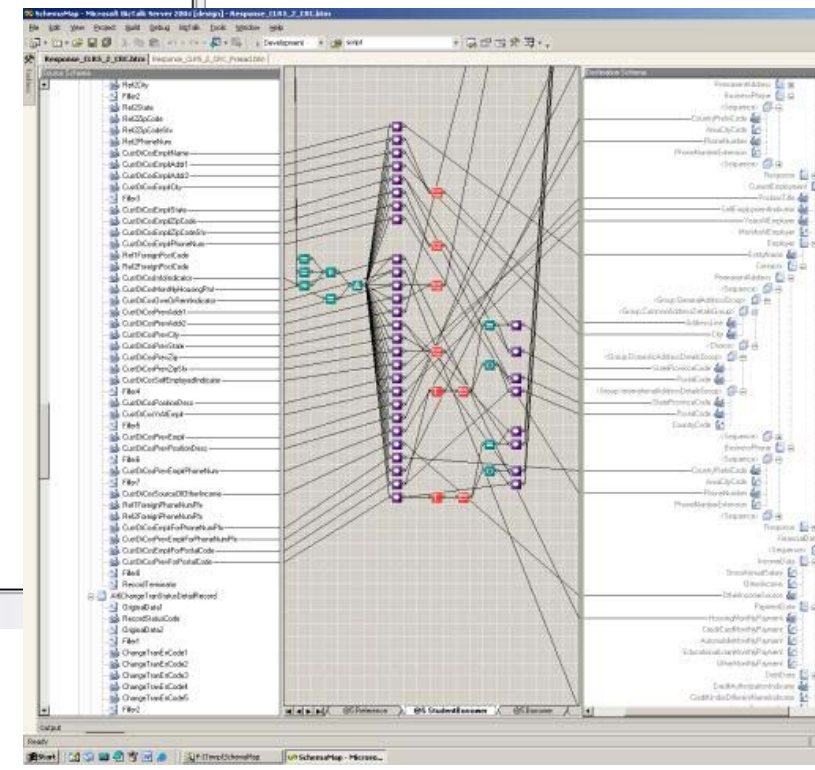
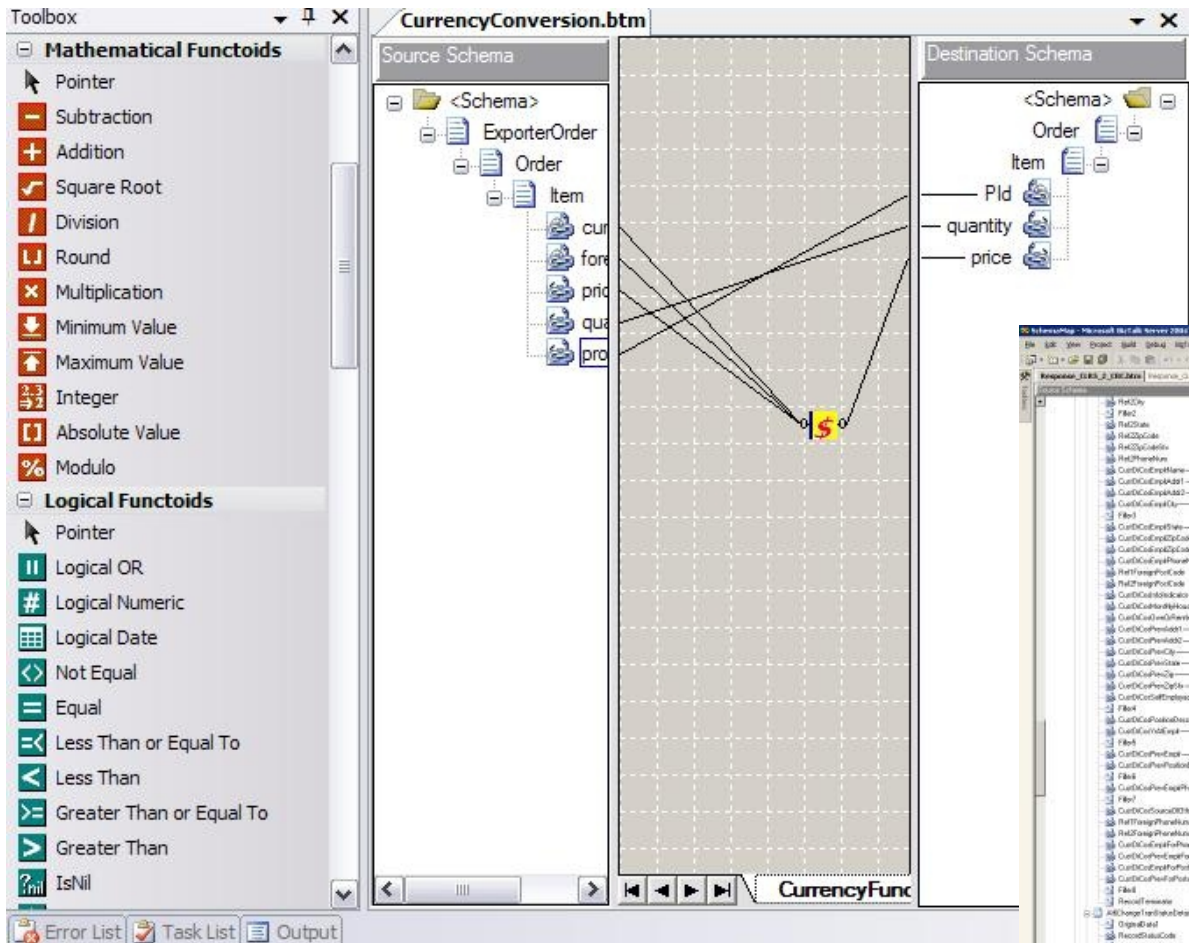
- Benchmarks
 - OAEI (Ontology Alignment Evaluation Initiative)
 - XBenchMatch
 - Thalia Benchmark
 - ...

