

Information Quality in Integrated Information Systems

Institute for Infocomm Research, Singapore

15.3.2005

Felix Naumann



Humboldt-Universität zu Berlin



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

2



Humboldt-Universität zu Berlin



- Wilhelm and Alexander von Humboldt
- Unity of teaching and research
- Freedom and independence of sciences
- 29 Nobel prize winners
 - Mommsen, Hertz, Koch, Hahn, Planck, Einstein,...
- 38,000 students, (1200 computer sciences)
- 560 professors (12 + 10 computer sciences)



15.3.2005

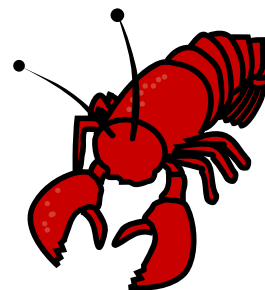
Felix Naumann - Humboldt-Universität zu Berlin

3

Research Group „Information Integration“



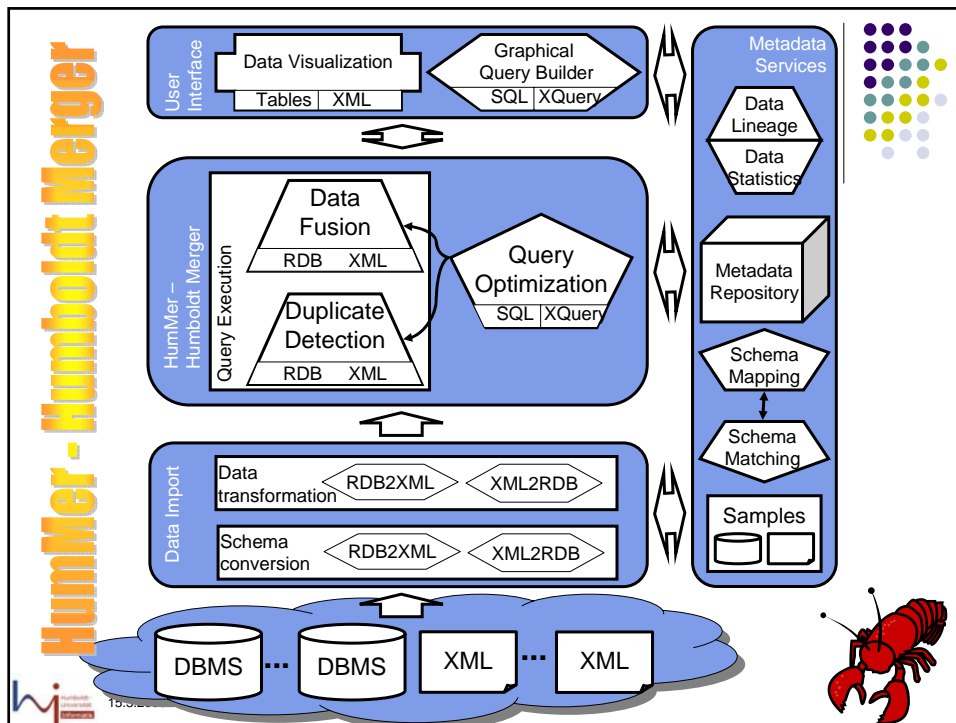
- Head: Felix Naumann (naumann@informatik.hu-berlin.de)
- PhD students
 - Jens Bleiholder (bleiho@informatik.hu-berlin.de)
 - Information fusion in relational data
 - Melanie Weis (mweis@informatik.hu-berlin.de)
 - Object identification in XML Data
- Affiliated PhD students
 - Armin Roth (aroth@informatik.hu-berlin.de)
 - Data quality in Peer-Data-Management-Systems
 - Alexander Bilke (bilke@cs.tu-berlin.de)
 - Schema Matching
- General Research Topics
 - Object identification
 - Information fusion
 - Optimization
 - Visualization
 - And IQ



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

4

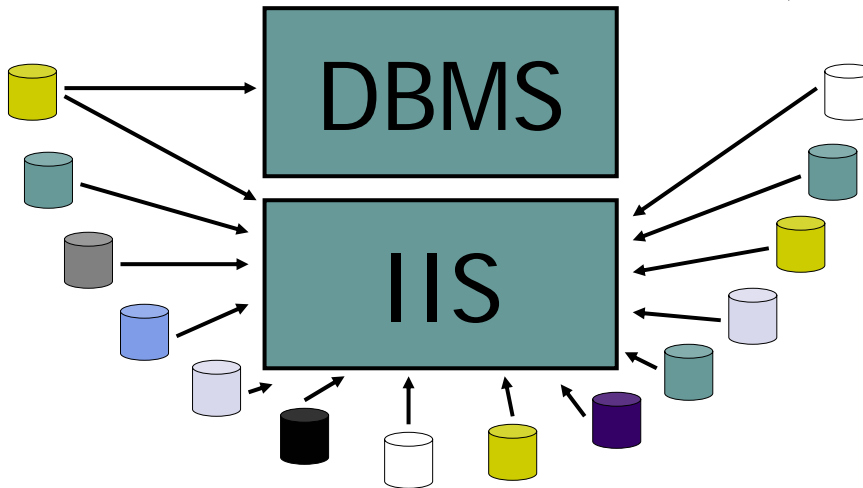


Outline

- ➔ ● Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration
- Further topics for discussion



