

Innovationslabor Semantische Integration von Webdaten

Workflow-basierte Daten- integration und Objekt-Matching

Dr. Andreas Thor

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

GEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung



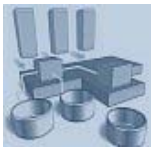
UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik



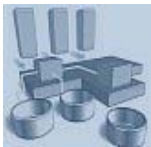
Workflow-basierte Datenintegration

- Ausgangspunkt / Problemstellung
 - Lösung eines konkreten Integrationsproblems erfordert eine koordinierte Ausführung mehrerer Teilschritte, u.a.
 - Anfragen an Suchmaschinen / Datenquellen
 - Abgleich der Daten (Objekt-Matching) und Metadaten (Schema-Matching)
 - bisherige Werkzeuge ungeeignet
 - Programmiersprachen (Java, C#, ...) sehr “low level” → viel Code
 - ETL-Werkzeuge zu statisch / unflexibel
- Ansatz: Erweiterung des Mashup-Ansatzes zur schnellen Realisierung von Datenintegrationsaufgaben
- Forschungsarbeiten seit 2005
 - Prototyp iFuice und Beispiel-Mashup OCS



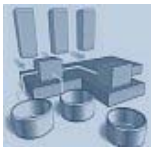
Prototyp iFuice: Features

- Workflow-artiger (programmatischer) Integrationsansatz
→ zur Lösung komplexer Integrationsaufgaben
- “Programmiersprache zur Datenintegration”
→ Operatoren für typische Aufgaben (query, match, ...)
- Nutzung bestehender Verknüpfungen (Mappings)
→ Wiederverwendung z.B. von Web-Links
- Verwendet ein flexibles Domänenmodell
→ typisierte Objekte mit Attribut-Wert-Paaren
- P2P-artige Kopplung von Datenquellen, kein zentrales Schema
→ Einfaches „Ankoppeln“ neuer Datenquellen „wo es am besten passt“ bzw. „am einfachsten ist“



Beispiel-Mashup: OCS

- Online Citation Service
 - Finden von Zitierungszahlen wissenschaftlicher Publikationen
 - Hochqualitative Zitierungsanalyse für Autoren, Konferenzen, ...
- Problemstellung / typisches Workflow-Muster
 - Finden intelligenter Suchanfragen an unterschiedliche Datenquellen, um eine Menge von Objekten (hier: Publikationen) zu finden
 - Extraktion der benötigten Informationen (hier: Zitierungszahl) aus den Ergebnissen
 - Zuordnung der gefundenen Objekte zu den Eingabeobjekten (Objekt-Matching)
- Analoge Anwendungsszenarien
 - Finden der Preise und/oder Bewertungen zu einer Menge von Produkten in anderen Datenquellen
 - Urlaubsbuchung mit abhängigen Komponenten
Flüge → “dazu passende” Hotels → “dazu passender” Mietwagen



Basket

Empty

Venues

- [1999 ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery 1999](#)
- [ACM SIGMOD Anthology 1999](#)
- [ACM SIGMOD Digital Review 1999](#)
- [ACM SIGMOD Digital Symposium Collection 1999](#)
- [SIGMOD Conference 1999](#)
- [SIGMOD Record 1999](#)

Publications of SIGMOD Conference 1999

 Step-By-Step

- Integration of Spatial Join Algorithms for Processing Multiple Inputs (SIGMOD Conference 1999)
Nikos Mamoulis, Dimitris Papadias
- Selectivity Estimation in Spatial Databases (SIGMOD Conference 1999)
Swarup Acharya, Viswanath Poosala, Sridhar Ramaswamy
- Efficient Concurrency Control in Multidimensional Access Methods (SIGMOD Conference 1999)
Koushik Chakrabarti, Sharad Mehrotra
- Snake Sandwiches: Optimal Clustering Strategies for a Data Warehouse (SIGMOD Conference 1999)
H. V. Jagannathan, V. S. Lakshmanan, Divesh Srivastava
- OPTICAL... (SIGMOD Conference 1999)
Sander...
- Logical Logging... (SIGMOD Conference 1999)
David B. Lomet, Mark R. Tuttle
- Efficient Concurrency Control for Broadcast Environments (SIGMOD Conference 1999)
Jayavel Shanmugasundaram, Arvind Nithrakashyap, Rajendran M. Sivasankaran, Krithi Ramasamy
- Update Propagation Protocols For Replicated Databases (SIGMOD Conference 1999)
Yuri Breitbart, Raghavan Komondoor, Rajeev Rastogi, S. Seshadri, Abraham Silberschatz
- Belief Reasoning in MLS Deductive Databases (SIGMOD Conference 1999)
Hasan M. Jamil
- A Multimedia Presentation Algebra (SIGMOD Conference 1999)
Shelagh M. Lee, Maria Lina Serrino, H. S. Subrahmanian