- Diverse Datensätze

- Maße
 - Recall / Precision / F-Measure
 - OverAll Maß (Cupid)

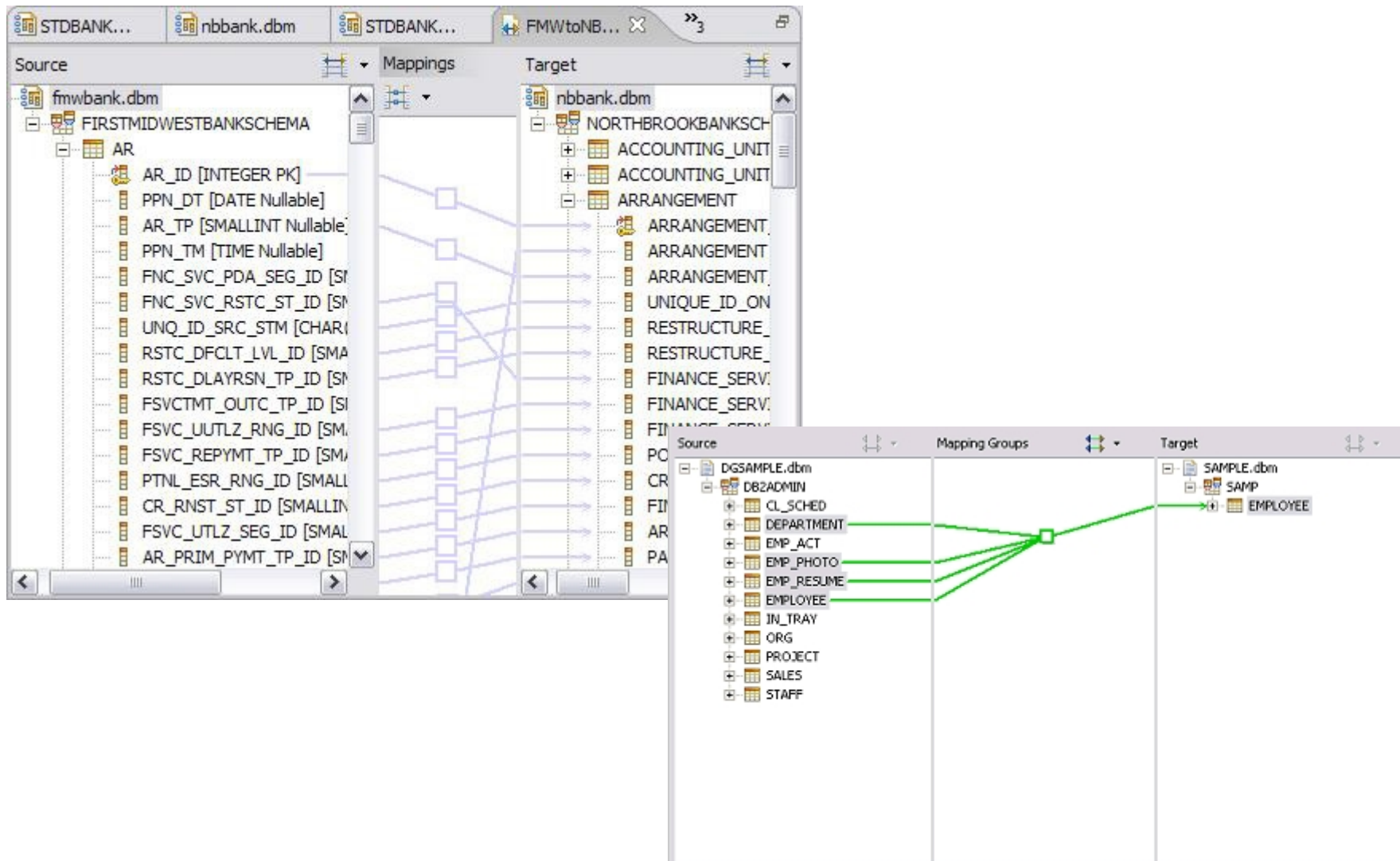
C.2. Microsoft BizTalk-Mapper

