

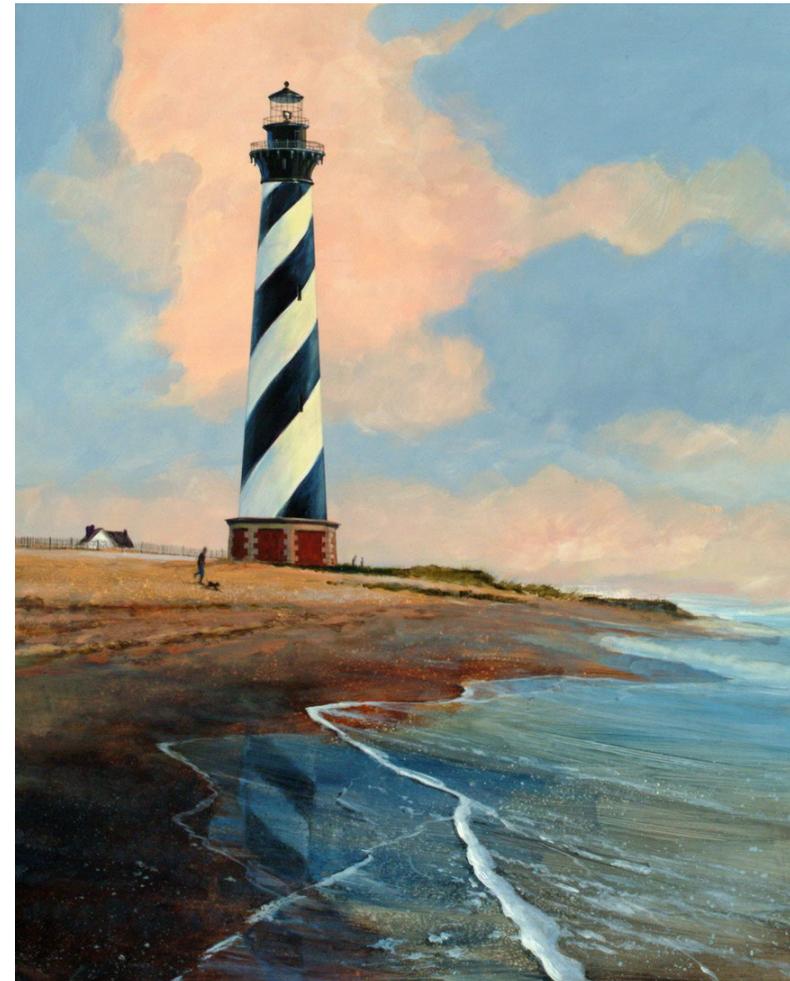


Information Integration  
Verteilung, Autonomie und Heterogenität

24.10.2019  
Felix Naumann

# Überblick

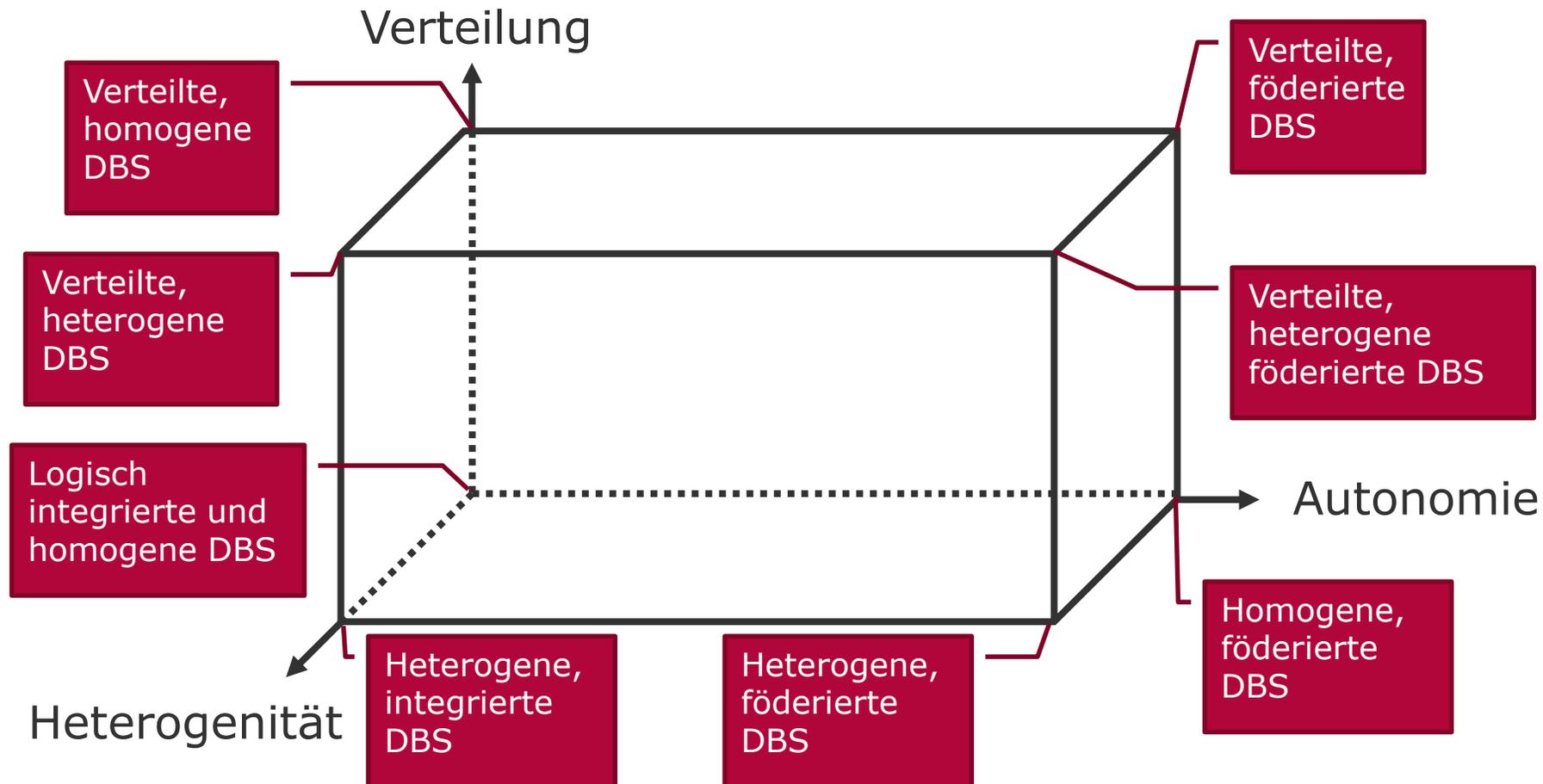
1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung



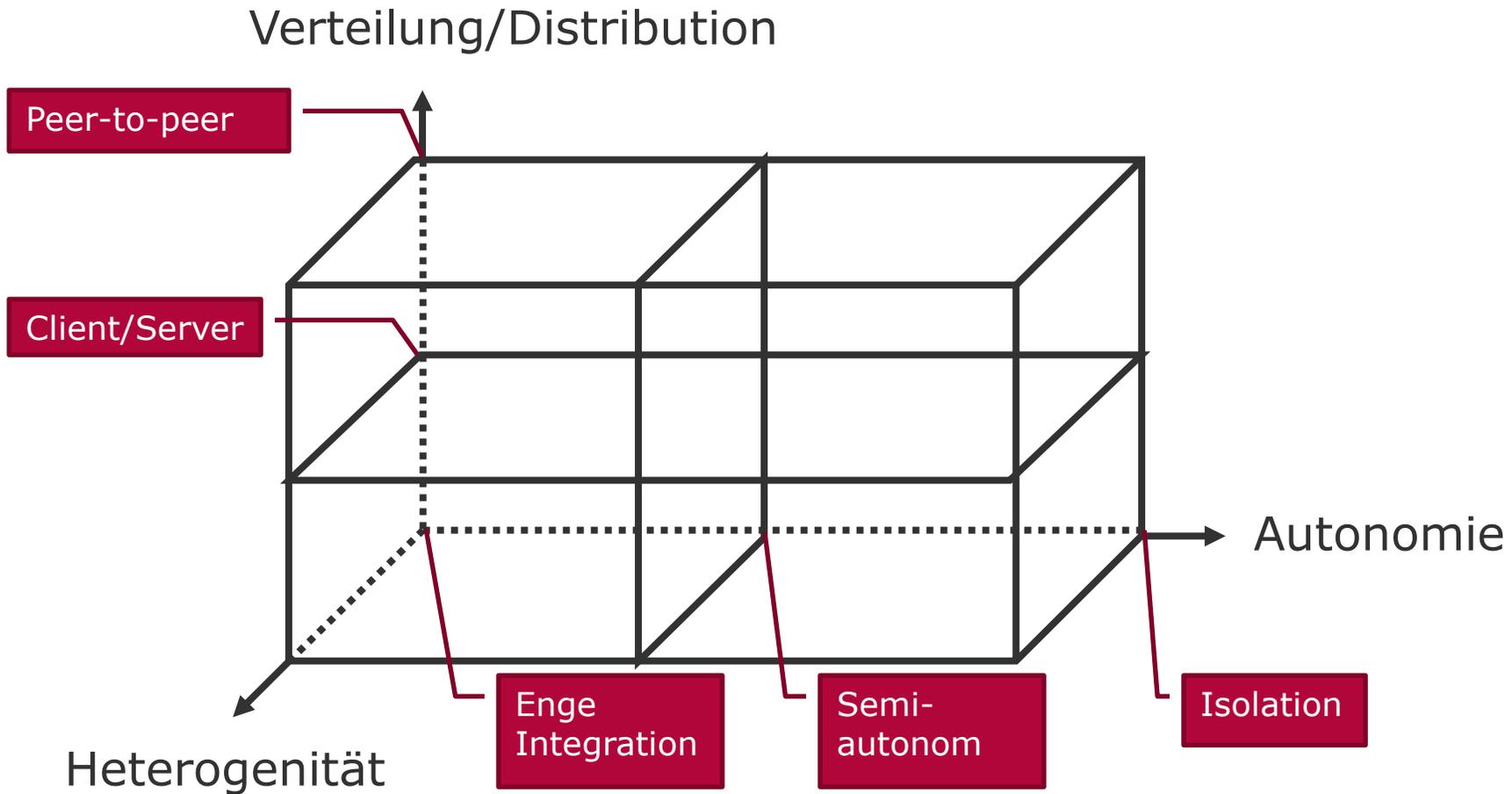
Felix Naumann  
Information Integration  
Winter 2019/20

