



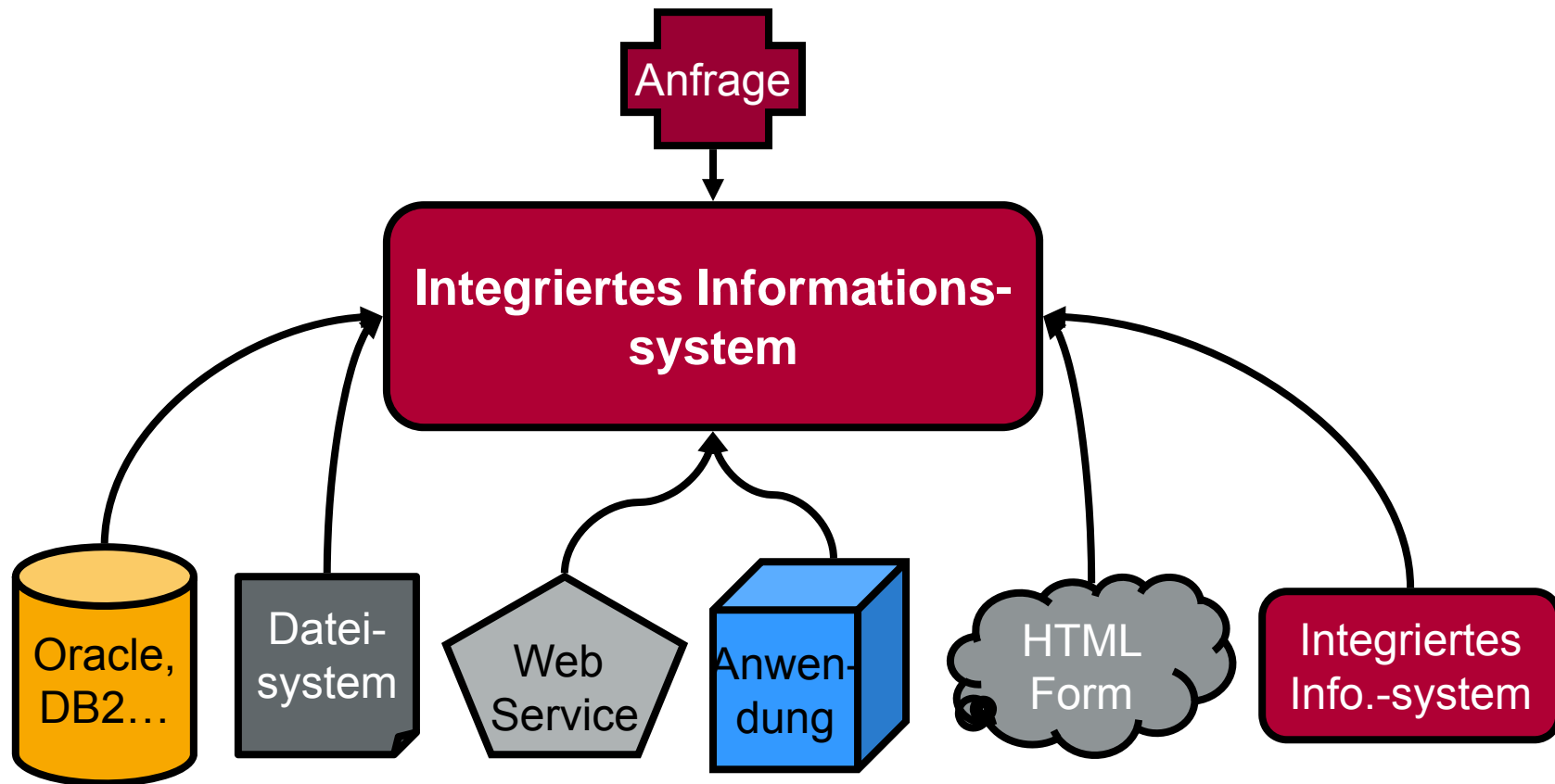
**Hasso  
Plattner  
Institut**

IT Systems Engineering | Universität Potsdam

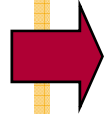
## Informationsintegration Einführung

14.4.2008

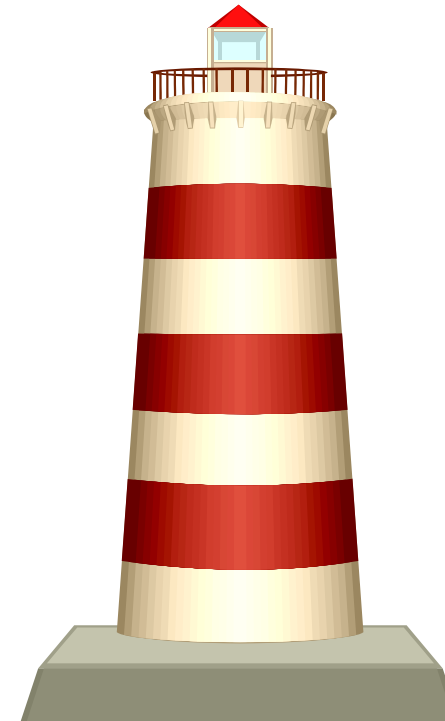
Felix Naumann



3



- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Informationssysteme
- Informationsintegration am Beispiel
- Ausblick auf das Semester



# Arbeitsgruppe Informationssysteme

4

project **ViCTOR**



Paul Führung



Patricia Hobro

**DQ Assessment**



Prof. Felix Naumann

**Information Integration**

**Information Quality**



Jens Bleiholder

**Data Fusion**

project **HumMer**

**Duplicate Detection**



Karsten Draba



Melanie Weis & Sascha Szott

**Data Cleaning**



Armin Roth

**Peer Data Management Systems**

**Matching**

**Data Integration for Life Science Data Sources**

**Duplicate Detection**

project **XClean**

project **System P**



Mohammed AbuJarour

**Service-Oriented Systems**

**Ontologies**



Frank Kaufer

project **Aladin**



Jana Bauckmann



Alexander Albrecht

**Personal Information Management**

**Data Profiling for Schema Management**

# Lehrveranstaltungen in diesem Semester

5

## Vorlesungen

- DBS II
- Informationsintegrati

## Seminare

- Bachelor: Beauty is our Business
- Bachelor: [www.ligageschichte.de](http://www.ligageschichte.de)
- Master: Duplikaterkennung
- Forschungsseminar



## 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;
- (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

CR Categories: 3.70, 3.73, 4.22, 4.28, 4.35, 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 cluster 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.

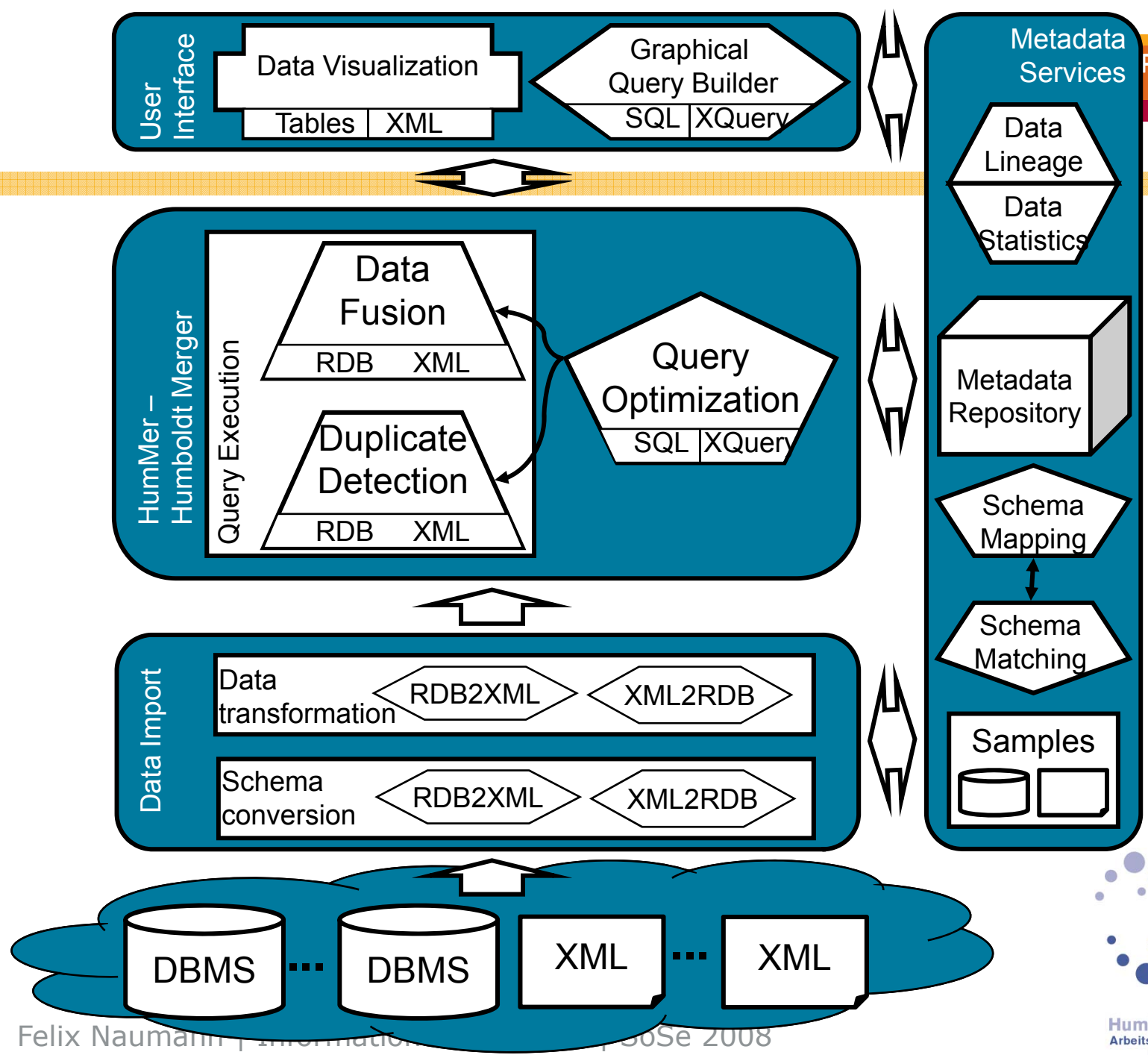
A version of this work was presented at the 1979 International Conference on Management of Data (SIGMOD), Boston, Mass., May 30–June 1, 1979.

Author's address: IBM Research Laboratory K01/282, 5609 Cottle Road, San Jose, CA 95193.

© 1979 ACM 0362-5915/79/1390-0387 \$00.75.

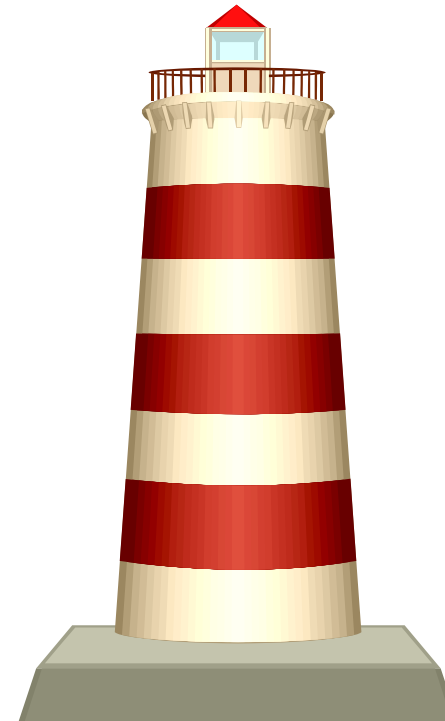
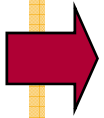
ACM Transactions on Database Systems, Vol. 4, No. 4, December 1979, Pages 397–424.