Database Management Systems vs. Integrated Information Systems

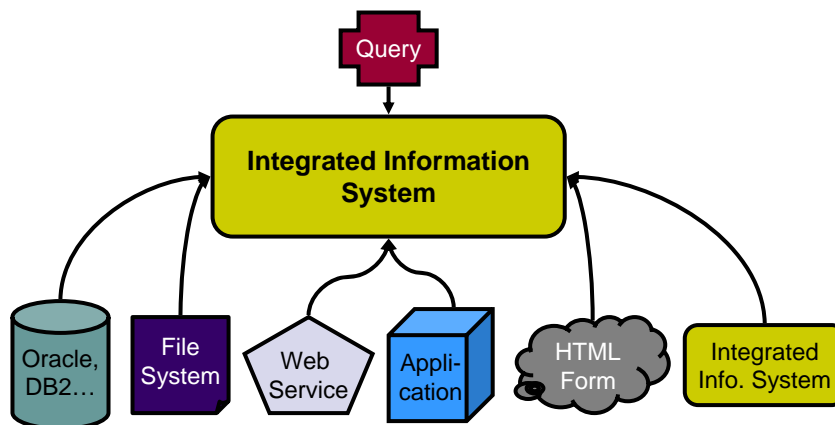


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

7

Integrated Information Systems

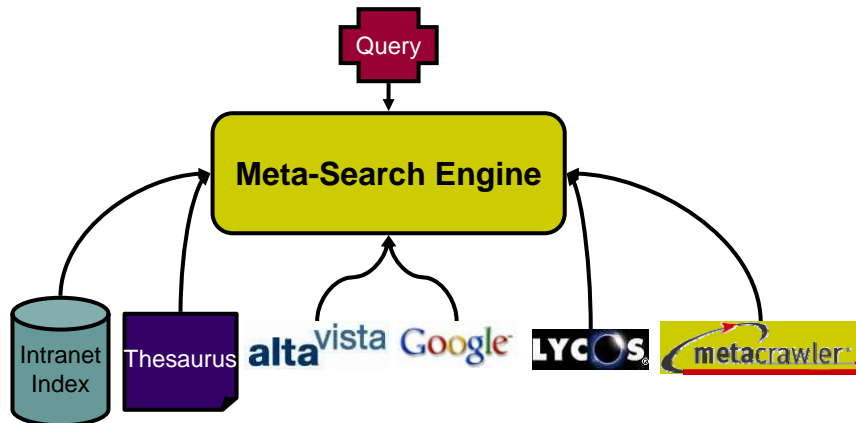


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

8

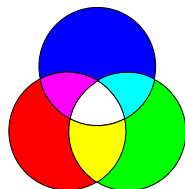
Integrated Search Engines



Redundancy



- Is good for you
 - More information
 - More detailed information
 - Verifiable information
 - Thus: **INTEGRATE!**
- But yields problems
 - More alternatives mean more **complexity**.
 - Technical and structural (and semantic) **heterogeneity**
 - Redundancy is only conceptual:
 - Object identity
 - Data conflicts
 - Mixed and low **quality**
 - New types of **errors** (duplicates)





DBMS vs. IIS

- Single source ↔ • Multiple sources
- Local ↔ • Distributed
- Controlled ↔ • Autonomous
- Structured ↔ • Semi-structured
- Data ↔ • Information
 - short
 - text, images, tables
- Professional user ↔ • Professional and casual user



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

11



DBMS Quality vs. IIS Quality

- Complete (assumed) ↔ • Incomplete
- Accurate ↔ • Inaccurate
- Trusted ↔ • Untrusted
- Fast ↔ • Slow
- Free ↔ • Possible cost

High expectations
(guarantees)
High quality

Low expectations
(no guarantees)
Low quality

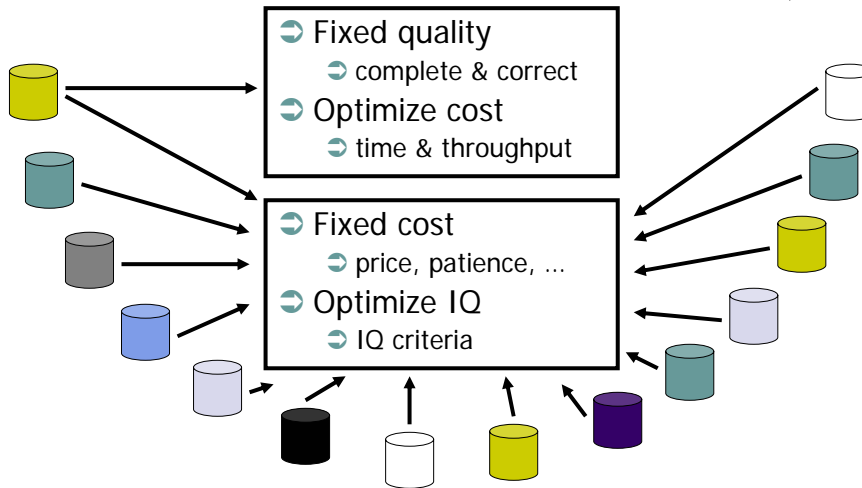


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

12

Optimize IQ!



