

Beyond Information Integration: Content Merging

Forschungskolloquium CIS/ISST

Felix Naumann

naumann@informatik.hu-berlin.de



Einige Untertitel

- Content Merging
- Objekt/Data Fusion
- Data Amalgamation
- Data Consolidation
- Intelligent Information Integration
- Data Cleansing
- Datenintegration
- Datenverschmelzung



Forschungsgruppe Informationsintegration

- Institut für Informatik
der Humboldt-Universität zu Berlin
- Juniorprofessor: Felix Naumann
- Mitarbeiter
 - Jens Bleiholder & Melanie Weis
- Forschungsthemen
 - Objektidentifikation
 - Informationsintegration
 - Optimierung
 - Visualisierung



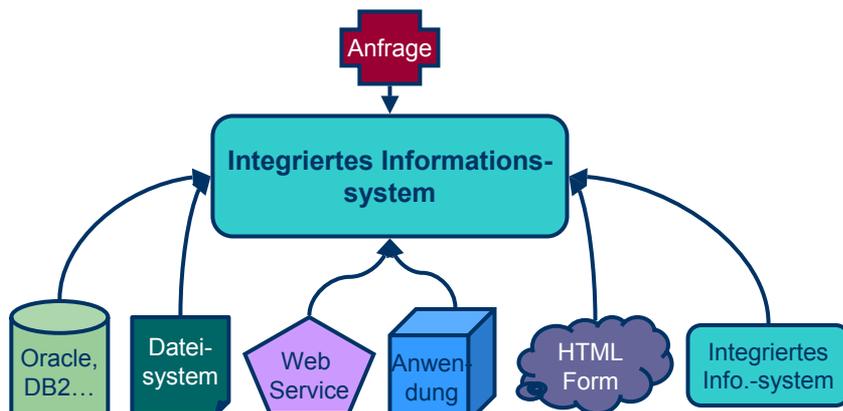
3



26.1.2004

Felix Naumann, HU Berlin

Integrierte Informationssysteme



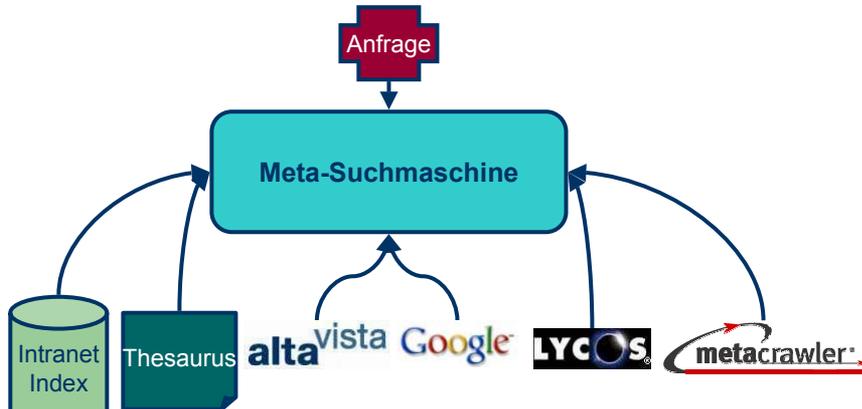
4



26.1.2004

Felix Naumann, HU Berlin

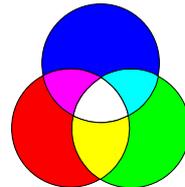
Integrierte Suchmaschinen



5

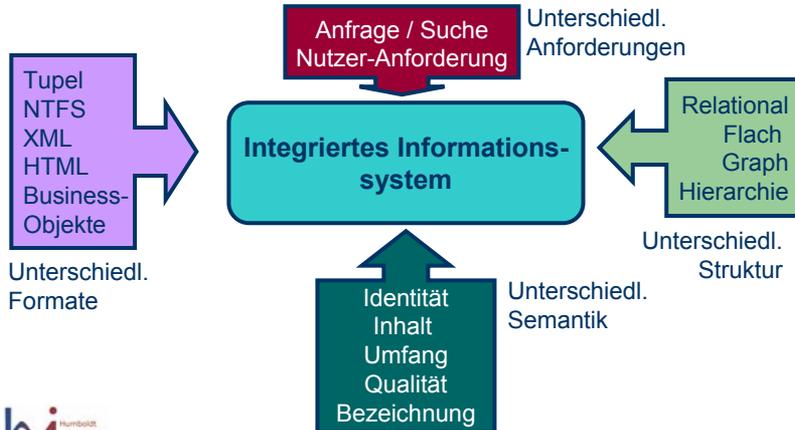
Redundanz – pro & contra

- Ist gut, denn sie bietet
 - mehr Informationen,
 - detaillierter Informationen,
 - und verifizierbare Informationen.
 - Deshalb sollten wir integrieren!
- Ist problematisch, denn
 - Redundanz herrscht nur konzeptionell.
 - Technische und strukturelle Schwierigkeiten
 - Konflikte und Alternativen
 - Deshalb ist Informationsintegration interessant.