Auswahl der zu analysierenden Publikationen (Eingabeobjekte)

Online Citation Service 2.0



	All	DBLP	GS	ACM	CS	Libra
# Publications	85	85	72	84	25	90
Σ Citations			3522	1516	940	2336
\emptyset Citations			41.4	17.8	11.1	27.5
H-Index			23	23	17	30

Title	<input type="range"/>	80
Auth	<input type="range"/>	50
Year	<input type="range"/>	-2
GS Citations (DESC)		<input type="button" value="v"/>

Unmark All | Mark All | Mark Selected | Expand All | Collapse All

	G		C	
<input checked="" type="checkbox"/> + Storing Semistructured Data with STORED (SIGMOD Conference 1999) <i>Alin Deutsch, Mary F. Fernández, Dan Suciu</i>	527	101	143	179
<input checked="" type="checkbox"/> + Bottom-Up Computation of Sparse and Iceberg CUBEs (SIGMOD Conference 1999) <i>Kevin S. Beyer, Raghu Ramakrishnan</i>	344	65	77	91
<input checked="" type="checkbox"/> + An Adaptive Query Execution System for Data Integration (SIGMOD Conference 1999) <i>Zachary G. Ives, Daniela Florescu, Manoj S. Mani, Alon Y. Levy, Daniel S. Weld</i>	301	73	80	135
<input checked="" type="checkbox"/> + Record-Boundary Discovery in Web Documents (SIGMOD Conference 1999) <i>David W. Embley, Y. S. Jiang, Yi Ma</i>	222		40	64
<input checked="" type="checkbox"/> + Online Association Rule Mining (SIGMOD Conference 1999) <i>Christian Hidber</i>	203	38	45	53
<input checked="" type="checkbox"/> + Self-tuning Histograms: Building Histograms for Data Mining (SIGMOD Conference 1999) <i>Ashraf Aboulnaga, Surajit Chaudhuri</i>	180	50		62
<input checked="" type="checkbox"/> + BOAT-Optimistic Decision Tree Construction (SIGMOD Conference 1999) <i>Johannes Gehrke, Venkatesh Ganti, Raghu Ramakrishnan, Wei-Yin Loh</i>	177	44		55
<input checked="" type="checkbox"/> + Selectivity Estimation in Spatial Databases (SIGMOD Conference 1999) <i>Swarup Acharya, Viswanath Poosala, Sridhar Ramaswamy</i>	160	47		66
<input checked="" type="checkbox"/> + On Random Sampling over Joins (SIGMOD Conference 1999) <i>Surajit Chaudhuri, Rajeev Motwani, Vivek R. Narasayya</i>	143	54		52
<input checked="" type="checkbox"/> + Query Optimization in the Presence of Limited Access Patterns (SIGMOD Conference 1999) <i>Daniela Florescu, Alon Y. Levy, Ioana Manolescu, Dan Suciu</i>	128	27	60	62
<input checked="" type="checkbox"/> + Query Rewriting for Semistructured Data (SIGMOD Conference 1999) <i>...</i>	118	22		45

"Runde 1": wenige, einfache Anfragen



Online Citation Service 2.0



	All	DBLP	GS	ACM	CS	Libra
# Publications	85	85	237	84	25	90
Σ Citations			7122	1516	940	2336
\emptyset Citations			83.8	17.8	11.1	27.5
H-Index			39	23	17	30

Title 80

Auth 50

Year -2

GS Citations (DESC)