6



7

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Informationssysteme
- Informationsintegration am Beispiel
- Ausblick auf das Semester



# Termine und Leistungserfassung

8

- Vorlesung
  - Montags 15:15 – 16:45
  - Donnerstags 13:30 – 15:00
- Praktikum / Übung
  - Ausgewählte Termine – WWW beachten
- Erste Vorlesung
  - 14.4.2008
- Letzte Vorlesung
  - 17.7.2008
- Feiertage
  - 1.5. Maifeiertag (Do)
  - 12.5. Pfingstmontag
- Sondertermine
  - 17.4. (Do): IBM Tag
  - 21.4.: Erste Übung
  - 24.4.: fällt leider aus ☹
- Prüfung
  - Mündlich, 30 Minuten
  - Erste Woche nach Vorlesungszeitraum
- 7 Übungstermine
  - Leitung: Alexander Albrecht
  - Theoretische und praktische Aufgaben
  - 2er Teams
- Voraussetzungen
  - Zur Teilnahme
    - ◇ Datenbankkenntnisse (DBS I)
  - Zur Prüfung
    - ◇ Besuch der Vorlesung
    - ◇ Aktive Teilnahme an den Praktikumsterminen
    - ◇ „Bestehen“ der Übungsblätter



# Feedback

9

- Evaluation am Ende des Semesters
- Fragen bitte jederzeit!
  - In der VL
  - Sprechstunde: Dienstags 15-16
  - Email: [naumann@hpi.uni-potsdam.de](mailto:naumann@hpi.uni-potsdam.de)
- Anregungen zur Verbesserung:
  - Z.B. zu
    - ◇ Gebrauch der Folien
    - ◇ Infos im WWW
  - Jeweils nach der VL oder in der Sprechstunde
  - Oder per Email: [naumann@hpi.uni-potsdam.de](mailto:naumann@hpi.uni-potsdam.de)

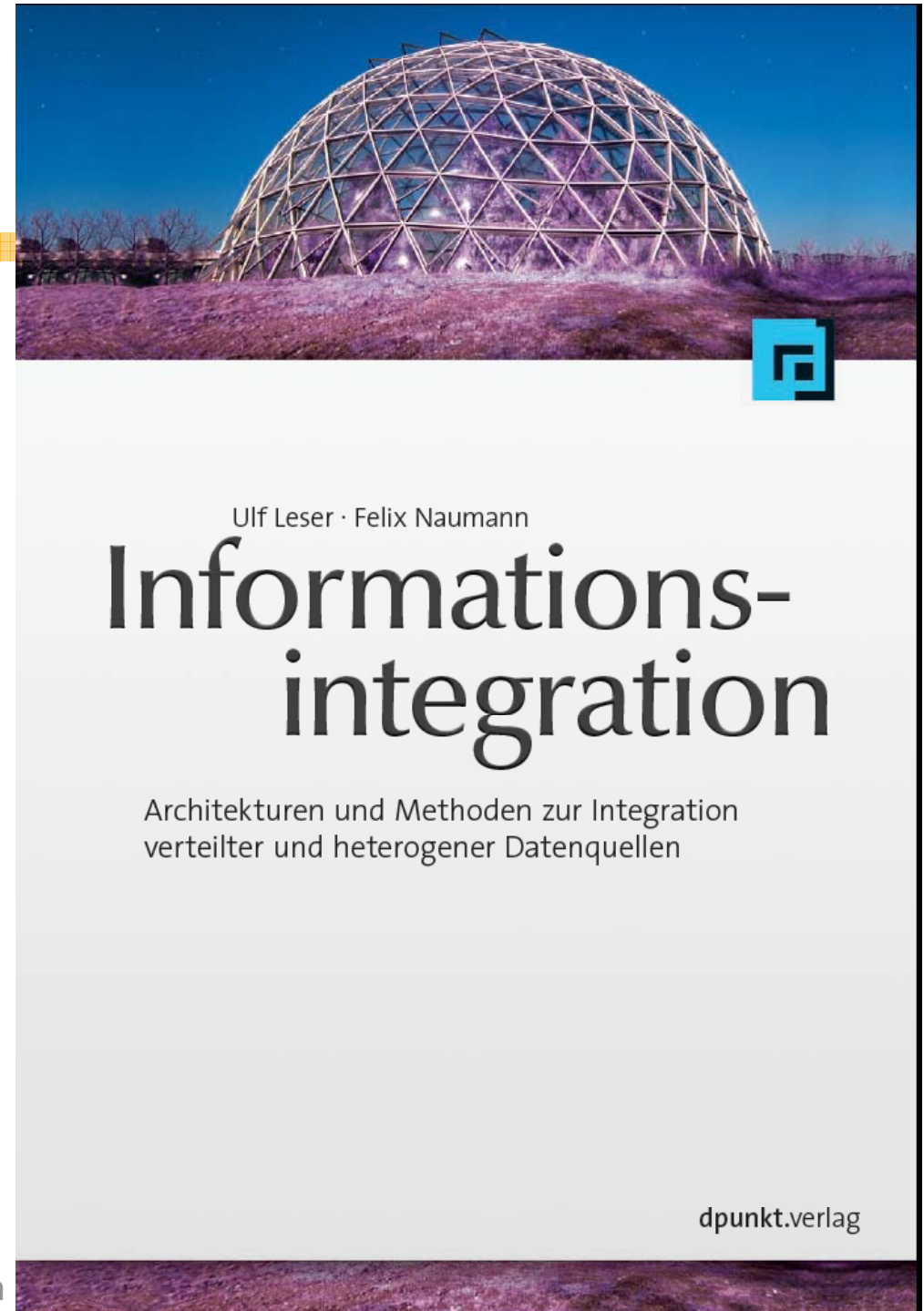
# Lehrbuch

10

- Informationsintegration
- Ulf Leser und Felix Naumann
  - dpunkt Verlag
- 42 Euro
- X-mal in Bibliothek
  
- Fehler gefunden =>



Felix Naumann | Informationsintegration



# Weitere Literatur

11

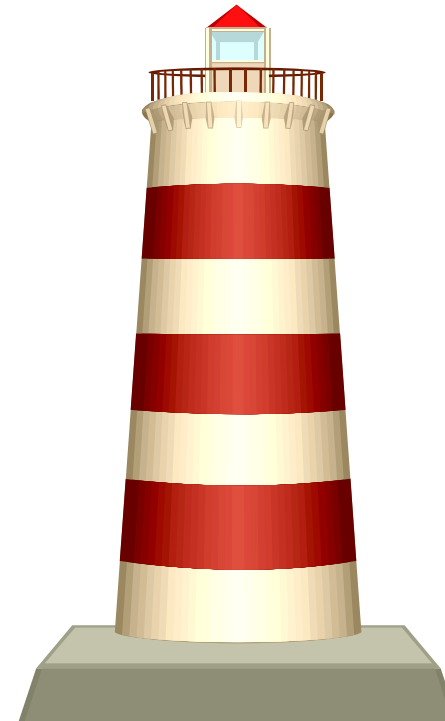
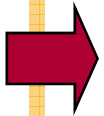
- Themen u.a. aus
  - Föderierte Datenbanksysteme. Konzepte der Datenintegration, Stefan Conrad, ISBN: 3540631763
  - Principles of Distributed Database Systems  
M. Tamer Özsu, Patrick Valduriez  
ISBN: 0136597076
- Jeweils Literaturhinweise in den Vorlesungen
- Alle genannten Artikel können von mir per Email angefragt werden. Oder:
  - Google Scholar: <http://scholar.google.com/>
  - DBLP: <http://www.informatik.uni-trier.de/~ley/db/index.html>
  - CiteSeer: <http://citeseer.ist.psu.edu/>
  - ACM Digital Library: [www.acm.org/dl/](http://www.acm.org/dl/)
  - Homepages der Autoren

# Vorstellung – Hörer

12

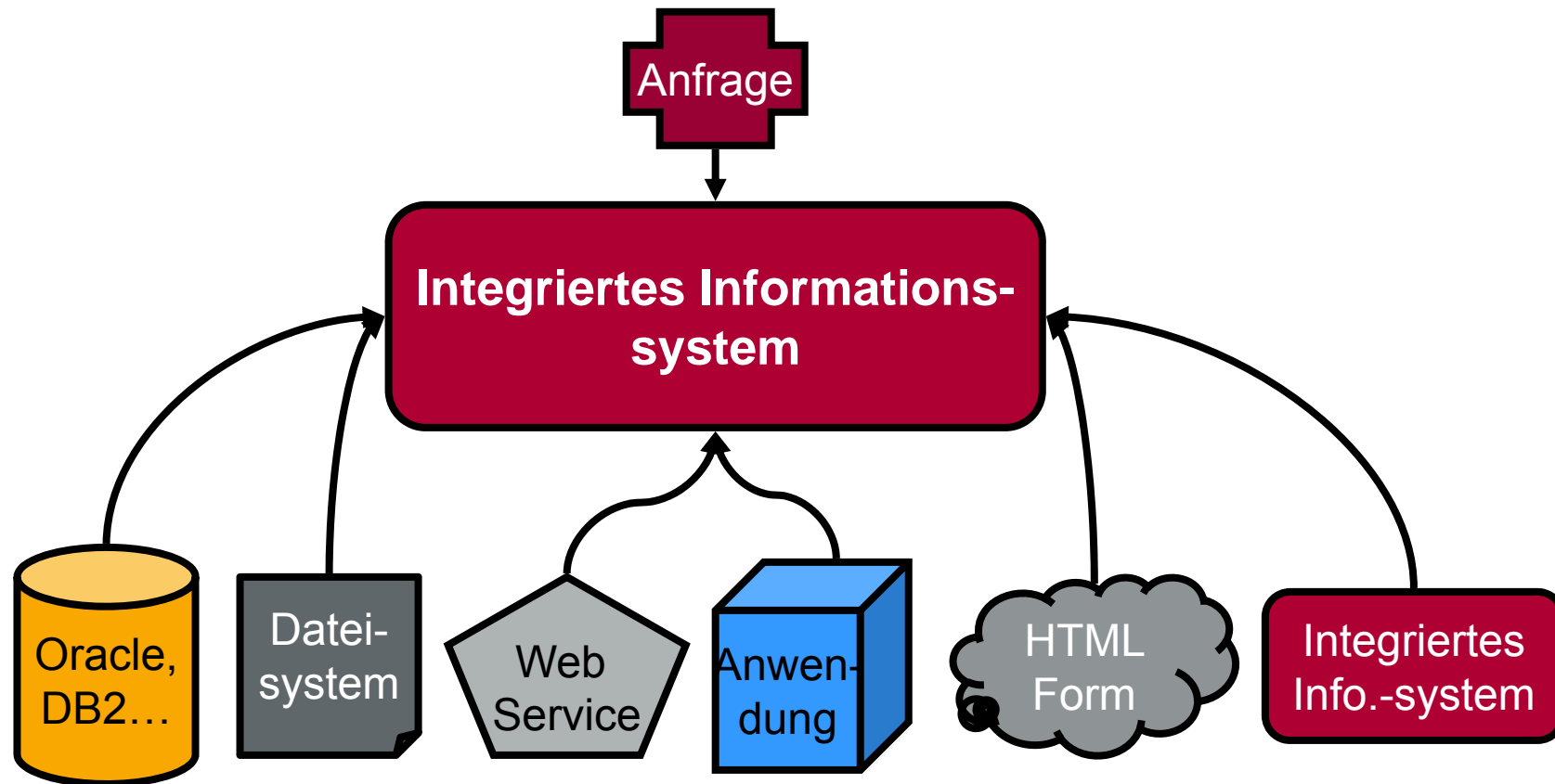
- Welches Semester?
- HPI oder IfI?
- Erasmus o.ä.?
- Datenbankkenntnisse?
  - Andere relevante Lehrveranstaltungen?
- Workshop/Seminar Duplikaterkennung?
- Ihre Motivation?
  - Schon mal integriert?
  - DWH?

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Informationssysteme
- Informationsintegration am Beispiel
- Ausblick auf das Semester



# Integrierte Informationssysteme

14



# Was ist Informationsintegration?

15

Informationsintegration ist die Zusammenführung von Daten und Inhalt verschiedener Quellen zu einer einheitlichen Informationsmenge.

Informationsintegration ist die **korrekte, vollständige** und **effiziente** Zusammenführung von Daten und Inhalt verschiedener, **heterogener** Quellen zu einer einheitlichen und **strukturierten** Informationsmenge zur effektiven **Interpretation** durch Nutzer und Anwendungen.