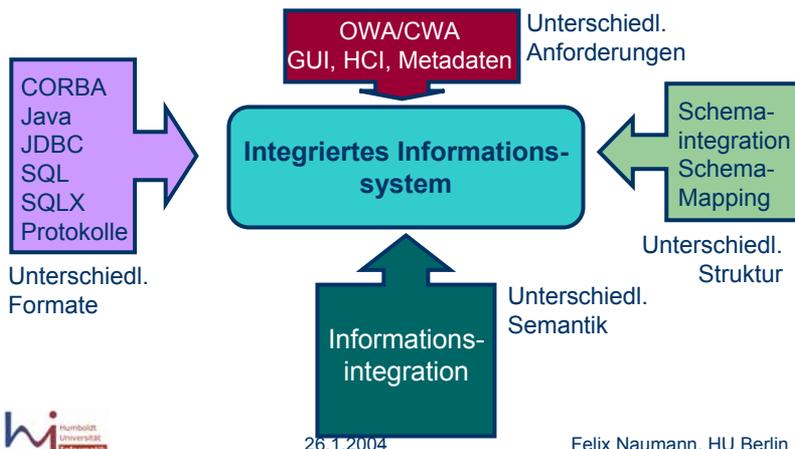


6

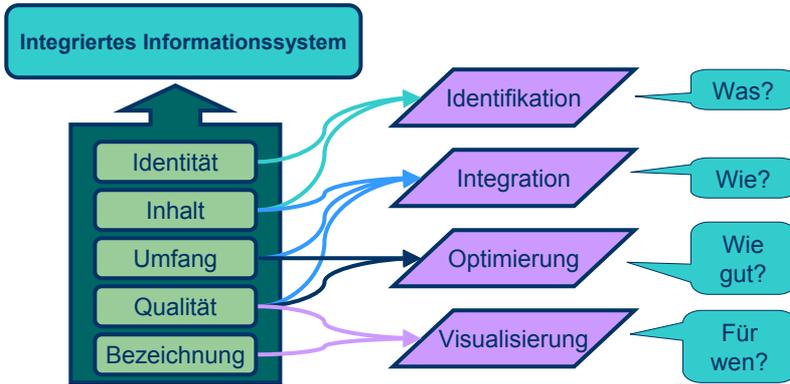
Integrierte Informationssysteme



Integrierte Informationssysteme



Lösung Unterschiedlicher Semantik



9

Beispiel der Informationsintegration



1sec.

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



5sec.

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

Identifikation

Integration

Optimierung

Visualisierung

10

Objektidentifikation

Identifikation Integration Optimierung Visualisierung

- Erkennung mehrere (ähnlicher) Darstellungen des gleichen Objekts
- Auch
 - Duplicate Detection
 - Data Cleansing
 - Record Linkage
- Domänenspezifische Algorithmen
 - Adressdaten
 - Mikrobiologische Daten
- Manchmal gibt es IDs
 - Bücher: ISBN
 - Personen: SSN / Personalausweisnummer
 - Webseiten: URL



11

Objektidentifikation

Identifikation Integration Optimierung Visualisierung

- Das klassische Problem
 - Finde Duplikate innerhalb einer Tabelle.
 - Sehr große Datenmenge
 - kein quadratischer Algorithmus
 - kein Hauptspeicher-Algorithmus
- Forschung
 - Merge/Purge Technik (Hernandez & Stolfo 1998)
 - Sorted-Neighborhood Methode
 - Record Linkage (Fellegi & Sunter, 1969)
 - u.a.
- Industrie
 - Trillium, Vality, ETI, et al.
 - Algorithmen sind gut gehütete Geheimnisse



12

Sorted Neighborhood

Identifikation Integration Optimierung Visualisierung

- Idee (Hernandez & Stolfo 1998)
 - Daten geschickt partitionieren.
 - Nur innerhalb dieser Partitionen Duplikate suchen.
- Algorithmus
 1. Create Key:
 - Schlüssel mittels relevanter Feldern erzeugen.
 2. Sort:
 - Daten nach dem Schlüssel sortieren.
 3. Merge:
 - Fenster (der Größe w) über sortierte Tupel schieben.
 - Nur Tupel innerhalb des Fensters miteinander vergleichen.

13



26.1.2004

Felix Naumann, HU Berlin

Objektidentifikation in XML Daten

Identifikation Integration Optimierung Visualisierung

- Data Warehouse Duplicates
 - Ausnutzung hierarchischer Daten
 - Nur Star-Schema
 - (Ananthakrishna, Chauduri & Ganti 2002)
- XML Duplikate
 - Kinder verschiedener Typen (Snowflake-Schema)
 - Was ist ein Duplikat?
 - Semistrukturierte Daten
 - Schemalose Daten

14



26.1.2004

Felix Naumann, HU Berlin

Objektidentifikation in XML Daten

Identifikation Integration Optimierung Visualisierung

```
- <author>
  <name>Bernd Amann</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
</author>
- <author>
  <name>Sophie Cluet</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Views in a Large Scale XML Repository.</title>
  <year>2001</year>
</publication>
</author>
```

- Vergleiche <author>
 - Mit Subelementen (<publication>)?
 - Wie tief?
- Vergleiche <publication>
 - Mit parallelen Elementen (<year>)?
 - Schema, oder Daten?
- Kurz: Was ist ein Duplikat?