38



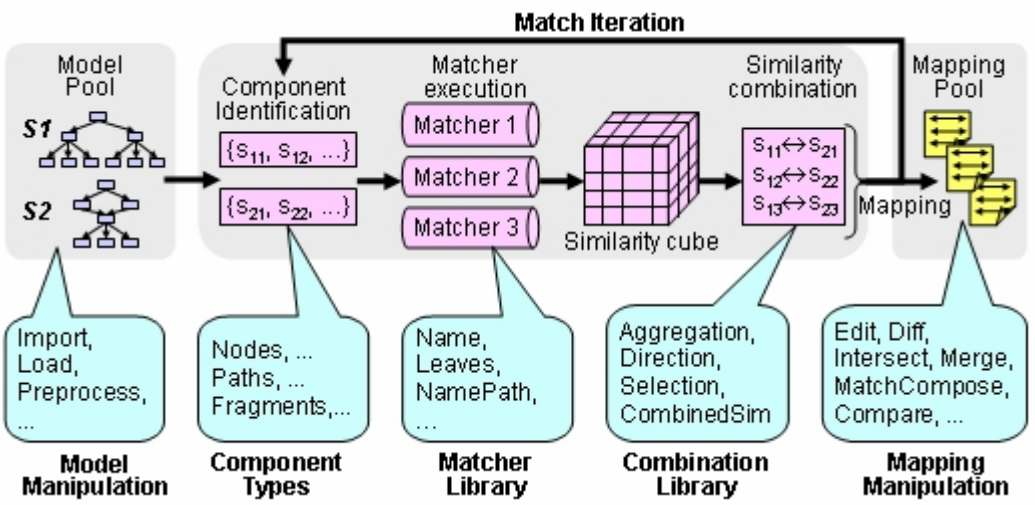
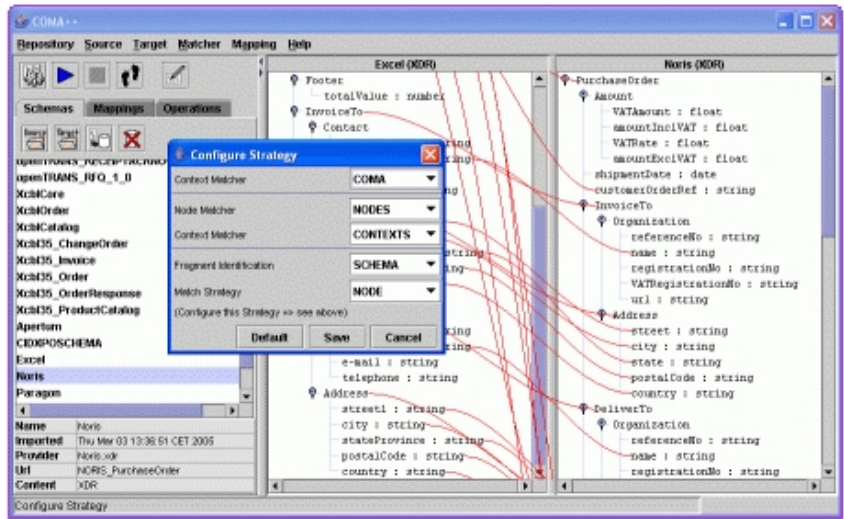
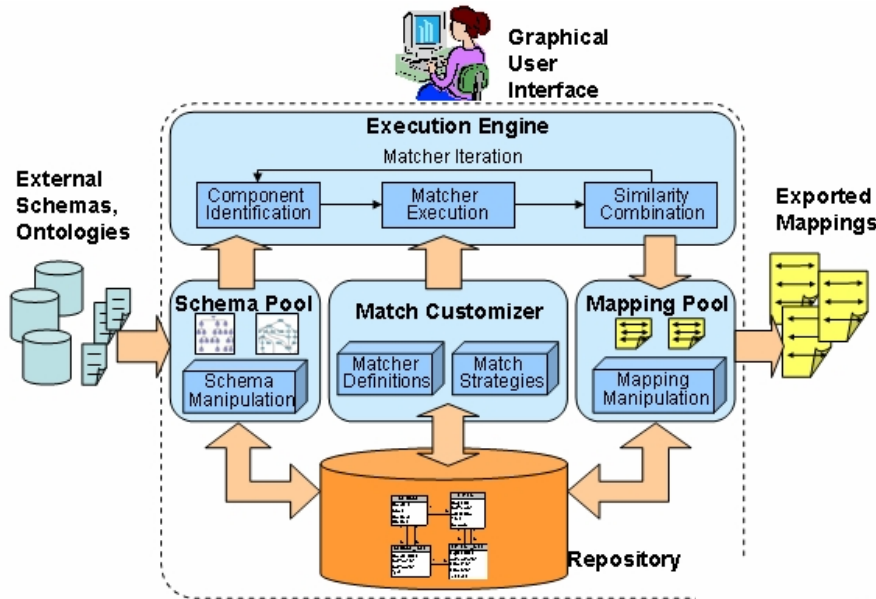
C.3. IBM Rational Data Architect

