

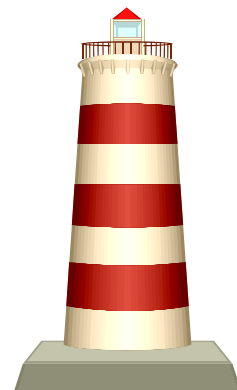
VL Informationsintegration
Verteilung, Autonomie und
Heterogenität

19.4.2007
Felix Naumann

Überblick

2

- ➔ ■ Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



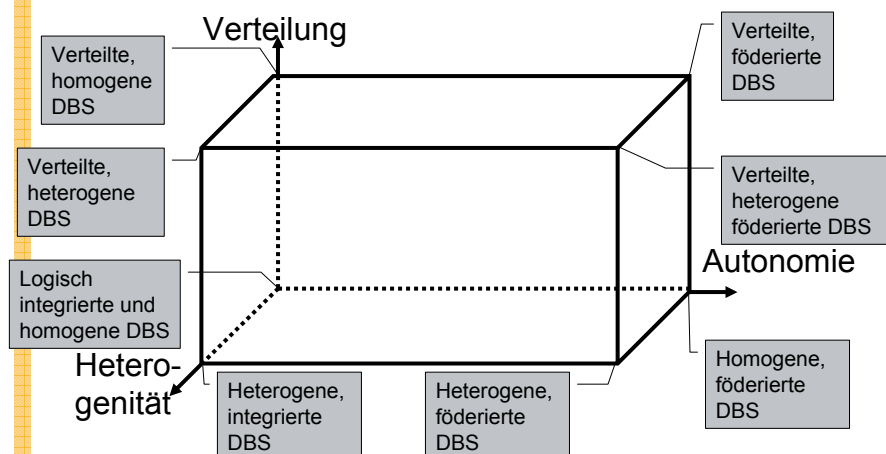
Klassifikation von Informationssystemen [ÖV99]

3

- Drei orthogonale Dimensionen
 - Verteilung
 - Autonomie
 - Heterogenität

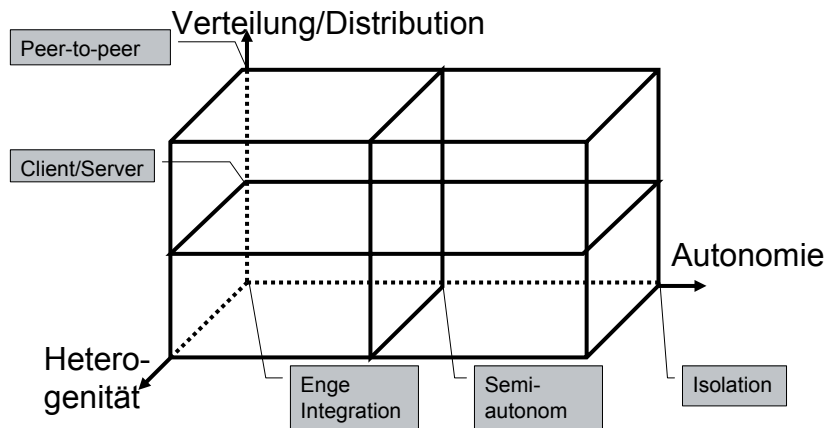
Klassifikation verteilter DBMS [ÖV91]

4



Klassifikation verteilter DBMS nach [ÖV99]

5



Felix Naumann | VL Informationsintegration | SS 2007

Zusammenhang mit Föderierten DBMS

6

- Verteilung führt zu Autonomie,
 - Intra-Organisation: Historisch
 - Inter-Organisation: Internet & WWW
- und Autonomie führt zu Heterogenität.
 - Verantwortung liegt bei lokalen Administratoren
 - Systempflege
 - Nutzbarkeit und Nützlichkeit
 - Erweiterungen am Informationssystem
 - Design
 - ...
- Diskussion
 - Historischer Entwicklung,
 - aber orthogonale Kriterien!

Felix Naumann | VL Informationsintegration | SS 2007

Verteilung (*Distribution*)

7

Ein verteiltes Informationssystem ist eine Sammlung mehrerer, logisch verknüpfter Informationssysteme, die über ein gemeinsames Netzwerk verteilt sind.

[ÖV91]

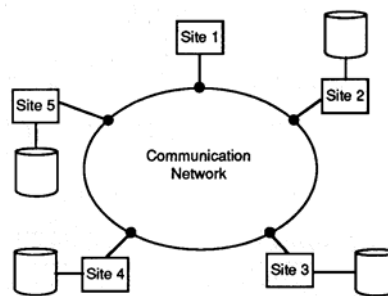


Figure 1.7 DDBS Environment

Felix Naumann | VL Informationsintegration | SS 2007

Physikalische Verteilung

8

- Motiviert durch Hardwareanforderungen (Hardwarebeschränkungen)
- Server stehen an unterschiedlichen Orten
 - Gleicher Raum, anderer Raum
 - Anderes Gebäude
 - Andere Stadt, anderes Land
- Shared Nothing
 - Server haben keine gemeinsamen, abhängigen Hardwarekapazitäten
 - Memory
 - Disk
 - CPU
 - Mit Ausnahme des Netzwerks
 - Im Gegensatz zu shared-disk und shared-memory

Felix Naumann | VL Informationsintegration | SS 2007

Logische Verteilung

9

- Motiviert durch Anwendungsanforderungen
 - Zuverlässigkeit
 - Bei Ausfall eines Servers
 - Verfügbarkeit
 - Bei Ausfall eines Netzwerkteils
 - Effizienz
- Redundanz
 - Replikation
 - Caching
- Partitionierung
 - Vertikal
 - Horizontal

Felix Naumann | VL Informationsintegration | SS 2007

Verteilung – Vor- und Nachteile

10

Vorteile aus Sicht der Quellen und des IIS

- Autonomie (gleich genauer)
- Performance: Kapazität dort, wo sie gebraucht wird
- Verfügbarkeit: Bei Ausfall eines Standorts
- Erweiterbarkeit
- Teilbarkeit (Verantwortung bei anderen Organisationseinheiten)

Nachteile aus Sicht des IIS

- Komplexität (Verwaltung, Optimierung)
- Kosten
- Sicherheit
- Autonomie

Felix Naumann | VL Informationsintegration | SS 2007

Verteilung – Techniken

11

HTTP, CORBA, ... nicht hier.

- Anwendungsentwicklung ohne Spezifikation der physikalischen Präsenz der Komponenten

Annahmen an Transparenz

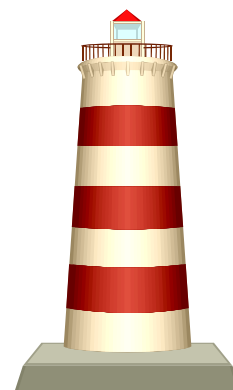
- Datenunabhängigkeit (jedes DBMS)
 - auch Speicherorttransparenz
- Netzwerktransparenz
- Replikationstransparenz
- Fragmentationstransparenz
 - auch Partitionierungstransparenz

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

12

- Verteilung
- ➔ ■ Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Felix Naumann | VL Informationsintegration | SS 2007

Autonomie (*Autonomy*)

13

Der Grad zu dem verschiedene DBMS unabhängig operieren können.
Bezieht sich auf Kontrolle, nicht auf Daten.

Klassen nach [ÖV99]

- Design-Autonomie
- Kommunikations-Autonomie
- Ausführungs-Autonomie

Design-Autonomie

14

- Auch: Entwurfsautonomie
- Freiheit des lokalen DBMS bezüglich
 - Datenmodell
 - Relational, hierarchisch, XML
 - Schema
 - Abdeckung der Domäne (*universe of discourse, miniworld*)
 - Grad der Normalisierung
 - Benennung
 - Transaktionsmanagement
 - Sperrprotokolle
- Freiheit dies jederzeit zu ändern.
 - Besonders problematisch!

Design-Autonomie – Beispiel

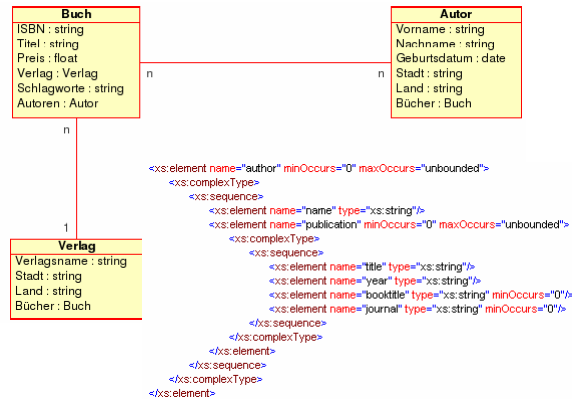
15

■ Schema und Datenmodell 1

- (Fast) relational
- Flach

■ Schema und Datenmodell 2

- XML
- hierarchisch



Felix Naumann | VL Informationsintegration | SS 2007

Kommunikations-Autonomie

16

■ DBMS frei bezüglich

- Wahl mit welchen Systemen kommuniziert wird
- Wahl wann mit anderen Systemen kommuniziert wird
 - Jederzeit Eintritt/Austritt aus integriertem System
- Wahl was (welcher Teil der Information) kommuniziert wird
- Wahl wie mit anderen Systemen kommuniziert wird
 - Anfragesprache
- Wahl welcher Teil der Anfragemöglichkeiten zur Verfügung gestellt werden
 - Prädikate
 - Sortierung
 - Write
 - ...

