



**Hasso  
Plattner  
Institut**

IT Systems Engineering | Universität Potsdam

Beauty is our Business

Organisatorisches und Einführung

23.4.2009

Felix Naumann

## Das Motto

2

... Wenn wir uns klarmachen, dass der Kampf gegen Chaos, Durcheinander, und unbeherrschte Kompliziertheit eine der größten Herausforderungen der Informatik ist, müssen wir zugestehen:

***„Beauty is our Business“.***

*Edsger W. Dijkstra, 1978*

Auch:

“Computer Science is no more about computers than astronomy is about telescopes.”

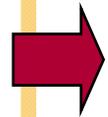
# Motivation

3

- Vorbereitung auf das Master Studium
  - Wissenschaftlich Arbeiten
- Vorbereitung auf den Beruf
  - Ideen „verkaufen“
- Interesse an den Themen

# Überblick

4



- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Ziele des Seminars
- Kurzvorstellung der Themen



5

Felix Naumann

- [naumann@hpi.uni-potsdam.de](mailto:naumann@hpi.uni-potsdam.de)
- Raum A-1.13 (über Frau Pamperin)
- HU, IBM Almaden, DFG/JP/HU, HPI

Universität Potsdam, Hasso-Plattner-Institut

- Fachgebiet „Informationssysteme“

Forschungsthemen (gleich mehr):

- Datenqualität
- Informationsintegration
- Peer Data Management
- Search
- <http://www.hpi.uni-potsdam.de/~naumann/>

# Information systems team

6

project **ViQTOR**



Paul Führung



Katrin Heinrich

**DQ Annotation & Assessment**



Prof. Felix Naumann

**Information Integration**



Jens Bleiholder

**Data Fusion**

project **HumMer**

project **fusem**



Christoph Böhm

**Data Profiling & Cleaning**

**Information Quality**



Armin Roth

**Peer Data Management Systems**

**Matching**

**Data Integration for Life Science Data Sources**



Alexander Albrecht

**ETL Management**

project **System P**



Mohammed AbuJarour

**Service-Oriented Systems**

**Ontologies, Profiling**

project **Aladin**



Jana Bauckmann

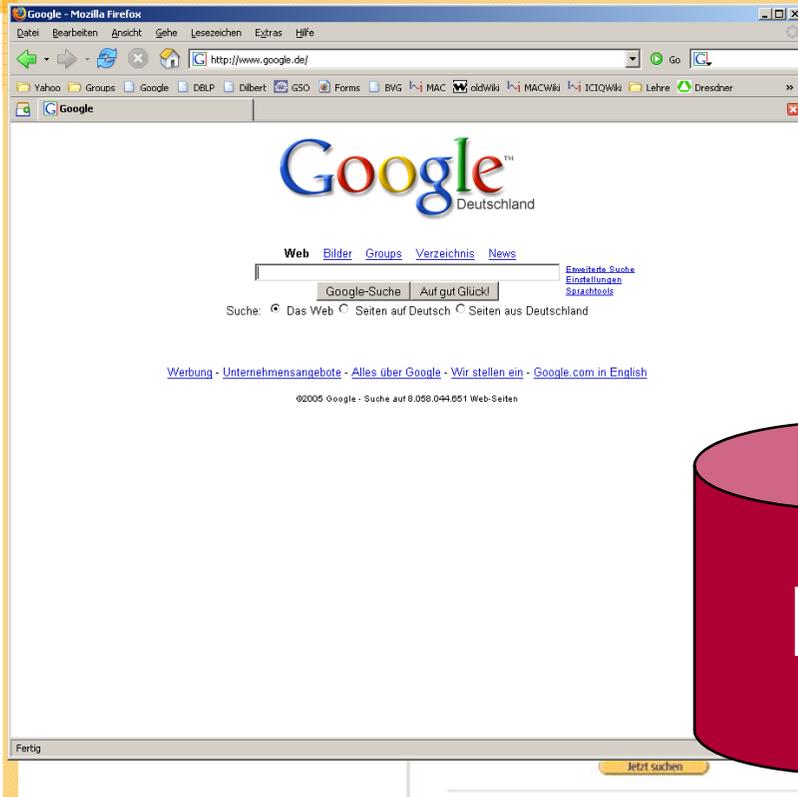
**Data Profiling for Schema Management**



