

Innovationslabor Semantische Integration von Webdaten

Prof. Dr. Erhard Rahm

<http://dbs.uni-leipzig.de/format>

GEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung



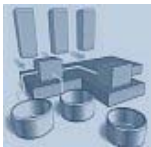
UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



Programmablauf

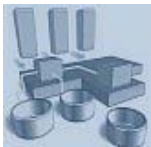
- Überblicksvortrag Prof. Rahm
- Feedback / Diskussion
- Vorstellung der Prototypen
 - Frau Maßmann
 - Dr. Thor
- Feedback / Diskussion
- Get-Together (Felix-Klein-Hörsaal)



BMBF: ForMaT



- Förderprogramm:
 - Forschung für den Markt im Team
 - Im Rahmen von Unternehmen Region (Innovationsinitiative für die Neuen Länder), <http://www.unternehmen-region.de>
- Ziel
 - Abstrakt: frühzeitiger Wissens- und Technologietransfer von der öffentlichen Forschung in die Wirtschaft
 - Konkret: Firmenausgründung nach Ende der Förderung
- Ansätze
 - Analyse der Forschungsuntersuchungen auf Verwertungseignung
 - Förderung der interdisziplinären Zusammenarbeit
- Förderung in mehreren Runden mit jeweils zwei Phasen
 - Phase 1: Potential-Screening, 6 Monate
 - Phase 2: Weiterentwicklung innerhalb von Innovationslabor, 2 Jahre
- Stand (Runde 2):
 - Derzeit 27 Initiativen in Phase 1 (ausgewählt aus ca. 100)



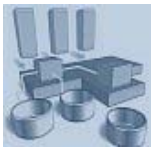
UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



Semantische Integration von Webdaten

- Initiative der Universität Leipzig, Abt. Datenbanken (Prof. Rahm)
- Aktuell noch in Phase 1, dem Potenzial-Screening
- Ziel: weitgehende Automatisierung von Datenintegrationsaufgaben mit hoher Datenqualität, insbesondere für Webdaten
 - Schnellere Umsetzung als mit Data-Warehouse-Ansätzen (Offline-Datenaufbereitung)
 - On-Demand Datengewinnung (fokussierte Suche) , Datenabgleich und Analyse
 - Gegenüber einfachen Mashup-Lösungen Anwendbarkeit auf großen Datenmengen und Unterstützung für hohe Datenqualität
- Basiert auf eigenen Forschungsergebnissen / Prototypen
 - *Workflow-basierte Datenintegration mit fokussierter Websuche: iFuice*
 - *Objekt-Matching, Dublettenbehandlung: MOMA und STEM*
 - *Schema und Ontologie-Matching: COMA und COMA++*

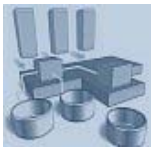
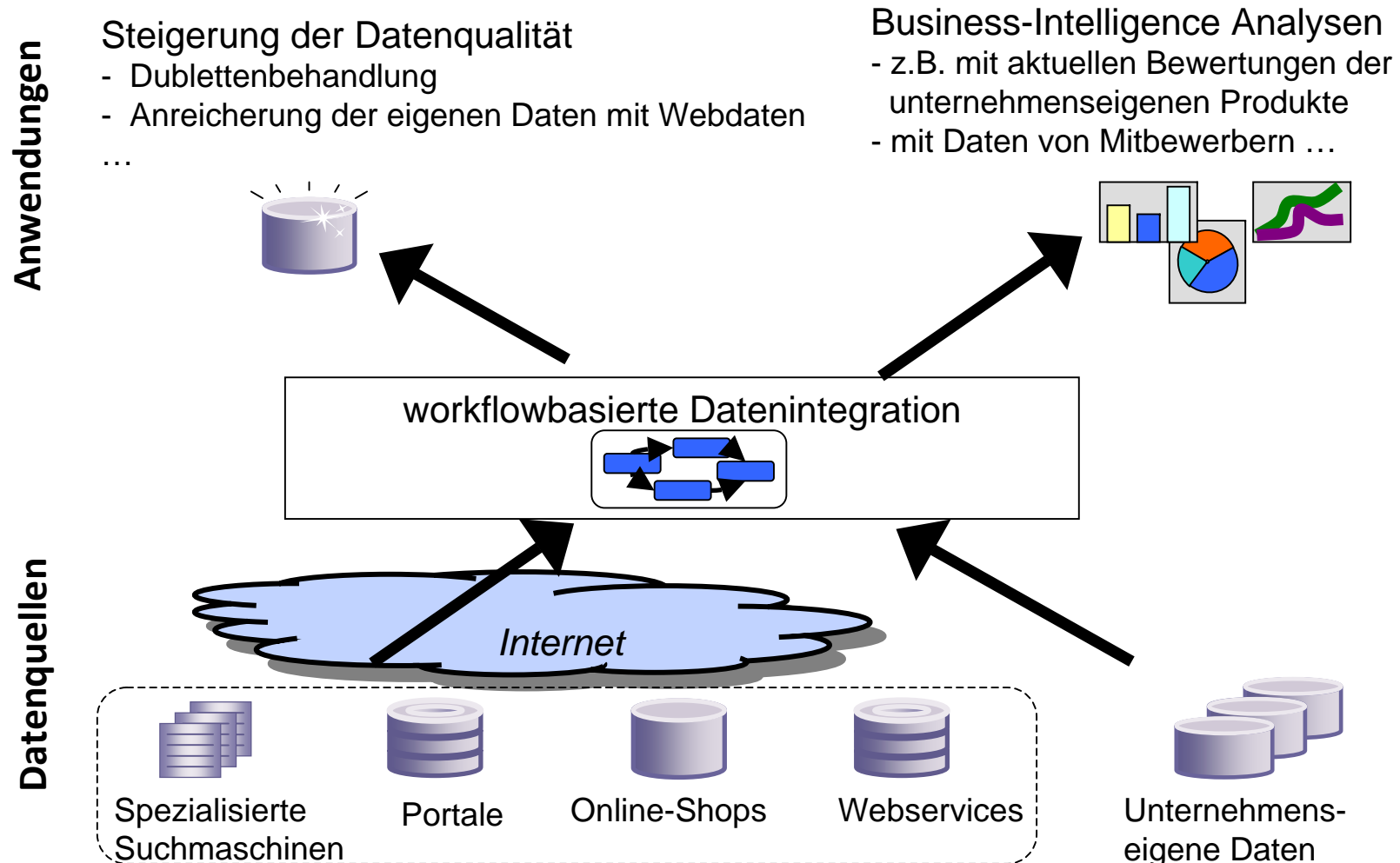


UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



Semantische Integration von Webdaten



Projektbeteiligte



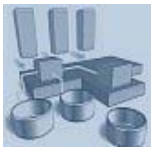
- Prof. Dr. E. Rahm
- Dipl.-Inf. S. Maßmann
- Dr. A. Thor