Unmark All Mark All Mark Selected Expand All Collapse All

<input checked="" type="checkbox"/>	+ OPTICS: Ordering Points To Identify the Clustering Structure (SIGMOD Conference 1999) <i>Mihael Ankerst, Markus M. Breunig, Hans-Peter Kriegel, Jörg Sander</i>	654	85		104
<input checked="" type="checkbox"/>	+ Storing Semistructured Data with ST... (SIGMOD Conference 1999) <i>Alin Deutsch, Mary F. Fernández, Dan...</i>	552	101	143	179
<input checked="" type="checkbox"/>	+ Fast Algorithms for Project... <i>Charu C. Aggarwal, Cec...</i>	409	87	61	91
<input checked="" type="checkbox"/>	+ An Adaptive Query Execution... <i>Zachary G. Ives, Daniela Florescu, Marc...</i>	351	73	80	135
<input checked="" type="checkbox"/>	+ Bottom-Up Computation of Overlap... (SIGMOD Conference 1999) <i>Kevin S. Beyer</i>	347	65	77	91
<input checked="" type="checkbox"/>	+ XML-Based Inf... <i>Chaitanya K. E... Papakonstanti...</i>	273	42	76	84
<input checked="" type="checkbox"/>	+ Approximate C... (SIGMOD Con... <i>Jeffrey Scott Vitter, Min Wang</i>	262	88		111
<input checked="" type="checkbox"/>	+ Record-Boundary Discover... <i>David W. Embley, Y. F...</i>	232		40	64
<input checked="" type="checkbox"/>	+ Online Association Rule... <i>Christian Hidber</i>	208	38	45	53
<input checked="" type="checkbox"/>	+ Join Synopses for Approximate Query Answering (SIGMOD Conference 1999) <i>Swarup Acharya, Phillip B. Gibbons, Viswanath Poosala, Sridhar Ramaswamy</i>	194	62	43	

"Runde 2": weitere, komplexere Anfragen

```
$Pubs := query (DBLP, "Konferenz = 'SIGMOD 1999'");  
$GS1 := query (GS, $Pubs, "allintitle:[[title]]");  
$GS2 := query (GS, $Pubs, "keywords:[[title]] [[author]]");  
$Res := attrMatch ($Pubs, $GS1+$GS2, "title", 0.8);
```

Definition der Programmlogik durch kurzes Skript

Online Citation Service 2.0



	All	DBLP	GS	ACM	CS	Libra
# Publications	85	85	237	84	25	90
Σ Citations			7122	1516	940	2336
\emptyset Citations			83.8	17.8	11.1	27.5
H-Index			39	23	17	30

Title

Auth

Year

GS Citations (DESC)

<input checked="" type="checkbox"/>	-	OPTICS: Ordering Points To Identify the Clustering Structure (SIGMOD Conference 1999) <i>Mihael Ankerst, Markus M. Breunig, Hans-Peter Kriegel, Jörg Sander</i>	654	85	104	
		OPTICS: Ordering Points To Identify the Clustering Structure				
		OPTICS: Ordering Points To Identify the Clustering Structure	639			
		OPTICS: Ordering points to identify the cluster structure	9			
		OPTICS: Ordering points to identify the clustering structure [A]. 1999 ACM-SIGMOD Int	1			
		OPTICS: Ordering Objects to Identify the Clustering Structure	1			
		Kriegel, H.-P. and Sander, J., 1999, OPTICS: Ordering points to identify the clustering structure	1			
		OPTICS: Ordering points to identify the clustering structure, In Proceedings of the ACM SIGMOD	1			
		OPTICS: Ordering Objects to Identify the Clustering Structure, Proc. ACM SIGMOD	1			
		JS OPTICS: Ordering Points To Identify the Clustering Structure Proc	1			
		OPTICS: ordering points to identify the clustering structure		85		
		OPTICS: Ordering Points To Identify the Clustering Structure			104	
<input checked="" type="checkbox"/>	+	Storing Semistructured Data with... <i>Alin Deutsch, Mary F. Fern</i>	552	101	143	179
<input checked="" type="checkbox"/>	+	Fast Algorithms for Projected... <i>Charu C. Aggarwal, Cecilia Magda</i>	309	87	61	91
<input checked="" type="checkbox"/>	+	An Adaptive Query Execution System for Data Integration (SIGMOD Conference 1999) <i>Zachary G. Ives, Daniela Florescu, Marc Friedman, Alon Y. Levy, Daniel S. Weld</i>	351	73	80	