# Wo herrscht Informationsintegration?

16

- Im weiteren Sinne
  - Business-Integration
  - Application-Integration
  - Prozess-Integration (Workflow-Integration)
- Im engeren Sinne
  - Datenbanken und Informationssysteme
    - ◇ Verteilt
    - ◇ Autonom
    - ◇ Heterogen



# Beispiele für Informationssysteme

17

- Dateisystem
  - WWW (HTML Dateien)
  - Desktop-Anwendungen (Textverarbeitung, etc.)
- Datei (Flat file)
  - Zeile /Token
  - Anfrage: Parser, File search, RegEx
  - Strukturiert / semi-strukturiert / unstrukturiert
  - Spezialfall: Markup-Datei
    - ◇ XML, HTML mit Struktur
    - ◇ Anfragespreache
  - Beispiele
    - ◇ Komma-delimited files
    - ◇ Annotated files
  - Einsatzgebiete
    - ◇ SwissProt

ID	RNGTPCHI	standard; RNA; ROD; 1016 BP.	<b>Molecule type</b>
XX			<b>Name</b>
DT	01-AUG-1991	(Rel. 28, Created)	<b>Date of creation and last update</b>
DT	04-MAR-2000	(Rel. 63, Last updated, Version 2)	
XX			
DE	Rat GTP cyclohydrolase I mRNA, complete cds.		<b>Free text description</b>
XX			
KW	GTP cyclohydrolase I.		<b>Keywords describing the molecule</b>
XX			
OS	Rattus norvegicus (Norway rat)		<b>Organism</b>
OC	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;		
OC	Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus.		
XX			
RN	[1]		<b>Article the sequence was published in</b>
RP	1-1016		
RX	MEDLINE; <a href="#">91093270</a> .		
RX	PUBMED; <a href="#">1985963</a> .		
RA	Hatakeyama K., Inoue Y., Harada T., Kagamiyama H.;		
RT	"Cloning and sequencing of cDNA encoding rat GTP cyclohydrolase I: The first enzyme of the tetrahydrobiopterin biosynthetic pathway";		
RL	J. Biol. Chem. 266(2):765-769(1991).		
XX			
FT	<a href="#">CDS</a>	128..853	<b>Structural annotation (coding sequence)</b>
FT		/codon_start=1	
FT		/db_xref="GOA: <a href="#">P22288</a> "	<b>Link to functional annotation of resulting protein</b>
FT		/db_xref="SWISS-PROT: <a href="#">P22288</a> "	
FT		/EC_number="3.5.4.16"	
FT		/gene="GTP cyclohydrolase I"	
FT		/product="GTP cyclohydrolase I"	
FT		/protein_id=" <a href="#">AAA41299.1</a> "	
FT		/translation="MEKPRGVRCINGFPERELPRPGASRPAEKSRPPEAKGAQPADAWK	<b>Translated protein sequence</b>
FT		AGRPRSEEDNELNLPNLAAAYSSILRSLGEDPQRQGLLKTPWRAATAMQFFTKGYQETI	
FT		SDVNLNDAIFDEHDHEMVIVKIDMFMSCEHHLVPFVGRVHIGYLPNKQVLGLSKLARIV	
FT		EIYSRRLQVQERLTKQIAVAITEALQPAGVGVVIEATHMCMVMRGVQKMNSKTVTSTML	
FT		GVFREDPKTREEFLTLIRS"	
FT			
SQ	Sequence 1016 BP; 236 A; 279 C; 291 G; 210 T; 0 other;		<b>Sequence of bases</b>
	gacttcgaac ctcattcggg gcagaactcc tgtcccgggtg acagccacag gtcacgggccc	60	
	ccggctaagc cgagccgcag cgcttggttag caccttaggg tgtctcggga gcaatcgcgc	120	
	cggttccatg gagaagccgc ggggtgtaag gtgcaccaat gggttccccg agcggggagct	180	
	...		
	catcaggagc tgaacttcg tgtgcgagcc cgggtttgca gacccccgct gaggccagcg	900	
	ttatctgtct cgattgtaca ttccagttcc agttggtata cttgtcaact ttatttctca	960	
	ccatgaattg tattaataa ttatttatag agatgtcaaa taaaggtgat caactt	1016	

//

# Beispiele für Informationssysteme

19

- Datenbank
  - Anfragesprachen, z.B. SQL
  - Relational, OO, Hierarchisch
  - XML DBMS
  - Data Warehouses
  - OLTP
- HTML Formular
  - Anfrage: Einfache Suche, komplexe Formulare
    - ◇ Radiobutton, dropdown-list, etc.
  - Struktur des Ergebnisses: wie Markup Datei: Flach, hierarchisch oder graph-basiert

Alle Kategorien ansehen

Suche

Bücher

LOS

Einkaufswagen

Bücher

Erweiterte Suche

Stöbern

Bestseller

Neuheiten

Hörbücher

Taschenbücher

Fachbücher

Sonderangebote

Bücher Verkaufen

Amazon.de  
KreditkarteSofort einsetzen **20€ sparen** und Amazon.de Punkte sammeln! [Jetzt 20€ sparen](#)

## Erweiterte Suche Bücher

Je mehr Felder Sie ausfüllen, desto zielgerichteter können wir suchen. Es reicht jedoch aus, nur eines der Felder auszufüllen.

Autor/in: Titel: Schlagwörter: ISBN: (10- oder 13-stellig, ohne Bindestriche) Verlag: 

Verfeinern Sie Ihre Suche, indem Sie nur nach bestimmten Buchformaten suchen.

Nur gebraucht: Format: Ordnen nach: Erscheinungsdatum: Suche in:  Deutsche Bücher  Englische Bücher

## Beispiele

- ◆ Geben Sie "fidelity" in das Titel- und "hornby" in das Autoren-Feld ein, erscheint *High Fidelity* von Nick Hornby.
- ◆ Indem Sie "grisham, john" in das Autor-Feld und "Taschenbücher" im Feld "Format" auswählen, finden Sie nur die Taschenbuchausgaben der Romane von John Grisham.

# Beispiele für Informationssysteme

21

- Services
  - Z.B. Web Services
    - ◇ XML Dokument
    - ◇ Anfrage als XML Dokument
  - Einsatzgebiete
    - ◇ Intra-organisatorische Workflows
    - ◇ E-Marketplaces
    - ◇ Datenaustausch
    - ◇ Mashups

**Recent Listings** [ [View the FULL LIST](#) ]

Publisher	Style		Service Name	Description	Implementation
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">Gale Group Web Domain Business Intelligence</a>	Based on a website domain get in-depth financial and corporate information	
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">Gale Group Business Intelligence</a>	Based on company name get in-depth financial and corporate information	
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">Gale Group Business Information</a>	Standard financial and corporate information for 440,000 U.S. and international companies	
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">IPligence Geo IP Location</a>	Comprehensive information on IP addresses	
mharvanek	DOC	<a href="#">Try It</a>	<a href="#">DuoShare Address Quality Integrator</a>	DuoShare Address Quality Integrator standardizes, corrects, validates (to the delivery point), and enhances U.S. addresses using a number of USPS® certified processes.	
CDYNE	DOC	<a href="#">Try It</a>	<a href="#">CDYNE Phone Notify!</a>	Calls any phone number and speaks text or sound file to the person. Also supports advanced Dialplans, Machine Detection and Incoming calls.	
CDYNE	DOC	<a href="#">Try It</a>	<a href="#">CDYNE SMS Notify!</a>	Send SMS Text Messages to your mobile telephone devices without the use developing your own software!	
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">The Sports Network Men's College Basketball</a>	Comprehensive sports data for NCAA Mens Division I Basketball	
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">Wall Street Horizon Interest Rate Calendar</a>	Access current and historical interest rate information	
EDIXMLOnline	DOC	<a href="#">Try It</a>	<a href="#">EDIXMLOnline</a>	Translate, generate & validate EDI files using X12/EDIFACT SEF files. Translate EDI to XML format.	MS .NET
StrikeIron	DOC	<a href="#">Try It</a>	<a href="#">Wall Street Horizon Real-Time Company Earnings</a>	Information to analyze and evaluate investments and future earning potential	
phi1281	RPC	<a href="#">Try It</a>	<a href="#">Weight Watchers Point Calculator</a>	Given the parameters of calories, fat grams, & fiber grams, computes the Weight Watchers point value for a given serving of food.	ColdFusion
SOATeader	RPC	<a href="#">Try It</a>	<a href="#">Captcha Web Service</a>	This Web service will create captcha	



- Business Users
- Enterprise IT Professionals
- Application Developers
- Partners

**STRIKEIRON:** Strikelron's Data Services give you access to live data to use now, integrate into applications or build into Web sites.

SEARCH WEB SERVICES

Search

All Categories

Browse

**LATEST NEWS**

**Strikelron Offers Gale Data on**  
Salesforce.com's AppExchange

**Strikelron Hits Home Run**  
with Major League Baseball Data from The Sports Network

**Strikelron Data Fuels**  
Enterprise Mashup Platforms



**EDITORS' CHOICE AWARD 2008**  
COMPANIES TO WATCH

**WHAT DO YOU WANT TO DO?**

**CONNECT**

- Find Data Services
- Browse Solutions
- Cleanse & Enhance Data

**CREATE**

- Build Mashups in Excel
- Create Business Solutions
- Integrate into Salesforce.com

**COLLABORATE**

- **New!** Developer Community
- Download Code and Clients
- Access Free Data

**DATA SPOTLIGHT**

**NEW WEB SERVICES**


- D&B Business Prospect
- Gale Business Intelligence
- Midnight Trader Financial News
- MapQuest Driving Directions
- TSN Men's College Basketball
- Wall Street Horizon Company Earnings

**HOT WEB SERVICES**

- US Address Verification
- Global SMS Pro
- Sales and Use Tax
- Email Verification
- Foreign Exchange Rates

**PARTNERS**

Strikelron is proud to partner with the following industry leaders:



Enabling the On Demand World



726 APIs

2946 Mashups

**Vertical Markets**

Social Platforms

Mapping

Video

Shopping

Government

Mobile+Telephony

**First time here?**  
take a tour

**Popular Directory Searches**

Celebrity Mashups

Video Mashups

Popular New Mashups

All Popular Mashups

Maps Mashups

Keeping you up to date with mashups, APIs and the Web as platform. [Learn more >>](#)

Find	Track	Do	Share
API Directory	Industry News	Learn	Sign-up
Mashup Directory	Market Trends	Search Code	Members
Most Popular	Major Players	Win Contests	Add Links
Mashup Matrix	Tag Cloud	Discuss	Talk to us

**Latest News Updates**

[more news >>](#)

**Google Chart API's New Schematic Maps >>**

Graphic, schematic maps can very useful in a wide range of charting scenarios and thanks to a recent upgrade, now you can use the [Google Chart API](#) to create them. >



Posted April 11, 2008. [Continue reading >](#)

**April 12, 2008**

**New Mashups >>**

- > 2008 Beijing Olympics Torch Relay Path
- > Svraka
- > Enomalism Elastic Computing
- > WebFOCUS for Google Maps

**Mashup of the Day >>**

2008 Beijing Olympics Torch Relay Path



**New APIs >>**

- > ThisNext
- > PilotOutlook
- > Summize Twitter Search
- > Enomalism
- > Mortgage Marvel