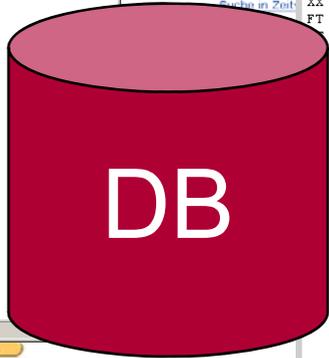
Frank Kaufer

# Was sind Informationssysteme?

7

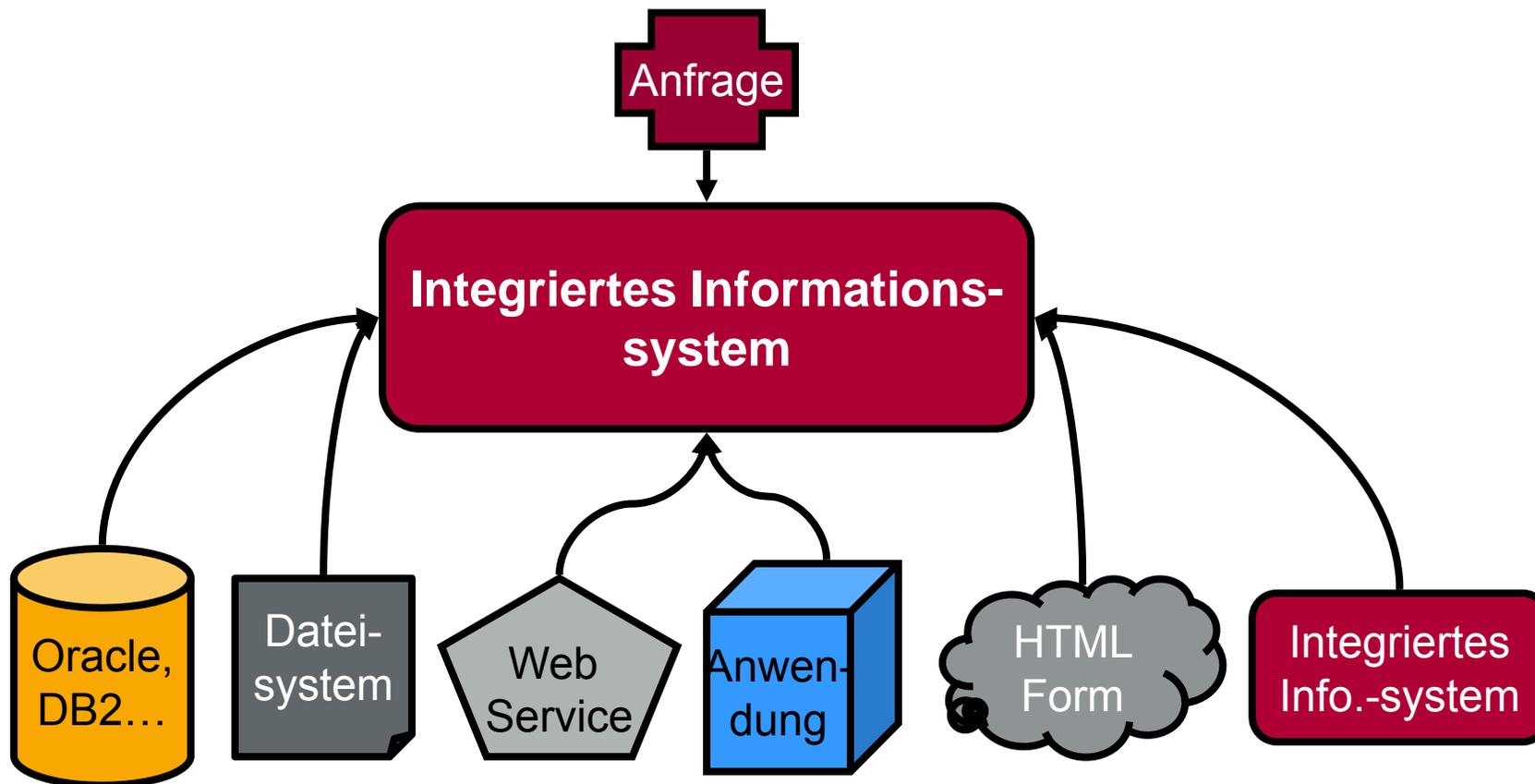


ID	RNGTPCHI	standard; RNA; ROD; 1016 BP.	Molecule type
XX			Name
DT	01-AUG-1991	(Rel. 28, Created)	Date of creation and last update
DT	04-MAR-2000	(Rel. 63, Last updated, Version 2)	
XX			
DE	Rat	GTP cyclohydrolase I mRNA, complete cds.	Free text description
XX			
KW	GTP cyclohydrolase I.		Keywords describing the molecule
XX			
OS	Rattus norvegicus	(Norway rat)	Organism
OC	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;		
OC	Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus.		
XX			
RN	[1]		Article the sequence was published in
RP	1-1016		
RX	MEDLINE; 91093270.		
RX	PUBMED; 1985963.		
RA	Hatakeyama K., Inoue Y., Harada 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(2):765-769(1991).		
XX			
FT	CDS	128..853	Structural annotation (coding sequence)
		/codon_start=1	
		/db_xref="GOA:P22288"	Link to functional annotation of resulting protein
		/db_xref="SWISS-PROT:P22288"	
		/EC_number="3.5.4.16"	
		/gene="GTP cyclohydrolase I"	
		/product="GTP cyclohydrolase I"	
		/protein_id="AAA41299.1"	