Outline

- Databases vs. Integrated Information systems
- ➔ • Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration
- Further topics for discussion



Motivating Examples



- Search
 - Meta Search Engines
- Data Fusion (Integration)
 - Data conflicts
 - Genome Databases
- Peer Data Management



Example: Meta Search Engines



- MetaGer



- Mamma



- MetaCrawler



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

MetaCrawler® Results | Search Query = humboldt+universit%at-4t -Netscape

MetaGer, Suche nach: Humboldt Universität

metacrawler® Search the Search Engines! Check Mail Tools & Tips

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

humboldt universität The Web 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](#)
 Englisch 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 ...](#)
 edoc - der Dokumenten- und Publikationsserver ist ein Service für alle Angehörigen der **Humboldt-Universität** zu Berlin zum elektronischen Publizieren ihrer ... [http://gochoad.zi.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.rewi.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.ib.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 Treffer
Excite.de:	22 Treffer
Uni-Hann-Harveste:	20 Treffer
Quali IL 1997-Teicus:	0 Treffer
Teicus:	20 Treffer
campus-search.de:	20 Treffer
Hitago:	10 Treffer
Quali 1997.ch:	20 Treffer
Quali IL 1997.ch:	0 Treffer
T-oni Int:	12 Treffer
Netchn:	2 Treffer
Vipppp:	10 Treffer
Vipppp:	10 Treffer
Getmanzahl:	156 Treffer

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

- Ihre Suchanfrage zu verfeinern, indem Sie weitere oder speziellere/treffendere Suchwörter eingeben, oder
- die Ergebnismenge zu verringern, indem Sie "Ausschlusswörter" 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

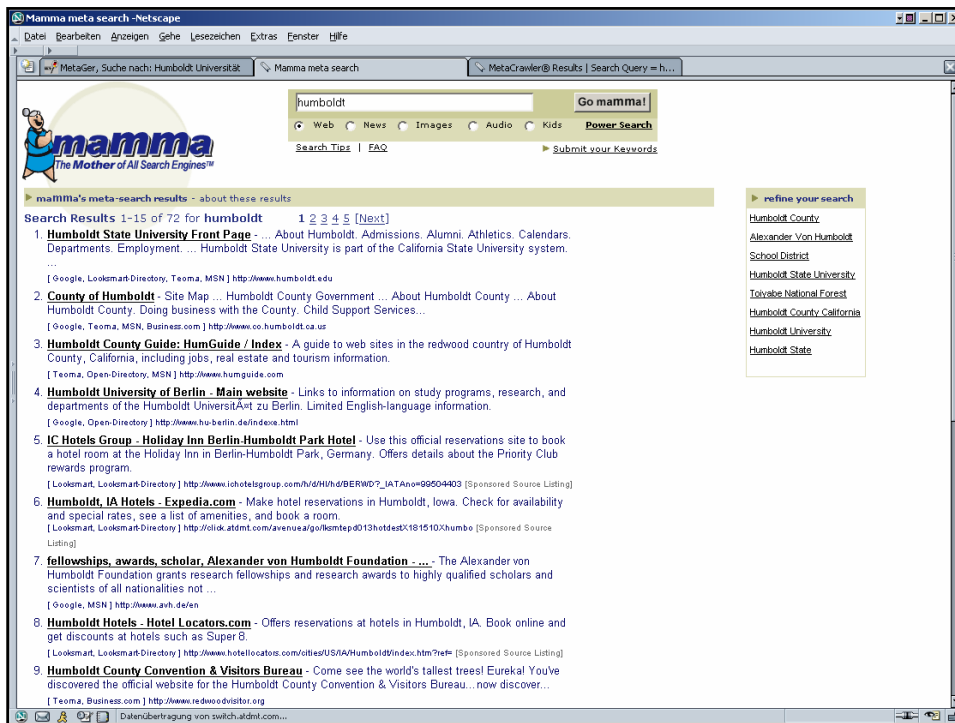
- (gefunden von [Teicus](#)) **HUMBOLDT UNIVERSITÄT** Berlin Mitte Bezirk. Mitte è wurde von 1748 bis 1766 von J. Boumann nach Entwürfen von Georg Wenzeslaus von Knobelsdorff als Palas für Prinz Heinrich erbaut è 1810 wird dieses Palas für die von Wilhelm von Humboldt gegründete

Volltreffer: Humboldt Universität in Berlin

[QCheck](#)

http://www.keichel.com/ausflug/berlin/humboldt_universitaet.html

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




HiQQ Meta Search Engine

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

[ICDT Home Page](#) [Altavista] *five search engines found this page*
[ICDT Home Page](#) [Northern Light]
[ICDT Home Page](#) [Fast]
[ICDT Home Page](#) [Hotbot] *the descriptions differ*
[ICDT Home Page](#) [Google]

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
data conflict IC 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]