# Klassifikation von Informationssystemen [ÖV99]

■ Drei orthogonale Dimensionen



# Klassifikation verteilter DBMS [ÖV99]



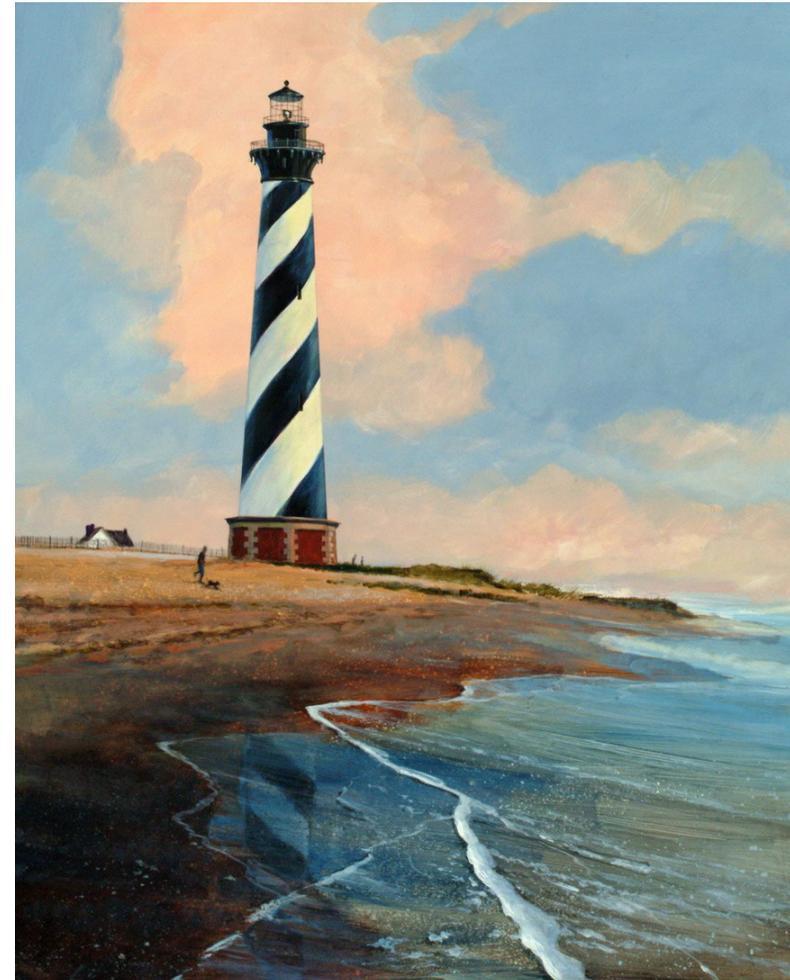
## Zusammenhang mit Föderierten DBMS

---

- 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!