ProgrammableWeb Sponsors

**Connect. Collect. Mashup. Everything!**

Get apps. Get paid. Userplane Money

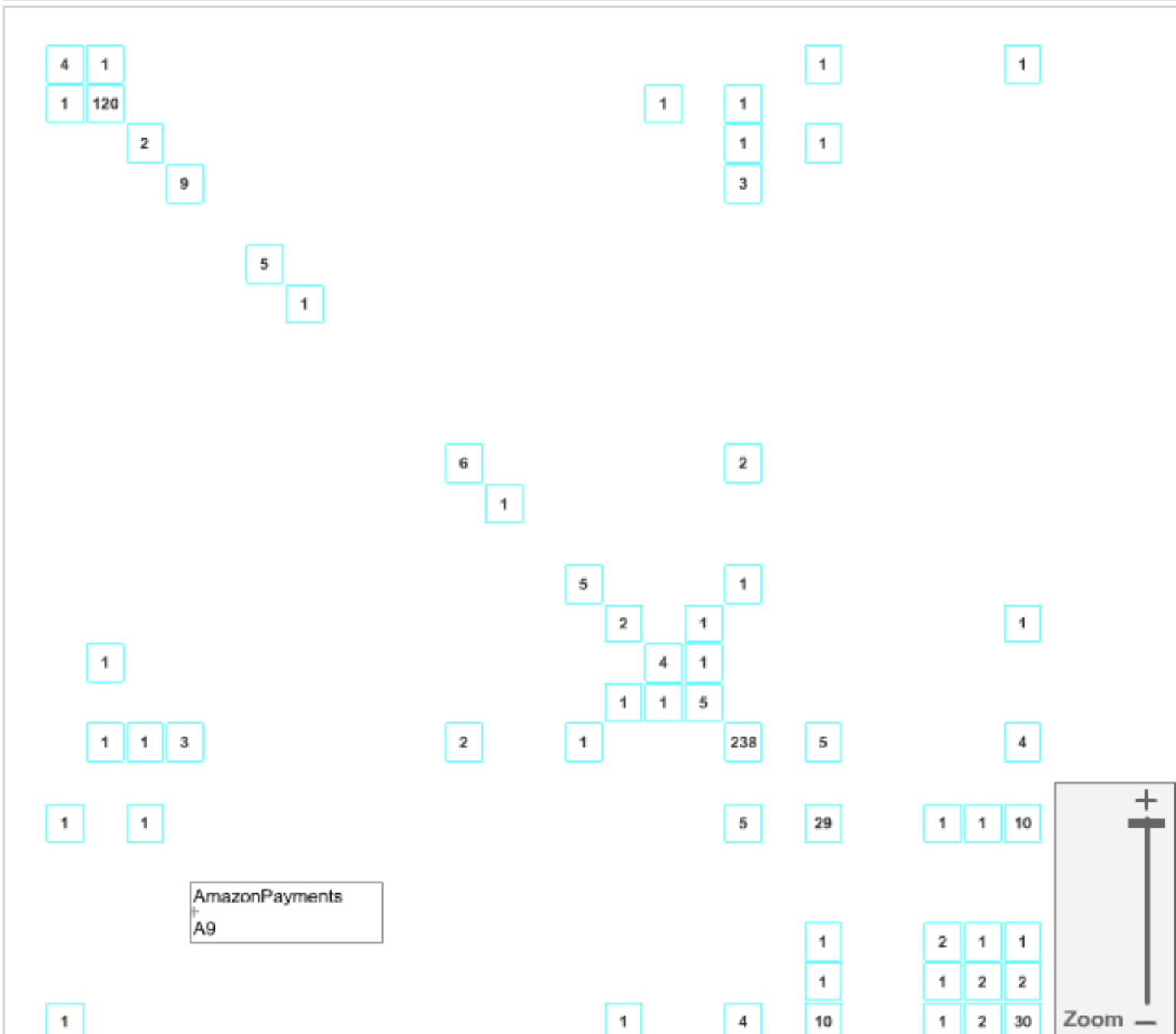
**Web21CSDK**  
Do Less : Achieve More  
[web21c.bt.com](http://web21c.bt.com)

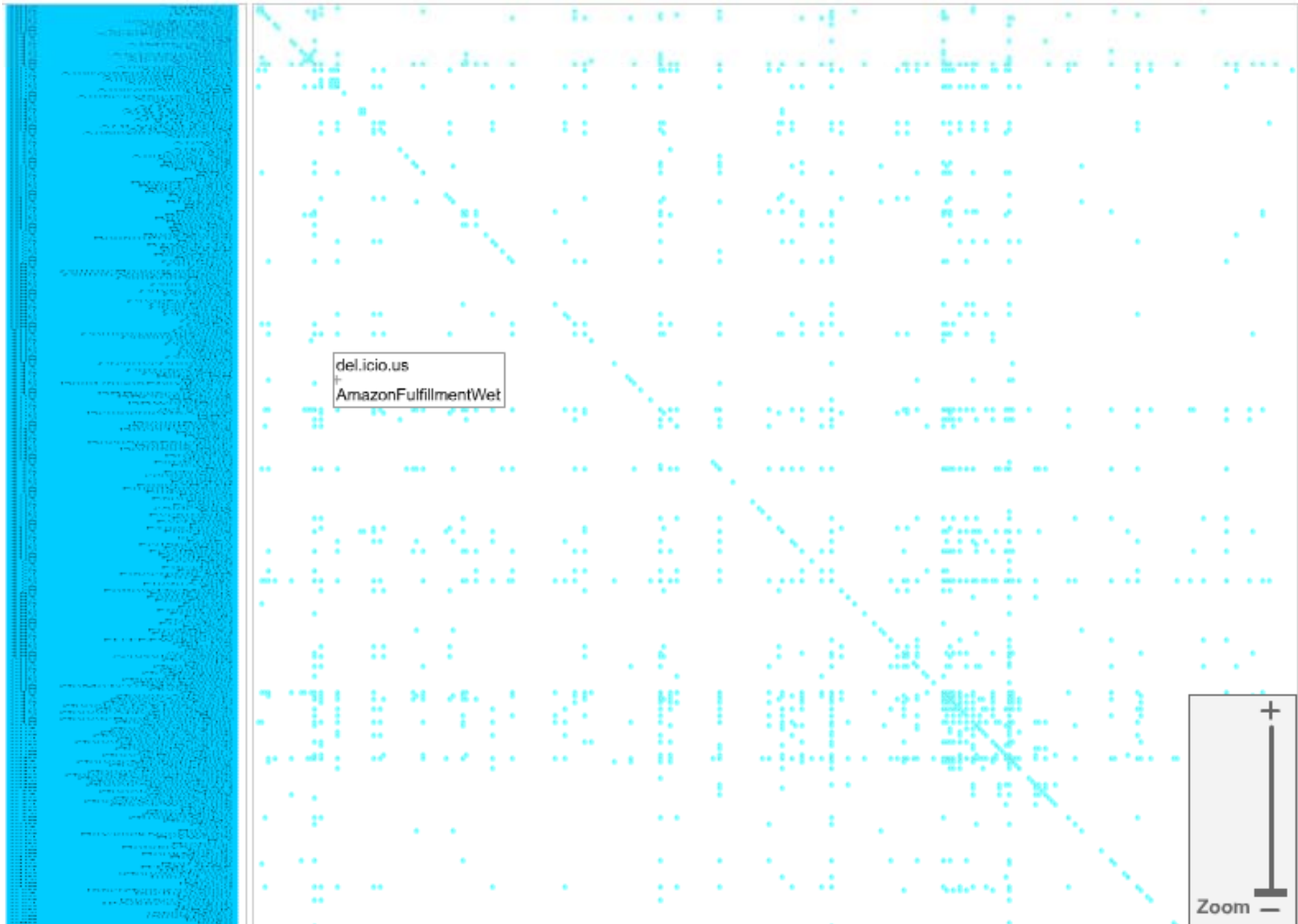
**Top Tags**

- mapping (1602)
- photo (438)
- shopping (389)
- search (348)
- video (293)
- travel (276)
- news (188)
- sports (176)



001	23
002	30Boxes
003	411Sync
004	43Things
005	A9
006	activeRenderer
007	AdobeOnAir
008	AdobeShare
009	AevumObscurum
010	AgentFactorTravel
011	AgentRank
012	AIM
013	AIMPhoneline
014	Akismet
015	Alexa
016	AlexaThumbnail
017	AlexaTopSites
018	AlexaWebInfo
019	Amazon
020	AmazonDevPay
021	AmazonEC2
022	AmazonFulfillmentWeb
023	AmazonHistorical
024	AmazonPayments
025	AmazonQueue
026	AmazonS3