15

20.1.2004

Felix Naumann, HU Berlin

Objektidentifikation in XML Daten

Identifikation Integration Optimierung Visualisierung

- Entwicklung eines Vergleichsmaßes für XML Objekte
 - Basierend auf IDFs
 - Symmetrisch
 - Einbeziehung beliebig tiefer Nachkommen
 - Effiziente Speicherung
 - Effizienter Vergleich
 - Prä-Selektion durch „korrekten“ Filter
- Future Work
 - Berücksichtigung von Schemainformationen
 - Annahme bisher: XML Daten gleichen Schemas
 - Jetzt: XML Daten heterogener Schemata

16

Objektidentifikation in XML Daten

Identifikation Integration Optimierung Visualisierung

```
- <author>
  <name>Bernd Amann</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <author>
  <name>Sophie Cluet</name>
- <publication>
  <title>XML Repository and Active Views Demonstration.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Active Views for Electronic Commerce.</title>
  <year>1999</year>
</publication>
- <publication>
  <title>Views in a Large Scale XML Repository.</title>
  <year>2001</year>
</publication>
</author>
- <inproceedings key="conf/vldb/AbiteboulAAACHHMMMSTV9"
  <author>Serge Abiteboul</author>
  <author>Vincent Aguilera</author>
  <author>Sébastien Ailleret</author>
  <author>Bernd Amann</author>
  <author>Sophie Cluet</author>
  <author>Brendan Hills</author>
  <author>Frédéric Hubert</author>
  <title>XML Repository and Active Views Demonstration.</t
  <pages>742-745</pages>
  <year>1999</year>
  <booktitle>VLDB</booktitle>
  <url>db/conf/vldb/vldb99.html#AbiteboulAAACHHMMMMS
  <crossref>conf/vldb/99</crossref>
  <ee>db/conf/vldb/AbiteboulAAACHHMMMSTV99.html</e
  <cdrom>VLDB99/P73.pdf</cdrom>
  <cite>conf/edbt/SantosAD94</cite>
  <cite>www/org/w3/dom</cite>
</inproceedings>
```

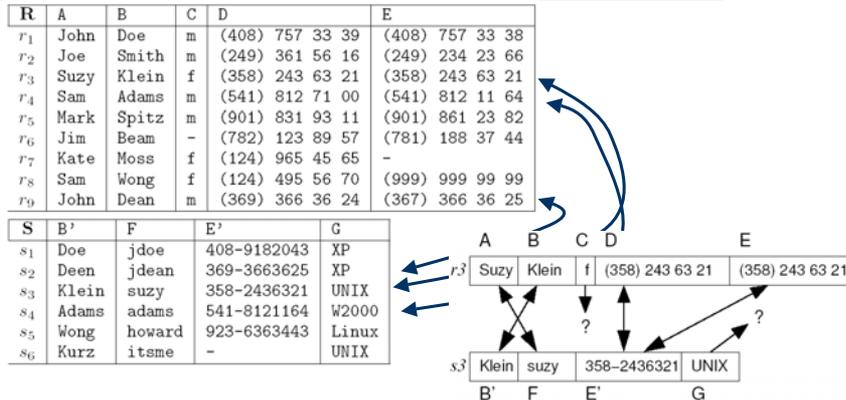
Duplikat-gesteuertes Schema Mapping

Identifikation Integration Optimierung Visualisierung

- Umgekehrte Idee
- Schema Matching
 - Klassisches Problem der Informationsintegration:
 - Finde Korrespondenzen zwischen zwei (heterogenen) Schemata
- Herkömmliche Lösungen
 - Auswertung von Attributnamen, Attributwerten und Struktur (Rahm & Bernstein 2001)
- Nun
 - Finde Duplikate (trotz mangelnder Schemata)
 - Korrespondenzen zwischen gleichen Attributwerten
 - In Kooperation mit Alexander Bilke

Duplikat-gesteuertes Schema Mapping

Identifikation Integration Optimierung Visualisierung



19



26.1.2004

Felix Naumann, HU Berlin

Objektidentifikation – Zusammenfassung

Identifikation Integration Optimierung Visualisierung

- Relationale Objekt Identifikation
 - Viele Algorithmen
 - No magic
- Identifikation in XML Daten
 - Wahl der Granularität
 - Nutzung struktureller Information
 - Hilfestellung für Schema Matching
- Agenda
 - Erweiterung existierender Algorithmen
 - Entwicklung neuer Ähnlichkeitsmaße und Algorithmen
 - Effiziente Ausführung
 - Datenstruktur
 - Algorithmen

Jetzt wissen wir, WAS integriert werden soll. Bloß WIE?