		/translation="MEKPRGVRCTNGFPERELPRPGASRPAEKSRPPEAKGQPADAWK	Translated protein sequence
		AGRPRSEEDNELNLPNLAAYSSILRSLGDPQRQGLLKPURAAATMQFFTKGYQETI	
		SDVLNDLIFDEHDHDMVIVKIDMFMSCEHHLVFPVGRVHIGYLPNKQVLGSLKLARIV	
		EIYSRRLQVQERLTKQIAVAITEALQPAAGVGVVIEATHMCMVMRGVQKMNKRTVSTML	
		GVFREDPKTREFLTLIRS"	
		Sequence 1016 BP; 236 A; 279 C; 291 G; 210 T; 0 other:	
		gacttgcgaac ctcattcggc gcagaactcc tgtcccgggt acagccacag gtcaggccgc	60
		ccggctaagc cgagccgcag cgtttgtag caacctaggg tgtctcggga gcaatcggc	120
		cggttcacg gagaagccgc ggggtgaag gtcaccaat gggttccccg agcgggagct	180
		...	
		catcaggagc tgaactccg tgtgcgagcc cgggtttgca gacccccgt gaggccagc	900
		ttatctgtct cgattgtaca ttccagttcc agttggtata ctgtccaact ttatttcca	960
		ccatgaattg tattaataa ttatttatag agatgtcaaa taaaggtgat caactt	1016



```
<buch>
  <isbn>0-201-318051</isbn>
  <titel>XML und Datenbanken</titel>
  <autor>Klettke/Meyer</autor>
</buch>
```





# Schematische und Daten-Heterogenität

9

## Variante 1

Männer	
Vorname	Nachname
Felix	Naumann
Jens	Bleiholder

Frauen	
Vorname	Nachname
Melanie	Weis
Jana	Bauckmann

## Variante 2

Personen			
Vorname	Nachname	Männl.	Weibl.
Felix	Naumann	Ja	Nein
Jens	Bleiholder	Ja	Nein
Melanie	Weis	Nein	Ja
Jana	Bauckmann	Nein	Ja