# Überblick

- 1. Verteilung**
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung

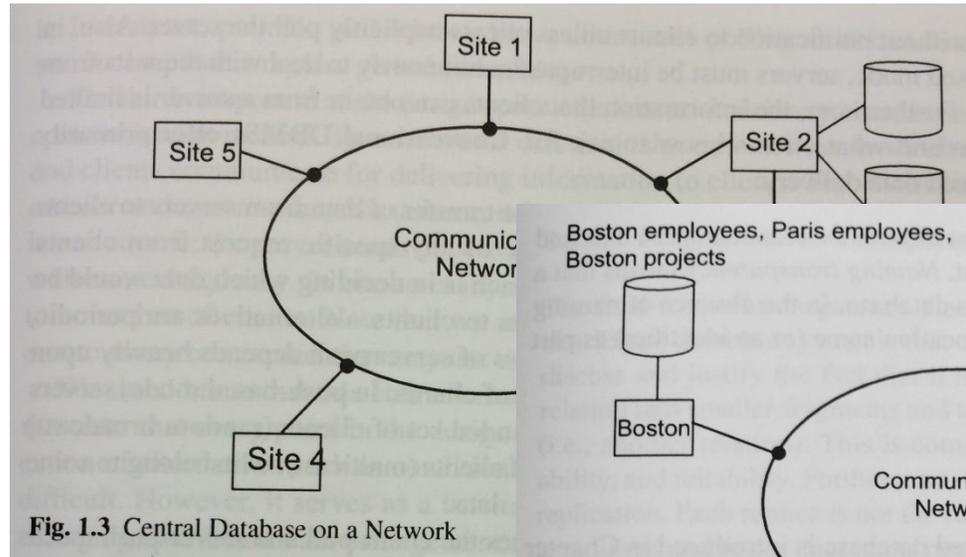


Felix Naumann  
Information Integration  
Winter 2019/20

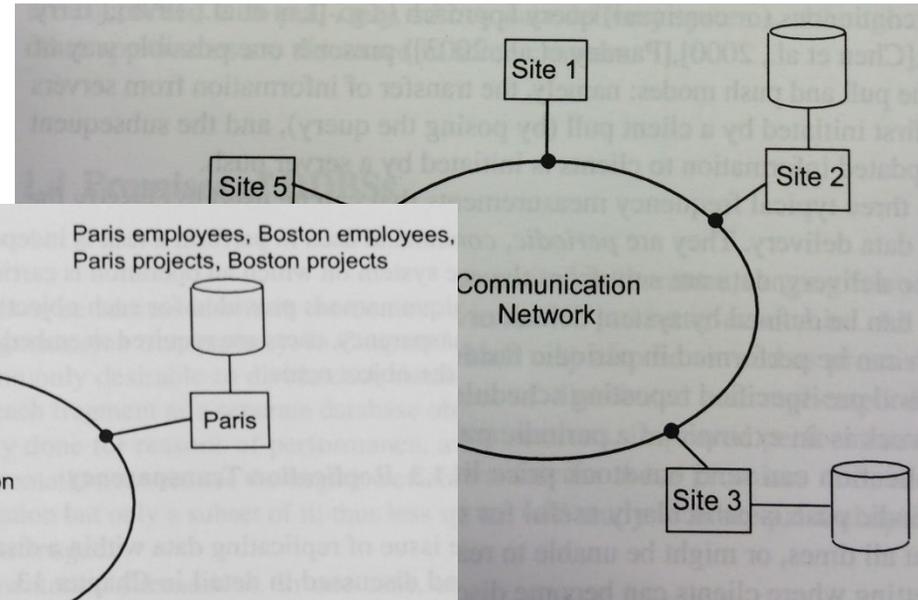
# Verteilung (Distribution)

A distributed database is a collection of multiple, logically interrelated databases distributed over a computer network.

[ÖV91]



**Fig. 1.3** Central Database on a Network



**Fig. 1.5** A Distributed Application

Felix Naumann  
Information Integration  
Winter 2019/20

## Physikalische Verteilung

---

- 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

## Logische Verteilung

---

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

## Verteilung – Vor- und Nachteile

---

### **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

## Verteilung – Techniken

- HTTP, CORBA, RPC, Grid, ... nicht hier.
  - Anwendungsentwicklung ohne Spezifikation der physikalischen Präsenz der Komponenten
- Annahmen an Transparenz
  - Datenunabhängigkeit (jedes DBMS)
    - auch Speicherorttransparenz
  - Netzwerktransparenz
  - Replikationstransparenz
  - Fragmentationstransparenz
    - auch Partitionierungstransparenz