20



26.1.2004

Felix Naumann, HU Berlin

Objektintegration

Identifikation Integration Optimierung Visualisierung



0766607194	H. Melville		\$3.98	
------------	-------------	--	--------	--



0766607194	Herman Melville	Moby Dick	\$5.99	
------------	-----------------	-----------	--------	--



21

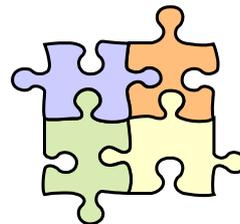
26.1.2004

Felix Naumann, HU Berlin

Relationale Objektintegration

Identifikation Integration Optimierung Visualisierung

- **Union** \cup
 - Duplikat-Eliminierung
- **Outer union** \oplus (Codd 1979)
 - Union bei heterogenen Schemata
- **Minimum union** \oplus (Ullmann 1989, Galindo-Legaria 1994)
 - Eliminierung subsummierter Tupel
- **Merge union**
 - Duplikatintegration
 - Konfliktlösung



Weitere Operatoren
-Generalized union
-... ??? JENS ???

22



26.1.2004

Felix Naumann, HU Berlin

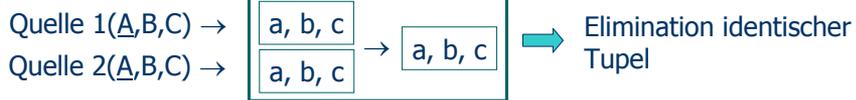
Wie Union integriert

Identifikation

Integration

Optimierung

Visualisierung



23



26.1.2004

Felix Naumann, HU Berlin

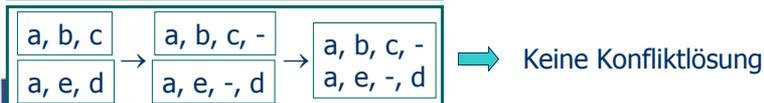
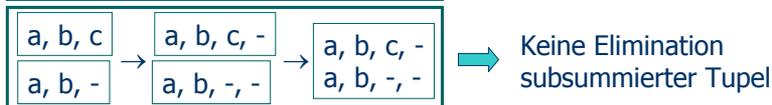
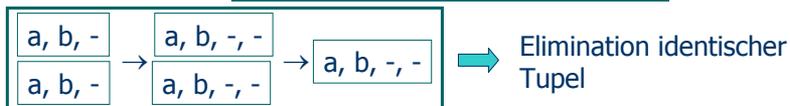
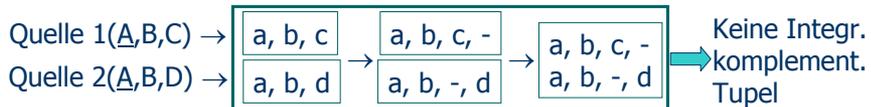
Wie Outer Union integriert

Identifikation

Integration

Optimierung

Visualisierung



24

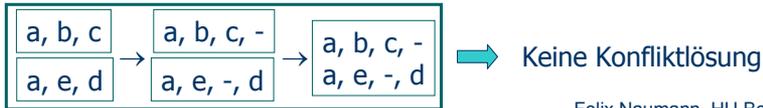
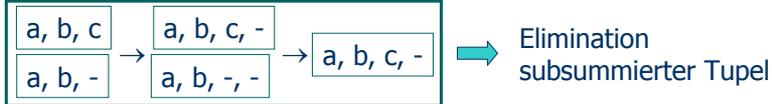
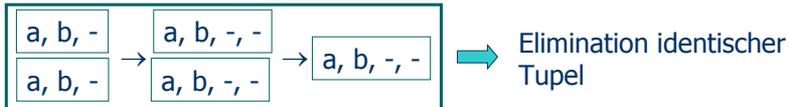
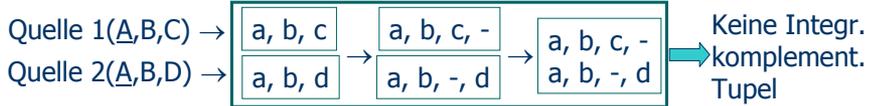


26.1.2004

Felix Naumann, HU Berlin

Wie Minimum Union integriert

Identifikation Integration Optimierung Visualisierung



25

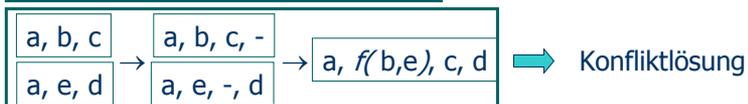
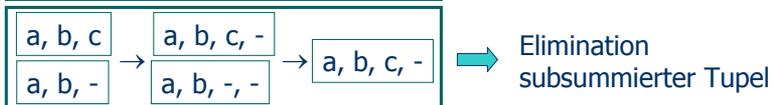
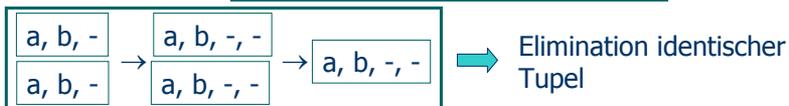
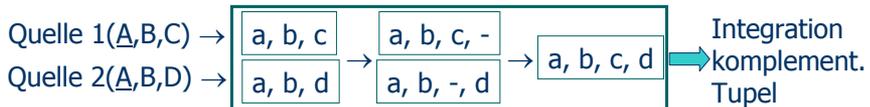
Informatica

20.1.2004

Felix Naumann, HU Berlin

Wie Merge Union integrieren sollte

Identifikation Integration Optimierung Visualisierung



26

Informatica

20.1.2004

Felix Naumann, HU Berlin

Konfliktlösung

Identifikation

Integration

Optimierung

Visualisierung

$$f(x, y) := \begin{cases} \perp & \text{if } x = \perp \text{ and } y = \perp \\ x & \text{if } y = \perp \text{ and } x \neq \perp \\ y & \text{if } x = \perp \text{ and } y \neq \perp \\ g(x, y) & \text{else} \end{cases}$$

null = unbekannt

interne
Konfliktlösungsfunktion

27



26.1.2004

Felix Naumann, HU Berlin

Konfliktlösungsfunktionen

Identifikation

Integration

Optimierung

Visualisierung

- Numerisch:

SUM, AVG, MAX, MIN, ...