39



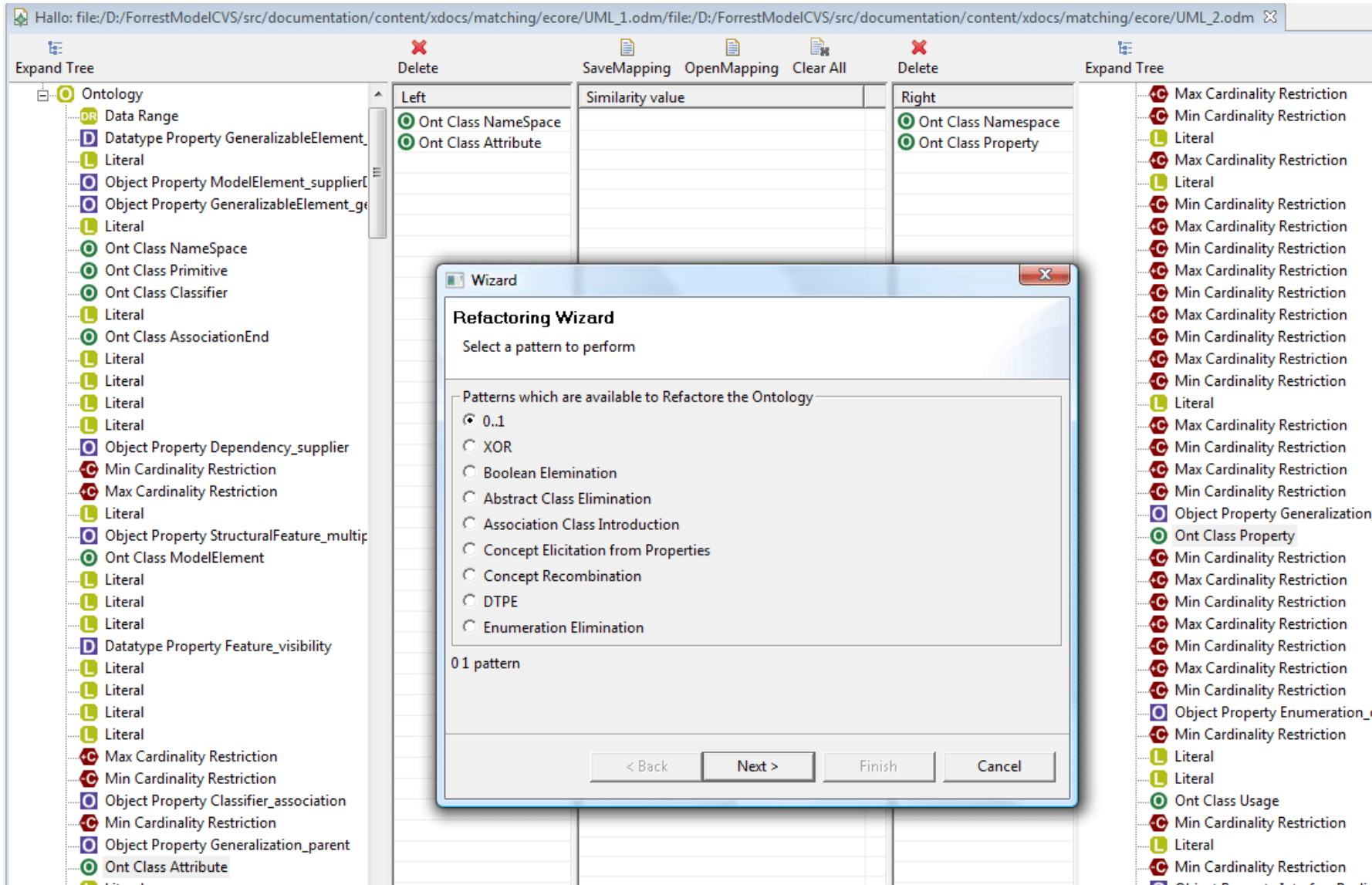
C.4. Coma++/QuickMig (Uni Leipzig/SAP)

40



C.4. Coma++/QuickMig Eclipse-Plugin Ontohammer

41



The screenshot displays the Eclipse IDE interface with the Ontohammer plugin. The main window shows a project named 'Hallo' with two UML files open. The 'Expand Tree' view on the left lists various ontology elements, including 'Ontology', 'Data Range', 'Datatype Property GeneralizableElement', 'Literal', 'Object Property ModelElement_supplier', 'Object Property GeneralizableElement_ge', 'Ont Class NameSpace', 'Ont Class Primitive', 'Ont Class Classifier', 'Ont Class AssociationEnd', 'Object Property Dependency_supplier', 'Min Cardinality Restriction', 'Max Cardinality Restriction', 'Object Property StructuralFeature_multip', 'Ont Class ModelElement', 'Datatype Property Feature_visibility', and 'Ont Class Attribute'. The central area contains a table with columns 'Left', 'Similarity value', and 'Right'. The 'Left' column lists 'Ont Class NameSpace' and 'Ont Class Attribute'. The 'Right' column lists 'Ont Class NameSpace' and 'Ont Class Property'. A 'Wizard' dialog box titled 'Refactoring Wizard' is open in the foreground, prompting the user to 'Select a pattern to perform'. The dialog lists several patterns available for refactoring the ontology: '0..1' (selected), 'XOR', 'Boolean Elimination', 'Abstract Class Elimination', 'Association Class Introduction', 'Concept Elicitation from Properties', 'Concept Recombination', 'DTPE', and 'Enumeration Elimination'. At the bottom of the dialog, it indicates '0 1 pattern' and provides navigation buttons: '< Back', 'Next >', 'Finish', and 'Cancel'.