## Variante 3

Personen		
Vorname	Nachname	Geschlecht
Felix	Naumann	Männlich
Jens	Bleiholder	Männlich
Melanie	Weis	Weiblich
Jana	Bauckmann	Weiblich

# Schematische und Daten-Heterogenität

10

### Variante 1

Männer	
Vorname	Nachname
Felix	Naumann
Jens	Bleiholder

### Frauen

Vorname	Nachname
Melanie	Weis
Jana	Bauckmann

### Variante 2

Personen			
FirstNa	Name	male	femal
Felix	Naumann	Ja	Nein
Jnes	Bleiho.	Ja	Nein
Melanie	Weiß	Nein	Ja
Jana	baukman	Nein	Ja

### Variante 3

Personen		
VN	NN	SEX
F.	Naumann	Männlich
J.	Bleiholder	Männlich
M.	Weis	Weiblich
J.	Bauckmann	Weiblich

# Schematische und Daten-Heterogenität

11

**Variante 1**

☛☛☛☛☛	
☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛☛

☛☛☛☛☛	
☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛
☛☛☛☛	☛☛☛☛☛

**Variante 2**

☛☛☛☛☛			
☛☛☛☛	☛☛☛☛	☛☛☛☛	☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛	☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛	☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛	☛☛☛☛

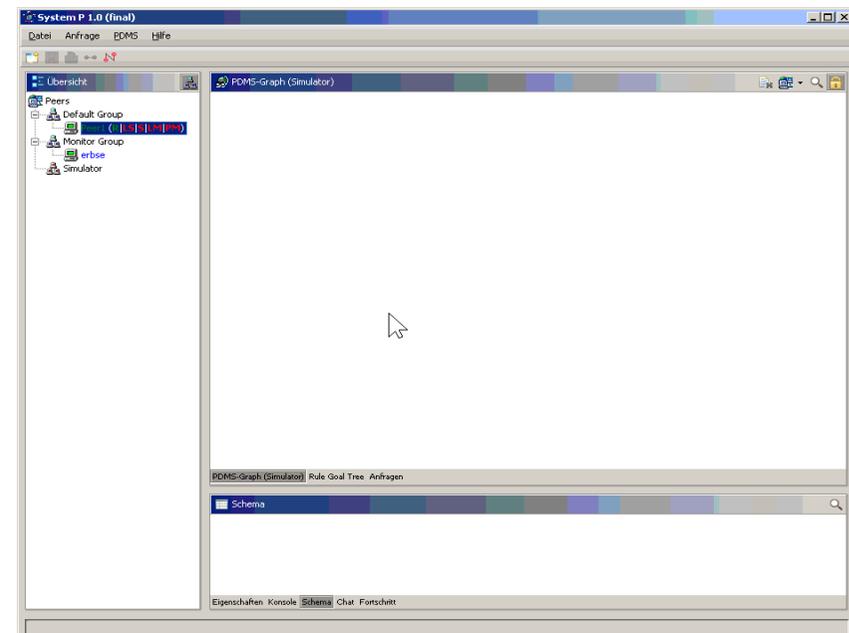
**Variante 3**

☛☛☛☛☛		
☛☛☛☛	☛☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛☛
☛☛☛☛	☛☛☛☛☛	☛☛☛☛☛

## Weitere Forschungsthemen

12

- Informationsintegration
  - ◇ Schema Matching
  - ◇ Duplikaterkennung
  - ◇ Datenfusion
- Datenqualität
- Peer Data Management
- Life Sciences: Aladin
- Search



# Lehrveranstaltungen in diesem Semester

13

## Vorlesungen

- DBS I
- Search engines

## Seminare

- Bachelor: Beauty is our Business
- Bachelor: Map/Reduce algorithms on Hadoop
- Master: Linked Data profiling
- Forschungsseminar

## Bachelorprojekt

- ETL Management



### Extending the Database Relational Model to Capture More Meaning

E. F. CODD  
IBM Research Laboratory

During the last three or four years several investigators have been exploring "semantic models" for formatted databases. The intent is to capture (in a more or less formal way) more of the meaning of the data so that database design can become more systematic and the database system itself can behave more intelligently. Two major thrusts are clear:

- (1) the search for meaningful units that are as small as possible—atomic semantics; atomics.
- (2) the search for meaningful units that are larger than the usual  $n$ -ary relation—molecular semantics.