- Prof. Dr. C. Lindemann
- Dipl.-Inf. C. Hoffmann

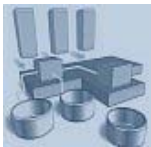


- Prof. Dr. H. Löbler
- Dipl.-Kfm. R. Müller



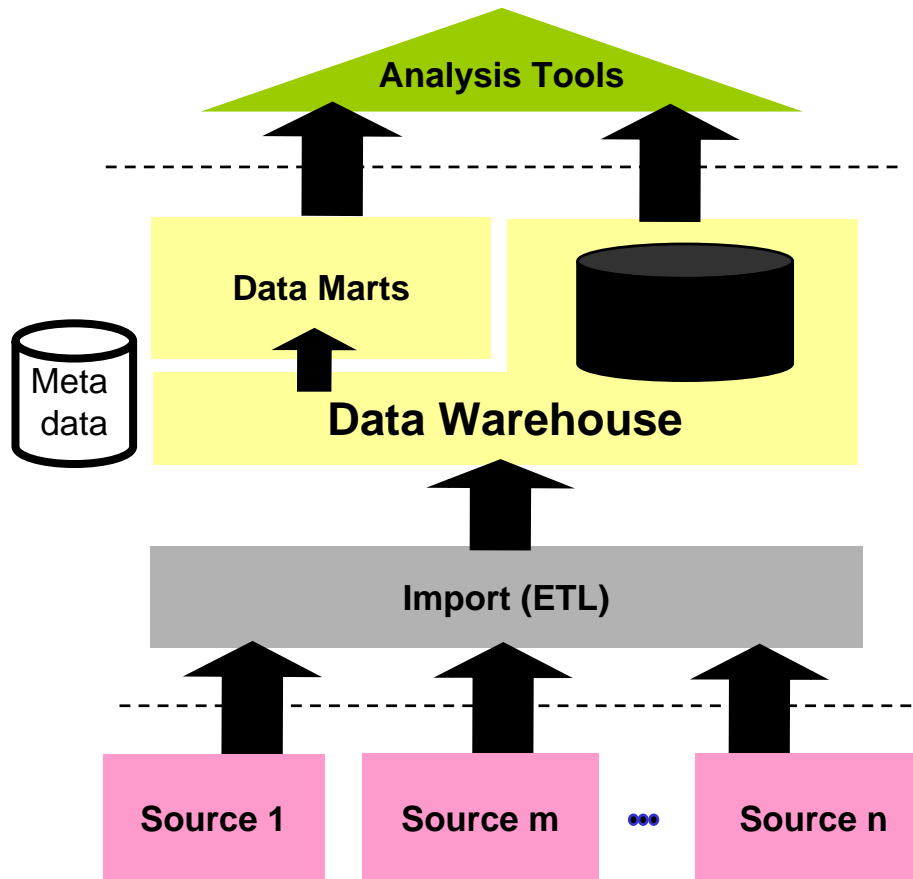
Vortragsinhalte

- Einleitung
- Derzeitige Ansätze zur Datenintegration
 - Warehousing, EII
 - Mashups
- Eigene Forschungsarbeiten / Prototypen
 - Schema/Ontologie Matching: COMA++
 - Objekt-Matching: MOMA, STEM
 - Workflowbasierte Integration: iFuice
- Feedback / Diskussion

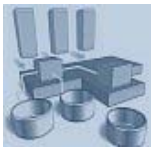
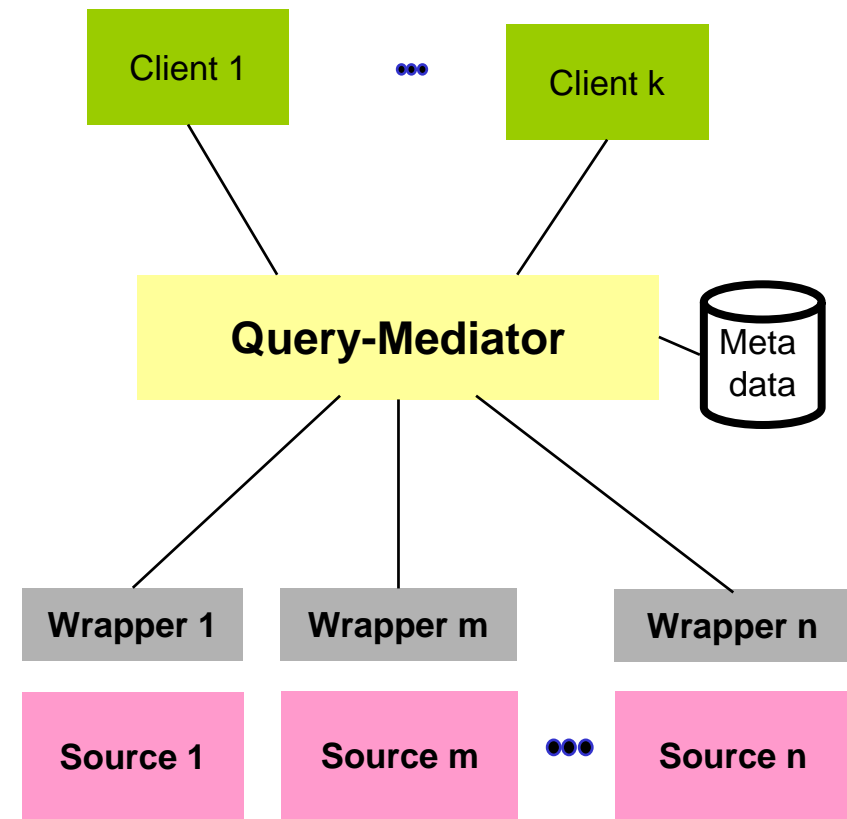


Physische vs. virtuelle Datenintegration

Data Warehousing, ETL (physische Datenintegration)

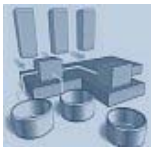


Enterprise Information Integration, EII (virtuelle Integration, Query-Mediatoren)