Objekt-Matching: Zuordnung der gefundenen Objekte zu den Eingabeobjekten

Online Citation Service 2.0



	All	DBLP	GS	ACM	CS	Libra
# Publications	85	85	109	71	37	83
Σ Citations			6058	1365	1146	2334
\emptyset Citations			71.3	16.1	13.5	27.5
H-Index			35	22	19	30

Title 100

Auth 50

Year -2

GS Citation (DESC)

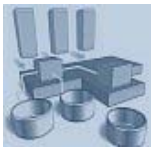
Unmark All Mark All Mark Selected Expand All Collapse All

		G	D	C	SR
<input checked="" type="checkbox"/> +	OPTICS: Ordering Points To Identify the Clustering Structure (SIGMOD Conference 1999) <i>Mihael Ankerst, Markus M. Breuninger, Anil K. Choudhury, Hans-Jörg Schuster</i>	648	85	76	104
<input checked="" type="checkbox"/> +	Storing Semistructured Data (SIGMOD Conference 1999) <i>Alin Deutsch, Mary Lou Sizer</i>	527	101	143	179
<input checked="" type="checkbox"/> +	Fast Algorithms for Projected Aggregate Queries (SIGMOD Conference 1999) <i>Charu C. Aggarwal, Cecilia Magdalena Procopiac, Victor V. Vassilvski, Yinpeng S. Yu, Jong Soo Park</i>	404	87	61	91
<input checked="" type="checkbox"/> +	An Adaptive Query Execution System for Data Integration (SIGMOD Conference 1999) <i>Zachary G. Ives, Daniela Florescu, Marc Friedman, Alon Y. Levy, Daniel S. Weld</i>	301	73	80	135
<input checked="" type="checkbox"/> +	Approximate Computation of Multidimensional Aggregates of Sparse Data Using Wavelets (SIGMOD Conference 1999) <i>Jeffrey Scott Vitter, Min Wang</i>	259	88		111
<input checked="" type="checkbox"/> +	XML-Based Information Mediation with MIX (SIGMOD Conference 1999) <i>Chaitanya K. Baru, Amarnath Gupta, Bertram Ludäscher, Richard Marciano, Yannis Papakonstantinou, Pavel Velikhov, Vincent Chu</i>	258	42	76	84
<input checked="" type="checkbox"/> +	Record-Boundary Discovery in Web Documents (SIGMOD Conference 1999) <i>David W. Embley, Y. S. Jiang, Yiu-Kai Ng</i>	225		40	64
<input checked="" type="checkbox"/> +	Online Association Rule Mining (SIGMOD Conference 1999) <i>Christian Hidber</i>	201	38	45	53
<input checked="" type="checkbox"/> +	Ripple Joins for Online Aggregation (SIGMOD Conference 1999) <i>Peter J. Haas, Joseph M. Hellerstein</i>	188	55	57	80
<input checked="" type="checkbox"/> +	Join Synopses for Approximate Query Answering (SIGMOD Conference 1999) <i>Swarup Acharya, Phillip B. Gibbons, Viswanath Poonala, Sridhar Ramaswamy</i>	186	62	43	