In this paper we propose extensions to the relational model to support certain atomic and molecular semantics. These extensions represent a synthesis of many ideas from the published work in semantic modeling plus the introduction of new rules for insertion, update, and deletion, as well as new algebraic operators.

Key Words and Phrases: relation, relational database, relational model, relational schema, database, data model, database schema, data semantics, semantic model, knowledge representation, knowledge base, concept model, conceptual schema, entity model  
C/ Categories: 3.70, 3.73, 4.22, 4.28, 4.30, 4.34, 4.39

#### 1. INTRODUCTION

The relational model for formatted databases [5] was conceived ten years ago, primarily as a tool to free users from the frustrations of having to deal with the clutter of storage representation details. This implementation independence coupled with the power of the algebraic operators on  $n$ -ary relations and the open questions concerning dependencies (functional, multivalued, and join) within and between relations have stimulated research in database management (see [30]). The relational model has also provided an architectural focus for the design of databases and some general-purpose database management systems such as MACAIMS [13], PRTV [38], RDMS(GM) [41], MAGNUM [19], INGRES [37], QBE [46], and System R [2].

During the last few years numerous investigations have been aimed at capturing permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

A version of this work was presented at the 1979 International Conference on Management of Data (SIGMOD), Boston, Mass., May 30–June 1, 1979.  
Author's address: IBM Research Laboratory K01/282, 5609 Cottle Road, San Jose, CA 95193.  
© 1979 ACM 0362-5915/79/1390-0387 \$00.75.

ACM Transactions on Database Systems, Vol. 4, No. 4, December 1979, Pages 397–424.

# Überblick

14

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Ziele des Seminars
- Kurzvorstellung der Themen



# Vorstellung Teilnehmer

15

- Konstantin Käfer
- Martin Kreichgauer
- Tobias Mohr
- Daniel Moritz
- Tobias Pape
- Philipp Tessenow

# Seminarleistungen

16

## Lesen

- Paper lesen und verstehen
- Verwandte Literatur lesen und verstehen
- Mindestens eine individuelle Besprechung mit mir
- Kurzvorstellung der Literatur (je 2 min)

1/5 Note

## Vortragen

- Mindestens eine Folien-Besprechung mit mir
  - Spätestens 1 Woche vor Vortrag
- 30 min. Vortrag am jeweiligen Termin
  - + 15 min Diskussion
- Aktive Teilnahme an anderen Vorträgen
  - Jeweils ein *designated diabolus*

2/5 Note

## Ausarbeitung

- Mindestens eine Gliederungs-Besprechung mit mir
  - Spätestens 2 Wochen vor Abgabetermin
- 8-10-seitige Ausarbeitung bis zum 14.8.
- Unter Verwendung der LaTeX-Vorlage im WWW

2/5 Note

**Aktive Teilnahme an sämtlichen gemeinsamen Terminen**

# Feedback

17

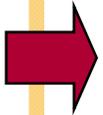
Fragen bitte jederzeit!

- Während des Seminars
- Während individueller Besprechungen
  - Termin bei Katrin Heinrich ausmachen
    - ◇ office-naumann@hpi.uni-potsdam.de
    - ◇ 0331 / 5509 280
  - Muss nicht zur Sprechstunde sein!
- Sprechstunde
  - Dienstags 15:00 – 16:00
  - Raum A-1.13
  - Am liebsten mit Anmeldung
- Email: [naumann@hpi.uni-potsdam.de](mailto:naumann@hpi.uni-potsdam.de)

# Überblick

18

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Ziele des Seminars
- Kurzvorstellung der Themen