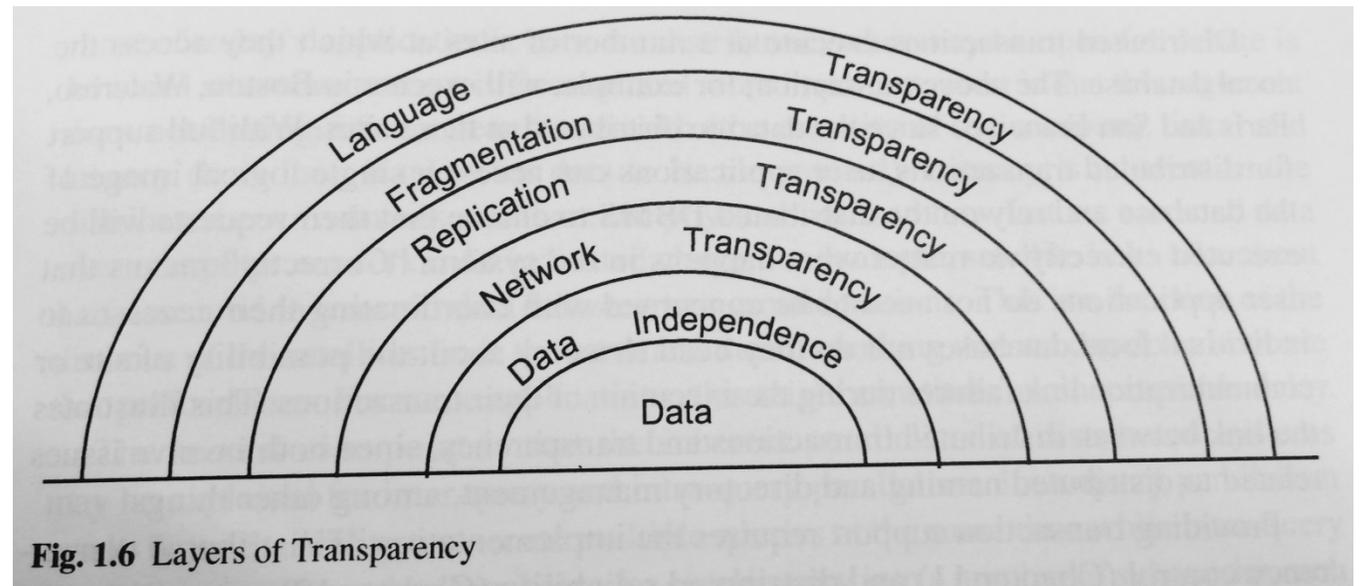
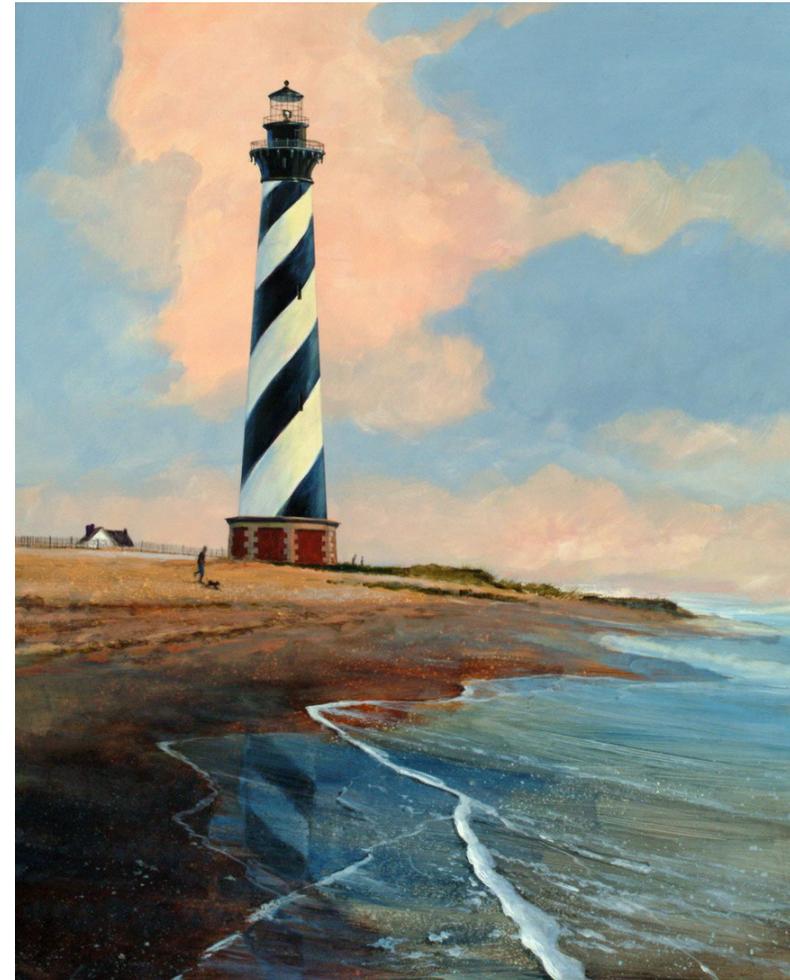


Fig. 1.6 Layers of Transparency

# Überblick

1. Verteilung
2. **Autonomie**
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung



Felix Naumann  
Information Integration  
Winter 2019/20

## Autonomie (Autonomy)

---

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

---

- 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!