- Nicht-numerisch:

MAXLENGTH, CONCAT, AnnCONCAT, ...

- Spezialisiert:

RANDOM, COUNT, CHOOSE, FAVOR, MaxIQ, ...

- Domänen-spezifisch

...

28



26.1.2004

Felix Naumann, HU Berlin

Implementierung von Merge Union

Identifikation Integration Optimierung Visualisierung

- Merging = Gruppierung & Aggregation
 - User-defined Grouping (zur Objektidentifikation)
 - Meist nicht möglich
 - Gruppen typischerweise sehr klein
 - User-defined Aggregation (zur Konfliktlösung)
 - Meist nicht möglich
 - Input nicht nur Werte des aggregierten Attributs
- Relationales Modell
 - Relationale Algebra / SQL
 - Programm
 - Erweiterung einer föderierten Datenbank
- XML Datenmodell
 - XML Query Algebra / XQuery
 - Programm
- Effizienz
- User Spezifikation / Interaktion

29



26.1.2004

Felix Naumann, HU Berlin

Integration - Zusammenfassung

Identifikation Integration Optimierung Visualisierung

- “Wahre” Integration auf Datenwert-Niveau
- Vereinen ist nicht Integrieren
 - Redundanzen
 - Widersprüche
 - Komplemente
- Agenda
 - Merge-Operatoren für relationale und XML Daten
 - Konfliktlösung
 - Implementierung

**Jetzt wissen wir, WIE integriert werden soll.
Aber WIE GUT schaffen wir das?**

30



26.1.2004

Felix Naumann, HU Berlin

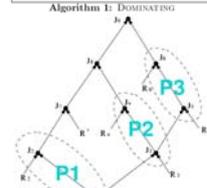
Optimierung - Ziele

Identifikation Integration Optimierung Visualisierung

1. Integrationsalgorithmen
 - Stand-alone Algorithmen zur schnellen Ausführung von Integrationsoperationen
2. Anfrageoptimierung
 - Strategien zur optimalen Ausführung komplexer Anfragen mit Integration
 - Redundanz \Rightarrow größerer Suchraum
 - *Merge union* Operator mit unklaren Kosten
3. Optimierungsziel
 - Paradigmenwechsel von *Zeit* zu *Inhalt*

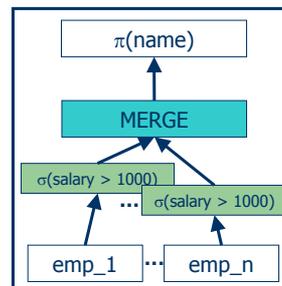
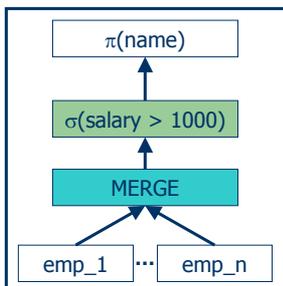
```

Input: Query Q, sources S = {s1, s2, ..., sn},
      costs {c1, c2, ..., cn}, limit L
Output: Result for Q
1: Answer ← []
2: R ← S; [remaining sources]
3: U ← 0; [used cost]
4: while U < (L - C) and R is not empty do
5:   greedy ← greedySelectCoverage(R);
6:   single ← singleLargestCoverage(R);
7:   if greedy > single then
8:     Next ← maxGreedySequence(R);
9:   else
10:    Next ← maxSource(R);
11:   R ← R - {Next};
12:   Execute Q at Next;
13:   if Next is available then
14:     Collect result into Answer;
15:     U ← U + cNext;
16: Return Answer;
    
```



Anfrageoptimierung mit Merging

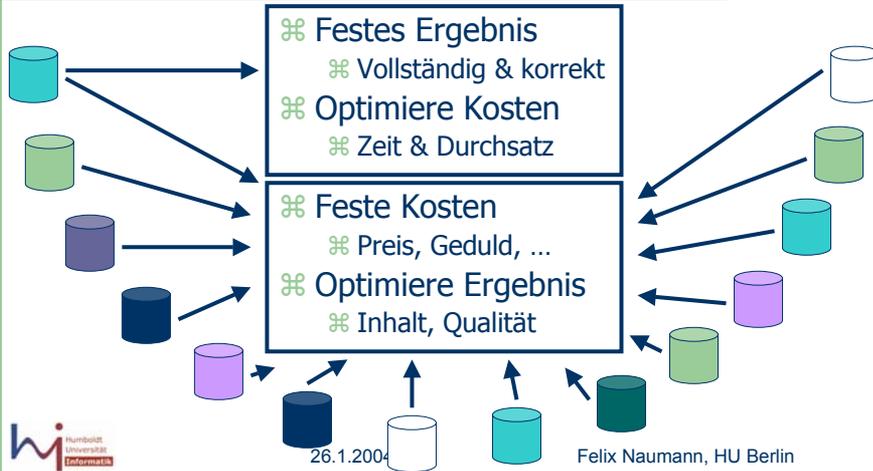
Identifikation Integration Optimierung Visualisierung



- Korrektheit?
- Vollständigkeit?
- Effizienz?