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



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

20

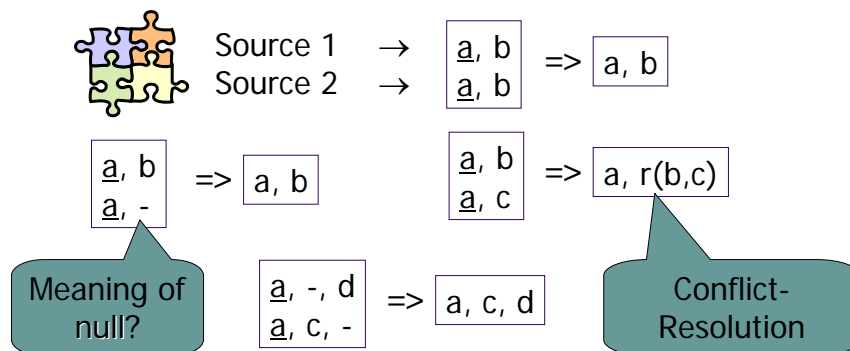
Motivating Examples



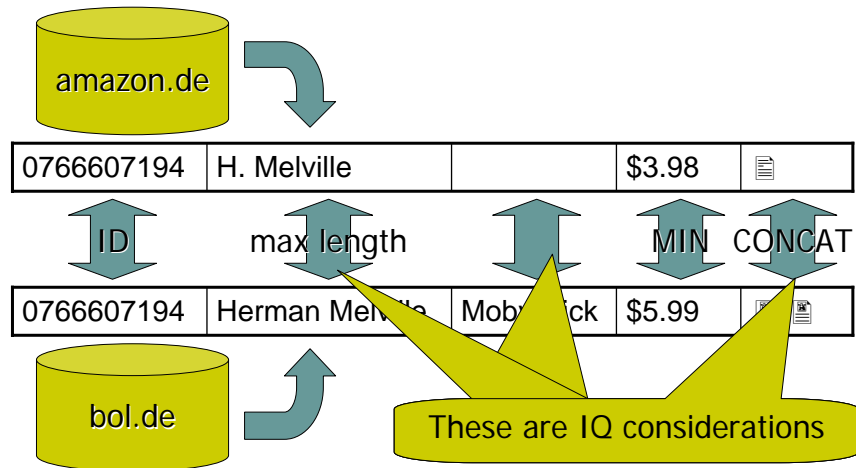
- Search
 - Meta Search Engines
- Data Fusion (Integration)
 - Data conflicts
 - Genome Databases
- Peer Data Management



Example: Data Fusion



Example: Data Fusion



Example: Data Fusion



- Genome Databases
 - Research labs produce many conceptually redundant results.
 - Quality of experimental results vary (poor IQ is intrinsic).
 - Quality of human (and machine) annotations vary.
 - Biologists are very opinionated about IQ of information sources.
 - ⇒ Perfect application for IQ-driven information integration

Example: Data Fusion

- Genome databases
 - Identification of redundancy very difficult (no standard representation)
 - No IQ metadata,
 - only human-readable annotation
 - (Some annotation on "methods used")

```

ID SNOTPCHI standard: RNA; BO#: 1016 BF.
XX
DT 01-ADG-1991 (Rel. 28, Created)
DT 04-MAR-2000 (Rel. 63, Last updated, Version 2)
XX
DE Rat GTP cyclohydrolase I mRNA, complete cds.
XX
KW GTP cyclohydrolase I,
XX
OC Mammalia; Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Eumetazoa; Mammalia;
OC Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus.
XX
RN [1]
RF 1-1016
RX MEFLINE; 91092270.
RX PUBMED; 1989962.
RA Satohyama H., Inoue Y., Nareda T., Kagamiyama H.
RT "Cloning and sequencing of cDNA encoding rat GTP cyclohydrolase I: The
RT first enzyme of the tetrahydrobiopterin biosynthetic pathway";
RL J. Biol. Chem. 266(8):765-769(1991).
XX
FT CDS
FT /code_start=1
FT /db_xref="GOA:G22288"
FT /db_xref="SWISS-PROT:G22288"
FT /EC_number="3.5.1.16"
FT /gene="GTP cyclohydrolase I"
FT /product="GTP cyclohydrolase I"
FT /protein_id="AA81109.1"
FT /translation="MEKFDVFCVTHGFFKELPPGASFPKESPPFAAGQADAVK
FT AGRPFEKINELMLPMLAALYSILBGLGEPQQLLETFWPAATANGFFTKYQKTI
FT SVGLDIFIDFQKREKIVDIDFYSKRELVFFQWRIQLFRVQLGSLKLVY
FT EIVSRQLQGERLTKQIAVITLQFAGVVVVEATHRCVRSQVQNSKVTVTSL
FT GVFFKSPKTRKFFLTLRS"
FT
SO Sequence 1016 BF: 236 A; 279 C; 194 G; 210 T; 0 other:
gaattgaaac ccaattcggc gcaaaactcc tgcctccggc aacagccacg gtaacggccg 60
ccggttaagc cagacccagc cgtctgtag caacttaggg tctctcggga caaatcgcgc 120
gggtccatg gagaagccgc ggggtctagg gtcacacat ggtctcccg agcggggcct 180
...
cattcagagc ccaacttcag tgcctcagc ccgctttaga gacccctcgt gtagccagcg 240
ctactcgtct caattgataa ttccagttcc agttgata ctgtcaact ctacttctca 300
ccatgattg taattataa ttatttatag agatgcataa taaagtgat caactt 360
1016
//
    
```

RNA sequence entry from EMBL



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

25

Motivating Examples

- Search
 - Meta Search Engines
 - Information Retrieval
- Data Fusion (Integration)
 - Data conflicts
 - Genome Databases
- Peer Data Management

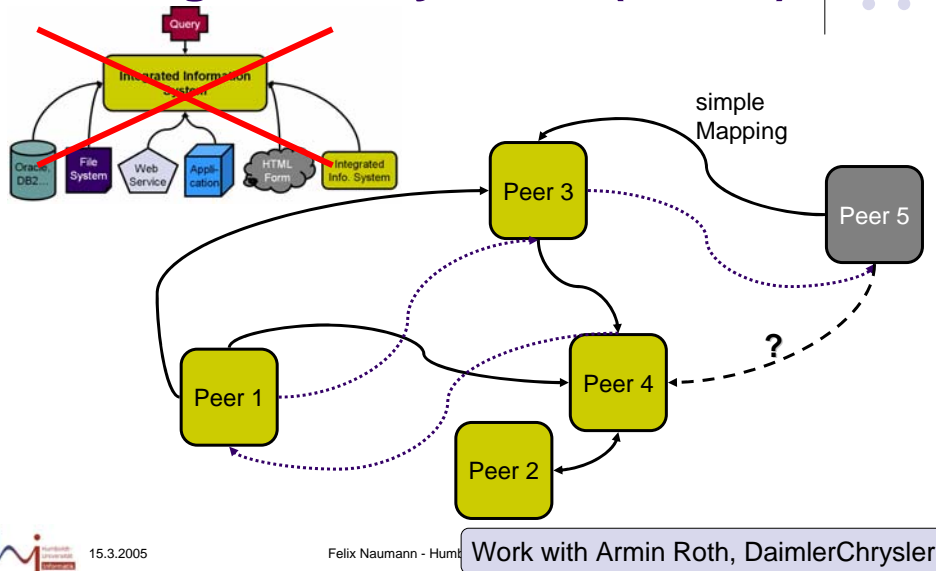


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

26

Example: Peer Data Management Systems (PDMS)