Einstellung des Objekt-Matchings durch Schwellwerte

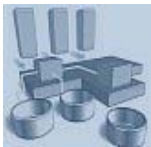
Objekt-Matching

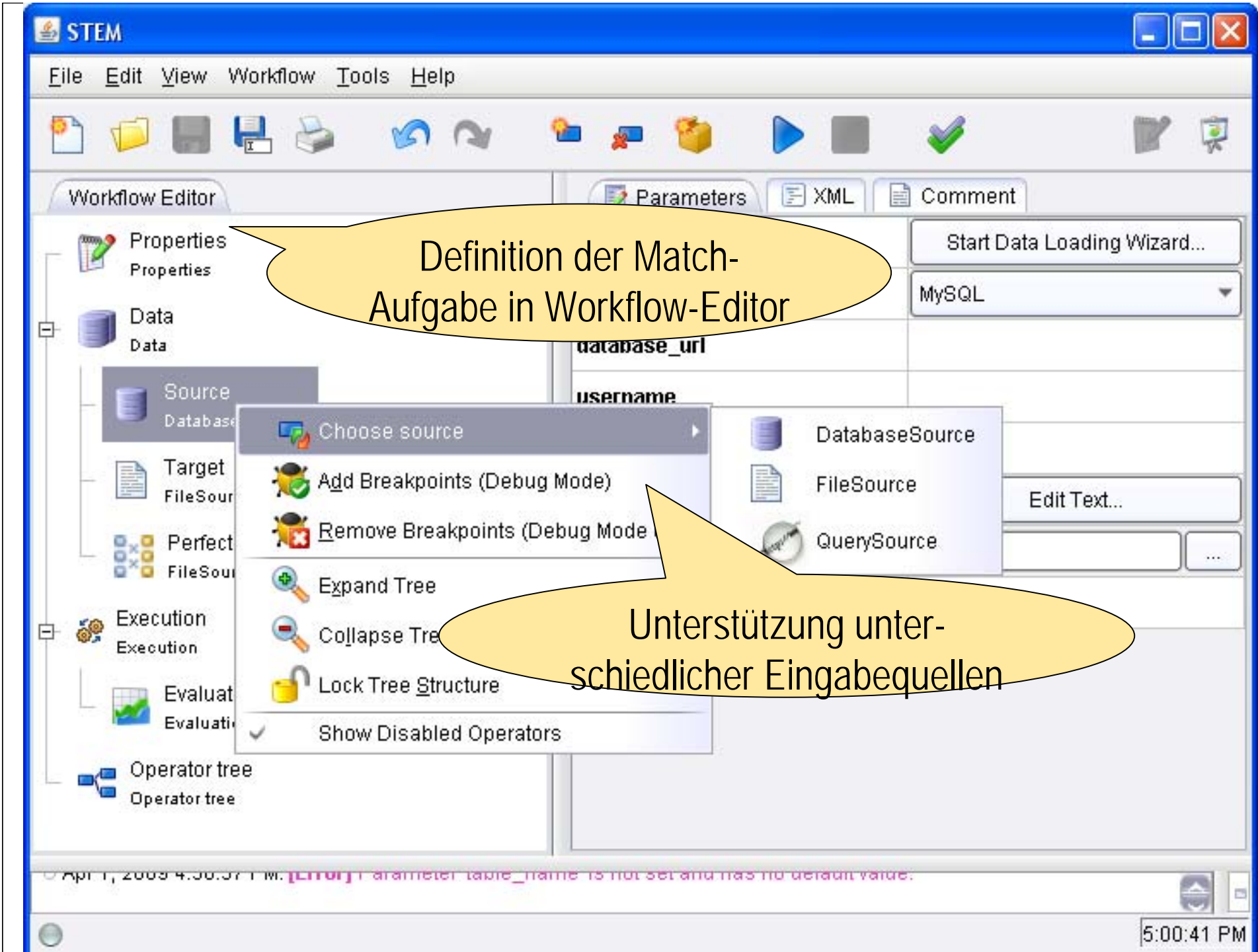
- Identifikation semantisch äquivalenter Objekte
 - innerhalb einer Datenquelle oder zwischen verschiedenen Quellen
 - um Objekte zu integrieren/mischen, zu vergleichen, Dubletten zu eliminieren, etc
- Anwendungsbereiche
 - Kunden-/Adressdaten
 - Produkte
 - Geografische Orte
 - ...
- Forschungsarbeiten seit 2006
 - Prototypen: MOMA und die Weiterentwicklung STEM



Prototypen MOMA und STEM: Features

- Erweiterbare Bibliothek von Match-Verfahren (Matcher)
 - Anpassbar an verschiedene Szenarien / Domänen
- Mapping-Kombination
 - Effiziente Berechnung
 - Qualitätsabsicherung / -steigerung
- Konstruktion von Match-Workflows
 - an Problem angepasste Lösungsstrategie
- Speichern von Mappings in Repository
 - Wiederverwendung von Match-Ergebnissen
- Automatische Einstellung relevanter Parameter
 - Verringerung des Konfigurationsaufwandes