# Wissenschaftliche Texte lesen

19

- Fachartikel
  - Journale, Konferenzen und Workshops
  - Entstehung:
    - ◇ Forschungsvorhaben
    - ◇ Begutachtungsprozess
- Struktur eines Artikels
  - Kritisches Lesen
  - Experimente
- Literatur
  - (Online-) Recherche
- Englisch

# Wissenschaftlichen Vortrag halten

20

- Gliederung
  - Die Kunst des Weglassens
- Foliengestaltung
  - Powerpoint
  - Overhead
- Zeit einhalten
- Techniken zur Vorbereitung
- Techniken während des Vortrags



# Ausarbeitung schreiben

21

- Gliederung
- Schreibstil
  - Objektives
  - Subjektives
- Plagiate
- LaTeX

# Zeitlicher Ablauf – Termine

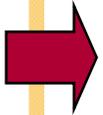
22

- Part 1: Lesen
    - Termin zur Literaturkritik
    - *Elevatorpitch*
  - Part 2: Vortragen (Zwei pro Termin)
    - Jeweils 30 Minuten Vortrag + Diskussion
  - Part 3: Schreiben
    - Ausarbeitung / Seminararbeit
- Vortragsreihenfolge
    - Siehe Web-Page
  - Dazwischen Vorträge von mir
    - Lesen wissenschaftlicher Texte
    - Vortragstechniken
    - LaTeX (bei Bedarf)
    - Tipps zur Gliederung und Ausarbeitung

# Überblick

23

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Ziele des Seminars
- Kurzvorstellung der Themen



# VLDB 10 Year Best Paper Awards

24

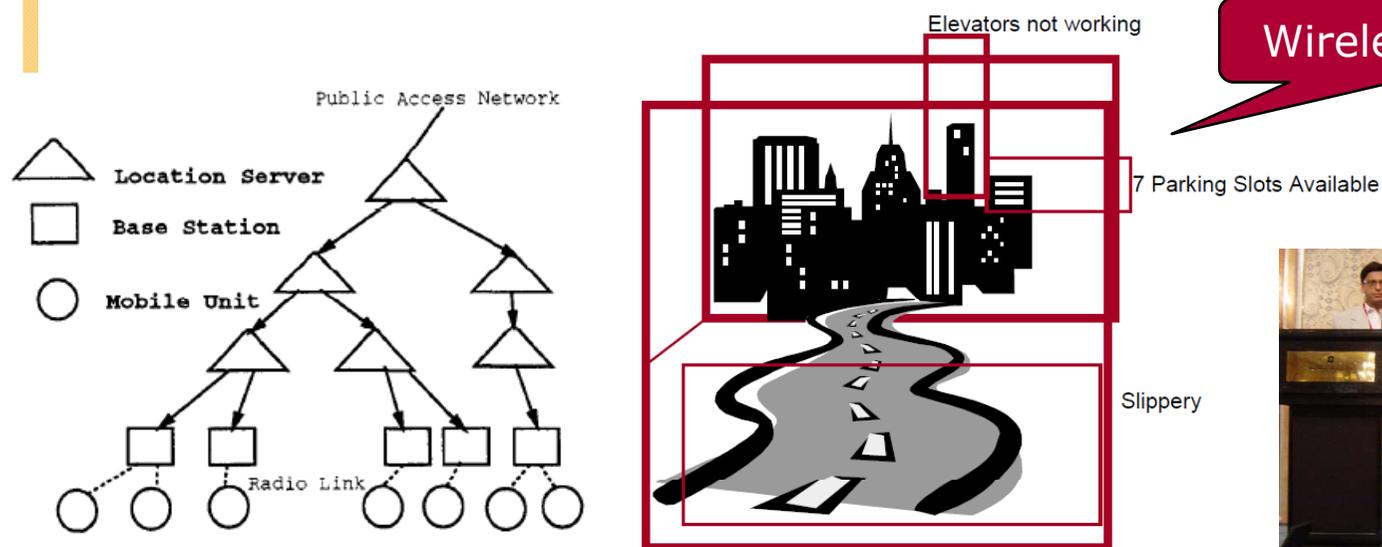
- Very Large Databases (VLDB)
  - Endowment: <http://www.vldb.org/>
  - VLDB Archives:  
<http://www.vldb.org/archives/public/10year/10year.html>
- This award is established for the author(s), whose paper appeared in the VLDB conference 10 years ago, that has the **most impact on database research** since then.
- Neues Auswahlgremium jedes Jahr
- Nominiert sind alle Paper der 10 Jahre zuvor stattgefundenen VLDB Konferenz