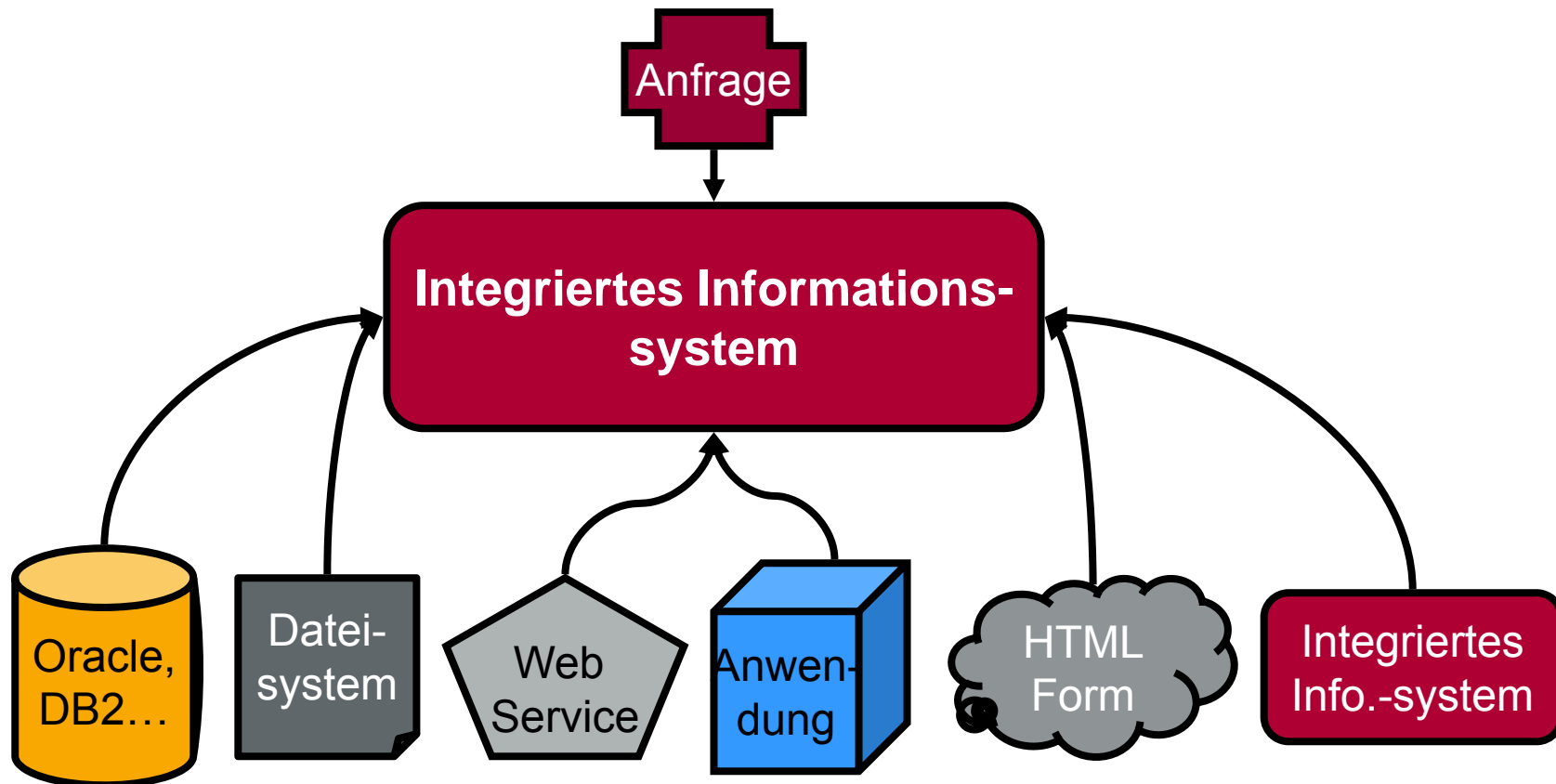




# Beispiele für Informationssysteme

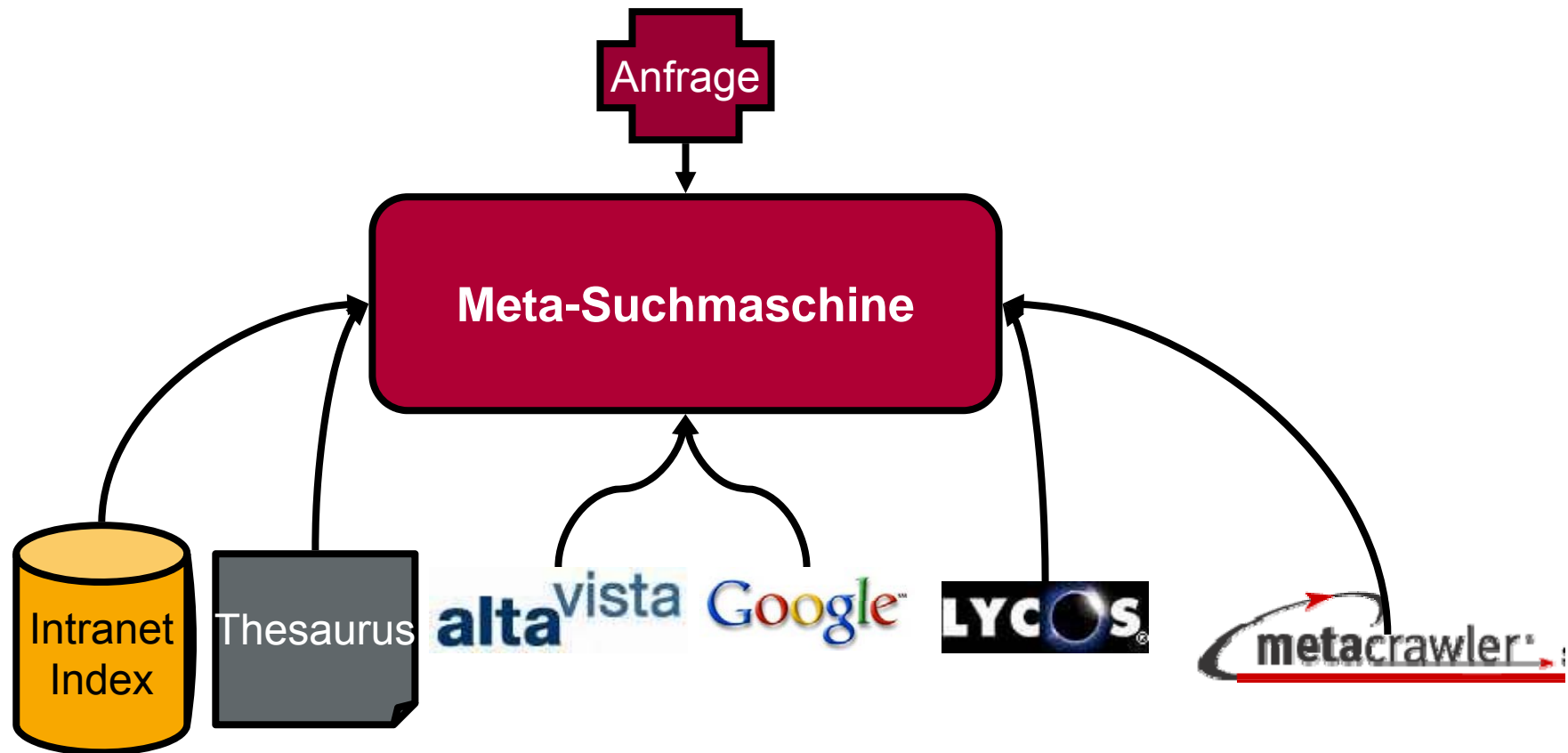
27

- Anwendung
  - Anfrage via Anwendungsschnittstelle oder GUI
  - Struktur: Objekt (Interface), Display (GUI)
  - Einsatzgebiete
    - ◇ Komplexe Analysen (Data Mining, Statistik)
- Integriertes Informationssystem
  - Verhält sich in Anfrage, Struktur und Informationseinheit je nach Design:
    - ◇ DBMS
    - ◇ HTML Formular
    - ◇ Web Service
  - Einsatzgebiete: Meta Search, Life Sciences, Int. Unternehmen, Intranets



# Integrierte Suchmaschinen

29



Weitere Beispiele?

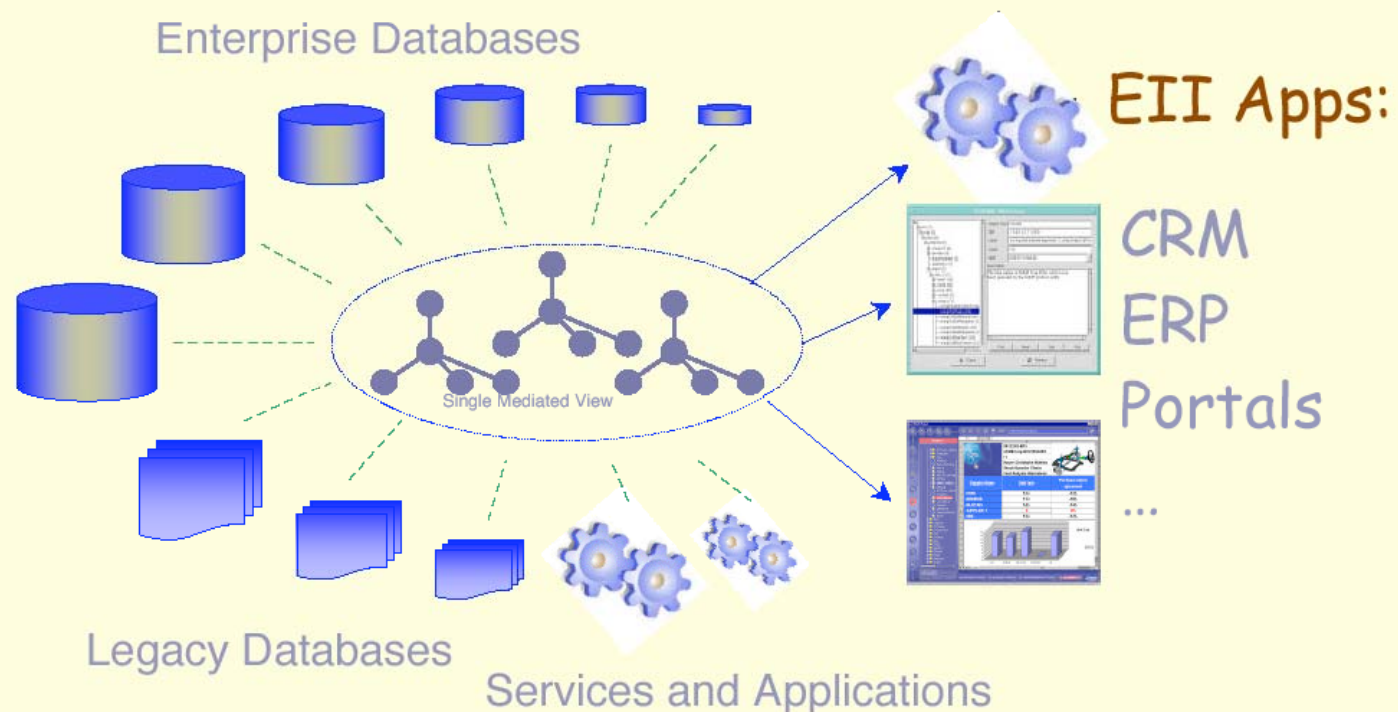
# Integration = Abstraktion

30

1. Logisches DB-Design abstrahiert von physischem DB-Design
  - Datenunabhängigkeit
  - Anfragen: Prozedural vs. deklarativ
2. Informationsintegration abstrahiert von logischen DB Design
  - Quellenunabhängigkeit (Speicherort)
  - Datenmodell- und Syntaxunabhängigkeit
  - Unabhängigkeit von semantischen Unterschieden (hoffentlich!)



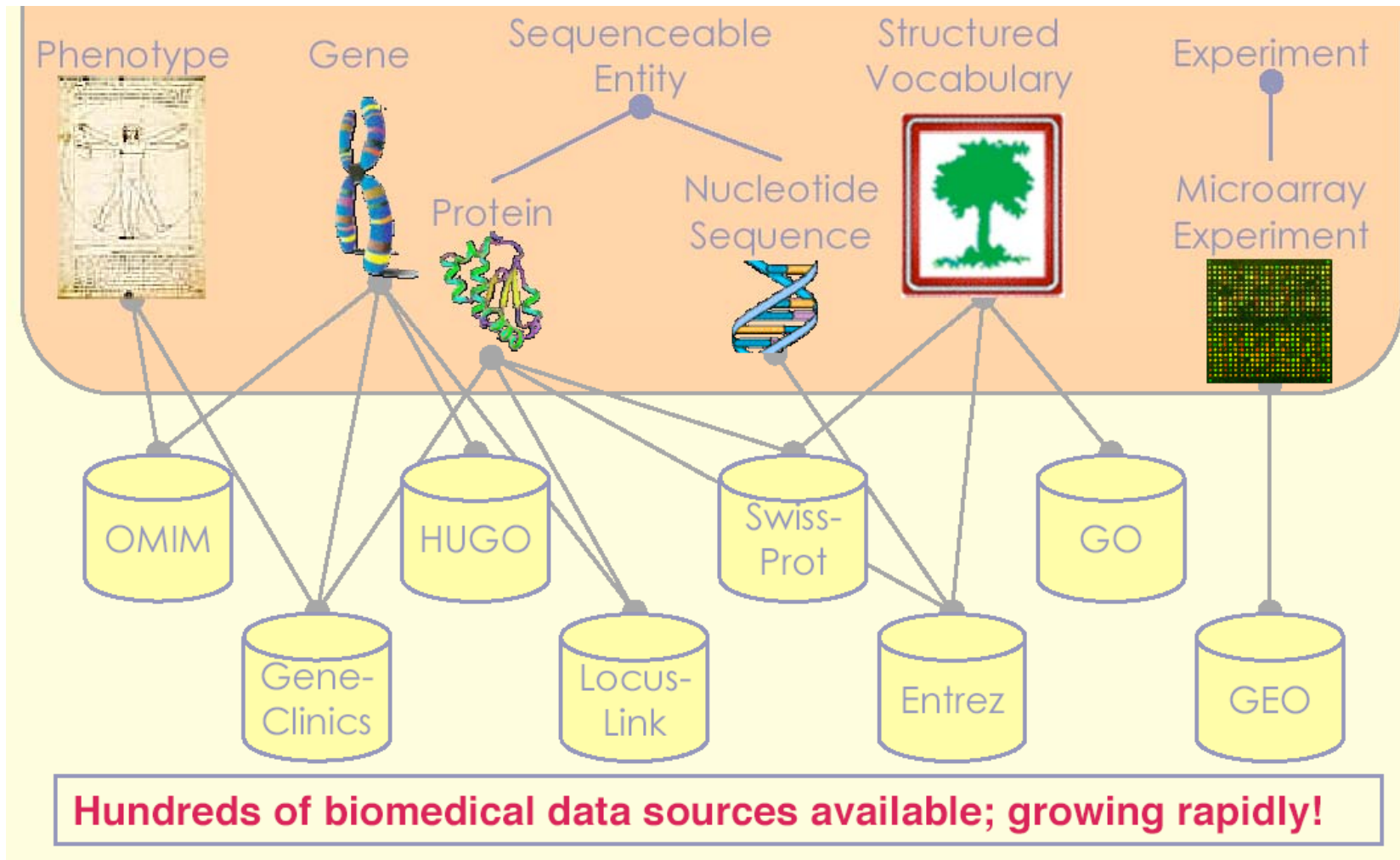
## Application Area 1: Business



# Anwendungsgebiet 2: Wissenschaft

[Halevy04]

32





## Application Area 3: The Web



# Informationsintegration: Ein altes Problem

34

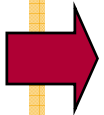
- Seit 50 Jahren auf der Forschungsagenda
- Frühe Systeme in den 70ern
- Integration per Hand natürlich noch früher
- Neue Probleme
  - Viele, viele Quellen
  - Heterogenität
  - Neue Arten von Daten (XML, GIS, OO,...)
  - Neue Arten von Anfragen (Search, UDFs,...)
  - Neue Arten von Ergebnissen (Ranking, Visualisierung, ...)
  - Neue Arten von Nutzern (Laien, Manager, Admins, ...)
- Alon Halevy: „It`s plain hard!“ [Halevy04]