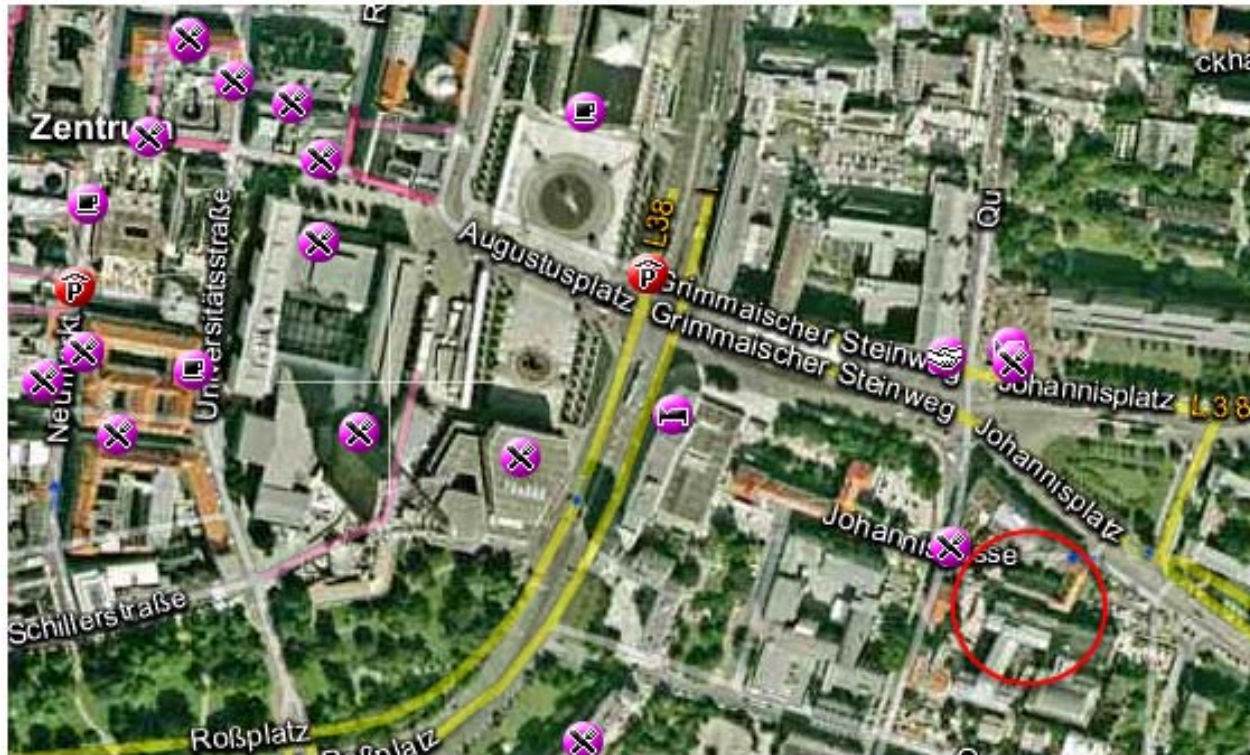
Probleme ETL / EII

- Hoher Realisierungsaufwand, lange Vorlaufzeiten
 - Einheitliche Sicht auf heterogene und verteilte Datenbestände
 - Aufwändige Erstellung und Wartung
 - für gemeinsames Schema sowie
 - für Mappings zu Quell-Schemas
- Unzureichende Unterstützung für Webdaten
 - Beschränkte Zugänglichkeit über Suchformulare bzw. Suchmaschinen
 - Größere Heterogenität der Daten
 - Geringe Datenqualität



Mashups – Datenintegration

Mashup = interaktive Web-Applikation zu Kombination von Daten mehrerer Online-Datenquellen im Rahmen eines neuen Dienstes



www.goyellow.de



UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



Mashup-Tools

- Basieren oft auf RSS-Feed
- Angepasst auf einige konkrete Datenquellen mit meist geringen Datenmengen
- Schnell zu erstellen
- Einfache Operationen mit Daten möglich z.B. Filtern, Sortieren, Teilen
- Z.B. Yahoo-Pipes (pipes.yahoo.com)

The screenshot shows the Yahoo Pipes editor interface. The pipe is titled "US population by state". The workflow consists of four modules:

- Fetch CSV:** URL: `http://www.census.gov/popest/st/`. Column separated by `.`. Skip the first 3 rows. Use rows 4 to 4 as column names.
- Filter:** Permit items that match all of the following. Rule: `item.col_1` Contains `.`
- Regex:** Use regular expression patterns here. Rule: In `item.col_1` replace `l` with `text`. Rule: In `item.col_1` replace `Washington` with `Washington State`.
- Truncate:** Truncate feed after 51 items.

The screenshot shows the output page for the pipe "US population by state". The page title is "US population by state" and it includes the author's name "Jonathan". The page description states: "Uses the census data from census.gov (in CSV format) to give a sorted list of population by state". The pipe web address is `http://pipes.yahoo.com/jonathan/_ja89ese3B6GM20YNivXUAA`. The page includes a "Map" view showing a map of the United States with red pins indicating population data for each state. The "List" view shows a sorted list of population by state. The page also includes a "View Source" link and a "Report abusive Pipe" link.

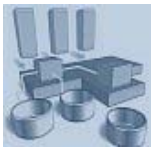


UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik

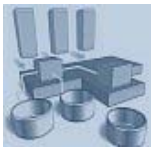
Mashups: Merkmale und Probleme

- Schnelle Entwicklung einfacher Workflows
 - zum Teil Erstellung über GUI; z.B. Yahoo Pipes
- Nutzung von Web-APIs / Web Services
- Einfache Datenintegration
 - z.B. über geographische Koordinaten
 - Kombination von Suchergebnissen zu Keyword
- Einschränkungen
 - Einfache Anfragen (keine Transformation für unterschiedliche Quellen)
 - Einfache Ergebnismachbearbeitung (Merge statt Match)
 - Unzureichende Vollständigkeit / Genauigkeit
 - Geringe Datenqualität (Dubletten)
 - Unternehmenstauglichkeit



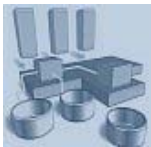
Vortragsinhalte

- Einleitung
- Derzeitige Ansätze zur Datenintegration
 - Warehousing, EII
 - Mashups
- Eigene Forschungsarbeiten / Prototypen
 - Schema/Ontologie Matching: COMA++
 - Objekt-Matching: MOMA, STEM
 - Workflowbasierte Integration: iFuice
- Feedback / Diskussion



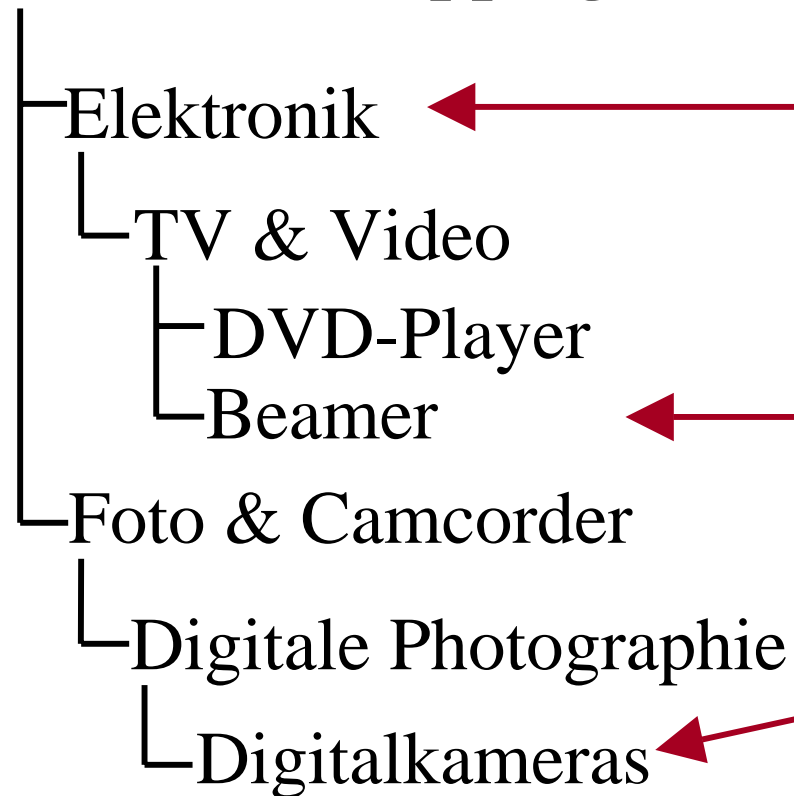
Schema / Ontologie-Matching

- Finden semantischer Korrespondenzen zwischen 2 Schemas bzw. Ontologien
 - DB-Schemas, XML-Schemas
 - Inhaltskategorisierungen, Produktkataloge, Suchformulare, ...
- Input: 2 Schemas (Ontologien) S_1 and S_2
 - Evtl. Dateninstanzen zu S_1 und S_2
 - *Hintergrundwissen*
- Output: *Mapping* zwischen S_1 und S_2
- Kritischer Schritt in zahlreichen Applikationen
 - Datenintegration: Mappings zwischen Datenquellen bzw. zwischen Datenquelle und globalem Schema
 - E-Business: XML-Datentransformation
 - Katalogintegration, Finden verwandter Inhalte ...