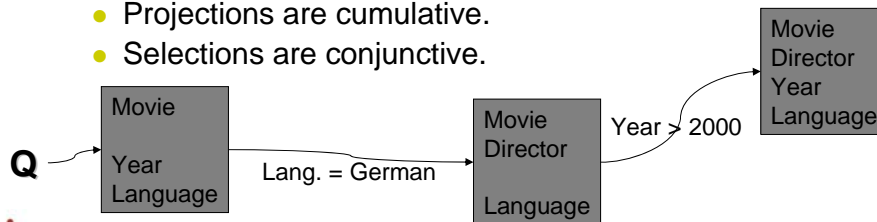


Example: Peer Data Management Systems (PDMS)



• IQ in PDMS

- You don't know who is answering your query.
- Mappings further degrade IQ:
 - Projections diminish density
 - Selections diminish completeness
- Problem compounded along mapping paths:
 - Projections are cumulative.
 - Selections are conjunctive.



Outline

- Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration
- Further topics for discussion

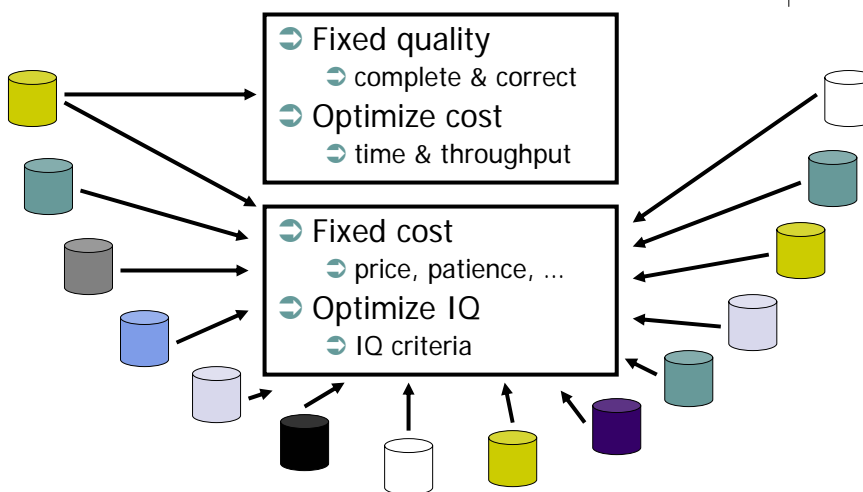


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

29

Optimize IQ!



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

30

A New Optimization Paradigm – Changes



- Cost criteria ↔ • Quality criteria
- Cost model ↔ • Quality model
- Optimization algorithm ↔ • Optimization algorithm
- + Information integration



15.3.2005

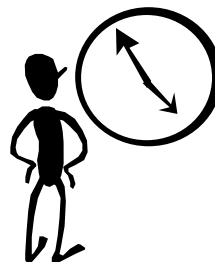
Felix Naumann - Humboldt-Universität zu Berlin

31

DB Cost Criteria



- Response time
- Execution time
- Latency
- Throughput
- Cardinality
- ...



Assessed through system parameters and statistics.



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

32

IIS Quality Criteria



IQ :=

*"Even though quality
cannot be defined, you
know what it is."*

Robert Pirsig



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

IIS Quality Criteria



IQ := {Understandability, Reputation,
Reliability, Timeliness,
Availability, Price, Conciseness
Consistency, Coverage,
Response time, Density,
Completeness, Amount,
Accuracy, Relevancy, ... }



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

34

IQ Criteria



- Accuracy
 - Definition:
 - Usually: Percentage of incorrect tuples
 - For integration: Percentage of incorrect data values
 - Assessment:
 - Domain and Constraint Testing
 - Lookup tables
 - Scientific measurements
 - Data-input experience
 - Improvement:
 - Often: Deletion
 - Better: "Data Scrubbing"
- Response Time
 - Definition:
 - Usually: Time until complete query result is received
 - For integration: Latency
 - Assessment:
 - "Cost Calibration"
 - Continuous assessment
 - Improvement:
 - Source selection
 - Classical optimization
 - Federated Optimization



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

35

IQ Criteria



- Completeness
 - Definition:
 - Coverage (extension): Number of real world objects represented
 - Density (intension): Number of attributes covered
 - NULL-values
 - Assessment:
 - Sampling
 - Existing Metadata
 - Improvement:
 - Source selection
 - "Best k" vs. "k best"
- Conciseness
 - Definition:
 - Extension: Number of real-world objects not represented doubly
 - Intension: Number of attributes not represented doubly
 - Assessment:
 - Duplicate detection
 - Schema Matching
 - Improvement:
 - Data Fusion



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

36

IIS Quality Criteria



IQ := {Understandability, Reputation,
Reliability, Timeliness,
Availability, Price, Conciseness
Consistency, Coverage,
Response time, Density,
Completeness, Amount,
Accuracy, Relevancy, ... }

Assessed in 3 classes...