Optimierung – Ein Paradigmenwechsel

Identifikation Integration Optimierung Visualisierung



Optimierung - Zusammenfassung

Identifikation Integration Optimierung Visualisierung

- Behandlung Integrierender Operatoren
 - Stand-alone
 - In Anfrageplänen
- Paradigmenwechsel der Optimierung
 - Nutzeranforderungen
- Agenda
 - Multi-dimensionale Optimierung
 - SQL, XQuery
 - Optimierung in DBMS

Jetzt wissen wir, WIE wir schnell und gut integrieren.
Aber FÜR WEN tun wir das?

Visualisierung – Ziele

Identifikation Integration Optimierung Visualisierung

- Darstellung integrierter Ergebnisse
 - Viele Quellen
 - Viele Operationen / Transformationen
- Nutzerfreundliche Darstellung
- Drill-Down
 - Metadaten
- Kein kleinster gemeinsamer Nenner!

35



26.1.2004

Felix Naumann, HU Berlin

Visualisierung – Metasuchmaschinen

Identifikation Integration Optimierung Visualisierung

- MetaGer
- Mamma
- MetaCrawler
 - Google: Title, Summary, Descr., Category, URL, Size
 - Fast: Title, Summary, Descr., URL, Size
 - Inktomi: Title, Summary, URL

36



26.1.2004

Felix Naumann, HU Berlin

MetaCrawler® Results | Search Query = humboldt-universit%4E4t-~netscape

MetaGer, Suche nach: Humboldt Universität

metacrawler® Search the Search Engines!

Web Pages Directory Listings Audio/MP3 Images Multimedia Shopping News Message Boards

humboldt universität The Web Search

any all phrase Yellow Pages White Pages Submit Your Site | Advanced Search

Are you looking for: [Humboldt](#) [Humboldt County](#) [Humboldt County California](#) [Humboldt University](#)
[Humboldt State University](#) [Humboldt State](#)

MetaCrawler searches these sites: Google • FAST • Ask Jeeves • Inktomi • About • Looksmart • FindWhat • Overture • Teoma

Meta-Search results for "humboldt universität" (1 - 20 of 51) page: 1 - 2 - 3 next

Search by: Relevance | [Source](#) [Send these results to a friend](#)

MetaCrawler Results [About Results](#)

- [Humboldt-Universität zu Berlin](#)
English version, **Humboldt-Universität** zu Berlin. Studium, Forschung, Angebote Zugang ... **Humboldt-Universität** zu Berlin. **Humboldt-Universität** ... [http://www.hu-berlin.de/](#) (Google, Fast, Inktomi, LookSmart Reviewed Sites) | [More like this](#)
- [Universitätsbibliothek der HU Berlin](#)
Willkommen auf der Hauptseite der Universitätsbibliothek der **Humboldt-Universität** zu Berlin: Kataloge, Dienstleistungen, Webinformationen. ... [http://www.ub.hu-berlin.de/](#) (Google, Fast, LookSmart Reviewed Sites) | [More like this](#)
- [Wirtschaftswissenschaftliche Fakultät der Humboldt-Universität Berlin \(Deutschland\)](#) [http://www.wwi.hu-berlin.de/](#) (Google, Fast, Inktomi) | [More like this](#)
- [Landwirtschaftlich-Gärtnerische Fakultät der HU zu Berlin](#)
Die Landwirtschaftlich-Gärtnerische Fakultät der **Humboldt-Universität** zu Berlin benutzt für ihr Webangebot die Frametechnologie. ... [http://www.agrar.hu-berlin.de/](#) (Google, Fast, Inktomi) | [More like this](#)
- [edoc - Dokumenten- und Publikationsserver der Humboldt-Universität](#)
edoc - der Dokumenten- und Publikationsserver ist ein Service für alle Angehörigen der **Humboldt-Universität** zu Berlin zum elektronischen Publizieren ihrer ... [http://edoc.rz.hu-berlin.de/](#) (Google, Fast, Inktomi) | [More like this](#)
- [Humboldt-Universität zu Berlin - Juristische Fakultät](#)
Das Gebäude der Juristischen Fakultät, **Humboldt-Universität** zu Berlin, Welcome Studium L&F Fakultät Studenten Alumni Online Service. FEHLER. ... [http://www.rwi.hu-berlin.de/](#) (Google, Fast) | [More like this](#)
- [Institut für Mathematik](#)
... Institut für Mathematik Mathematisch-Naturwissenschaftliche Fakultät II **Humboldt-Universität** zu Berlin, Uni Homepage. Allgemeines: Unser Institut. ... [http://www.mathematik.hu-berlin.de/](#) (Google, Fast) | [More like this](#)
- [Humboldt-Universität zu Berlin, Institut für ...](#)
... Institut für Bibliothekswissenschaft der **Humboldt-Universität** zu Berlin Geschäftsführender Direktor Prof. Dr. Konrad Umlauf. ... [http://www.lib.hu-berlin.de/](#) (Google, Fast) | [More like this](#)

MetaGer, Suche nach: Humboldt Universität -Netscape

MetaGer, Suche nach: Humboldt Universität

Für detaillierte Anfragen empfehlen wir Ihnen die direkte Benutzung dieser Suchdienste.