Definition der Match-Aufgabe in Workflow-Editor

Unterstützung unterschiedlicher Eingabequellen

The screenshot shows a software interface for a workflow editor. On the left is a tree view with categories like Properties, Data, Source, Target, Execution, Evaluation, and Operator tree. The Operator tree contains various join operators such as Merge, EdJoin, SortedNeighborhood, FuzzyJoin, and CrossJoin. A central menu is open, showing options like 'New Operator', 'Add Breakpoints', 'Remove Breakpoints', 'Expand Tree', 'Collapse Tree', 'Lock Tree Structure', and 'Show Disabled Operators'. A sub-menu for 'Mapping' is also open, listing 'Blocking', 'Combination', 'Learning', 'SimilarityJoin', and 'TrainSelect'. A further sub-menu for 'SimilarityJoin' lists 'EdJoin', 'FuzzyJoin', 'PPJoinPlus', and 'StringMap'. A top toolbar contains icons for file operations and workflow execution. The main workspace shows a table with columns for 'attribute name' and 'source'.

Workflow Editor

Parameters XML Comment

attribute name source

attrib

tau

Vielzahl von Operatoren,
die flexibel kombiniert
werden können

New Operator

Entities

Mapping

Blocking

Combination

Learning

SimilarityJoin

TrainSelect

EdJoin

FuzzyJoin

PPJoinPlus

StringMap

Add Breakpoints (Debug Mode)

Remove Breakpoints (Debug Mode off)

Expand Tree

Collapse Tree

Lock Tree Structure

Show Disabled Operators

Operatorbaum =
strukturierte Darstellung
des Match-Workflows

Properties
Properties

Data
Data

Source
DatabaseSource

Target
FileSource

Perfect Mapping
FileSourceMapping

Execution
Execution

Evaluation
Evaluation

Operator tree
Operator tree

Merge
Merge

EdJoin
EdJoin

SortedNeighborhood
SortedNeighborhood

FuzzyJoin
FuzzyJoin

CrossJoin
CrossJoin



Mapping

Meta Data View Data View Plot View

474 correspondences

Select source attributes:

- Select target attributes:
- id
 - name
 - price
 - availability
 - description
 - upc
 - Approximate Weight:
 - Features:
 - Specifications:
 - regularprice
 - yourprice
 - Approximate Dimensions:
 - id
 - name
 - listprice
 - ourprice
 - shipping
 - instock
 - Manufacturer
 - Mfg Part#
 - UPC
 - Buy.com Sku
 - Item#
 - Buy.com Sales Rank
 - description
 - Features
 - Tech Specs
 - price
 - Format