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

37

IQ Assessment



Information source:	S_1	S_2	S_3
Understandability	5	7	7
Reputation	5	5	7
Reliability	2	6	4
Age	30	30	2
Availability	99	99	60
Completeness	1	1	0.5
Response Time	0.2	0.2	0.2
Accuracy	99.9	99.9	99.8
Relevancy	60	80	90

- User
 - Questionnaires
- Information Source
 - Parsing
 - Metadata
- Query Process
 - Time Measurement



15.3.2005

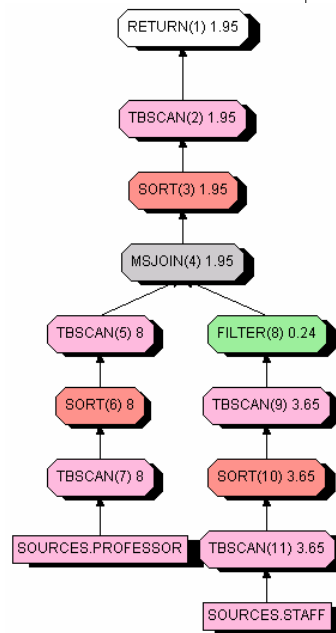
Felix Naumann - Humboldt-Universität zu Berlin

38

DB Cost Models

- Operators
 - + (add)
 - max
 - X (multiply)

Easy ;-)



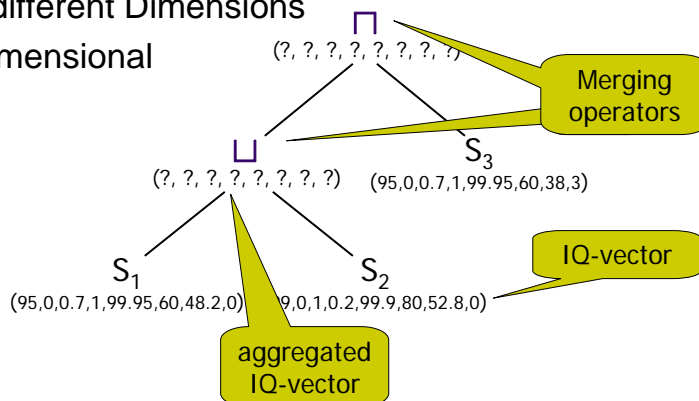
15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

39

IIS Quality Model

- 2 Problems
 - Many different Dimensions
 - Multidimensional



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

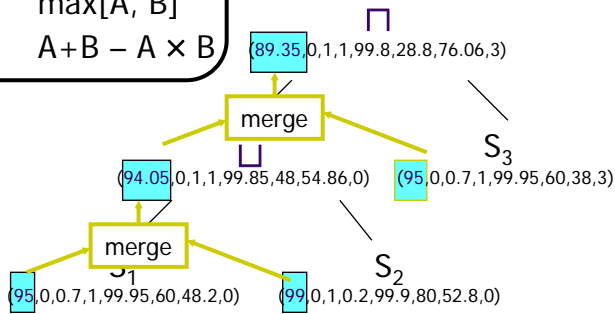
40

Merge IQ in many Dimensions



IQ Merge Functions

Availability: $A \times B$
 Price: $A + B$
 Response Time: $\max[A, B]$
 Coverage: $A+B - A \times B$



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

41

Multidimensional IQ



(89.35, 0, 1, 1, 99.8, 28.8, 76.06, 3) > (82.35, 0, 2, 1.5, 95, 32, 71.77, 2) ?

- IQ-criteria have
 - Different units
 - Different ranges
 - Different importance
- So...
 - convert
 - scale
 - weight

MADM methods:
 SAW, TOPSIS, ELECTRE, AHP, DEA



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

42

Outline

- Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration
- Further topics for discussion



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

43

DB-type Optimization

- Goal
 - Minimize response time (single user)
 - Maximize throughput (multi-user)
- Restrictions
 - Complete
 - Correct (not just accurate: filter conditions...)

Find best plan!



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

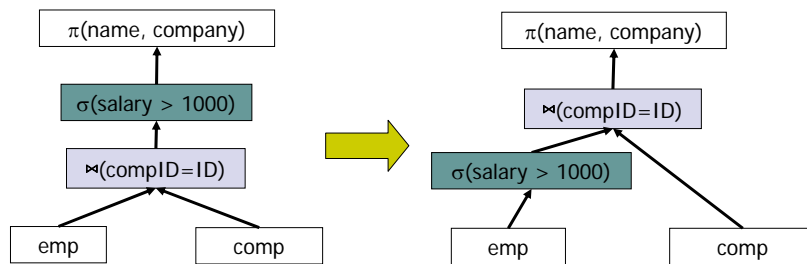
44



DB-type Optimization



SELECT name, company FROM emp, comp
WHERE emp.compID = comp.ID AND emp.salary > 1000