Felix Naumann | VL Informationsintegration | SS 2007

17

- Extrem 1: Voller SQL Zugang
 - z.B. via JDBC
 - Transaktionen
 - Optimierung
 - Lesend (und Schreibend?)
 - Schemaveränderungen?
 - Antwort als Ergebnisrelation
- Extrem 2: HTML Formular
 - Nur ein (oder mehr) Suchfelder
 - Antwort als HTML Text
 - Nur Teile der Daten (public area)



Felix Naumann | VL Informationsintegration | SS 2007

18

- DBMS frei bezüglich
 - Wahl wann Anfragen ausgeführt werden
 - Wahl wie Anfragen ausgeführt werden
 - Wahl der Scheduling-Strategien
 - Wahl Optimierungs-Strategien
 - Wahl ob globale Transaktionen unterstützt werden

Felix Naumann | VL Informationsintegration | SS 2007

Ausführungs-Autonomie – Beispiel

19

- Optimierung und Scheduling
 - Behandlung externer vs. lokaler Anfragen
 - *Golden customers*
 - Garantierte Antwortzeiten

- Transaktionen
 - Dirty-read egal?



Felix Naumann | VL Informationsintegration | SS 2007

Autonomie → Heterogenität

20

- Verteilung als „Ursache“ für Autonomie
- Autonomie als Ursache für Heterogenität:
 - Autonome Systeme
 - ⇒ Gestaltungsfreiheit
 - ⇒ Unterschiedliche Entscheidungen
 - ⇒ Heterogenität

Felix Naumann | VL Informationsintegration | SS 2007

Heterogenität (*Heterogeneity*)

21

Heterogenität herrscht, wenn sich zwei miteinander verbundene Informationssysteme syntaktisch, strukturell oder inhaltliche unterscheiden.

- Syntaktische Heterogenität
 - Auch: „Technische Heterogenität“
- Strukturelle Heterogenität
- Semantische Heterogenität

Heterogenitäten zu überbrücken ist die Kernaufgabe der Informationsintegration.

Felix Naumann | VL Informationsintegration | SS 2007

Heterogenitätsklassen

22

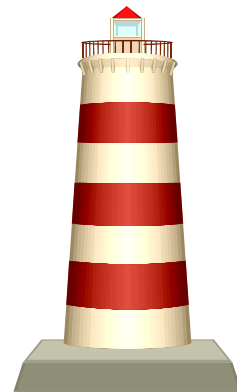
- Auch andere Klassifikationen möglich, z.B. [BKLW99]
 - Syntaktische Heterogenität
 - Datenmodell Heterogenität
 - Logische Heterogenität
- Oder nach [SPD92]
 - Semantische Konflikte
 - Beschreibungskonflikte
 - Heterogenitätskonflikte
 - Strukturelle Konflikte
 - Datenkonflikte

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

23

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung

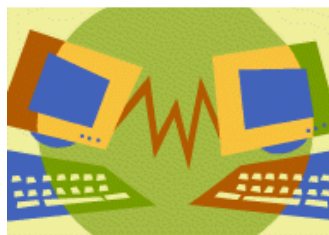


Felix Naumann | VL Informationsintegration | SS 2007

Syntaktische Heterogenität

24

- Hardware-Heterogenität
- Software-Heterogenität
- Schnittstellen-Heterogenität

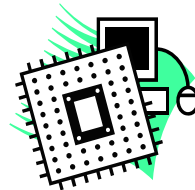


Felix Naumann | VL Informationsintegration | SS 2007

Hardware Heterogenität

25

- Bandbreite
- Hauptspeicher
- CPU
 - Art
 - Geschwindigkeit



Nicht hier

Felix Naumann | VL Informationsintegration | SS 2007

Software Heterogenität

26

- Betriebssystem
- Dateisystem
- Protokolle
 - HTTP, ODBC, Java API, CORBA, etc.
- Zustandsbehaftet vs. zustandsfrei
- Sicherheit
 - Security level
 - Log-on Prozedur



Felix Naumann | VL Informationsintegration | SS 2007

Software Heterogenität – Beispiel

27

```
String sqlQuery = „SELECT Name, Strasse FROM Hersteller  
WHERE PLZ = 69115“;  
...  
Connection jdbcCon = DriverManager.getConnection(dbURL, ...);  
Statement stmt = jdbcCon.createStatement();  
ResultSet table = stmt.executeQuery(sqlQuery);  
...
```

```
String webQuery = „plz=69115“;  
...  
URL url = new URL(„http://www.system2.de/cgi-bin  
/search.cgi“ + „?“ + webQuery);  
URLConnection urlCon = new url.openConnection();  
InputStreamReader reader = new InputStreamReader(  
urlCon.openStream());  
...
```

Nicht hier

Quelle: VL: Föderierte
Datenbanksysteme
Peter Tomczyk, FZI &
Uni Karlsruhe

Schnittstellen Heterogenität

28

Schnittstellen von Informationssystemen sind im wesentlichen deren
Anfragensprache:

- HTML Formular,
- „Google“-Sprache (+, - , ...),
- SQL,
- XQuery,
- etc.

Jetzt hier!

Schnittstellen Heterogenität

29

- Negation vs. keine Negation
 - Oft zu teuer
- Gleichheit / Ungleichheit
 - „=" oder auch „>, <, ≥, ≤“
- Konjunktion (UND)
 - oder auch Disjunktion (ODER)
- Prädikate nur mit Konstanten (author = „Melville“)
 - Oder auch mit anderen Variablen (ResidenceCountry = Nationality)
- Gebundene und freie Variablen [RSU95,LC00,YLGU99]
 - später
- Andere Einschränkungen
 - Joins über maximal 3 Relationen
 - z.B. Prädikate nur über eine Auswahl von Werten

Felix Naumann | VL Informationsintegration | SS 2007

Schnittstellen Heterogenität - Beispiel

30

The screenshot shows a webmail interface with a search window titled 'Nachrichten durchsuchen'. The search criteria are:

- Betreff: enthält integration
- Betreff: enthält nicht Humboldt
- Betreff: enthält (dropdown menu is open showing options: enthält, enthält nicht, gleich, ungleich, beginnt mit, endet mit)

 Annotations with yellow arrows point to:

- 'Suche' pointing to the search window title.
- 'Konjunktion/Disjunktion' pointing to the search criteria list.
- 'gleich/ungleich' pointing to the dropdown menu options.

Felix Naumann | VL Informationsintegration | SS 2007

Schnittstellen-Heterogenität – Beispiel

31

Gebundene Variablen

Prädikat nur mit
Auswahl von Werten

amazon.de

Heiße Tage, heiße Nächte! Musik zum Verführen und mehr (Erotik)

Erweiterte Suche Bücher

Sie können auch nur eines der Felder ausfüllen.
Bitte geben Sie in die unten stehenden Suchfelder ein oder mehrere Suchbegriffe ein und klicken Sie auf "Jetzt suchen".

Autor/in:

Titel:

Schlagwörter:

ISBN:

Verlag:

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

Nur gebraucht:

Format:

Ordnen nach:

Errscheinungsdatum:

Suche in:

Felix Naumann | VL Informationsintegration | SS 2007

Schnittstellen Heterogenität

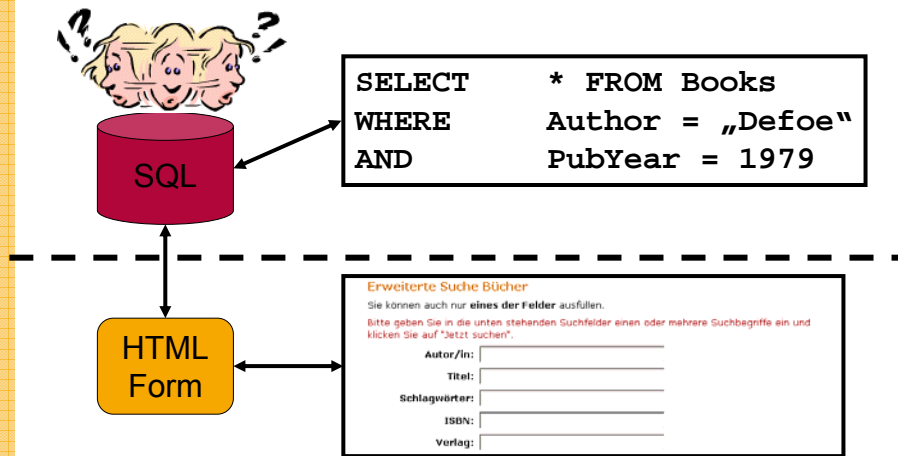
32

- In einzelnen Systemen kein Problem
- Probleme für integrierte Systeme
 1. Globale Anfragesprache ist mächtiger als lokale Anfragesprache
 - Anfragen eventuell nicht ausführbar
 - Oder globales System muss kompensieren
 2. Lokale Anfragesprache ist mächtiger als globale Anfragesprache
 - Verpasste Chance, lokale (effiziente) Ausführung auszunutzen
 3. Gebundene und freie Variablen sind inkompatibel
 - Anfragen eventuell nicht ausführbar