Yahoo.de	10	TestEffec
Searchtopportal	22	TestEffec
HotMan-Motors	20	TestEffec
Quali11NK-Telex	0	TestEffec
Telex	20	TestEffec
clustest-search.de	20	TestEffec
Mirago	10	TestEffec
Quali11NK.ch	20	TestEffec
Quali11NK.ch	0	TestEffec
T-Oni line	12	TestEffec
winch	2	TestEffec
zipppp	10	TestEffec
zippp	10	TestEffec
testwanahl	100	TestEffec

HINWEIS: Sie haben SEHR viele Ergebnisse erhalten. Möglicherweise ist es sinnvoll:

- Ihre Suchanfrage zu verfeinern, indem Sie weitere oder speziellere/treffendere Suchworte eingeben, oder
- die Ergebnismenge zu verringern, indem Sie "Ausschlussworte" vorgeben. Klicken Sie hierzu die [MetaGer-Tips](#) an, und lesen dort ggf. die Ziffer 3. oder
- Sie fassen Ergebnisse zusammen, indem Sie auf der MetaGer-Startseite anklicken: "Ausgabe ... clustern" und ggf. zusätzlich "nur Kompakt-Darstellung ausgeben".

Enthalten die Suchergebnisse wirklich das, was Sie suchten?
 -> Klicken Sie auf "QCheck" (QuickCheck, Schnellprüfung) links vor dem jeweiligen Treffer ...

Volltreffer: Humboldt Universität

[QCheck](#)

[http://www.berlinfotos.de/humboldt_universitaet.htm](#)

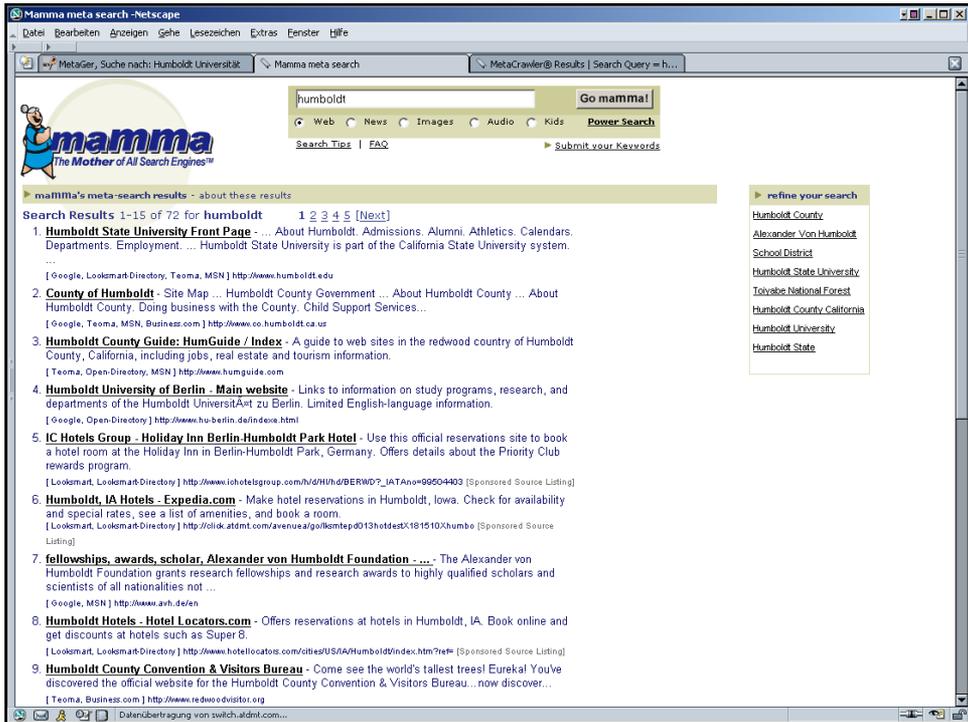
- (gefunden von [Tricus](#)) **HUMBOLDT UNIVERSITÄT** Berlin Mitte Bezirk: Mitte die wurde von 1748 bis 1766 von J. Boumann nach Entwürfen von Georg Wenzeslaus von Knobelsdorff als Palais für Prinz Heinrich erbaut ist 1810 wird dieses Palais für die von Wilhelm von Humboldt gegründete

Volltreffer: Humboldt Universität in Berlin

[QCheck](#)

[http://www.keichel.com/austzug/berlin/humboldt_universitaet.html](#)