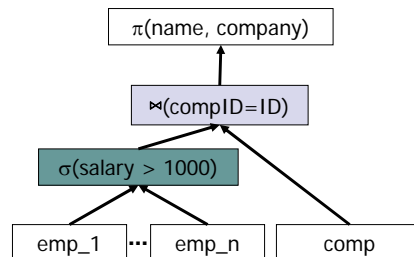


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

45

DB-type Optimization

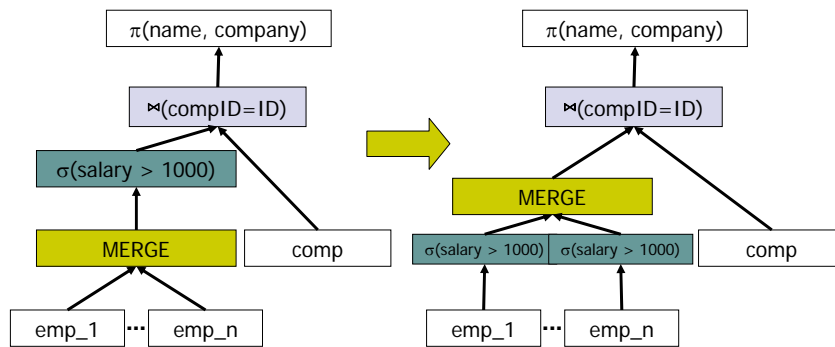


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

46

Integration-type Optimization

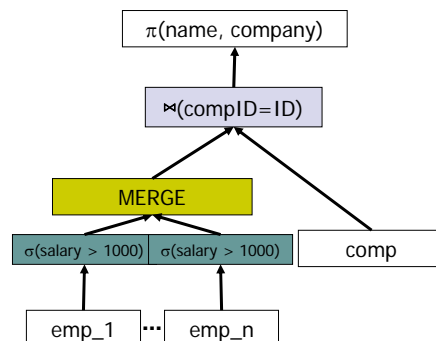


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

47

Integration-type Optimization



- Change is efficient
- But:
 - Result can be incomplete.
 - Depending on "Merge function"
- Preferences? Quality!



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

48

IIS-type Optimization

- Goal
 - Maximize information quality
 - Maximize completeness
- Restrictions
 - Price
 - Bandwidth
 - Time (user patience)



Find K best sources – Find best K sources

IIS-type Optimization

- K best sources
 - (or K best conjunctive plans)
 - Simple IQ model, but
 - Sources may not **complement** each other
 - at tuple level (replication)
 - at attribute level
- Best K sources
 - (or best disjunctive normal form plan)
 - Finds optimal query result
 - Uses IQ merging

Outline

- Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
- ➔
 - IQ-driven integration
- Further topics for discussion



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

51

Data Integration in DBMS

- UNION, UNION, UNION
- OUTER UNION,
- GROUP BY

- Cannot compute **similarity**.
 - For object identification
- Cannot deal with **inconsistencies**.
 - For conflict resolution



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

52

Information Integration in IS



- WHAT?
 - **Problem:** Identify different representations of same real-world object.
 - **Solution:** Object identification techniques
- HOW?
 - **Problem:** Merge conflicting data about same real-world object.
 - **Solution:** Conflict resolution techniques



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

53

Object Identification



- Also
 - Duplicate Detection
 - Record Linkage
 - Data Reconciliation/Consolidation
- Domain-specific
 - Address data
 - Microbiological data
 - ...
- IQ weighting can help

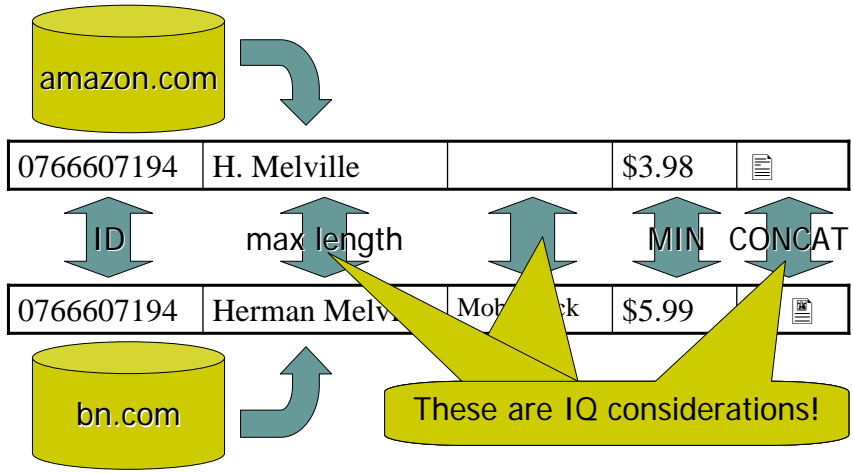


15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

54

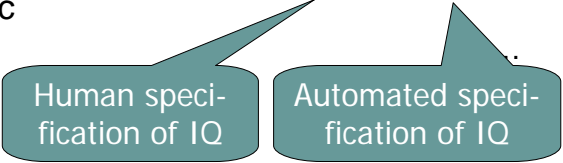
Conflict Resolution



Conflict Resolution



- Numerical: SUM, AVG, MAX, MIN, ...
- Non-numerical: MAXLENGTH, CONCAT, AnnCONCAT, ...
- Special: RANDOM, COUNT, CHOOSE, FAVOR, MaxIQ, ...
- Domain-specific



Outline

- Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration
- ➔ • Further topics for discussion



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

60

Further topics for discussion

- One slide each on ongoing work:
 - IQ in Off-the-Shelf-DBMS
 - Detecting regular errors in DBMS using Data Mining
 - Duplicate detection in XML
 - Duplicate detection without a schema
 - Data Fusion – Completeness and Conciseness
- 10th Int. Conference on Information Quality



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

61

IQ in Off-the-Shelf-DBMS



1. Most quality-critical data is stored in DBMS.
 - Data Warehouses
 - Web back-ends
2. DBMS are carefully designed to preserve high data quality.
 - Decades of development
 - ACID
3. Sometimes DBMS even improve data quality.
 - Aggregation
 - Integration

Work with Mary Roth, IBM



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

62

Detecting regular errors in DBMS using Data Mining



ID	ORG.	SEX	COLOR	SIZE
1	Frog	F	Green	20,4
2	Toad	F	Green	25,1
3	Newt	M	Grey	30,4
4	Frog	F	Blue	19,6
5	Toad	M	Blue	20,8
6	Newt	M	Grey	30,1
7	Newt	F	Grey	28,9
8	Frog	M	Green	20,3