# 2002: Mobility

Querying in highly mobile distributed environments

26

- Neue Hardware: Wireless personal devices
- Probleme
  - Location update volume: „Do not tell me - if I do not want to know“
  - Distributed location-queries: “I do not know but can find out“

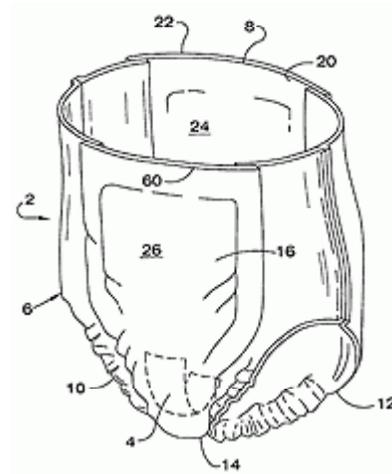


# 2004: Data Mining

## Fast Algorithms for Mining Association Rules

28

- Bahnbrechendes Papier
- Setzte intensive Forschung zu Data Mining in Gang
- Zwei Algorithmen zur schnellen Entdeckung von Assoziationsregeln
  - Apriori
  - AprioriTid



Rakesh Agrawal

# 2005: Web-Anfragen

W3QS: A Query System for the World-Wide Web

29

- Bisher: Suche basierend auf Indizes und Navigation mittels Browser
- Aber: WWW wächst ständig; Daten sind heterogen; strukturelle Anfragen unmöglich
- Idee: WWW aus DB-Perspektive
  - Graph-basiertes Datenmodell
  - SQL-ähnliche Anfragespache
  - Content-queries
  - Structure-queries
- Features
  - Nutzt vorhandene Indizes (Suchmaschinen)
  - Füllt Formulare aus (Deep-Web)
  - Deklarativ, deshalb optimierbar