C.5. Model Management

42

Ziele

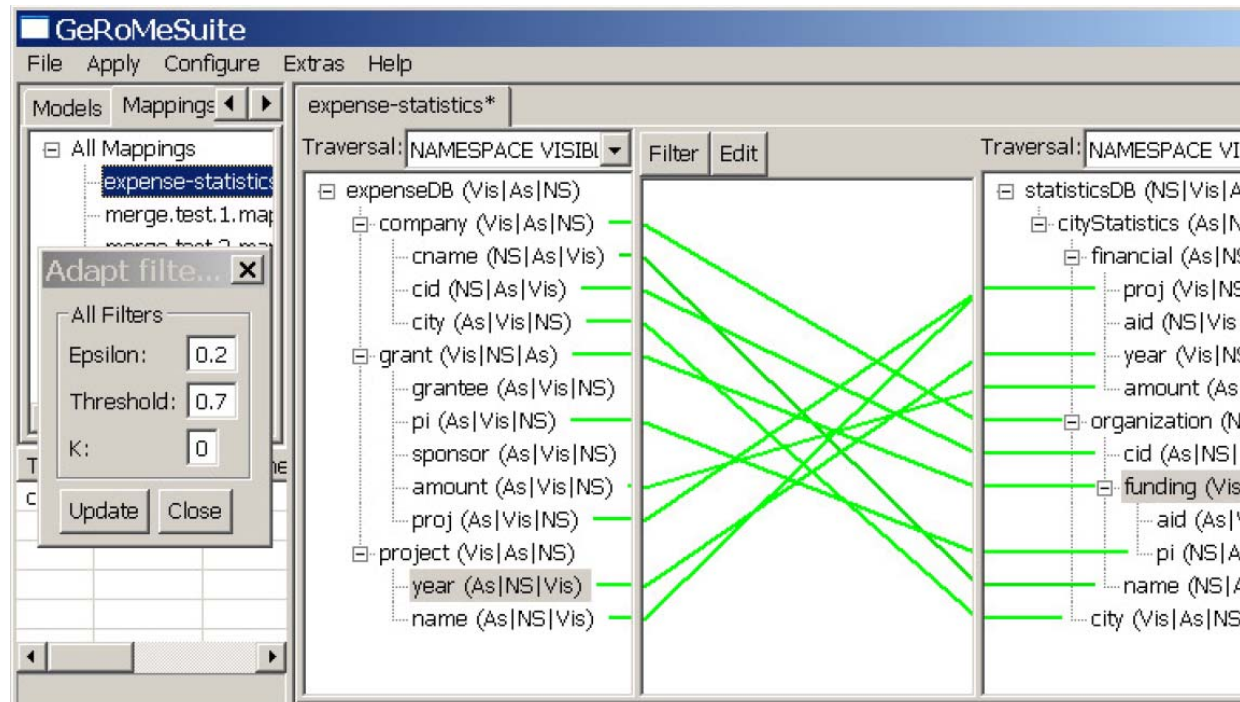
- Data Programmability
- Modellevolution
- Generische Operatoren für Datenmodelle

Operatoren

- Match
- Merge
- Compose
- Inverse
- Diff
- ModelGen
- ...