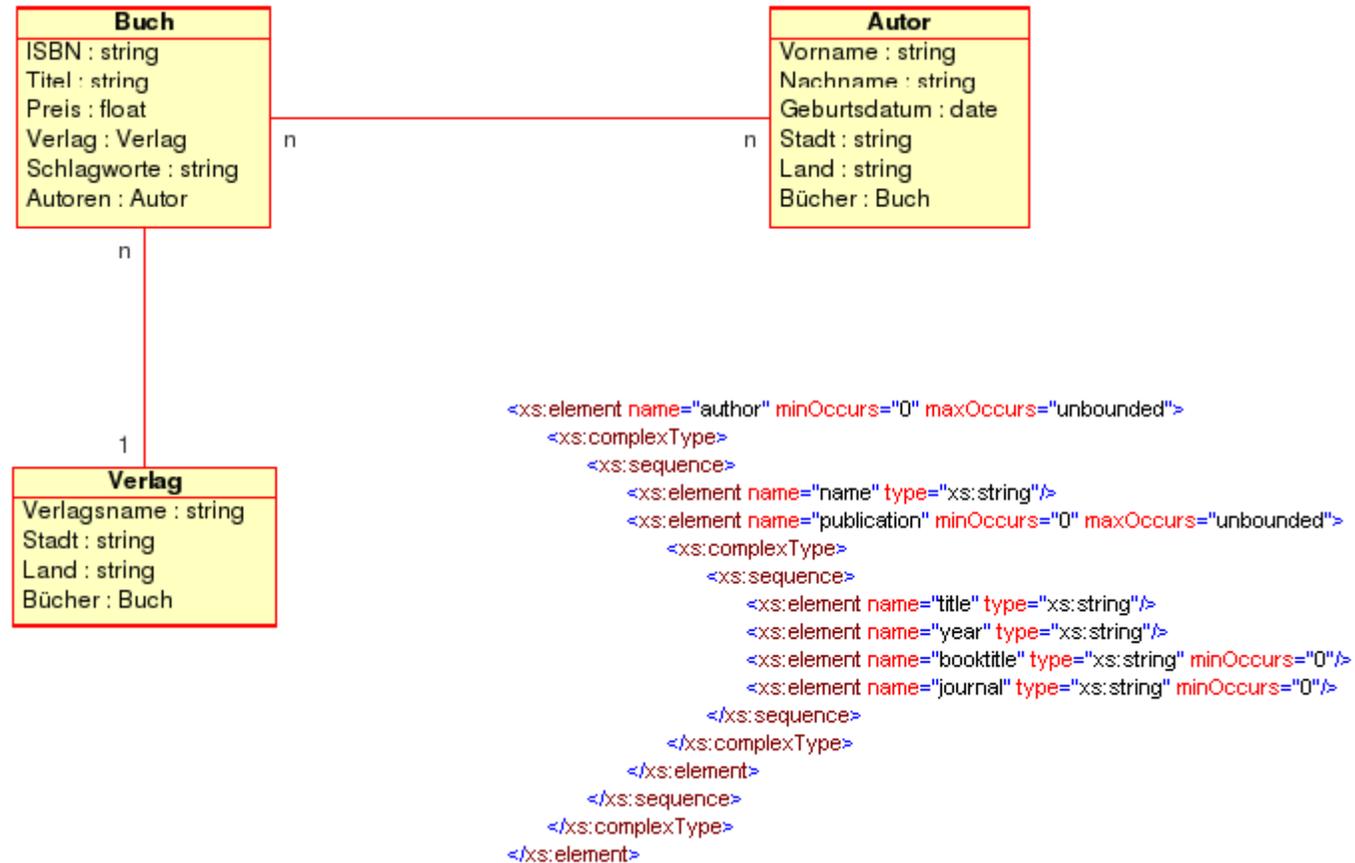
# Design-Autonomie – Beispiel

- Schema und Datenmodell 1

- (Fast) relational
- Flach

- Schema und Datenmodell 2

- XML
- hierarchisch



Felix Naumann  
Information Integration  
Winter 2019/20

## Kommunikations-Autonomie

---

- 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
    - ...

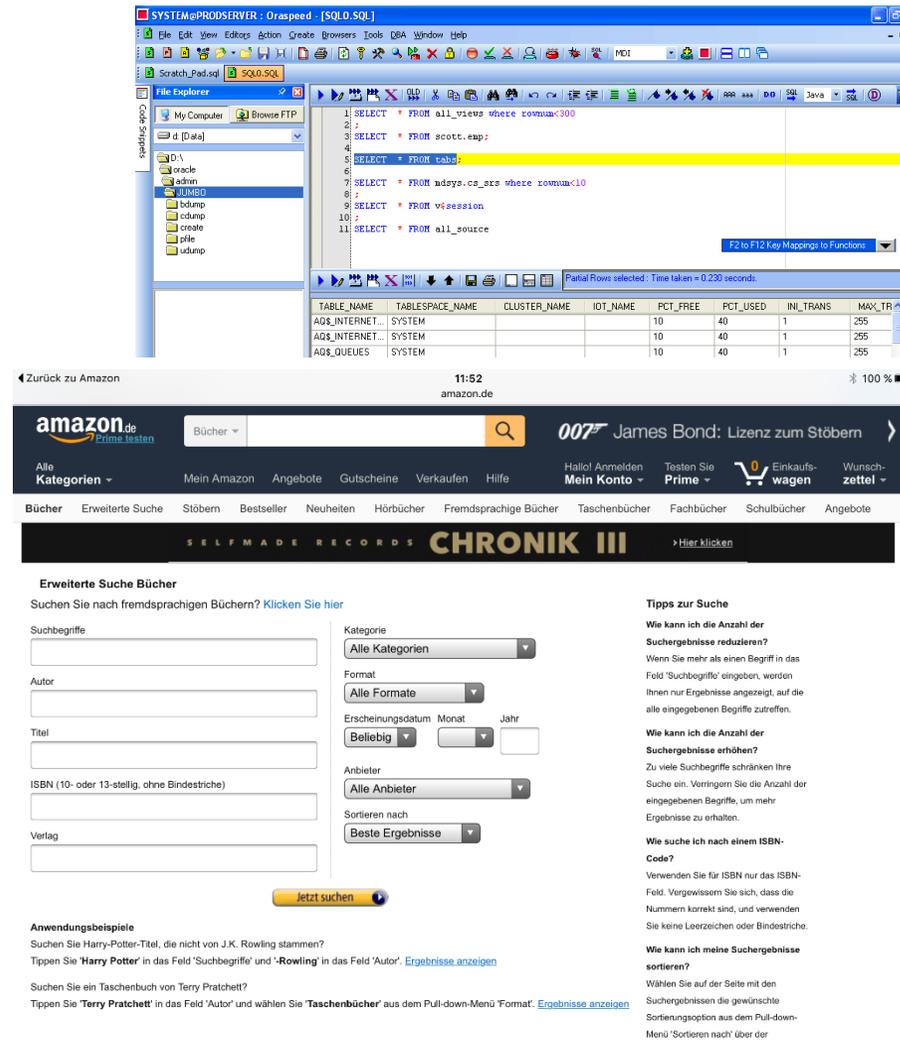
# Kommunikations-Autonomie – Beispiel

## ■ 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)



The top screenshot shows a SQL IDE window titled 'SYSTEM@PRODSEVER: Oraspeed [SQL.O.SQL]'. It contains a query with 11 lines of SQL code. The results table below shows columns: TABLE\_NAME, TABLESPACE\_NAME, CLUSTER\_NAME, IOT\_NAME, PCT\_FREE, PCT\_USED, INL\_TRANS, and MAX\_TR. The data rows are:

TABLE_NAME	TABLESPACE_NAME	CLUSTER_NAME	IOT_NAME	PCT_FREE	PCT_USED	INL_TRANS	MAX_TR
AG\$_INTERNET...	SYSTEM			10	40	1	255
AG\$_INTERNET...	SYSTEM			10	40	1	255
AG\$_QUEUES	SYSTEM			10	40	1	255

The bottom screenshot shows the Amazon.de website with an 'Erweiterte Suche Bücher' form. The form includes fields for 'Suchbegriffe', 'Autor', 'Titel', and 'ISBN (10- oder 13-stellig, ohne Bindestriche)'. It also has dropdown menus for 'Kategorie', 'Format', 'Anbieter', and 'Sortieren nach'. A 'Jetzt suchen' button is at the bottom. To the right of the form are 'Tipps zur Suche' instructions.

Felix Naumann  
Information Integration  
Winter 2019/20

## Ausführungs-Autonomie

---

- Informationssystem 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

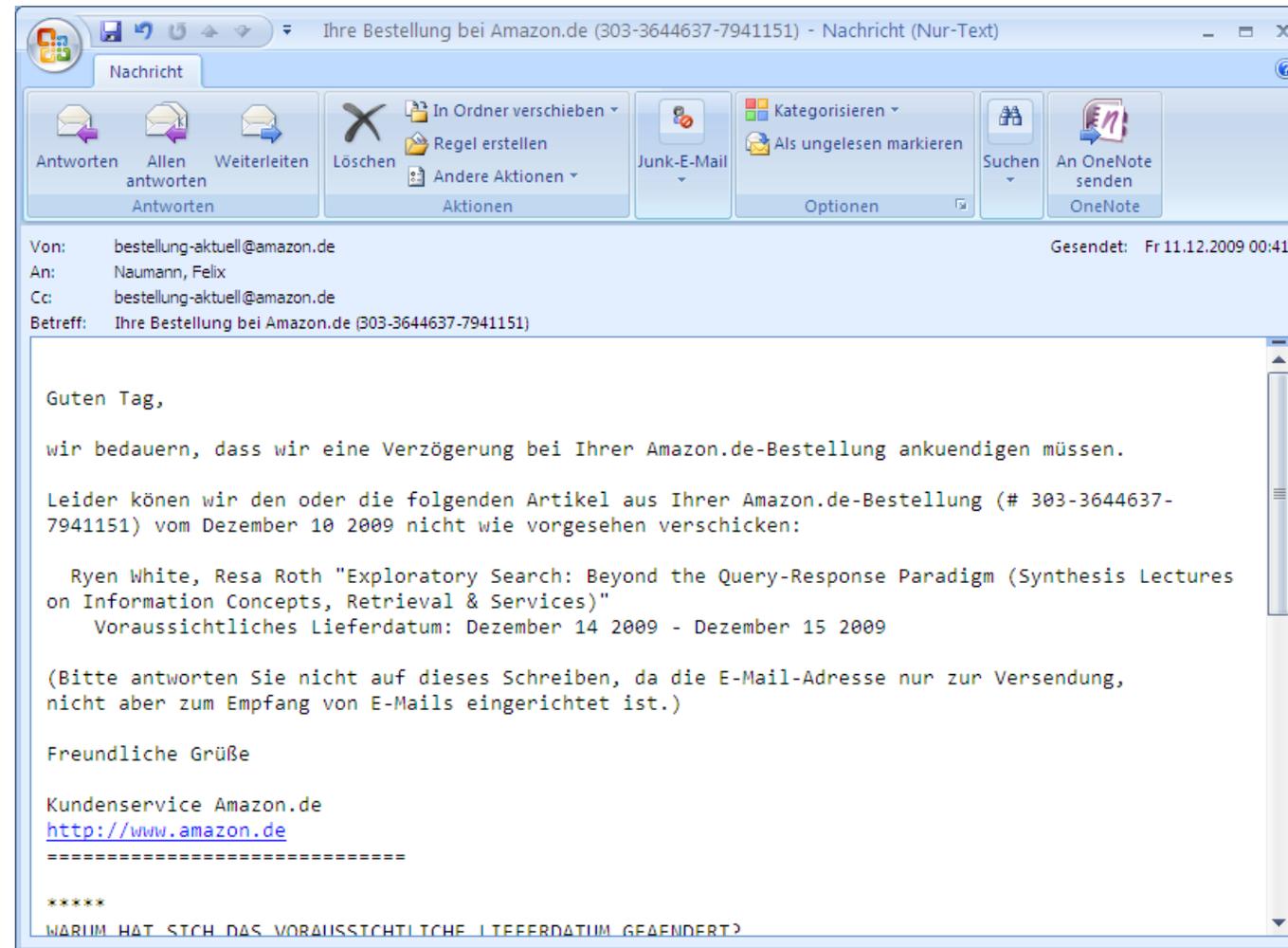
# Ausführungs-Autonomie – Beispiel

- Optimierung und Scheduling

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

- Transaktionen

- Dirty-read egal?
- Amazon-Beispiel



Felix Naumann  
 Information Integration  
 Winter 2019/20