Felix Naumann | VL Informationsintegration | SS 2007

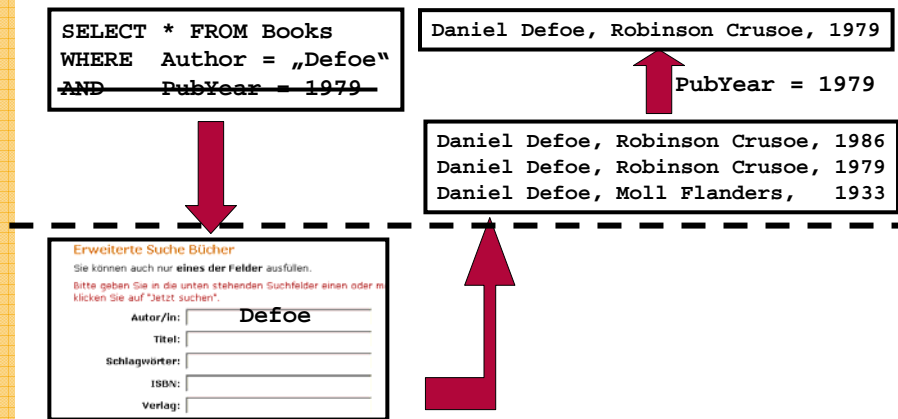
Mächtige globale Anfragesprache

33



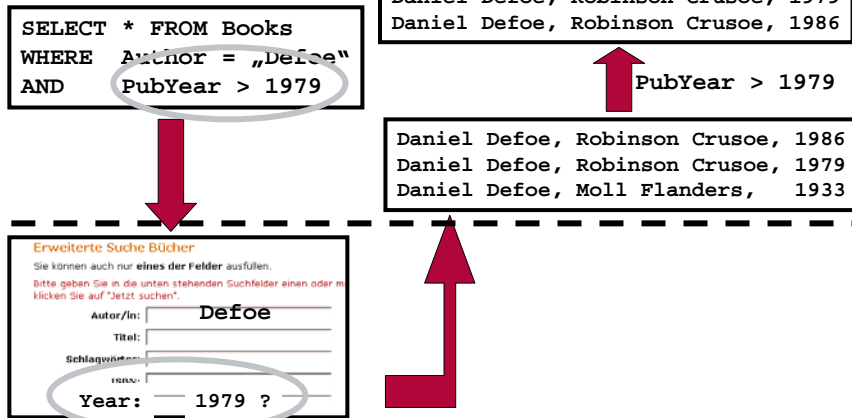
Mächtige globale Anfragesprache

34



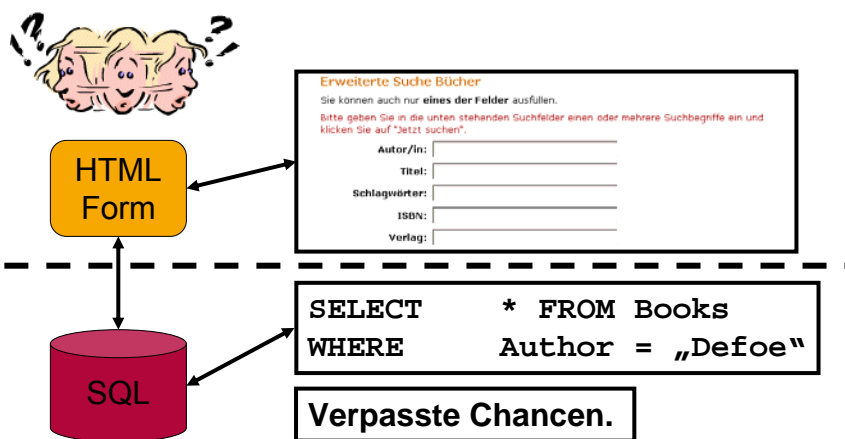
Mächtige globale Anfragesprache

35



Mächtige lokale Anfragesprache

36



Gebundene & Freie Variablen

37

- **Gebundene Variablen** müssen bei einer Anfrage gebunden werden.
 - z.B.: „Search“-Feld bei Google
- **Freie Variablen** müssen nicht gebunden werden.
 - z.B. „Autor“-Feld bei Amazon.de, falls Titel gebunden ist.

Gebundene & Freie Variablen – Beispiel & Ausblick

38

SONGS	Song	CD
	Friends	Life
	Friends	Love

CDs	CD	Künstler	Preis
	Love	Lucy	15
	Story	Snoopy	14

Künstler	CD	Künstler	Preis
	Story	Lucy	13
	Love	Snoopy	10
	Life	Charlie	8

Bastelaufgabe 1:
Wie teuer ist die billigste CD mit einem Song namens "Friends"?

Quelle: [LC00]

Gebundene & Freie Variablen – Beispiel & Ausblick

39

SONGS	<u>Song</u>	CD
	Friends	Life
	Friends	Love

CDs	<u>CD</u>	Künstler	Preis
	Love	Lucy	15
	Story	Snoopy	14

Künstler	CD	<u>Künstler</u>	Preis
	Story	Lucy	13
	Love	Snoopy	10
	Life	Charlie	8

Unterstrichen
= gebundene
Variable

Bastelaufgabe 2:
Welches ist die billigste CD mit einem Song
namens "Friends", *die Sie anfragen können?*

Mehr später...

Felix Naumann | VL Informationsintegration | SS 2007

Syntaktische Heterogenität - Zusammenfassung

40

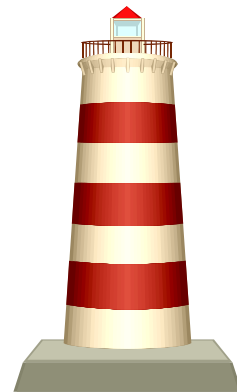
- Hardware Heterogenität
 - Bandbreite, CPU, ...
- Software Heterogenität
 - Protokolle, Sicherheit, ...
- Schnittstellen Heterogenität
 - Mächtigkeit der Anfragesprachen
 - Gebundene & freie Variablen

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

41

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Strukturelle Heterogenität

42

- Datenmodell-Heterogenität
 - Unterschiedliche Semantik
 - Unterschiedliche Struktur
- Schematische Heterogenität
 - Integritätsbedingungen, Schlüssel, Fremdschlüssel, etc.
 - Schema (Attribut vs. Relation etc.)
 - Struktur (Gruppierung in Tabellen)

Datenmodell-Heterogenität

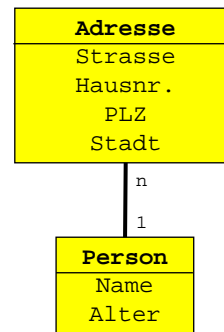
43

Datenmodelle

- Relationales Modell
- XML Modell
- OO Modell
- Hierarchisches Modell

```
Adresse(PersonId, Strasse,
        Hausnr., PLZ, Stadt)
```

```
Person( Id, Name, Alter)
```



Schematische Heterogenität

44

Sprachregelung hier:

Schemas statt Schemata

Alt-Griechisch wird ignoriert...



Schematische Heterogenität

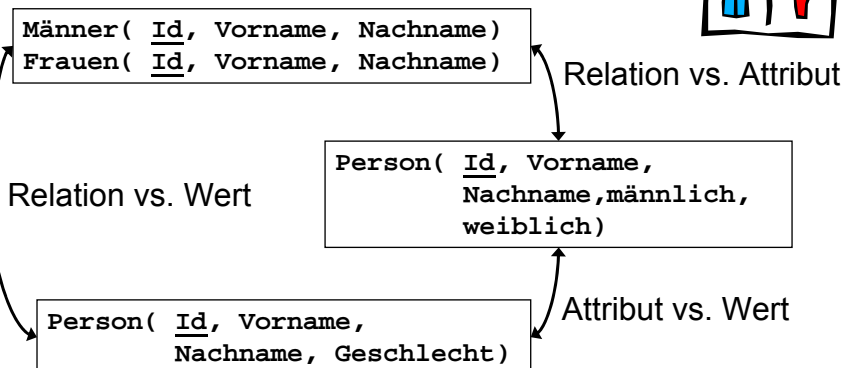
45

- Struktur
 - Modellierung
 - Relation vs. Attribut
 - Attribut vs. Wert
 - Relation vs. Wert
 - Benennung
 - Relationen
 - Attribute
 - Jeweils Homonyme und Synonyme
 - Normalisiert vs. Denormalisiert
 - Geschachtelt vs. Fremdschlüssel
- Diese Probleme sogar bei gleichem Datenmodell!

Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität

46



Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität

47

Tabellen-Tabellen Konflikte

- Namenskonflikte
 - Semantisch gleiche Tabellen mit verschiedenen Namen (Synonym)
 - Verschiedene Tabellen mit gleichem Namen (Homonym)
- Strukturkonflikte
 - fehlende Attribute
 - fehlende, aber ableitbare Attribute
- IC-Konflikte (Integrity-Constraint = Integritätsbedingungen)

Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität – Beispiel

48

mitarbeiter

p_id	vorname	nachname	funktion
1	Peter	Müller	Sachbearb.
5	Petra	Weger	Sekr.
...

mitarbeiter (leitend)

p_id	vorname	name
2	Stefanie	Meier
2	Petra	Weger
2	Andreas	Zwickel
...

Homonym

Fehlendes
(ableitbares)
Attribut

IC Konflikt (Eindeutigkeit)

Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität

49

Attribut-Attribut Konflikte

- Namenskonflikte
 - Verschiedene Namen für gleiche Attribute (Synonyme)
 - Gleiche Namen für verschiedene Attribute (Homonymie)
- Default-Wert-Konflikte
- IC-Konflikte
 - Datentypkonflikte
 - Bedingungskonflikte

Schematische Heterogenität – Beispiel

50

mitarbeiter			
p_id	Vorname VARCHAR(35)	nachname	alter
1	Wolfgang	Meyer	33
5	Klaus	Schmidt	NULL
...

IC: alter > 18

mitarbeiter			
p_id	Vorname VARCHAR(20)	name	alter
1	Peter	Müller	0
5	Petra	Weger	17
...

Synonym

Default Werte

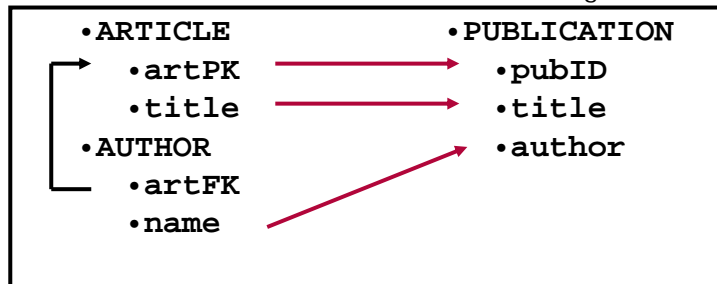
Datentypkonflikt

Schematische Heterogenität - Beispiel

51

Normalisiert vs. Denormalisiert

- 1:1 Assoziationen zwischen Werten wird unterschiedlich dargestellt
 - Durch Vorkommen im gleichen Tupel
 - Durch Schlüssel-Fremdschlüssel Beziehung



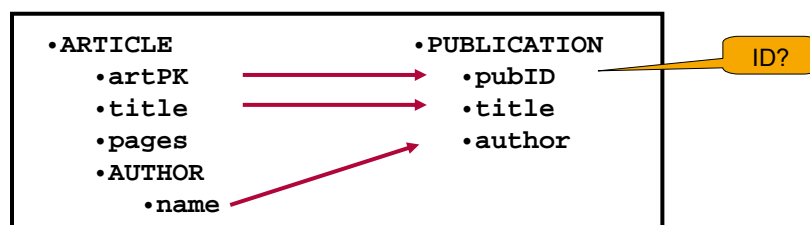
Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität - Beispiel

52

Geschachtelt vs. Flach

- 1:n Assoziationen werden unterschiedlich dargestellt
 - Als geschachtelte Elemente
 - Als Schlüssel-Fremdschlüssel Beziehung



Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität - Lösungen

53

- Problem
 - Einheitlich auf beide Schemas zugreifen
 - Auf Schemaebene: Schema Mapping und Schema-Sprachen
 - Auf Datenebene: Virtuelle Integration
 - Beide Schemas in eine gemeinsames neues Schema integrieren
 - Auf Schemaebene: Schemaintegration
 - Auf Datenebene: Materialisierte Integration
- Für die materialisierte Integration
 - Schemaintegration
 - ETL
- Für die virtuelle Integration
 - Schema-Sprachen
 - Z.B. SchemaSQL, MSQL, CPL
 - Schema Mapping
 - Z.B. Clio, RONDO, u.a.

Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität – Lösungen (Ausblick)

54

SchemaSQL [LSS96]

- Erweiterung von SQL
- Daten und Metadaten werden gleich behandelt
- Umstrukturierungen innerhalb der Anfrage
- Dynamische Sicht-Definition
- Horizontale Aggregation

```
SELECT RelA
FROM uniA->RelA, uniA::RelA A, uniB::grundgehalt B
WHERE RelA = B.institut
AND A.Kategorie = „Student“
AND A.grundgehalt > B.Student
```

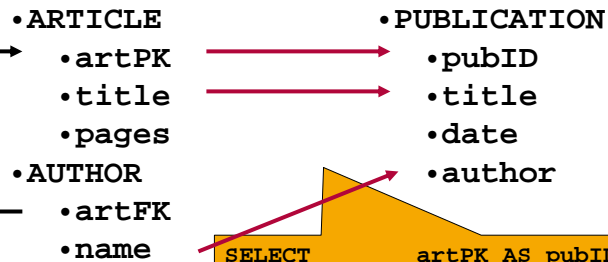
High-order Join

Felix Naumann | VL Informationsintegration | SS 2007

Schematische Heterogenität – Lösungen (Ausblick)

55

Schema Mapping



```
SELECT artPK AS pubID
       title AS title
       null AS date
       name AS author
FROM ARTICLE, AUTHOR
WHERE ARTICLE.artPK = AUTHOR.artFK
```

Felix Naumann | VL Informationsintegration | SS 2007

Zusammenfassung

56

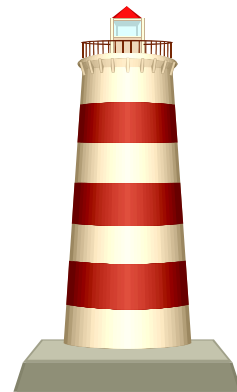
- Verteilung
- Autonomie
 - Design-Autonomie
 - Kommunikations-Autonomie
 - Ausführungs-Autonomie
- Heterogenität
 - Syntaktische Heterogenität
 - Hardware Heterogenität
 - Software Heterogenität
 - Schnittstellen Heterogenität
 - Strukturelle Heterogenität
 - Datenmodell-Heterogenität
 - Schematische Heterogenität

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

57

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Felix Naumann | VL Informationsintegration | SS 2007

Semantik

58

Fremdwörterduden "Semantik":

- Teilgebiet der Linguistik, das sich mit den Bedeutungen sprachlicher Zeichen und Zeichenfolgen befasst
- Bedeutung, Inhalt eines Wortes, Satzes oder Textes

„Semantische Heterogenität ist ein überladener Begriff ohne klare Definition. Er bezeichnet die Unterschiede in Bedeutung, Interpretation und Art der Nutzung.“

[ÖV91]

Felix Naumann | VL Informationsintegration | SS 2007

Semantik vs. Struktur

59

Strukturelle Heterogenität

- Betrifft Schemas
- Bedeutung der Labels im Schema egal
- Annahme bisher: Gleiche Label -> Gleiche Semantik

Männer(<u>Id</u> , Vorname, Nachname)	A(<u>Id</u> , X, Y)
Frauen(<u>Id</u> , Vorname, Nachname)	B(<u>Id</u> , X, Y)
Person(<u>Id</u> , Vorname, Nachname, Männlich, weiblich)	P(<u>Id</u> , X, Y, a, b)

Semantische Heterogenität

- Betrifft Daten
- Betrifft „Bedeutung“

Felix Naumann | VL Informationsintegration | SS 2007

Unterschiedliche Namen

60

- Die Probleme
 - Konzept (z.B. Gen)
 - Definition des Konzepts
 - Synonyme (z.B. surname vs. last name)
 - Homonyme (z.B. biweekly)
 - Einheiten (z.B. cm vs. inch)
 - Werte (z.B. „manager“)
- Eher auf Schema Ebene

Felix Naumann | VL Informationsintegration | SS 2007

Konzept

61

- Definition eines Konzepts
 - Noch nicht einmal hier sind sich immer alle einig.
 - Gen, Transaktion, Bestellung, Mitarbeiter
- Semantisch überlappende Weltausschnitte mit einander entsprechenden Klassen
- Korrespondenzarten zwischen Klassenextensionen:
 - $A=B$ Äquivalenz
 - $A \subseteq B$ Inklusion
 - $A \cap B$ Überlappung
 - $A \neq B$ Disjunktion

Felix Naumann | VL Informationsintegration | SS 2007

Konzept

62

„Wie viele Mitarbeiter hat IBM?“

- Definition Mitarbeiter:
 - temporäre MA
 - Diplomanden
 - Berater
 - Studentische Mitarbeiter
 - Stellen oder Köpfe?
- Definition IBM
 - Welche Region? Welcher Geschäftsbereich?
 - Informix?
 - PWC?
- Welcher Zeitpunkt?
- Definition der Zählung:
 - Doppelte Zählung bei mehreren Anstellungen?