ID	ORG.	SEX	COLOR	SIZE
1	Frog	W	Olive	20
2	Frog	W	Green	25,1
3	Newt	M	Grey-Spotted	30,4
4	Frog	W	Blue	19,6
5	Frog	M	Blue	20,8
6	Newt	M	Grey-Spotted	30,1
7	Newt	W	Grey&Yellow	28,9
8	Frog	M	Olive	20

ID	C[ORGANISM]			C[SEX]			C[COLOR]			C[SIZE]		
	R ₁ -ORG	R ₂ -ORG	C ₁	R ₁ -SEX	R ₂ -SEX	C ₂	R ₁ -COLOR	R ₂ -COLOR	C ₃	R ₁ -SIZE	R ₂ -SIZE	C ₄
1	Frog	Frog	0	F	W	1	Green	Olive	1	20,4	20	1
2	Toad	Frog	1	F	W	1	Green	Green	0	25,1	25,1	0
3	Newt	Newt	0	M	M	0	Grey	Grey-Spotted	1	30,4	30,4	0
4	Frog	Frog	0	F	W	1	Blue	Blue	0	19,6	19,6	0
5	Toad	Frog	1	M	M	0	Blue	Blue	0	20,8	20,8	0
6	Newt	Newt	0	M	M	0	Grey	Grey-Spotted	1	30,1	30,1	0
7	Newt	Newt	0	F	W	1	Grey	Grey&Yellow	1	28,9	28,9	0
8	Frog	Frog	0	M	M	0	Green	Olive	1	20,3	20	1

Work with Heiko Müller, Humboldt



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

Duplicate detection in XML



```

- <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>
  
```

- Compare <author>
 - With Subelements (<publication>)?
 - How deep?
- Compare <publication>
 - With siblings (<year>)?
 - Schema, or Data?
- In short: What is a duplicate?



15.3.2005

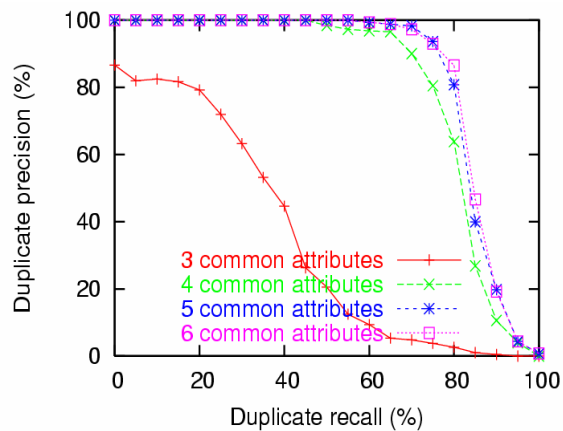
Felix Naumann - Humboldt-Universität zu Berlin

Work with Melanie Weis, Humboldt

Duplicate detection without a schema



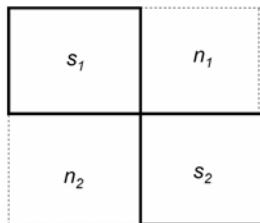
- Opaque attribute names
- Only some intensional overlap
- Similarity measure based on tokenization of tuple and TFIDF
- Whirl algorithm to efficiently find top duplicates



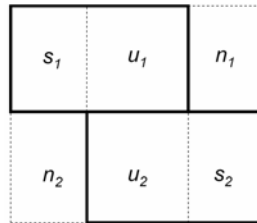
15.3.2005

Work with Alexander Bilke, University of Technology Berlin

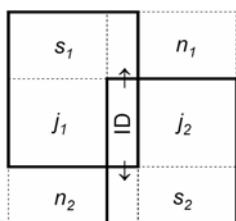
Data Fusion – Completeness and Conciseness



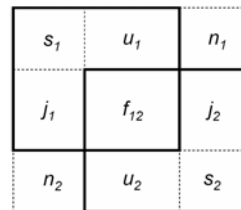
„Union Join“



Outer Union



Outer Join



Data Fusion



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

Work with Jens Bleiholder, Humboldt-Universität zu Berlin

CfP: 10th ICIQ



- International Conference on Information Quality
- Nov. 4 to 6 @ MIT in Cambridge
- Paper deadline: June 30 !
- Suggested topics:
 - IQ Concepts, Tools, Metrics, Measures, Models, and Methodologies
 - IQ Policies and Standards
 - IQ Assessment
 - IQ Practices: Case Studies and Experience Reports
 - IQ Product Experience Reports
 - Cost/Benefit Analysis of IQ and IQ Improvement
- Information Product Implementation, Delivery, and Management
- IQ in Databases, the Web, and e-Business
- Data Warehouses and Data Mining
- Corporate Household Data
- IQ in Scientific Data Management
- Data Cleansing and Reconciliation
- IQ Education and Curriculum Development
- Trust, Knowledge, and Society in the IQ Context
- Academic, research-in-progress, or practice-oriented paper.



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

67

Questions

- Databases vs. Integrated Information systems
- Motivating Examples
- Components for IQ-driven query processing
 - IQ measure + IQ model
 - IQ-driven optimization – A new paradigm
 - IQ-driven integration



www.hiqiq.de

[www.informatik.hu-berlin.de/mac/
naumann@informatik.hu-berlin.de](http://www.informatik.hu-berlin.de/mac/naumann@informatik.hu-berlin.de)



15.3.2005

Felix Naumann - Humboldt-Universität zu Berlin

68