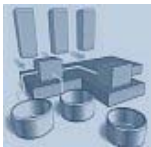
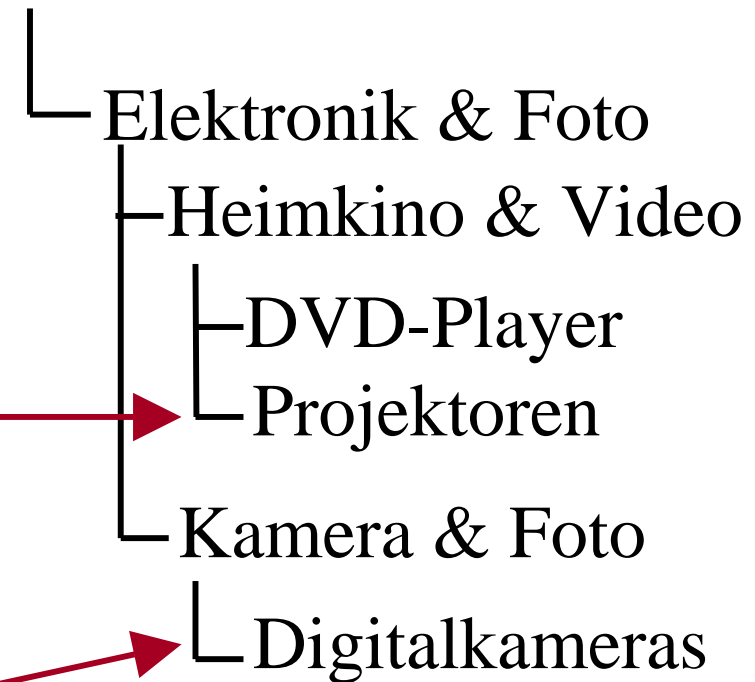


Match-Beispiel: Produktkataloge

Yahoo.de Shopping

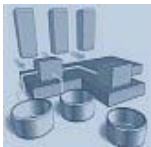


Amazon.de



Derzeitige Situation

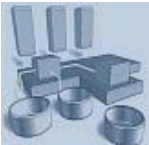
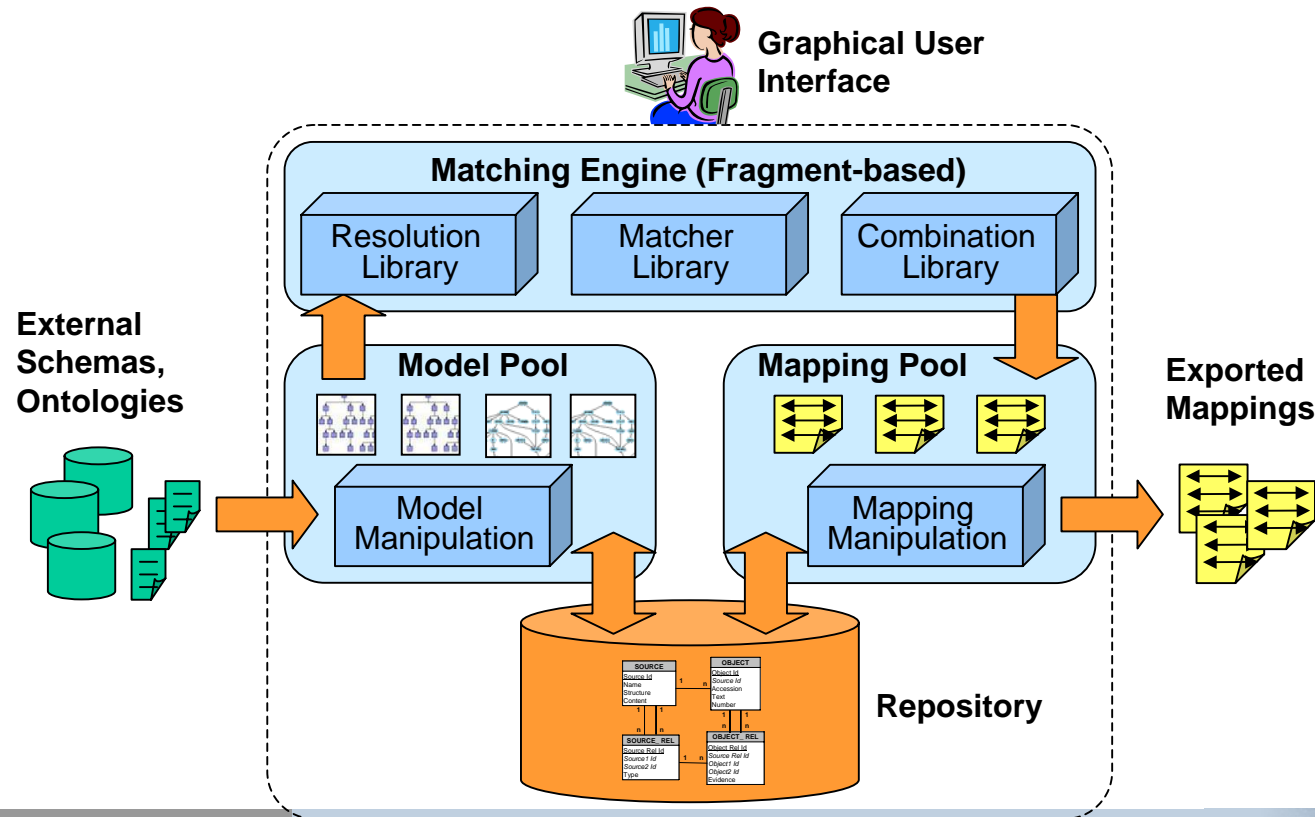
- Finden der Mappings ist Flaschenhals
 - Weitgehend manuell
 - Aufwändig und fehleranfällig
- Zunehmende Dringlichkeit
 - Wachsende Nutzung verteilter Daten
 - Starke Zunahme an XML-Daten und -Schemas
 - Zunahme an Ontologien / semantischen Kategorisierungen
- Skalierbarkeit erfordert semi-automatische Tools
 - Vollautomatische Lösungen aufgrund semantischer Heterogenität nicht perfekt
 - Namensproblematik (Synonyme, Homonyme)
 - begrenzte Mächtigkeit von Metadaten / Schemasprachen