# Warum ist es so schwer? [Halevy04]

35

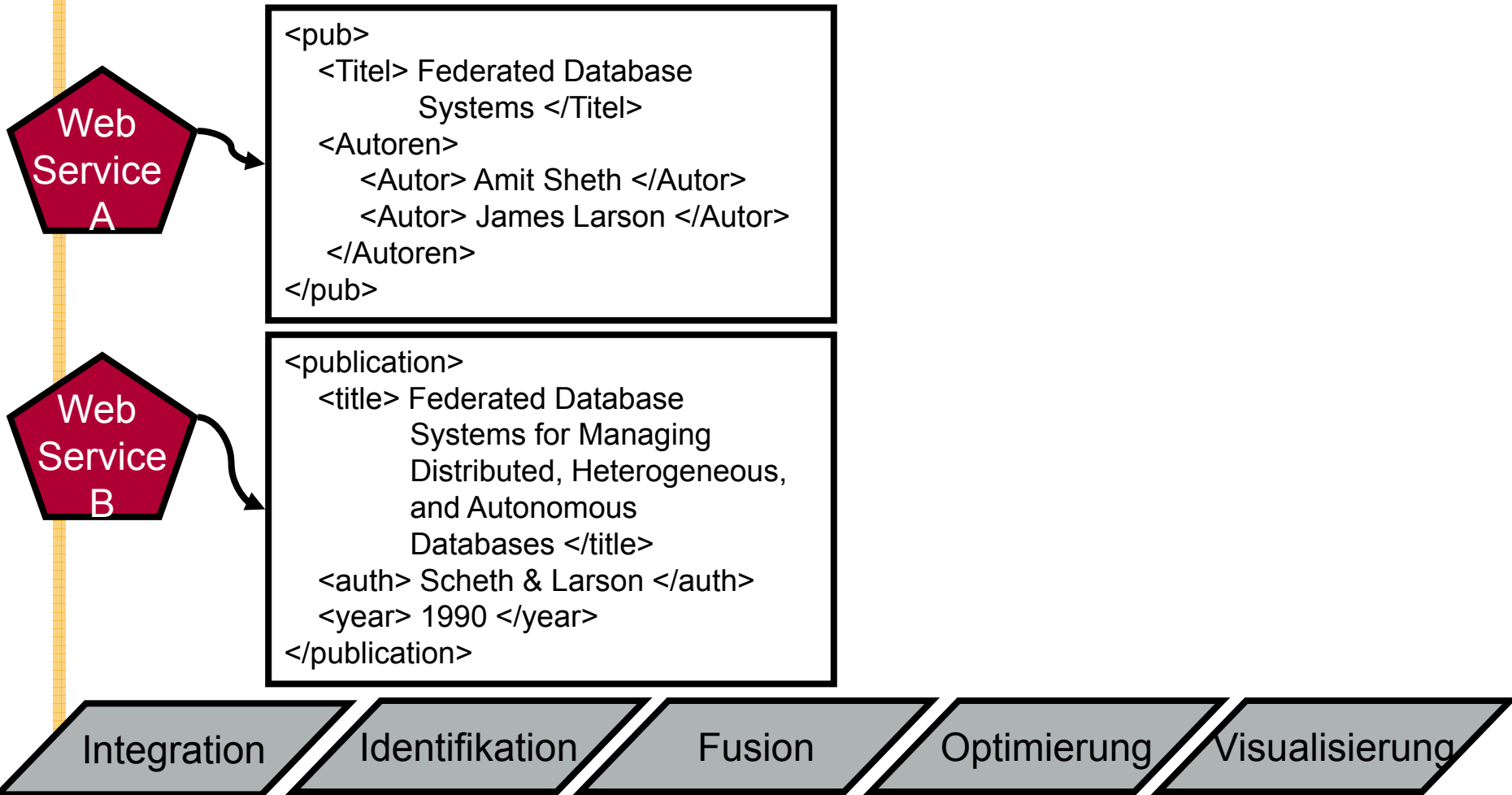
- System-bedingte Gründe
  - Verschiedene Plattformen
  - Anfragebearbeitung über mehrere Systeme
- Soziale Gründe
  - Finden relevanter Daten in Unternehmen
  - Beschaffen relevanter Daten in Unternehmen
  - Menschen zur Zusammenarbeit überreden
- Logik-bedingte Gründe
  - Schema- und Datenheterogenität
  - Dies ist unabhängig von der jeweiligen Integrationsarchitektur.

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Informationssysteme
- Informationsintegration am Beispiel
- Ausblick auf das Semester



# Informationsintegration

37





## Beispiel – Web Service A Output

39

```
<xs:element name="pub" maxOccurs="unbounded">
  <xs:complexType>
    <xs:all>
      <xs:element name="Titel" type="xs:string" nillable="true"/>
      <xs:element name="Autoren">
        <xs:complexType>
          <xs:sequence maxOccurs="unbounded">
            <xs:element name="Autor" type="xs:string" nillable="false"/>
          </xs:sequence>
        </xs:complexType>
      </xs:element>
    </xs:all>
  </xs:complexType>
</xs:element>
```

## Beispiel – Web Service B

40

- Standort: HPI
- Operation: myPubs(Autor, Jahr)
- Struktur:





## Beispiel – Web Service B Output

41

```
<xs:element name="publication" maxOccurs="unbounded">
  <xs:complexType>
    <xs:all>
      <xs:element name="Title" type="xs:string" nillable="true"/>
      <xs:element name="Auth" type="xs:string" nillable="false"/>
      <xs:element name="Year" type="xs:string" nillable="false"/>
    </xs:all>
  </xs:complexType>
</xs:element>
```

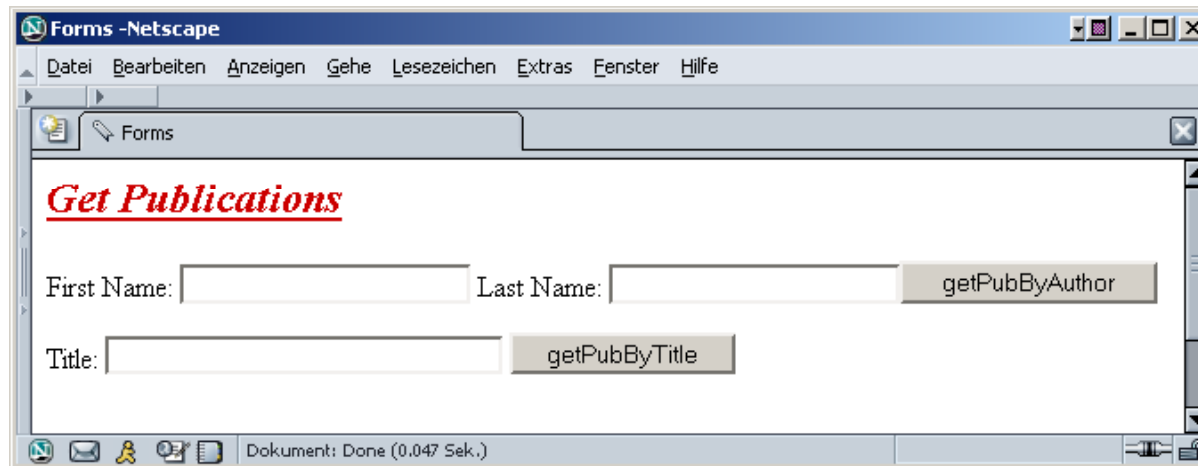
# Integration von Web Services A & B

42

1. Nutzerschnittstelle
2. Schema Integration / Schema Mapping
3. Anfrage-Umwandlung
4. Zeit abschätzen (Optimierung)
5. Requests an beide Services abschicken
6. Antworten einholen
7. Objektidentifikation
8. Integrationsschritte
  1. Konfliktlösung etc.
  2. Entscheidung kleinster gemeinsamer Nenner?
  3. Durchführung (deklarativ, prozedural)
9. Anzeige beim Nutzer