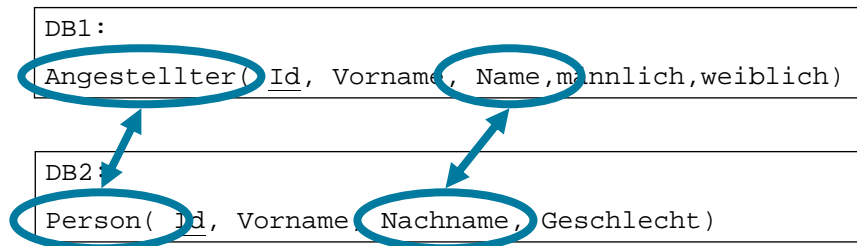
„Wieviele Hardware haben wir ans HPI verkauft?“

Felix Naumann | VL Informationsintegration | SS 2007

Synonyme

63

- Verschiedene Worte mit gleicher Bedeutung
- Im Kontext der zu integrierenden Datenbanken

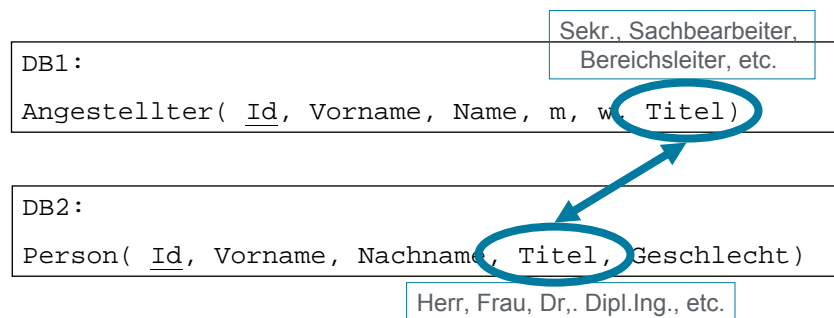


Felix Naumann | VL Informationsintegration | SS 2007

Homonymie

64

- Gleiche Worte verschiedener Bedeutung
- Andere Domäne
- Andere Bedeutung



Felix Naumann | VL Informationsintegration | SS 2007

-nym Wörter

65

- **Synonym**
 - Verschiedene Wörter, gleiche Semantik
- **Homonym**
 - Gleiche Wörter, verschiedene Semantik
- **Antonym**
 - Verschiedene Wörter, gegenteilige Semantik
- **Auto-Antonym:**
 - Gleiche Wörter, gegenteilige Semantik
 - Transparenz
 - Overlook
- **Heteronym**
 - Gleiche Schreibung, verschiedene Aussprache, verschiedene Semantik
- **Autonym (selbstbeschreibend, Wort = Semantik, „Substantiv“)**
- **Pseudonym u.v.a.m.**

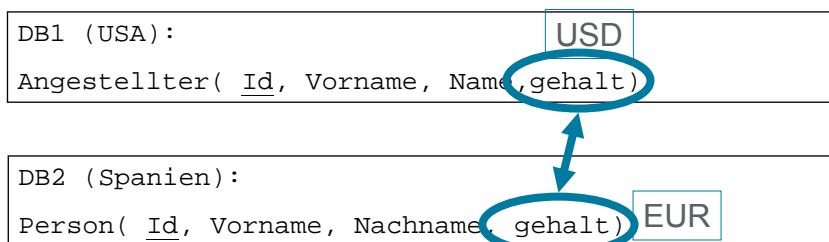
- http://www.fun-with-words.com/nym_words.html

Felix Naumann | VL Informationsintegration | SS 2007

Einheiten

66

- Gleiche „Bedeutung“ aber anderes Maß.
- Werden auch als Homonym bezeichnet, da anderes Maß eine andere Bedeutung erzeugt.

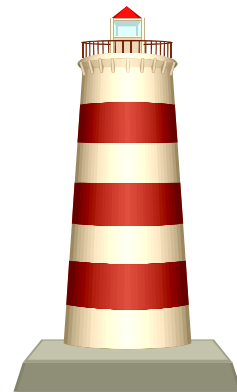


Felix Naumann | VL Informationsintegration | SS 2007

Überblick

67

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung

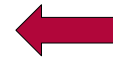


Felix Naumann | VL Informationsintegration | SS 2007

Identität

68

- Drei zentrale Fragen
 - Was ist ein Objekt?
 - XML: Über mehrere Schachtelungsebenen hinweg
 - Relationales Modell: Über mehrere Relationen hinweg
 - Repräsentiert Objekt A die gleiche Entität wie Objekt B?
 - Wie finde ich effizient gleiche Repräsentationen?
- Namen des Problems
 - Duplikaterkennung
 - Objektidentifikation
 - Record Linkage
 - Data Cleansing
 - ...
- Auf Datenebene



Felix Naumann | VL Informationsintegration | SS 2007

Typische Anwendungen

69

- Personen- und Adressdaten
 - Volkszählungen
 - Werbeaktionen
 - Kundenpflege
- Molekularbiologische Daten
- Bibliographische Daten
 - Zentrale Register
- Typische Merkmale zur Entstehung:
 - Gleiches Objekt mehrfach beobachtet
 - Manuelle Erfassung der Daten
 - Objekt ändert Eigenschaften von Zeit zu Zeit
 - Keine global konsistente ID
 - ISBN, IBAN, URL, ISO, EAN, SSN, etc.



Felix Naumann | VL Informationsintegration | SS 2007

Duplikaterkennung

70

- Duplikate in Relationen
 - Zwei Tupel, die das gleiche real-world Objekt repräsentieren
 - Semantik!
 - Attributwerte dürfen sich unterscheiden.
- Formales Problem
 - Eine Tabelle (der Größe N), potentiell mit Duplikaten
 - Erzeuge für jedes Tupel einen Identifier, so dass Duplikate gleiche Identifier erhalten
- Problemerweiterungen
 - Zwei Tabellen mit unterschiedlichem Schema
 - Ein XML Dokument mit Duplikaten

Felix Naumann | VL Informationsintegration | SS 2007

Duplikaterkennung

71

- Praktisches Problem
 - Wie entscheide ich, ob zwei Tupel das gleiche Objekt repräsentieren?
 - Ähnlichkeitsmaße und Klassifikation
 - Edit-Distance
 - N-grams
 - IDs
 - Wahrscheinlichkeitstheoretische Ansätze
 - Maschinelles Lernen
 - Augenschein

Felix Naumann | VL Informationsintegration | SS 2007

Duplikaterkennung

72

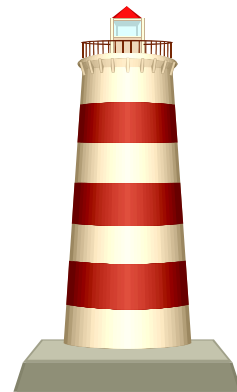
- Praktisches Problem
 - Sehr große Datenmenge
 - Millionen Tupel
 - Kein quadratischer Algorithmus
 - Kein Hauptspeicher-Algorithmus
- Als SQL Anfrage
 - Sei R die Relation mit Duplikaten
 - `SELECT C1.*, genID(C1,C2)` ← ID Erzeugung
`FROM R as C1, R as C2` ← Kreuzprodukt
`WHERE M(C1,C2)` ← Ähnlichkeit
 - Schwieriger als normaler Join
 - Ähnlichkeitsmaß ist nicht nur Gleichheit
 - Algorithmen zur Objektidentifikation in VL „Duplikaterkennung“

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

73

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte

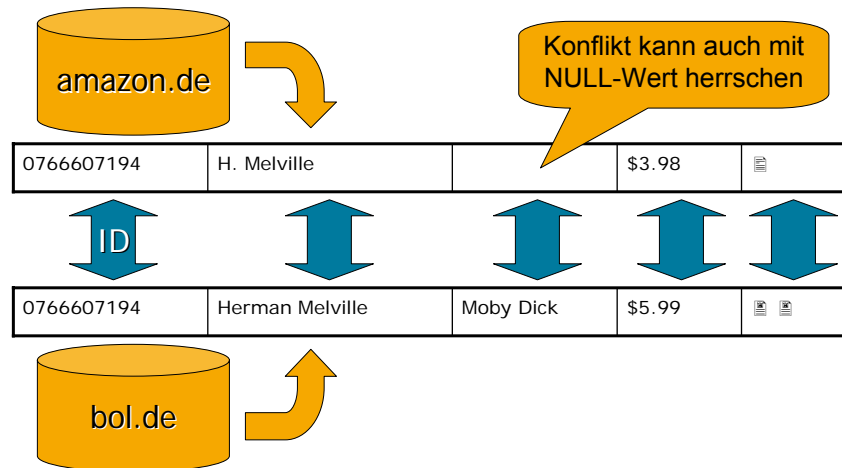
74

- Datenkonflikt:
 - Zwei Duplikate haben unterschiedliche Attributwerte für ein semantisch gleiches Attribut.
 - Im Gegensatz zu Konflikten mit Integritätsbedingungen
- Datenkonflikte entstehen
 - innerhalb eines Informationssystems (intra-source) und
 - bei der Integration mehrerer Informationssysteme (inter-source).
- Voraussetzung:
 - Duplikat!
 - d.h. Identität schon festgestellt.

Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte - Beispiel

75



Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte – Entstehung

76

Innerhalb eines Informationssystems

- Mangels Integritätsbedingungen oder Konsistenz-Checks
- Bei redundanten Schemata
- Bei Entstehung von Duplikaten
- Nicht korrekte Einträge
 - Tippfehler, Übertragungsfehler
 - Falsche Rechenergebnisse
- obsolete Einträge
 - div. Aktualisierungszeitpunkte
 - ausreichende Aktualität einer Quelle
 - verzögerte Aktualisierung
 - vergessene Aktualisierung

Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte – Entstehung

77

Innerhalb eines Informationssystems

- bei div. Datentypen (mit/ohne Codierung)
 - 1,2,...,5 bzw. "sehr gut", "gut", ..., mangelhaft"
- bei gleichem Datentyp
 - Schreibvarianten
 - Kantstr. Kantstrasse Kant Str. Kant Strasse
 - Kolmogorov Kolmogoroff Kolmogorow
 - Typische Verwechslungen U<->V,
O<->o, usw. (OCR)

Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte – Behebung

78

- Referenztabelle für exakte Wertabbildung
 - Z.B. Städte, Länder, Produktnamen, Codes...
- Ähnlichkeitsmaße
 - bei Tippfehlern
 - bei Sprachvarianten (Meier, Mayer,...)
- Standardisieren und transformieren
- Nutzung von Hintergrundwissen (Metadaten)
 - bzgl. von Konventionen (landestypische Schreibweisen)
 - Ontologien zur Behandlung von Zusammenhängen
 - Thesauri, Wörterbücher zur Behandlung von Homonymen, Synonymen, ...

Felix Naumann | VL Informationsintegration | SS 2007

Datenkonflikte – Entstehung

79

Bei der Integration von Informationssystemen

- Lokal konsistent aber global inkonsistent
- Duplikate (extensionale Redundanz)
- Andere Datentypen
- Lokale Schreibweisen/Konventionen

Datenkonflikte – Behebung

80

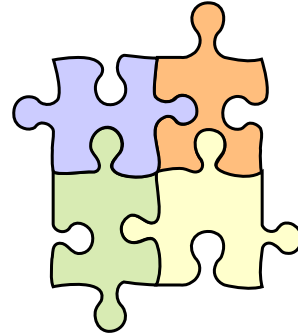
- Präferenzordnung über Datenquellen
 - nach Aktualität, Trust (Vertrauen), Öffnungszeiten usw.
- Informationsqualität
- Konfliktlösungsfunktionen

- Wie implementieren?

Relationale Objektintegration

81

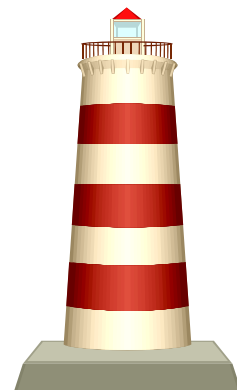
- Union (Vereinigung)
 - Duplikat-Eliminierung
 - Minimum Union
 - Eliminierung sub-summierter Tupel
 - ...
 - Aber keine
 - Duplikatintegration
 - Konfliktlösung
- Mehr dazu in VL „Datenfusion“



Überblick

82

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Gebundene & Freie Variablen

83

Gebundene Variablen müssen bei einer Anfrage spezifiziert werden.

- z.B.: „Search“-Feld bei Google

Freie Variablen müssen nicht gebunden werden.

- z.B. „Autor“-Feld bei Amazon.de, falls Titel gebunden ist.

Einordnung:

- Heterogenität
 - Syntaktische Heterogenität
 - Schnittstellenheterogenität

Gebundene und Freie Variablen – Adornments

84

- Jede Quelle exportiert eine oder mehrere relationale Sichten.
- IIS erlaubt Anfragen auf diese Sichten mittels Join, Union, Selektion und Projektion.



Quelle: [YLGU99]

Gebundene und Freie Variablen – Adornments

Z gebunden

85

Beispiel Quelle 1:

$R_1(X,Y,Z)$

Daten:

(x_1,y_1,z_1)

(x_1,y_2,z_1)

(x_2,y_2,z_2)

Beispiel Anfrage 1:

$Q_1(X,Y,z_1)$

Beispiel Anfrage 2:

$Q_2(X,y_1,Z)$

Beispiel Ergebnis:

(x_1,y_1,z_1)

(x_1,y_2,z_1)

Beispiel Ergebnis:

(x_1,y_1,z_1)

5 Quellen (für später):

$R_1(X,Y,Z)$

$R_2(X,Y,Z)$

$R_3(X,Y,Z)$

$R_4(Z,U)$

$R_5(U,V,W)$

Quelle: [YLGU99]

Gebundene und Freie Variablen – Adornments

86

- Anfragefähigkeiten der Quellen als templates
 - Wie ein WWW Formular
 - Templates bestehen aus einem adornment für jedes Attribut
- Anhänge (adornments = Verzierungen) an Attribute schränken ein:
 - f: free
 - Frei: Kann in Anfrage spezifiziert werden, muss aber nicht.
 - u: unspecifiable
 - Unbestimmbar: Kann nicht spezifiziert werden.
 - Ist aber Teil des Ergebnisses.
 - b: bound
 - Gebunden: Muss spezifiziert werden.
 - c[s]: constant
 - Auswahl aus einer Menge s von Konstanten
 - Implizit bound: muss spezifiziert werden
 - o[s]: optional
 - Auswahl aus einer Menge s von Konstanten
 - Implizit free: Muss nicht spezifiziert werden.

Quelle: [YLGU99]

Adornments - Beispiele

87

Beispiel Quelle 1:
 $R_1(X,Y,Z)$

Anfragemöglichkeit 1:
X muss spezifiziert werden
Y kann nicht spezifiziert werden
Z kann spezifiziert werden

Template:
buf

Anfragemöglichkeit 2:
X kann nicht spezifiziert werden
Y kann spezifiziert werden
Z ist entweder z_1 oder z_2

Template:
ufc[z_1, z_2]

Felix Naumann | VL Informationsintegration | SS 2007

Adornments – Anfragebearbeitung

88

Anfragebearbeitung

- $R_1(X,Y,Z)$: bff, ffb
- $R_2(X,Y,Z)$: fbf
- Sei $M = R_1 \cup R_2$ eine integrierte Sicht des IIS, gegen die man Anfragen stellen kann.
- Annahme über Anfragebearbeitung:
 - Anfragen werden übersetzt in je eine Anfrage pro Quelle (gebundene Variablen werden weitergereicht)
 - Ergebnisse werden entsprechen der Sicht verknüpft (hier \cup)
- Frage: Was ist das Template der Sicht M?

$$\begin{array}{c} \mathbf{bff} \\ \cup \\ \mathbf{fbf} \\ = \end{array}$$

$$\begin{array}{c} \mathbf{ffb} \\ \cup \\ \mathbf{fbf} \\ = \end{array}$$

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

3 Sichten und deren Adornments:

$R_1(X,Y,Z)$: bff, ffb

$R_2(X,Y,Z)$: fbf

$R_3(X,Y,Z)$: ffc[s₁], c[s₂]ff

$R_1 \cup R_2$:

$\text{bff} \cup \text{fbf} = \text{bbf}$

$\text{ffb} \cup \text{fbf} = \text{fbb}$

$(R_1 \cup R_2) \cup R_3$:

$\text{bbf} \cup \text{ffc}[s_1] = \text{bbc}[s_1]$ usw.

	f	o[s ₃]	b	c[s ₄]	u
f	f	o[s ₃]	b	c[s ₄]	u
o[s ₁]	o[s ₁]	o[s ₁ ∩ s ₃]	c[s ₁]	c[s ₁ ∩ s ₄]	u
b	b	c[s ₃]	b	c[s ₄]	-
c[s ₂]	c[s ₂]	c[s ₂ ∩ s ₃]	c[s ₂]	c[s ₂ ∩ s ₄]	-
u	u	u	-	-	u

Quelle: [YLGU99]

■ Unterschied zu UNION

- Nicht jedes Attribut der integrierten Sicht ist auch Attribut jeder beteiligten Quelle.
- Beispiel: $R_1(X,Y,Z)$ und $R_4(Z,U)$
- Sicht: $M(X,Y,Z,U) = R_1(X,Y,Z) \bowtie R_4(Z,U)$

■ Berechnung des Templates der Sicht

- Adornments der nicht-Join-Attribute werden kopiert.
- Adornments der Join-Attribute werden gemäß der UNION Tabelle vereint.