Generische Match-Plattform

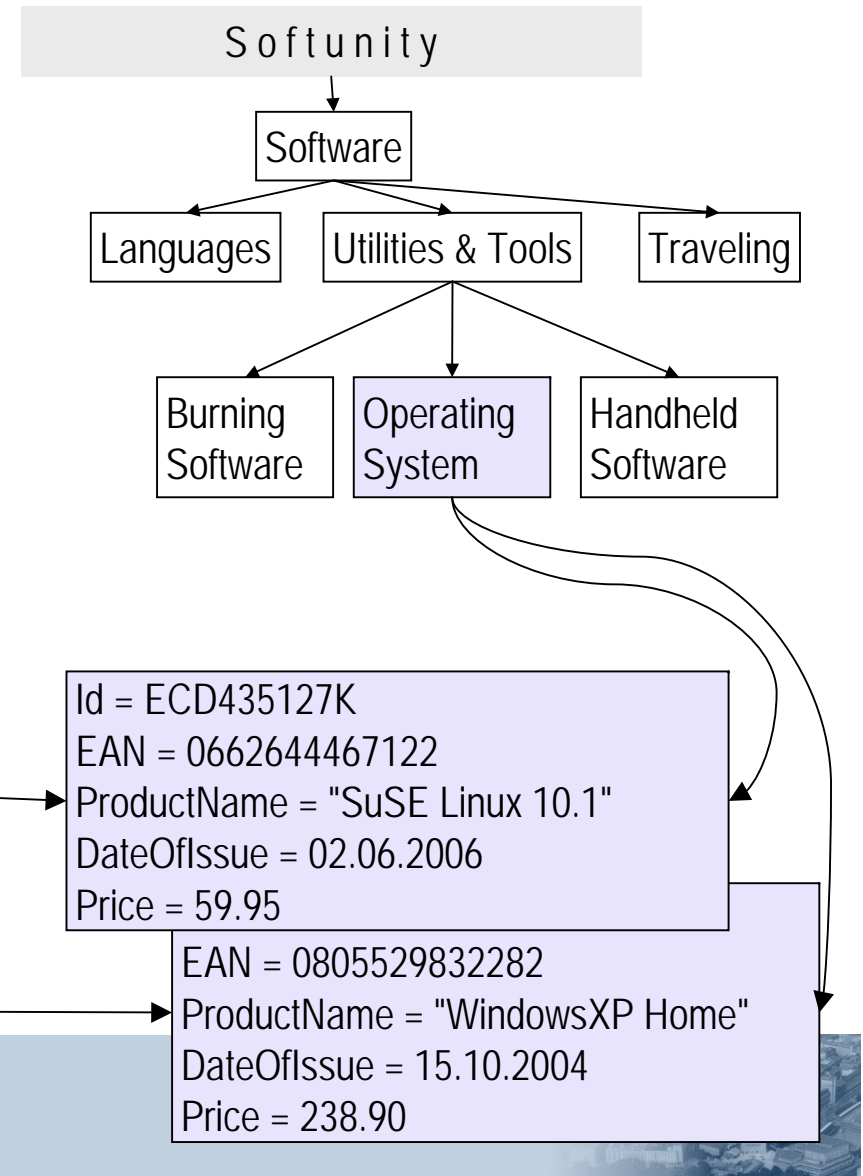
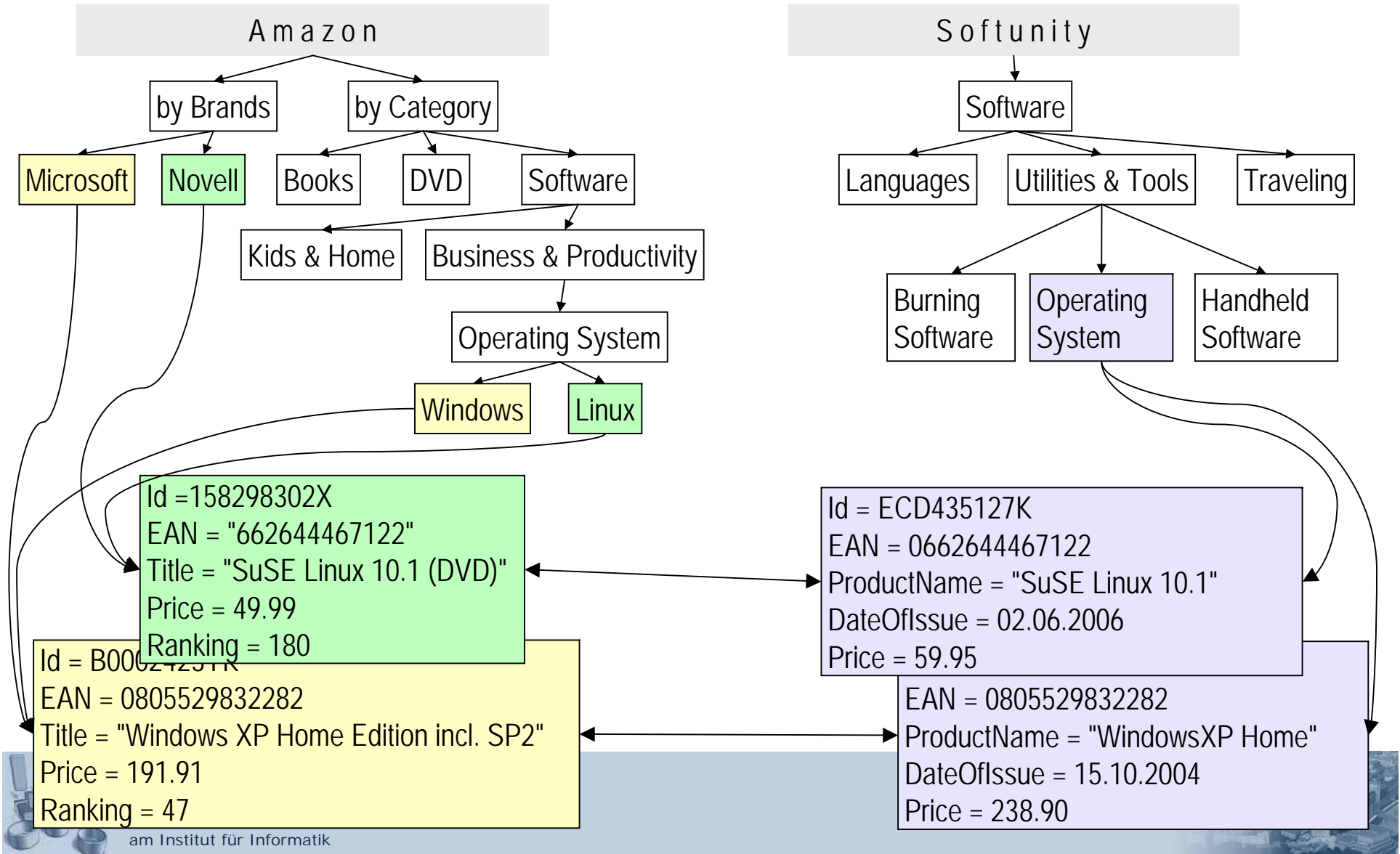
- Unterstützt relationale und XML-Schemas, OWL-Ontologien
- Flexible Kombination mehrerer Match-Algorithmen
- Match-Strategien für große Schemas
- Wiederverwendung vorhandener Match-Ergebnisse



UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik

Instanz-basiertes Ontologie-Matching



Id = 158298302X
 EAN = "662644467122"
 Title = "SuSE Linux 10.1 (DVD)"
 Price = 49.99
 Ranking = 180

Id = ECD435127K
 EAN = 0662644467122
 ProductName = "SuSE Linux 10.1"
 DateOfIssue = 02.06.2006
 Price = 59.95

Id = B00021251K
 EAN = 0805529832282
 Title = "Windows XP Home Edition incl. SP2"
 Price = 191.91
 Ranking = 47

EAN = 0805529832282
 ProductName = "WindowsXP Home"
 DateOfIssue = 15.10.2004
 Price = 238.90

Objekt-Matching-Problem

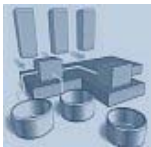
- Identifikation semantisch äquivalenter Objekte
 - innerhalb einer Datenquelle oder zwischen verschiedenen Quellen
 - um Objekte zu integrieren/mischen, zu vergleichen, Dubletten zu eliminieren, etc.

Quelle1: Kontakt

<i>KID</i>	<i>Name</i>	<i>Strasse</i>	<i>Stadt</i>	<i>Frau</i>
11	Kristen Schmid	Hanse Pl 2	Berlin	1
24	Christian Schmied	Hanse Str 2	Berlin	0

Quelle2: Kunde

<i>Knr</i>	<i>Nachname</i>	<i>Vorname</i>	<i>Geschl</i>	<i>Adresse</i>	<i>Telefon</i>
11	Schmid	Chris	M	Hansestr. 2, 18182 Bentwich	
493	Schmid	Kris L.	W	Hansa-Platz 2, 10557 Berlin	030-9627621



Objekt-Matching-Problem für Webdaten

Viele heterogene Dubletten

- unterschiedliche Namen
- nur teilstrukturiert
- unterschiedliche Formate
- heterogene Inhalte ...