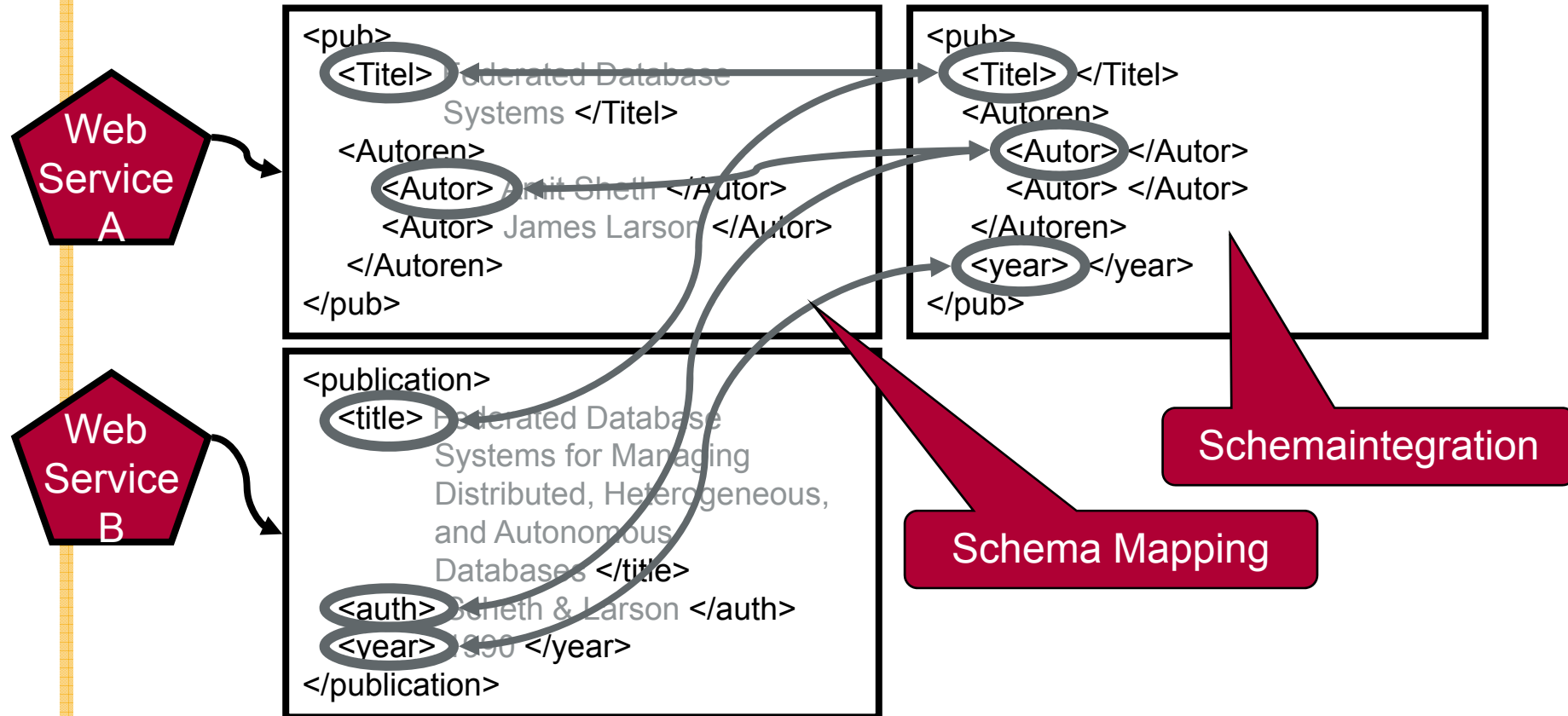
# Nutzerschnittstellen

43



# Informationsintegration

44



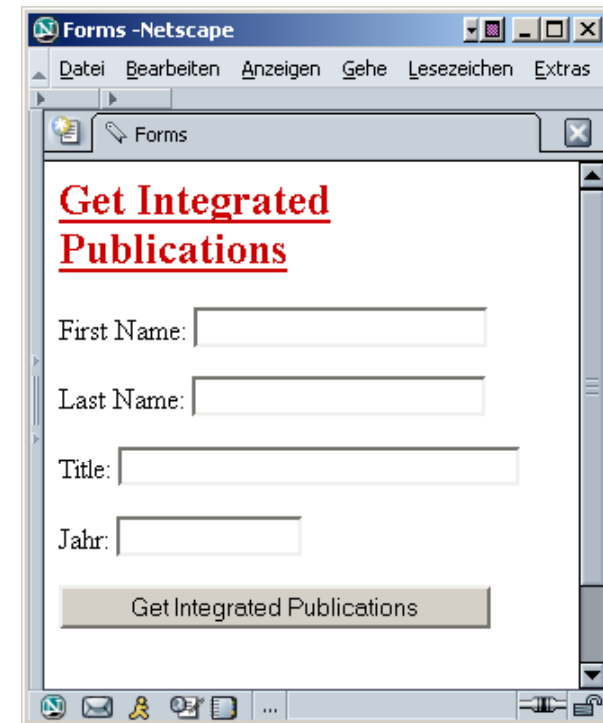
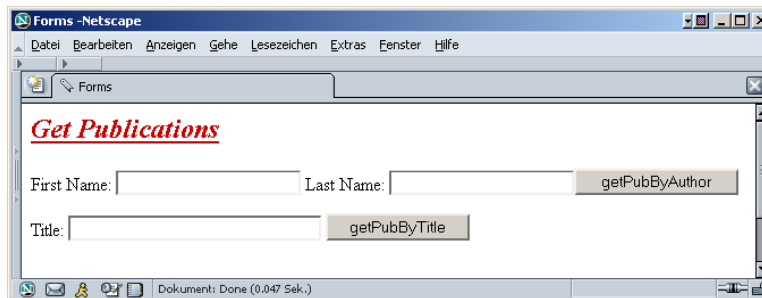
Modellierung durch eine Menge von Anfragen (Views)

# Anfrage Umwandlung

45

Integration der Anfrage durch Mediator:

- Integrierte Schnittstelle
- Z.B.  
Concat(First Name, Last Name)  
= Autor



# Anfrageoptimierung

46

- Was ist besser: Eine schnelle Antwort oder vollständige Antwort?
  - Web Service A in Trier (remote)
  - Web Service B in Griebnitzsee (local)
  - Web Service A hat mehr Attribute und mehr Objekte.
  - Web Service B hat weniger Attribute.
- Außerdem:
  - Eine Suche nach „year“ kann nur durch Web Service B beantwortet werden.
  - Transformationen können teuer sein.

# Zwei Resultate

47

## Web Service A

```
<?xml version="1.0" encoding="UTF-8" ?>
<!-- edited with XMLSPY v2004 rel. 2 U (http://www.xmlspy.com) by Felix Naumann
Universität zu Berlin) -->
<!-- Sample XML file generated by XMLSPY v2004 rel. 2 U (http://www.xmlspy.com)
- <getPub xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="C:\Dokumente und Einstellungen\naumann\Eigene
  Dateien\Lehre\Winter03_04\VL InfoInt I WS03_04\getPub.xsd">
- <pub>
  <Titel>Real-world Data is Dirty: The Merge/Purge Problem</Titel>
  - <Autoren>
    <Autor>Mauricio Hernandez</Autor>
    <Autor>Salvatore Stolfo</Autor>
  </Autoren>
</pub>
- <pub>
  <Titel>MAC: Merging Autonomous Content</Titel>
  - <Autoren>
    <Autor>Felix Naumann</Autor>
    <Autor>Jens Bleiholder</Autor>
  </Autoren>
</pub>
</getPub>
```

```
<?xml version="1.0" encoding="UTF-8" ?>
<!-- Sample XML file generated by XMLSPY v2004 rel. 2 U (http://www.xmlspy.com)
- <myPub xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="C:\Dokumente und Einstellungen\naumann\Eigene
  Dateien\Lehre\Winter03_04\VL InfoInt I WS03_04\myPubs.xsd">
- <publication>
  <Title>Merging Autonomous Content</Title>
  <Auth>Naumann</Auth>
  <Year>2003</Year>
</publication>
- <publication>
  <Title>Object Mathcing for Information Integration</Title>
  <Auth>Doan</Auth>
  <Year>1999</Year>
</publication>
</myPub>
```

## Web Service B

# Schema Matching

48

```

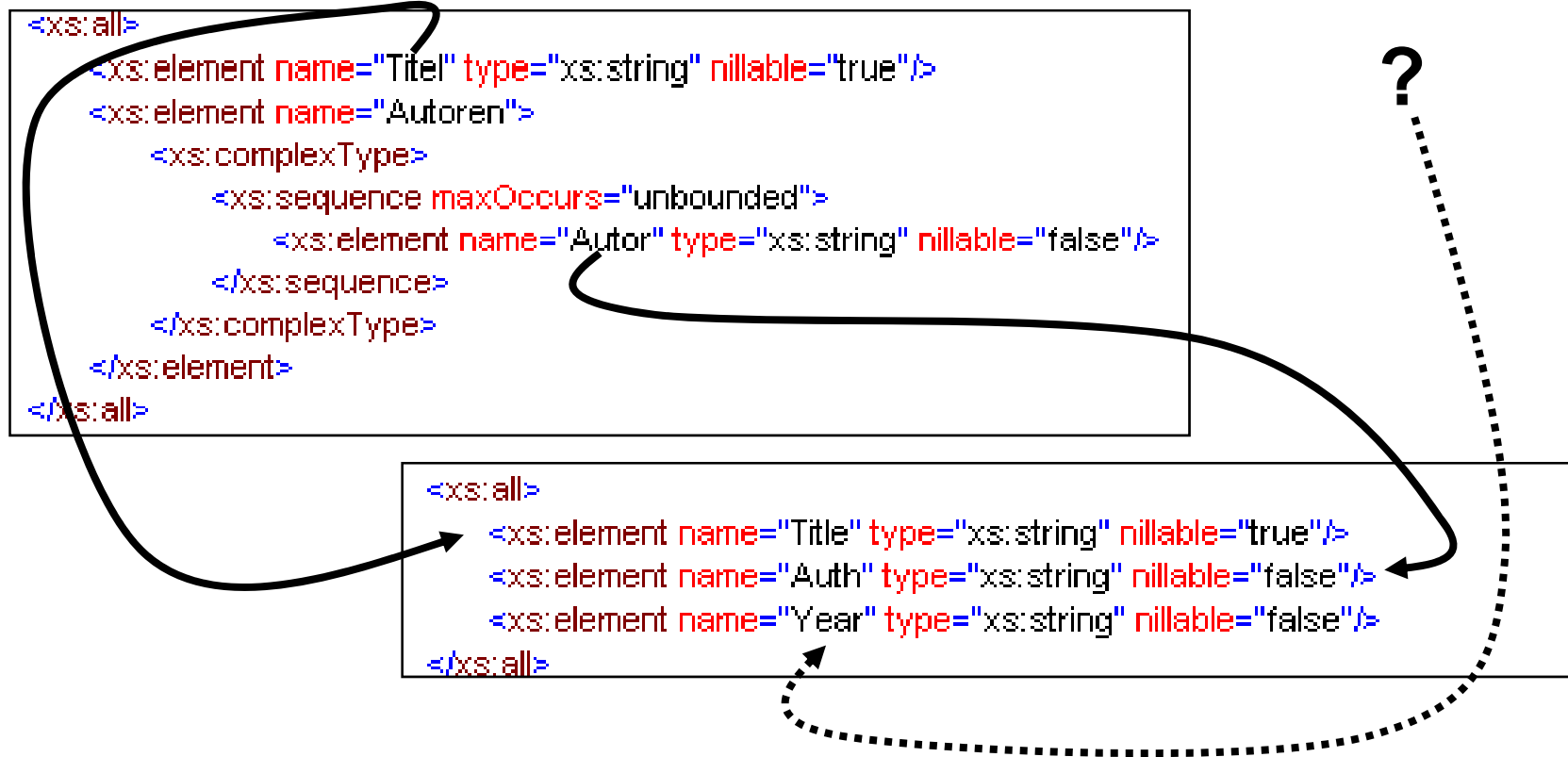
<xs:all>
  <xs:element name="Titel" type="xs:string" nillable="true"/>
  <xs:element name="Autoren">
    <xs:complexType>
      <xs:sequence maxOccurs="unbounded">
        <xs:element name="Autor" type="xs:string" nillable="false"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>
</xs:all>

```

```

<xs:all>
  <xs:element name="Title" type="xs:string" nillable="true"/>
  <xs:element name="Auth" type="xs:string" nillable="false"/>
  <xs:element name="Year" type="xs:string" nillable="false"/>
</xs:all>

```





# Objektidentifikation

49

```

- <pub>
  <Titel>Real-world Data is Dirty: The Merge/Purge Problem</Titel>
  - <Autoren>
    <Autor>Mauricio Hernandez</Autor>
    <Autor>Salvatore Stolfo</Autor>
  </Autoren>
</pub>
- <pub>
  <Titel>MAC: Merging Autonomous Content</Titel>
  - <Autoren>
    <Autor>Felix Naumann</Autor>
    <Autor>Jens Bleiholder</Autor>
  </Autoren>
</pub>
</getPub>

```

```

- <publication>
  <Title>Merging Autonomous Content</Title>
  <Auth>Naumann</Auth>
  <Year>2003</Year>
</publication>
- <publication>
  <Title>Object Mathcing for Information Integration</Title>
  <Auth>Doan</Auth>
  <Year>1999</Year>
</publication>
</myPub>

```

# Objektidentifikation

50

```
- <pub>  
  <Titel>MAC: Merging Autonomous Content</Titel>  
  - <Autoren>  
    <Autor>Felix Naumann</Autor>  
    <Autor>Jens Bleiholder</Autor>  
  </Autoren>  
</pub>  
</getPub>
```

Edit-distance: 5 }  
Edit-distance: 6 } Zusammen?

```
- <publication>  
  <Title>Merging Autonomous Content</Title>  
  <Auth>Naumann</Auth>  
  <Year>2003</Year>  
</publication>
```

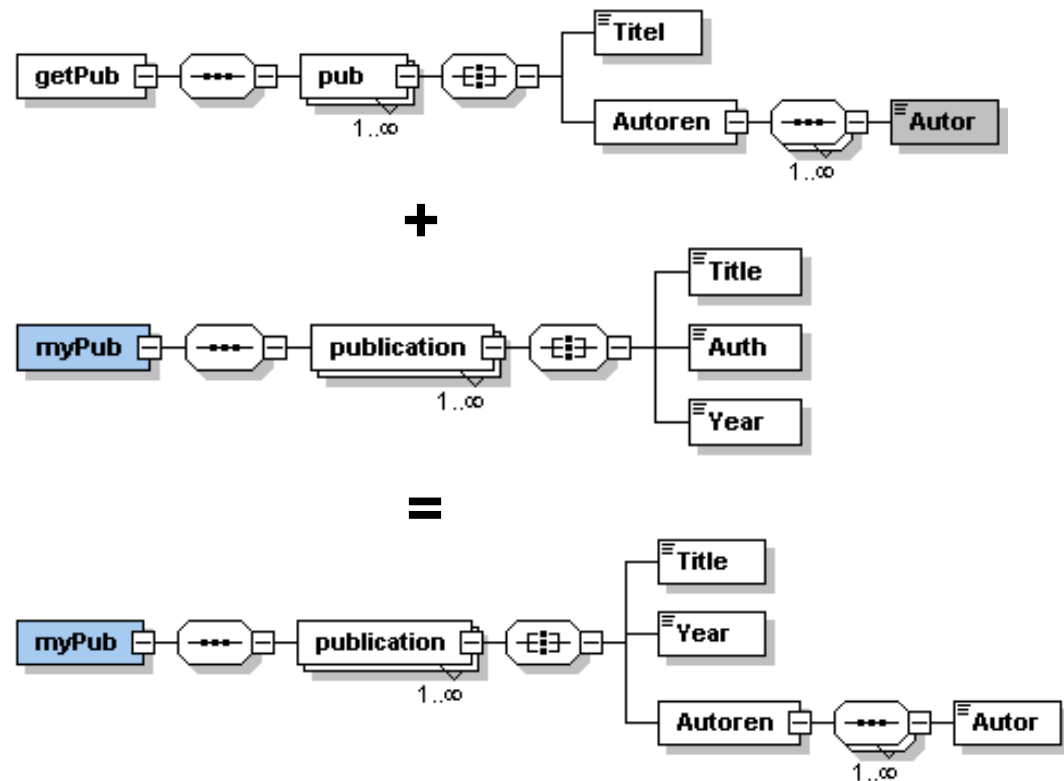
# Stand der Dinge

51

- Wir haben die heterogenen Informationen.
- Wir wissen, was wir integrieren wollen.
- Aber noch nicht wie:
  - Integriertes Schema
  - Integrierte Daten