Adornments – Selektion und Projektion

91

- Selektion
 - Sicht im IIS selektiert mit Prädikaten.
 - $X = \text{'Test'}$ oder $U > 1999$
 - Prädikate werden auf Ergebnisse der Quellen angewandt.
 - Deshalb: Kein Einfluss auf adornments
- Projektion
 - Einfach projizierte Attribute weglassen.
 - Aber: Falls Attribut mit b oder c adornment durch Projektion wegfallen soll => Sicht des IIS nicht ausführbar
 - Sonst: Adornments bleiben erhalten

Adornments – Anfragebearbeitung

92

Problem

- UNION-Matrix zu restriktiv

	f	o[s₃]	b	c[s₄]	u
f	f	o[s ₃]	b	c[s ₄]	u
o[s₁]	o[s ₁]	o[s ₁ ∩ s ₃]	c[s ₁]	c[s ₁ ∩ s ₄]	u
b	b	c[s ₃]	b	c[s ₄]	-
c[s₂]	c[s ₂]	c[s ₂ ∩ s ₃]	c[s ₂]	c[s ₂ ∩ s ₄]	-
u	u	u	-	-	u

Idee: Erhöhung der Menge beantwortbarer Anfragen

- durch Post-Processing
- durch Passing Bindings

93

$R_1(X,Y,Z)$: bfu
 $R_2(X,Y,Z)$: buf
 $R_1 \cup R_2 = buu$

Anfrage 1: (x_1, Y, Z) beantwortbar?
 Anfrage 2: (x_1, y_1, z_1) beantwortbar?

Idee: (x_1, y_1, Z) an R_1
 (x_1, Y, z_1) an R_2
 Dann im Mediator filtern:
 $Z=z_1$ bzw. $Y=y_1$

Was ist neu?

$u = f$: durch nachträgliches Filtern (postprocessing)
 $o[s] = f$: falls Bindung nicht in s, weglassen und später

Filtern

Zusammen: $R_1 \cup R_2 = bff$

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

94

	f	o[s₃]	b	c[s₄]	u
Vorher:	f	o[s ₃]	b	c[s ₄]	u
o[s₁]	o[s ₁]	o[s ₁ ∩ s ₃]	c[s ₁]	c[s ₁ ∩ s ₄]	u
b	b	c[s ₃]	b	c[s ₄]	-
c[s₂]	c[s ₂]	c[s ₂ ∩ s ₃]	c[s ₂]	c[s ₂ ∩ s ₄]	-
u	u	u	-	-	u
Nachher:	f	f	b	c[s ₄]	f
o[s₁]	f	f	b	c[s ₄]	f
b	b	b	b	c[s ₄]	b
c[s₂]	c[s ₂]	c[s ₂]	c[s ₂]	c[s ₂ ∩ s ₄]	c[s ₂]
u	f	f	b	c[s ₄]	f

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

JOIN über templates ohne *passing bindings*

$R_1(X,Y,Z) : \text{fbf}$
 $R_5(Z,U) : \text{bf}$
 $R_1 \bowtie R_2 = \text{fbff}$

Anfrage 1: (X,y_1,z_1,U) beantwortbar?
 Anfrage 2: (X,y_1,Z, U) beantwortbar?

Idee: (X,y_1,Z) an R_1
 $(z_1,U) \dots (z_n,U)$ an R_5

Passing Bindings: Ergebnisse einer Sicht werden vom Mediator in die gebundene Variable der nächsten Sicht eingetragen.

JOIN über templates mit *passing bindings*:
 $R_1 \triangleright \triangleleft R_5 = \text{fbff}$

Quelle: [YLGU99]

Vorher:

	f	o[s₃]	b	c[s₄]	u
f	f	o[s ₃]	b	c[s ₄]	u
o[s₁]	o[s ₁]	o[s ₁ ∩ s ₃]	c[s ₁]	c[s ₁ ∩ s ₄]	u
b	b	c[s ₃]	b	c[s ₄]	-
c[s₂]	c[s ₂]	c[s ₂ ∩ s ₃]	c[s ₂]	c[s ₂ ∩ s ₄]	-
u	u	u	-	-	u

Zweite Quelle

Nachher:

	f	o[s₃]	b	c[s₄]	u
f	f	f	f	c[s ₄]	f
o[s₁]	f	f	f	c[s ₄]	f
b	b	b	b	c[s ₄]	b
c[s₂]	c[s ₂]	c[s ₂]	c[s ₂]	c[s ₂ ∩ s ₄]	c[s ₂]
u	f	f	f	c[s ₄]	f

Erste Quelle

Quelle: [YLGU99]

Variante 1: $R_1(X,Y,Z), X < x_1$: bfu

$Q(x_2, Y, Z)$ beantwortbar?
 $Q(x_2, Y, z_1)$ beantwortbar?
 $Q(X, y_1, z_1)$ beantwortbar?

bfu wird zu bff mit postprocessing

Variante 2: $R_1(X,Y,Z), X = x_1$: bfu

$Q(X, y_1, z_1)$ beantwortbar?

$Q(X, y_1, z_1) = Q(x_1, y_1, z_1)$
 wegen Prädikat
 bfu wird zu bff wird zu fff

Quelle: [YLGU99]

Vorher Nachher

Base View Adornment	Sel. Attribute Adornment
f	f
$o[s_1]$	f
b	f or b
$c[s_1]$	f or $c[s_1]$
u	f

Quelle: [YLGU99]

Viele Templates

99

Problem: Quellen exportieren oft mehrere templates

- Beispiel: Amazon (Autor Titel, Schlagwort, ISBN, Verlag)
- bffff, fbfff, ffbff, fffb, fffb
- Beispiel: Verlage (Verlag, Ort)
- bf, fb
- Sicht im IIS: Amazon \bowtie Verlag Verlage
- Templates der Sicht aus jeder Kombination:
 - bffff, fbfff, ffbff, fffb, fffb
 - bffffb, fbffb, ffbfb, fffb, fffb
 - + fffffb (ffffb \bowtie fb mit passing binding)

Erweiterte Suche Bücher

Sie können auch nur eines der Felder ausfüllen.
Bitte geben Sie in die unten stehenden Suchfelder einen oder mehrere Suchbegriffe ein und klicken Sie auf "Suchen".

Autor/Se: _____

Titel: _____

Schlagwörter: _____

ISBN: _____

Verlag: _____

Lösung:

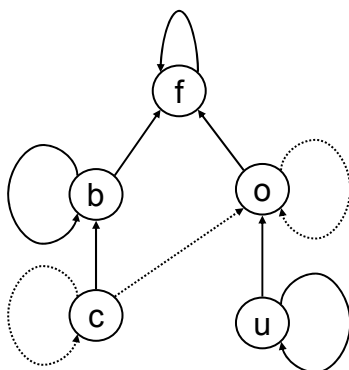
- Einige templates sind redundant

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

Redundanz in Templates

100



- Weniger restriktiv
- Weniger restriktiv falls Auswahllisten Teilmengen sind

- bffff, fbfff, ffbff, fffb, fffb
- bffffb, fbffb, ffbfb, fffb, fffb

Algorithmus zur Entfernung redundanter templates.

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

Adornments – Fallbeispiel

101

- Amazon
 - Formular 1: Mindestens eine Spezifikation aus author, title, subject, format (format aus Auswahlliste)
 - Formular 2: ISBN spezifizieren
 - Formular 3: Mindestens eine Spezifikation aus keyword, publisher, date
 - Antwortrelation: author, title, ISBN, publisher, date, format, price, shipping info
- Barnes & Noble
 - Formular 1: Mindestens eine Spezifikation aus author, title, keywords; optionale Spezifikation in format subject, price, age (alles aus Auswahllisten)
 - Formular 2: ISBN spezifizieren

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

Adornments - Fallbeispiel

102

Amazon

author	title	format	subject	KW	ISBN	pub	date	price	ship
b	f	o	f'	u'	u	u	u	u	u
f	b	o	f'	u'	u	u	u	u	u
f	f	c	f'	u'	u	u	u	u	u
f	f	o	b'	u'	u	u	u	u	u
u	u	u	u'	u'	b	u	u	u	u
u	u	u	u'	f'	u	b	f	u	u
u	u	u	u'	b'	u	f	f	u	u
u	u	u	u'	f'	u	f	b	u	u

Barnes & Noble

author	title	format	subject	KW	ISBN	pub	date	price	ship	age
f	b	o	o'	f'	u	u	u	o	u	o'
b	f	o	o'	f'	u	u	u	o	u	o'
f	f	o	o'	b'	u	u	u	o	u	o'
u	u	u	u'	u'	b	u	u	u	u	u'

IIS

author	title	format	subject	KW	ISBN	pub	date	price	ship	age
f	f	f	u'	u'	b	f	f	f	f	u'
f	f	f	u'	b'	f	f	f	f	f	u'
b	f	f	o'	u'	f	f	f	f	f	u'
f	b	f	o'	u'	f	f	f	f	f	u'
b	f	f	u'	f'	f	b	f	f	f	u'
f	b	f	u'	f'	f	b	f	f	f	u'
f	b	f	u'	f'	f	f	b	f	f	u'
b	f	f	u'	f'	f	f	b	f	f	u'

Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

Adornments - Fallbeispiel

103

	author	title	format	subject	KW	ISBN	pub	date	price	ship	age	
IIS	f	f	f	u'	u'	b	f	f	f	f	u'	
	f	f	f	u'	b'	f	f	f	f	f	u'	
	b	f	f	o'	u'	f	f	f	f	f	u'	
	f	b	f	o'	u'	f	f	f	f	f	u'	
	b	f	f	u'	f'	f	b	f	f	f	u'	
	f	b	f	u'	f'	f	f	f	f	f	u'	
	f	b	f	u'	f'	f	f	b	f	f	u'	
	b	f	f	u'	f'	f	f	b	f	f	u'	
		f	f	f	u'	f'	f	f	f	f	f	u'