CANON VIXIA HV30

CANON VIXIA HV30 _zustand: Neu - price inkl.
[Zur Einkaufsliste hinzufügen](#)

€656,90

www.ATLANTIV.com



Canon VIXIA HV30 High Definition Camcorder (PAL)

Canon VIXIA HV30 High Definition Camcorder **Canon HV30** inherits the HV20's **Canon** HD Camera System – the unique combination of a genuine **Canon** HD Video Lens, ...
★★★★★ [285 Beurteilungen](#) - [Zur Einkaufsliste hinzufügen](#)

€814,48

[DSLR Cameras Division - Phones Corporation Ltd.](#)



Canon VIXIA HV30 DV Camcorder

Canon VIXIA HV30 DV Camcorder Digital, DV, NTSC, Bis zu 3.15 Megapixels Einzelbild, 10x Optischer Zoom, 0,5 kg Mit dem HD-Camcorder **HV30** präsentiert **Canon** ...
[Zur Einkaufsliste hinzufügen](#)

€749,00

komponentenfabrik.de



CANON VIXIA HV30 Ladegeräte, VIXIA HV30 Ladegeräte

CANON VIXIA HV30 Ladegeräte, neu Ladegeräte für **CANON VIXIA HV30**, 1 Jahr Garantie, Akku lokalisiert in Deutschland.
[Zur Einkaufsliste hinzufügen](#)

€19,01

kaufen-akkus.de



Canon VIXIA HV30 Camcorder

The **Canon VIXIA HV30** Camcorder is a versatile High Definition camcorder which ... The **Canon VIXIA HV30** lets moviemakers adapt the camcorder to their needs ...
[Zur Einkaufsliste hinzufügen](#)

€1.139,22

DigitalRev.com
[3 Händlerbewertungen](#)



UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



Canon HF100 Camcorder - 3.31 MP - 12 x opt. .



von €532 bis €873 bei [76 Händl](#)

Flash card, SD Memory Card, SDHC-Speicher
 Der attraktive HF100 speichert ganz bequem 1.9 dabei einfach atemberaubend brillante HD Video
 Der HF100 speichert HD-Videos direkt auf SDH teilen. Die Karten lassen sich nämlich prima ver
 Der HF100 ist ein Camcorder im superkompakte Speicherung ohne mechanisch bewegliche Korr
 Der HF100 zeichnet in 1.920 x... [mehr >](#)

[Zur Einkaufsliste hinzufügen](#)

[Preise vergleichen](#)

Technische Spezifikationen

[Verw](#)

[Allgemein](#) - [Hauptmerkmale](#) - [Speicher / Speicher](#) - [Objektivsystem](#) - [Zusätzl](#)
[Batterie](#) - [Universal Product Identifiers](#)

Allgemein

Produkttyp	Camcorder
Breite	7.3 cm
Tiefe	12.9 cm
Höhe	6.4 cm
Gewicht	380 g

Hauptmerkmale

Sensorauflösung Camcorder	3.31 Mpix
Effektive Videoauflösung Camcorder	2.07 Mpix
Effektive Fotoauflösung Camcorder	2.76 Mpix
Widescreen-Videoaufzeichnung	Videoaufnahme im Breitbildformat
Medientyp	Flash card
Farbunterstützung	Farb
Typ des optischen Sensors	CMOS
Größe des optischen Sensors	1/3.2"
Mindestbeleuchtung	0.2 Lux
Digitales Zoom	200 x
Aufnahmegeschwindigkeit	

Informationsfusion

Canon HF100 Produkt vergleichen

Camcorder

- Video-System: HD-Video
- Zoom: 12x optisch, 200x digital
- Brennweite: 4,80 mm - 57,60 mm
- Bild-Sensor 1/3,20"
- 3.310.000 Pixel Bild-Sensor
- 2,70" LCD-Monitor
- 211.000 Pixel LCD-Monitor Auflösung
- Standbildaufnahme
- Effektiv 2.760.000 MegaPixel

► **Meinung Ø:** ★★★★★ (15 Meinungen)
 ► **Testnote Ø:** 1,8 (16 Testberichte)

Hersteller: [Canon](#)

Alle Angaben ohne Gewähr

PREIS
SUCHMASCHINE

Abbildung ähnlich

Canon HF100 günstig bei [ebay](#) ersteinern

Preisvergleich ab € 553,99 **Meinungen** Preisverlauf

Meinung zu Canon HF100

Klasse Cam 09.09.2008

Hallo, nach Vergleich mit einer Panasonic 3CCD (HDC-SD9) muss ich sagen, das im Low-Light-Bereich die Canon mehrere Nasen Vorsprung vor der 3CCD-Panasonic hat. Die Aufnahmen sind selbst bei wenig Licht sehr gut. Die Panasonic hat hier eine wesentlich schlechtere Bildqualität. Bei ausreichend Licht sind beide Cams in etwa gleichwertig. Wenn man eine 3-Chip-Cam mit 1-Chip-Cam vergleicht, so sollte die 3-Chip-Variante deutlich teurer sein. Bei Preisgleichstand gewinnt meist die 1-Chip-Cam. Insbesondere weil die Canon einen empfindlichen C-MOS Chip hat. Positiv bei der Canon: Größe + Gewicht Haptik Schärfe/Bildqualität Tonqualität Bildqualität von Fotos Negativ bei der Canon: Automatischer Weißabgleich Akkulaufzeit mitgelieferte Software Bei der Canon sollte man die vordefinierten Weissableichinstellungen, oder den manuellen Weissabgleich bevorzugen (wenn möglich). Ich werde die Canon behalten und die Panasonic zurückgeben. Für das Geld gibts z.Zt. nichts besseres.

Meine Wertung: ★★★★★

Autor: [skeptiker](#)

Canon HF100 Test (Camcorder)

[Merken und Vergleichen](#)

Testberichte.de



GUT
1,6

586,95 Euro bei [amazon.de](#)

Google-Anzeige:

[Canon HF100](#)

Canon Camcorder im Vergleich! Testberichte lesen & sparen.
[www.ciao.de/Canon](#)

Produktdaten: Typ: Speicherkarten - Camcorder;
 Sonderfunktionen: Digitalkamerafunktion; Optischer Zoom: 12x;
 Gewicht: 424g; Monitorgröße: 2,7"; Digitaler Zoom: 200x;
 Aufzeichnungsformat: AVCHD; Speichermedien: SDHC ... [mehr Infos](#)

↶ Bild 1/2 ↷

Testberichte (14)

Preise (ca. 587 €) Meinungen (2) Datenblatt

videofilmen

Ausgabe 2/2009
 Platz 3/14 im Test

„sehr gut“ (146 von 200 Punkten)

„... Bei der Bildschärfe liegen die HF-Modelle fast an der Grenze des technisch Machbaren bei der HD-Auflösung. Das schafft Canon auch ohne Tricks wie die elektronische Kantenerhebung, bei der auffällige Schärfungssäume auftreten. ...“

[Im Testbericht lesen](#)

[Download \(3,00 €\)](#)

[586,95 bei Amazon.de](#) [Camcorder suchen bei Otto.de](#)

video

Ausgabe 8/2008
 Platz 1/3 im Test

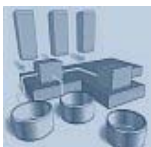
„sehr gut“ (75 von 100 Punkten)