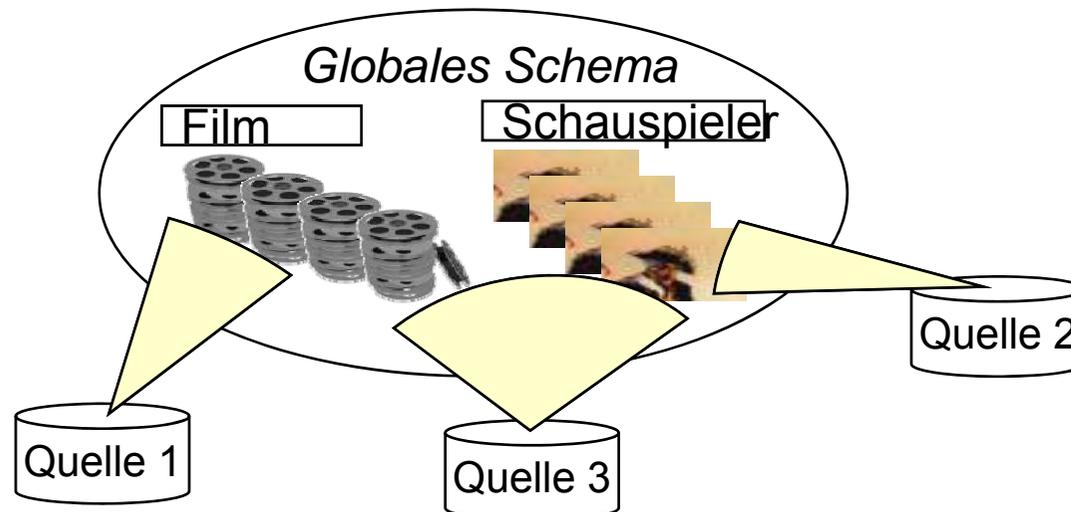
Oded Shmueli

# 2006: Information Manifold

Querying Heterogeneous Information Sources Using Source Descriptions

30

- Virtuelle Integration heterogener Datenquellen
- Modellierung der Fähigkeiten der Datenquellen als VIEW
  - Local-as-View (LaV)
  - Einfacher, da Quelle-orientiert; andere Quellen egal
  - Präziser: Einschränkungen der Quelle können leicht modelliert werden
- Anfragebearbeitung: Bucket-Algorithmus



Alon Halevy

# 2007: Self-Tuning Databases

An Efficient, Cost-Driven Index Selection Tool for Microsoft SQL Server

31

- Gegeben eine Workload ist die optimale Wahl geeigneter Indizes sehr schwierig
  - Job des DBA („an artform“)
- Self-CHOP
  - Configuring – Healing – Optimizing – Protecting
- Suchraum ist enorm
  - Instrumentierung des Optimierers
- AutoAdmin Projekt bei Microsoft immer noch aktiv
- Evaluation mittels TPC Benchmarks



Surajit Chaudhuri

# 2009: Top-K selection queries

## Evaluating Top- $k$ Selection Queries

33

- SELECT \* FROM R  
WHERE  $A_1=v_1$  AND ... AND  $A_n=v_n$   
ORDER  $k$  BY Score
- SELECT \* FROM R  
WHERE  $(a_1 \leq A_1 \leq b_1)$  AND ... AND  $(a_n \leq A_n \leq b_n)$

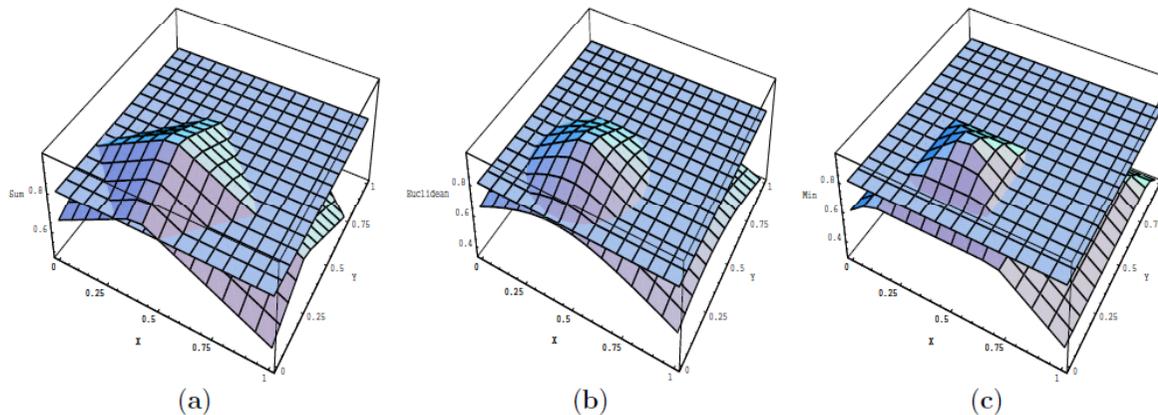


Figure 1: The scores ( $z$  axis) for query  $q = (0.4, 0.3)$  for the different  $(x, y)$  pairs and scoring functions *Sum* (a), *Euclidean* (b), and *Min* (c).



Surajit Chaudhuri

# Nächste Schritte

34

- Sie: E-Mail an mich
  - Mit Wunschliste (Platz 1 – 3)
    - ◇ Daraus: Themenvergabe
    - ◇ Daraus: Vergabe des Diskutantenplatzes
  - Bis Mitternacht
- Ich: Themenvergabe bald
  - E-Mail Benachrichtigung
  - [www.vldb.org](http://www.vldb.org)
  - <http://www.informatik.uni-trier.de/~ley/db/index.html>
- Ich: Vortrag zum „Lesen wissenschaftlicher Texte“