Ableiten von 4 Formularen im IIS nach [YLGU99]

- Spezifikation der ISBN (template 1)
- Spezifikation des keyword (template 2)
- Mindestens author oder title spezifizieren (templates 3 und 4)
- Mindestens author oder title und mindestens publisher oder date spezifizieren (templates 5-8)

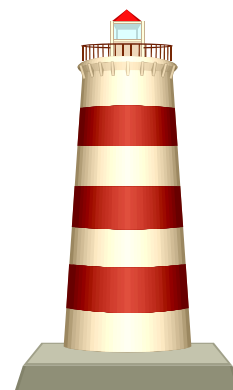
Quelle: [YLGU99]

Felix Naumann | VL Informationsintegration | SS 2007

Überblick

104

- Verteilung
- Autonomie
- Syntaktische Heterogenität
- Strukturelle Heterogenität
- Semantische Heterogenität
 - Namenskonflikte
 - Identität
 - Datenkonflikte
- Gebundene und Freie Variablen
 - Adornments
 - Anfrageplanung



Felix Naumann | VL Informationsintegration | SS 2007

Gebundene & Freie Variablen – Beispiel

105

Sources	View Schemas	Must Bind
1	v1(<u>Song</u> , CD)	Song
2	v2(CD, Artist, Price)	CD
3	v3(CD, <u>Artist</u> , Price)	Artist

Query: “Find the prices of CDs containing a song titled *Friends*.”

v1(*Friends*, CD) ⋈ v2(CD, Artist, Price)
v1(*Friends*, CD) ⋈ v3(CD, Artist, Price)

Quelle: [LC00]

Gebundene & Freie Variablen – Beispiel

106



v1(Song, CD)
<Friends, Life>
<Friends, Love>

v2(CD, Artist, Price)
<Love, Lucy, \$15>
<Story, Snoopy, \$14>

Bastelaufgabe 1:
Wie teuer ist die billigste CD mit einem Song namens “Friends”?

v3(CD, Artist, Price)
<Story, Lucy, \$13>
<Love, Snoopy, \$10>
<Life, Charlie, \$8>

Quelle: [LC00]

Gebundene & Freie Variablen – Beispiel

107

$v_1(\text{Song}, \text{CD})$
 <Friends, Life>
 <Friends, Love>

$v_2(\text{CD}, \text{Artist}, \text{Price})$
 <Love, Lucy, \$15>
 <Story, Snoopy, \$14>

Bastelaufgabe 2:
 Welches ist die billigste CD mit einem Song namens "Friends", die Sie anfragen können?

$v_3(\text{CD}, \text{Artist}, \text{Price})$
 <Story, Lucy, \$13>
 <Love, Snoopy, \$10>
 <Life, Charlie, \$8>

Gebundene & Freie Variablen – Beispiel

108

$v_1 \bowtie v_2: \{\$15\}$

$v_1(\text{Song}, \text{CD})$
 <Friends, Life>
 <Friends, Love>

$v_2(\text{CD}, \text{Artist}, \text{Price})$
 <Love, Lucy, \$15>
 <Story, Snoopy, \$14>

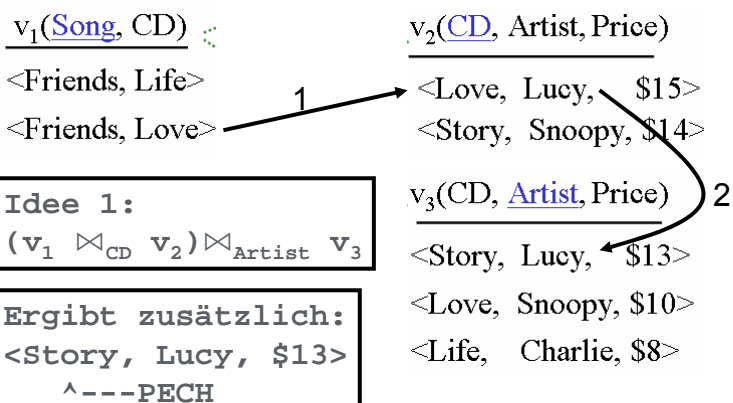
$v_1 \bowtie v_3$: empty, no binding for Artist.

$v_3(\text{CD}, \text{Artist}, \text{Price})$
 <Story, Lucy, \$13>
 <Love, Snoopy, \$10>
 <Life, Charlie, \$8>

Quelle: [LC00]

Gebundene & Freie Variablen – Beispiel

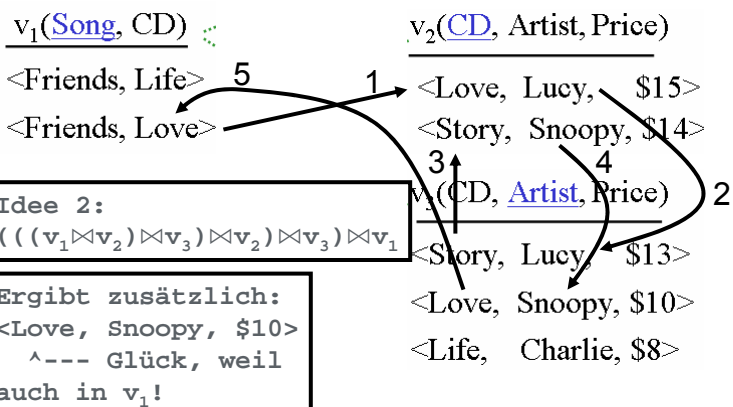
109



Felix Naumann | VL Informationsintegration | SS 2007

Gebundene & Freie Variablen – Beispiel

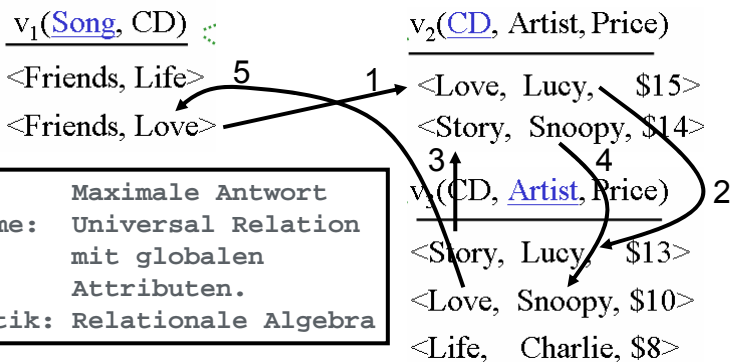
110



Felix Naumann | VL Informationsintegration | SS 2007

Gebundene & Freie Variablen – Beispiel: Semantik

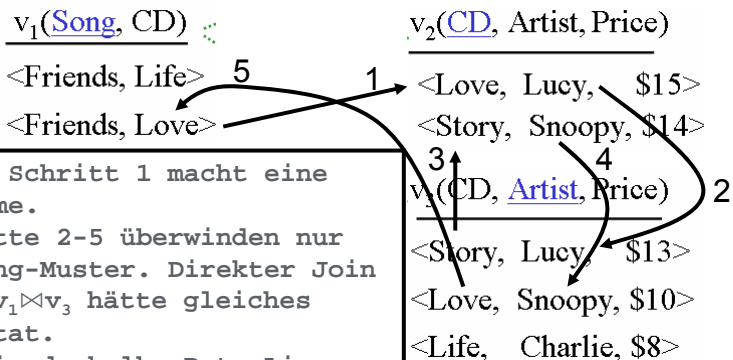
11



Ziel: Maximale Antwort
Annahme: Universal Relation mit globalen Attributen.
Semantik: Relationale Algebra

Gebundene & Freie Variablen – Beispiel: Semantik

12



Schon Schritt 1 macht eine Annahme.
 Schritte 2-5 überwinden nur Binding-Muster. Direkter Join über $v_1 \bowtie v_3$ hätte gleiches Resultat.
 Wichtig deshalb: Data Lineage und Visualisierung

Wichtigste Literatur für heute

- [BKLW99] Busse, Kutsche, Leser, Weber, Federated Information Systems: Concepts, Terminology and Architectures. Forschungsbericht 99-9 des FB Informatik der TU Berlin, 1999.
Online: http://www.informatik.hu-berlin.de/~leser/publications/tr_terminology.ps
- [ÖV99] Principles of Distributed Database Systems
M. Tamer Özsu, Patrick Valduriez, Prentice Hall, (1991/1999).
Kapitel 1 und 4
- [YLGU99] Ramana Yerneni, Chen Li, Hector Garcia-Molina, Jeffrey D. Ullman, „Computing Capabilities of Mediators“, SIGMOD 1999

Weitere Literatur

- [Con97] Föderierte Datenbanksysteme. Konzepte der Datenintegration
Stefan Conrad, Springer Verlag, 1997
- [LC00] Chen Li, Edward Chang „Query Planning with Limited Source Capabilities“, ICDE 2000