Preis/Leistung: "hervorragend", „Testsieger“

„Plus: schärfstes Bild im Test; manuelle Tonaussteuerung; Mikro- und Kopfhörerbuchse; optischer Bildstabilisator. Minus: Zubehörschuh nur für Canon-Equipment.“

[Download \(1,20 €\)](#)

[586,95 bei Amazon.de](#) [Camcorder suchen bei Otto.de](#)



UNIVERSITÄT LEIPZIG

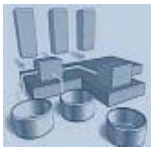
Abteilung Datenbanken
 am Institut für Informatik



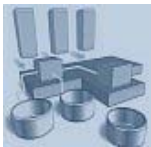
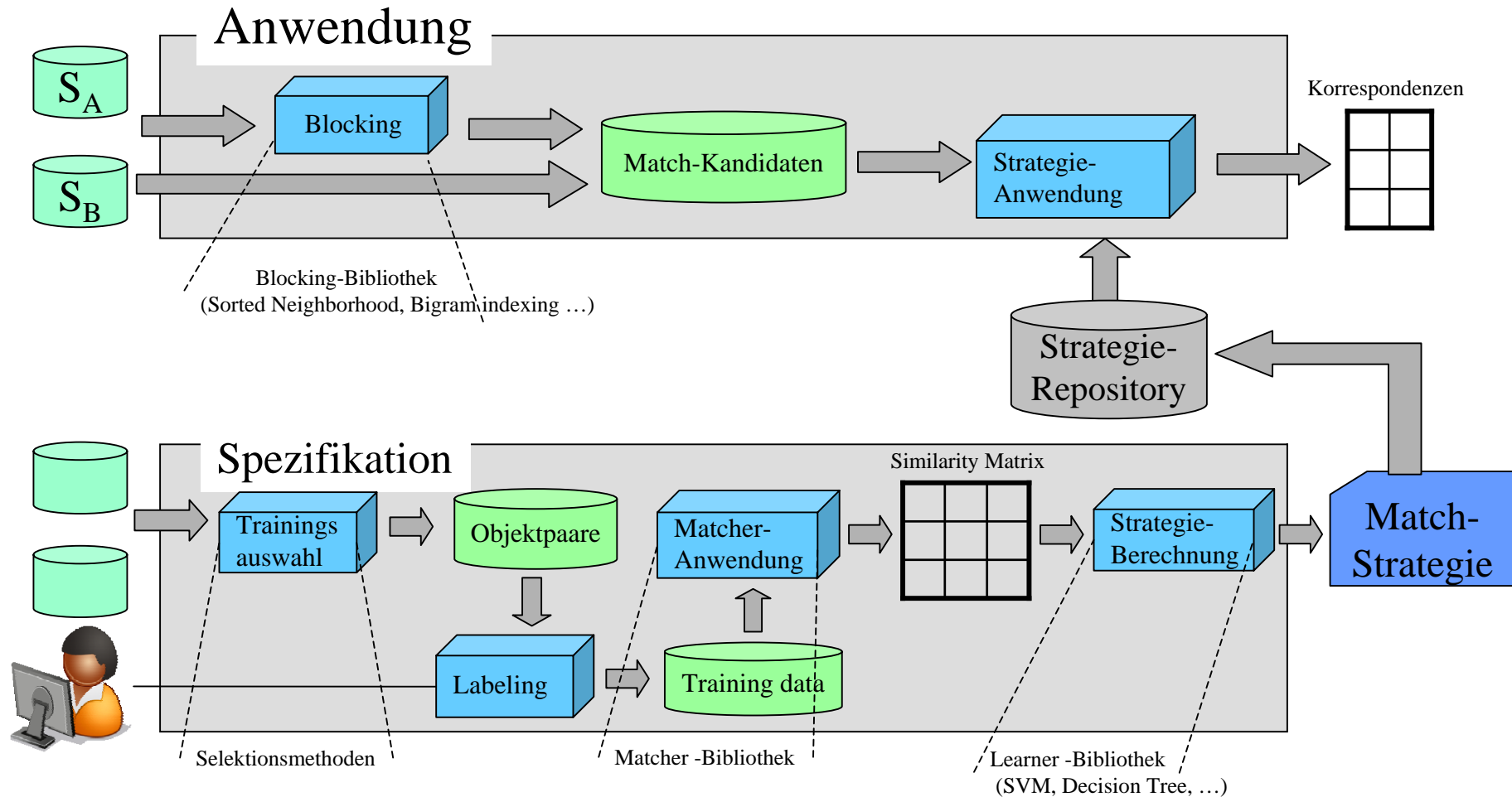
Prototypen zum Objekt-Matching

- Aufgabe: Erkennung äquivalenter Objekte (Dubletten) in Eingabedaten
 - Ergebnis: Mapping mit Korrespondenzen
- MOMA = Mapping based Object Matching
- STEM = Self-Tuning Entity Matching
- Kombinierte Verwendung mehrerer Match-Algorithmen
 - String-Vergleich zwischen Attributwerten
 - Auswertung von Nachbarschaftsbeziehungen
- STEM: Nutzung von *Machine Learning-Verfahren* zur automatischen Parameteroptimierung (z.B. Schwellwerte für Ähnlichkeiten)

LDS _A	LDS _{A'}	Sim
a ₁	a' ₁	1
a ₂	a' ₁	0.9
a ₃	a' ₃	0.8



STEM Architektur



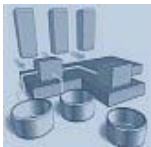
UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik

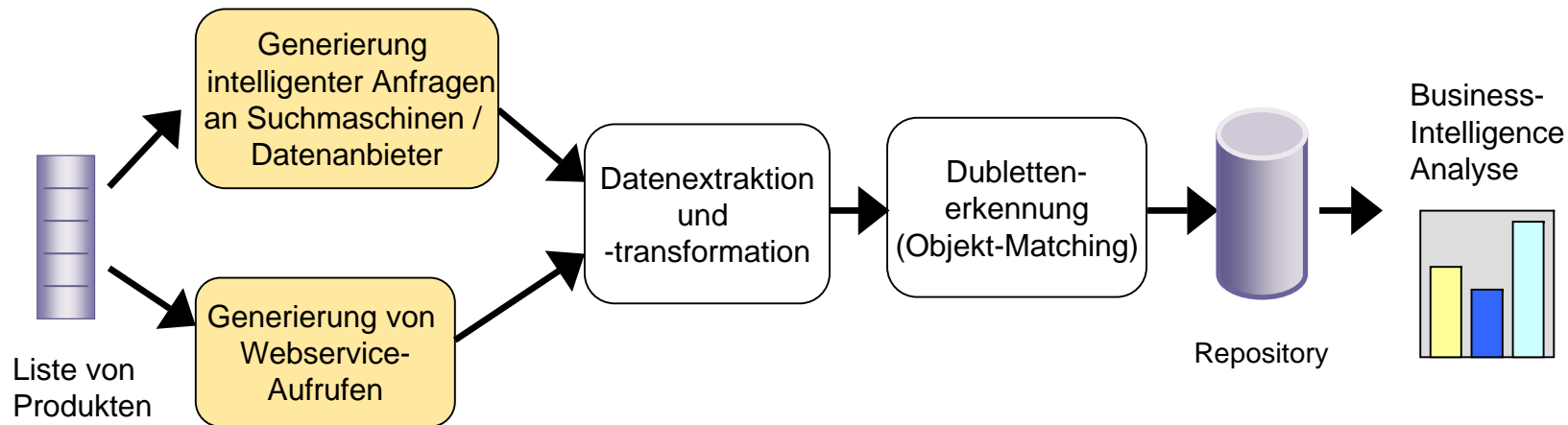


Workflow-artige Datenintegration mit iFuice

- iFuice-Framework: Erweiterung des Mashup-Ansatzes zur schnellen Realisierung von Datenintegrationsaufgaben
- Weitgehende Nutzung existierender Services, bereits integrierter Datensammlungen bzw. Portale
- Programmierung von Workflows durch Skript-Programme mit mächtigen mengenorientierten Operatoren (Anfragen, Matching, Mapping-Traversierung, Mengenverarbeitung etc.)
- Iterative Verbesserung von Suchanfragen zur Ergebnisverbesserung (Vollständigkeit, Genauigkeit)