- (gefunden von [Tricus](#)) Sehenwürdigkeiten in und um Berlin **HUMBOLDT UNIVERSITÄT** in Berlin Die **HUMBOLDT UNIVERSITÄT** in Berlin wurde 1748 1766 als Palais für Prinz Heinrich, einen Bruder Friedrichs des Großen gebaut. 1809 wurde das Gebäude der Universität übergeben, die 1949



HiIQ Meta Search Engine

Identifikation Integration Optimierung Visualisierung

<http://www.icdt.org/> (14 [Altavista] 3 [Northern Light] 1 [Fast] 1 [Hotbot] 3 [Google])

ICDT Home Page [Altavista]
 ICDT Home Page [Northern Light]
 ICDT Home Page [Fast]
 ICDT Home Page [Hotbot]
 ICDT Home Page [Google]

five search engines found this page

the descriptions differ

description Download your FREE evaluation copy at www.auscomp.com. [Altavista]
 Islamic Centre for Development of Trade. Islamic Centre for Development of Trade web site was built to facilitate your business and trade with OIC. Trade / [Northern Light]
 BusinessOpportunities|EconomicOperators|TradeGuides|Events&Fairs|Indicators|Statistics|Publications
 Islamic Centre for Development of Trade has built this web site to facilitate your business and trade with OIC member states. For Members and NFP [Fast]
 BusinessOpportunities|EconomicOperators|TradeGuides|Events&Fairs|Indicators|Statistics|Publications
 Islamic Centre for Development of Trade has built this web site to facilitate your business and trade with OIC member states. For Members and NFP [Hotbot]
 BusinessOpportunities|EconomicOperators|TradeGuides|Events&Fairs|Indicators|Statistics|Publications
 Islamic Centre for Development of Trade has built this web site to ... [Google]

data conflict

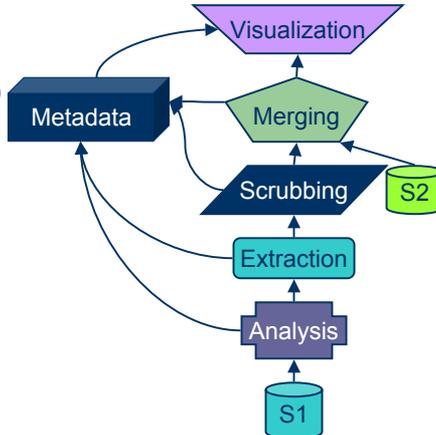
date 07/19/1999 [Altavista] 07/19/1999 [Northern Light] 1/18/2000 [Hotbot] size 8k [Altavista] 8k [Google] language English [Altavista] category Non-profit site [Northern Light] rating 94% [Northern Light]

different attributes from different engines

Metadaten der Integration

Identifikation / Integration / Optimierung / Visualisierung

- Herkunft
 - Data lineage
 - z.B. (Cui & Widom 2001)
 - Durch Operatoren hindurch
- Transformationen
 - Aggregation / Integration
- Konflikte
- Datenqualität
- Ranking



41

Visualisierung – Zusammenfassung

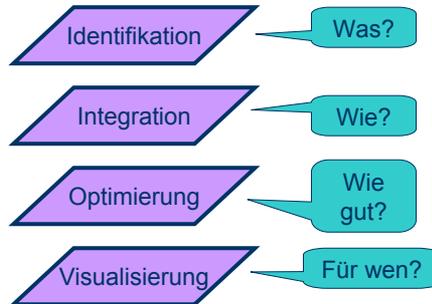
Identifikation / Integration / Optimierung / Visualisierung

- Zeige mehr, nicht weniger
- Metasuchmaschinen
- Metadaten
 - Data Lineage
 - Data Quality
 - Transformationen
- Ziele
 - Spezifikation notwendiger Metadaten
 - Automatisierte Erfassung
 - Metadaten-Verwaltung
 - Drill-Down Visualisierung

Jetzt kennen wir alle notwendigen Komponenten zur effektiven Informationsintegration.

42

Informationsintegration - Zusammenfassung



www.informatik.hu-berlin.de/macnaumann@informatik.hu-berlin.de

43

Literatur

- Fellegi & Sunter 1969
 - *A theory of record linkage*. Journal of the American Statistical Association, 64: 1183-1210.
- Codd 1979
 - *Extending the Relational Database Model to Capture more Meaning*. ACM Transactions on Database Systems (TODS), 4(4): 397-434.
- Galindo-Legaria 1994
 - *Outerjoins as Disjunctions*. ACM SIGMOD Conference, 348-358.
- Hernandez & Stolfo 1998
 - *Real-world data is dirty: Data cleansing and the merge/purge problem*. Data Mining and Knowledge Discovery, 2(1): 9-37.
- Rahm & Bernstein 2001
 - *A survey of approaches to automatic schema matching*. VLDB Journal 10(4), 334-350.
- Cui & Widom 2001
 - *Lineage tracing for general data warehouse transformations*. VLDB Journal 12(1): 41-58.
- Ananthakrishna, Chauduri & Ganti 2002
 - *Eliminating Fuzzy Duplicates in Data Warehouses*. VLDB Conference: 586-597.

44