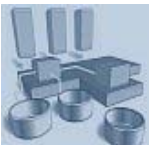
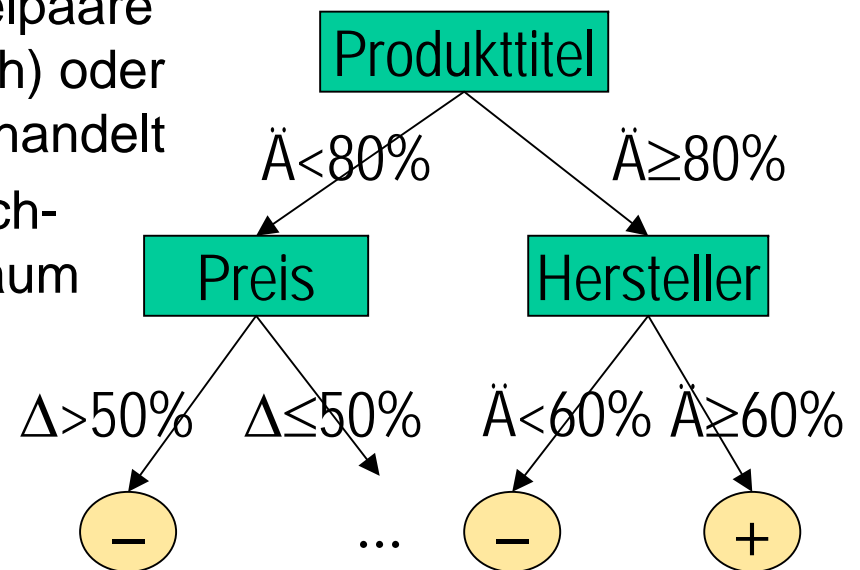
id	name			name
10101	Omnimount Wall Speaker Mount - 20WLBK			Omnimount Universal Wall Speaker Mounting Kit - 20.0 WA
10102	Omnimount Wall Speaker Mount - 20WLWH	0.5		Omnimount Universal Wall Speaker Mounting Kit - 20.0 WA
11801	Escort Passport Radar And Laser Detector - Black Finish - 8500	0.5		Escort Passport 9500ix Radar/Laser Detector
13155	Peerless Wall TV Mounts In Black - PM1327BK	0.462	1	Peerless Adjustable TV Wall Mount - PM1327
13202	Panasonic Laser Toner Cartridge - KXFA83	0.517	1	Panasonic Black Toner Cartridge - KX-FA83
13213	Sony Compact Disc Player/Recorder - RCDW500C		1	Sony Compact Disc Player/Recorder - RCDW500C
13700	Sanus Euro Foundations Satellite Speaker Stand - EFSATB	0.507	1	Sanus Euro Foundation Series III Speaker Stand - EFSATE
13701	Sanus Euro Foundations Satellite Speaker Stand - EFSATS	0.529	1	Sanus Euro Foundation Series III Speaker Stand - EFSATS
13945	Onkyo 6 Disc CD Player Changer - DXC390B	0.509	1	ONKYO 6-DISC CD CHANGER W/VLSC, BLACK NIC - DX-
13954	Panasonic Corded Phone - KXTS3282B	0.761	1	Panasonic KX-TS3282B Corded Phone
13996	Escort Cordless Solo Radar Detector - S2E	0.575	1	Cordless Solo S2 Radar/Laser Detector - 010ES20A
	Garmin Deluxe Carrying Case - Black Finish - 0101023101	0.500	1	Garmin Canvas Deluxe Carry Case - 010-10231-01
14061		0.614	1	Kenwood KDC-C669 Car CD Changer - KDCC669
14563				Sharp 1100 Watt Over the Counter Microwave
				Maytag 2.0 Cu. Ft. Over-the-Range Microwave Oven
				Maytag 2.0 Cu. Ft. Over-the-Range Microwave Oven
1630				Sony SCD-CE595 CD Player - SCDCE595
16668	Omnimount TV Top Shelf Mount - CCH1P		1	Omnimount TV Top Shelf Mount - CCH1P
16669	Omnimount TV Top Shelf Mount - CCH1B	0.697	1	Omnimount TV Top Shelf Mount - CCH1P
16741	Delonghi Twenty Four Seven Coffee Maker In Black - DC50B	0.510	1	DeLonghi DC50B Twenty Four Seven Drip Coffee Maker, I
16758	Sanus Silver LCD Television Turntable - TVLCDS	0.627	1	Sanus Television Turntable - TV/LCDS
16765	Frigidaire 27" Electric Stack Washer Dryer Combo - FEX831WH	0.539	1	"Frigidaire 27"" Electric Stack Washer Dryer Combo - FEXI
16877	Sony Digital Photo Printer Paper 40 Pack - SVMF40P	0.501	1	Sony Print Paper - SVMF40P
17067	Canon Rechargeable Battery - 9763A001	0.897	1	Canon NB-4L Rechargeable Camera Battery - 9763A001

Darstellung des Match-Ergebnisses inkl. Qualität

Auswahl darzustellender Attribute zur manuellen Inspektion

Automatische Match-Konfiguration (Self-Tuning)

- Richtige Match-Strategie: Fragestellungen
 - Welche (relevanten) Attribute sollen verglichen?
 - Wie sollen die Ähnlichkeitswerte kombiniert / verrechnet werden?
- Automatische Konfiguration mittels Lernverfahren
 - Nutzer gibt für (wenige!) Beispielpaare an, ob es sich um gleiche (match) oder ungleiche (non-match) Objekte handelt
 - System errechnet optimale Match-Strategie, z.B. Entscheidungsbaum



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