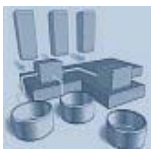


Beispiel-Workflows



Anwendungsbeispiele:

- Suche nach Zusatzinformationen (z.B. Nutzerbewertungen) in verschiedenen Quellen zu bestimmten Produkten
- Dynamische Preis/Verfügbarkeitsanalyse z.B. im Rahmen einer Metasuchmaschine (Hotels, Autos, etc.) auf Basis einer zuvor erstellten Präferenzliste
- Kopplung mehrerer Metasuchmaschinen: Flug, Hotel, Leihwagen
- Interaktive Recherche nach lokalen Anbietern bzw. Lieferanten bestimmter Dienstleistungen / Produkte

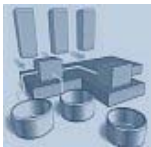
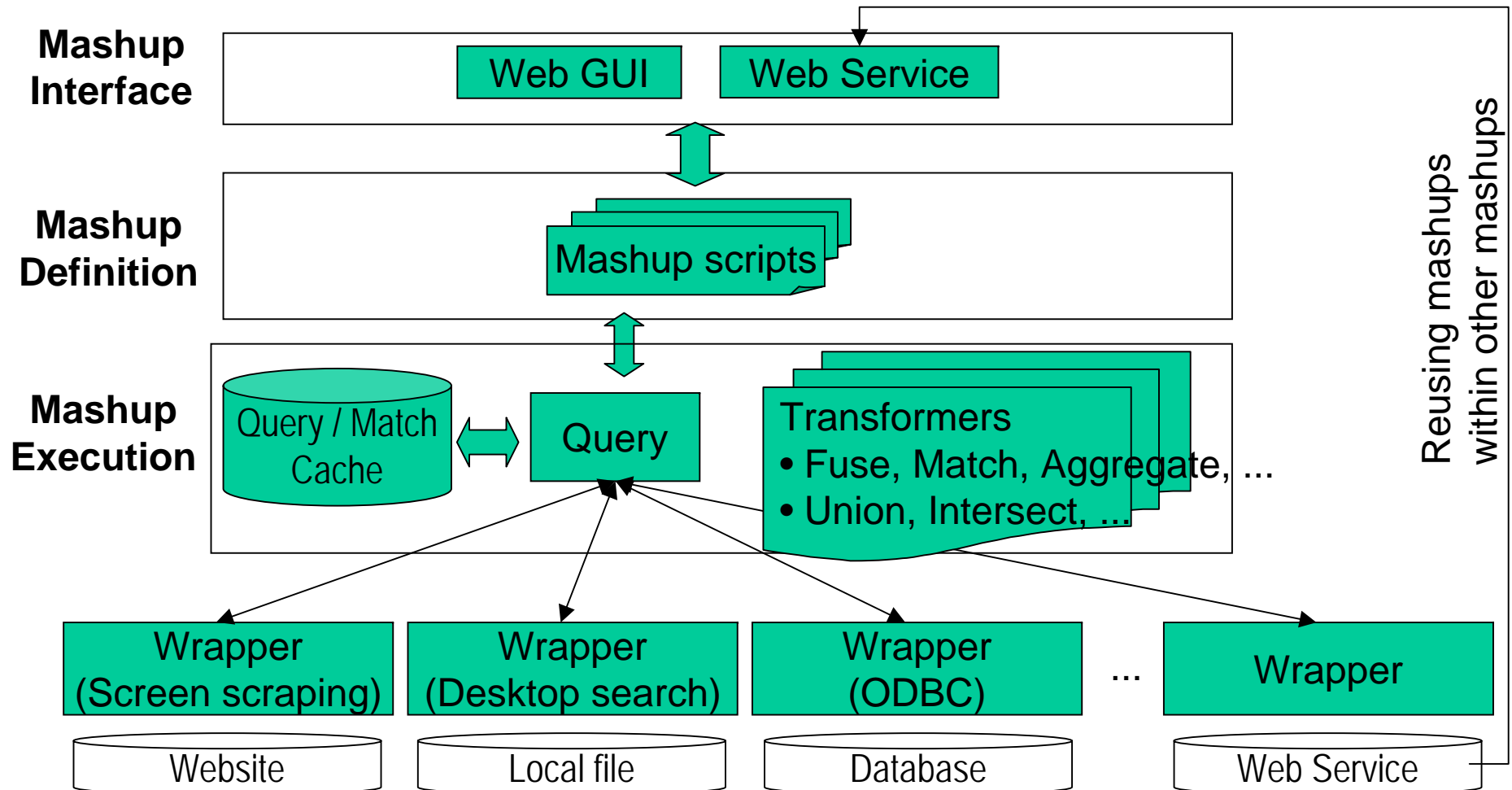


UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik

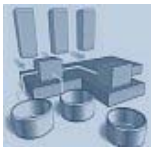


iFuice-Architektur



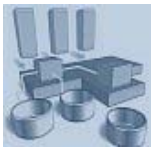
Vortragsinhalte

- Einleitung
 - Derzeitige Ansätze zur Datenintegration
 - Warehousing, EII
 - Mashups
 - Eigene Forschungsarbeiten / Prototypen
 - Schema/Ontologie Matching: COMA++
 - Objekt-Matching: MOMA, STEM
 - Workflowbasierte Integration: iFuice
- Feedback / Diskussion



Innovationslabor

- an Univ. Leipzig angesiedelt
- Erprobung / Weiterentwicklung der Techniken für Praxis/Unternehmenseinsatz
 - Workflow-basierte Datenintegration
 - Objekt-Matching / Dublettenbehandlung
 - Schema/Ontologie-Matching
- Kopplung mit verfügbaren Werkzeugen, z.B. von Datenbank/BI-Anbietern bzw. Open-Source-Lösungen
- Vorbereitung eventueller Ausgründungen



Feedback / Diskussion

- Wie beurteilen Sie die vorgestellten Verfahren (Stärken, Schwachstellen, Erweiterungsbedarf) ?
- Sehen Sie Anwendungsmöglichkeiten für die vorgestellten Verfahren?
 - zur Lösung konkreter Integrationsaufgaben
 - zur Kombination mit eigenen Entwicklungen/Produktangeboten
- Wie schätzen Sie die Marktrelevanz der Lösungsansätze ein?
- Welche Kooperationsformen mit dem Innovationslabor können Sie sich vorstellen, z.B.
 - gemeinsam betreute studentische Arbeiten
 - gemeinsame Projekte, z.B. zur Pilotanwendung der Werkzeuge
 - Beratung zu Datenintegrationsthemen