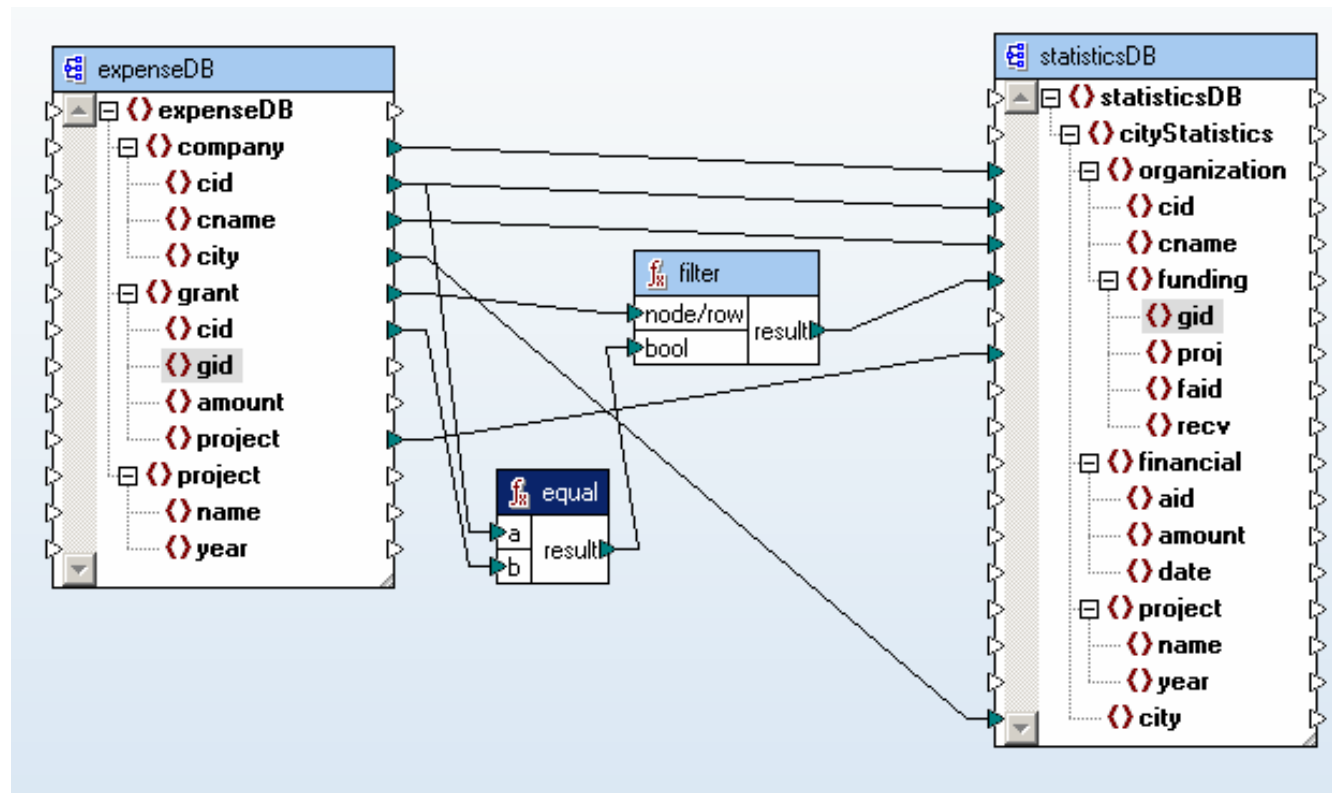
Systeme

- Rondo (Microsoft Research)
- GeRoMeSuite (RWTH Aachen)



C.6. Altova MapForce

43



Themenübersicht (Betreuer)

A. Modellierung (Frank Kaufer)

B. Matchingtechniken

1. Intrinsisches Elementmatching (Melanie Weis)
2. Extrinsisches Elementmatching (Frank Kaufer)
3. Extensionales Matching (Melanie Weis)
4. Strukturbasiertes Matching (Frank Kaufer)
5. Globales Matching (Frank Kaufer)
6. Maschinelles Lernen (Frank Kaufer)
7. Semantisches Matching (Frank Kaufer)

C. Matchingsysteme & Evaluation

1. Evaluationsverfahren (Felix Naumann)
2. Microsoft BizTalk-Mapper (Felix Naumann)
3. IBM Rational Data Architect (Felix Naumann)
4. Coma++/QuickMig (SAP/Uni-Leipzig) (Felix Naumann)
5. Model Management (Alexander Albrecht)
6. Altova MapForce (Alexander Albrecht)



Weitere Themen bei Bedarf

Organisation



Ablauf

- Mittwoch, 9:15 – 10:45: interne Treffen
- Vorträge en bloc im Dezember und/oder Februar
- Betreuungstreffen: individuell per Email anfragen
- Mailingliste
 - seminar-sm0708@hpi.uni-potsdam.de
 - <http://lists.hpi.uni-potsdam.de/listinfo/seminar-sm0708>
 - Fachliche & organisatorische Diskussion
 - Archiv
- Wikizugang nach Anmeldung
 - <http://www.hpi.uni-potsdam.de/naumann/a/matching>

Themenvergabe

- 3 Wunschthemen
- Email an Frank Käufer
- Vergabe nach
 - Priorität
 - Kohärenz
- Bewerbungsfrist: Donnerstag, 25.10.2007, 23:59 Uhr
- Anmeldefrist: 06.11.2007

Leistungserfassungsprozess

- Gesamtbewertung von
 - Wikiarbeit (Deadline 01.03.2008)
 - Vortrag 30 min + Diskussion
 - Praktische Arbeit

- Bonus
 - Wiki -> Wikipedia
 - Mailingliste
 - Gemeinsame, kombinierte Implementierungen

Fragen!