# Angestrebtes Integrationsergebnis

52



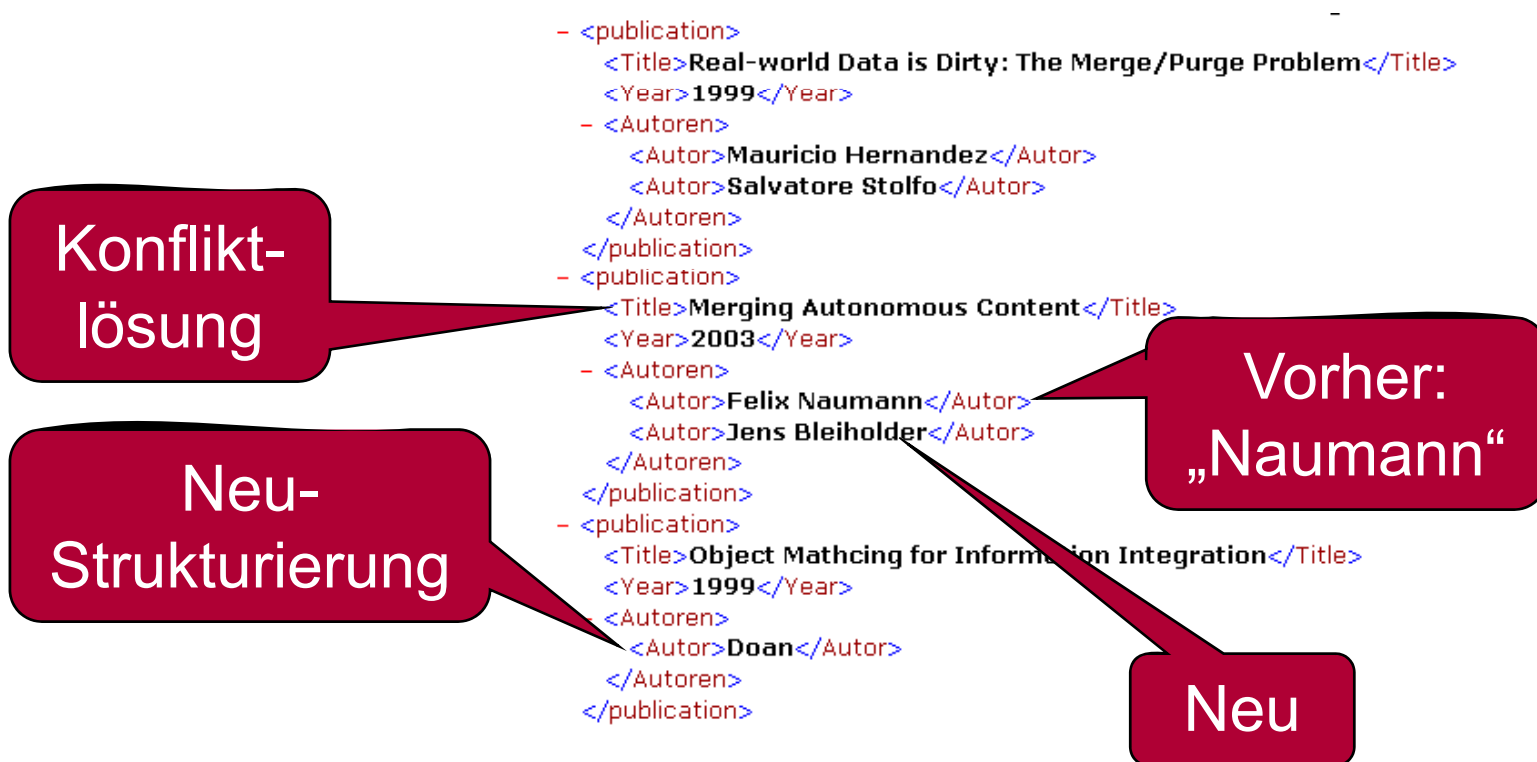
Integriertes Schema:

## Integrierte Daten:

- <publication>  
  <Title>Real-world Data is Dirty: The Merge/Purge Problem</Title>  
  <Year>1999</Year>
- <Autoren>  
  <Autor>Mauricio Hernandez</Autor>  
  <Autor>Salvatore Stolfo</Autor>  
</Autoren>
- </publication>
- <publication>  
  <Title>Merging Autonomous Content</Title>  
  <Year>2003</Year>
- <Autoren>  
  <Autor>Felix Naumann</Autor>  
  <Autor>Jens Bleiholder</Autor>  
</Autoren>
- </publication>
- <publication>  
  <Title>Object Matching for Information Integration</Title>  
  <Year>1999</Year>
- <Autoren>  
  <Autor>Doan</Autor>  
</Autoren>
- </publication>

# Integrierte Daten – was ist passiert?

54



# Implementierung

55

- Auf Folien ist alles klar, aber wie implementieren?
- Deklarativ?
  - SQL, XQuery, XSLT
  - Oft nicht alles möglich
  - Langsam
- Prozedural?
  - Java, C++
  - Schlecht wartbar
  - Schnell

# Anzeige beim Nutzer

56

Konflikt-  
lösung

Visualisierung der

- Datenherkunft
- Qualität
- veränderten Daten
- Operationen

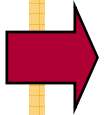
```
<publication>  
<Title>Real-world Data is Dirty: The Merge/Purge Problem</Title>  
<Year>1999</Year>  
- <Autoren>  
  <Autor>Mauricio Hernandez</Autor>  
  <Autor>Salvatore Stolfo</Autor>  
</Autoren>  
</publication>  
<publication>  
<Title>Merging Autonomous Content</Title>  
<Year>2003</Year>  
- <Autoren>  
  <Autor>Felix Naumann</Autor>  
  <Autor>Jens Bleiholder</Autor>  
</Autoren>  
</publication>  
- <publication>  
<Title>Object Mathcing for Information Integration</Title>  
<Year>1999</Year>  
- <Autoren>  
  <Autor>Doan</Autor>  
</Autoren>  
</publication>
```

Vorher:  
„Naumann“



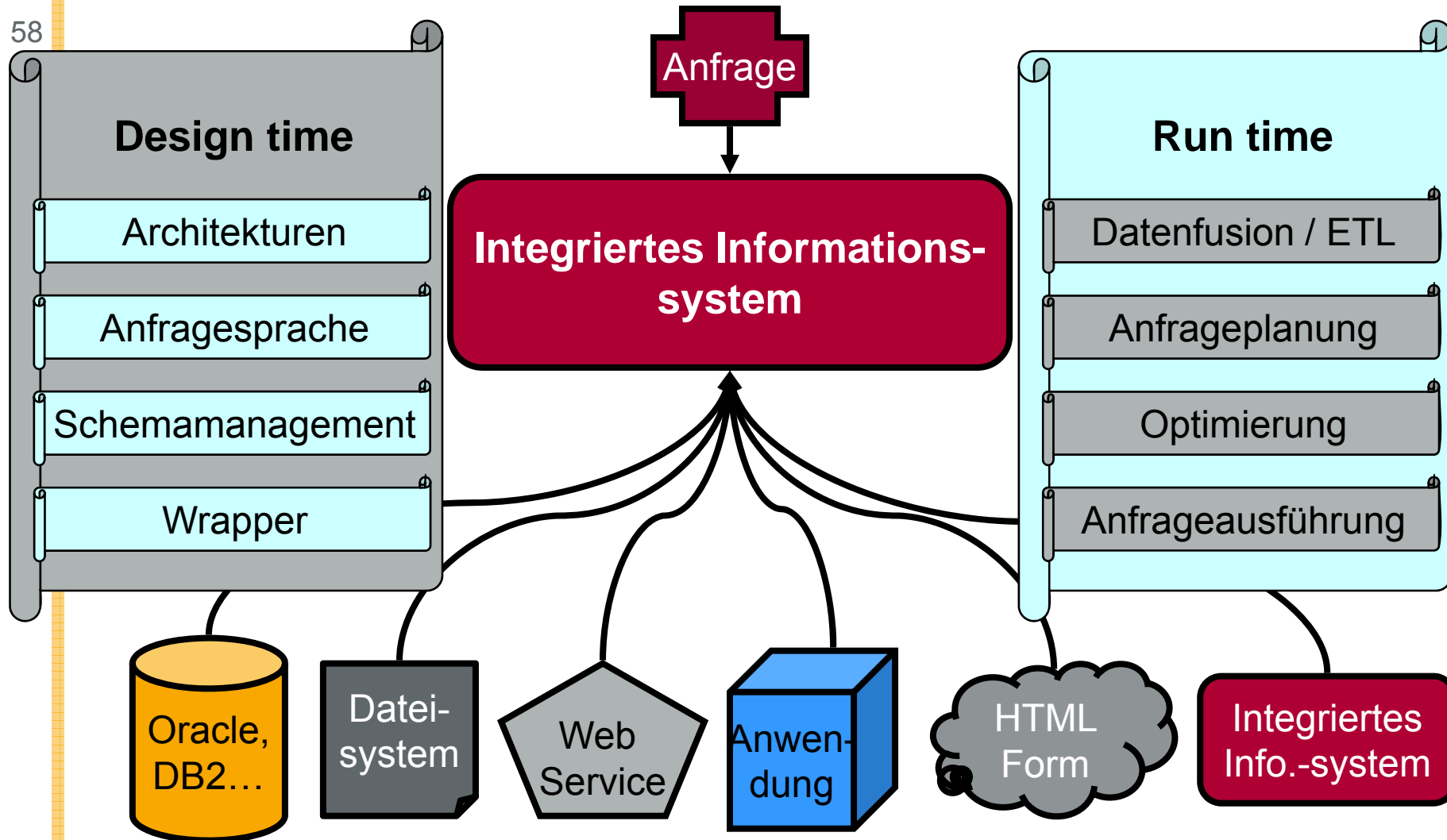
57

- Vorstellung der Arbeitsgruppe
- Organisatorisches
- Informationssysteme
- Informationsintegration am Beispiel
- Ausblick auf das Semester



# Integrierte Informationssysteme

58



- Einführung in die Informationsintegration

- Szenarien der Informationsintegration

- Verteilung und Autonomie

- Heterogenität

- Materialisierte und virtuelle Integration

- 5-Schichten Architektur

- Mediator/Wrapper-Architektur / PDMS

- Schema Mapping

- Schema Matching

- SchemaSQL

- Global-as-View und Lokal-as-View Modellierung

- Global-as-View Anfragebearbeitung

- Containment & Local-as-View Anfragebearbeitung
- Bucket Algorithmen
- Verteilte Anfragebearbeitung

## Anfragen

- Duplikaterkennung
- Datenfusion – Union & Co.
- DWH, ETL & Data lineage
- Informationsqualität

## Datenintegration

- Hidden Web
- Semantic Web

## Anwendungen

# Fragen, Wünsche und Vorstellungen

61

- Jetzt, oder...
- Raum: A.1-13
- Sprechstunden: Dienstags 15-16 Uhr  
oder n.V.
- Email: [naumann@hpi.uni-potsdam.de](mailto:naumann@hpi.uni-potsdam.de)
- Telefon: (0331) 5509 280

*The end.*