## Autonomie → Heterogenität

---

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

## Heterogenität (*Heterogeneity*)

---

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

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

**Heterogenitäten zu überbrücken ist die Kernaufgabe der Informationsintegration.**

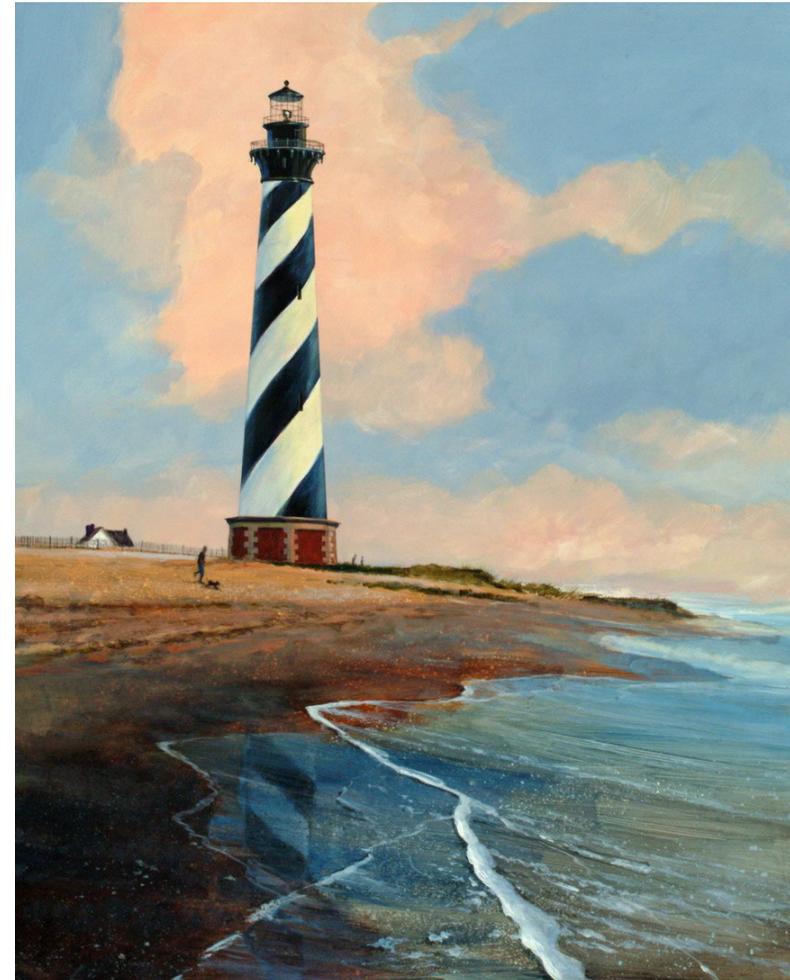
## Weitere Heterogenitätsklassifikationen

---

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

## Überblick

1. Verteilung
2. Autonomie
- 3. Syntaktische Heterogenität**
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung

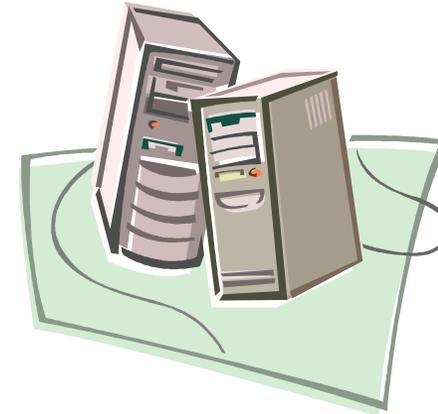
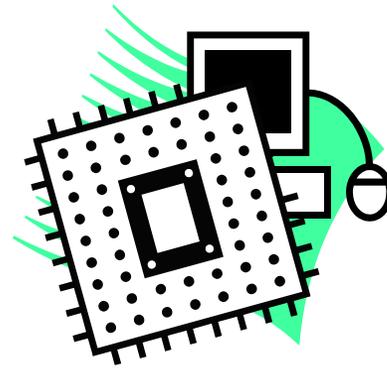


Felix Naumann  
Information Integration  
Winter 2019/20

# Syntaktische Heterogenität

- Hardware-Heterogenität

- Bandbreite
- Hauptspeicher
- CPU
  - Art
  - Geschwindigkeit



**Nicht hier**

- Software-Heterogenität

- Schnittstellen-Heterogenität

# Software Heterogenität

---

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



## Software Heterogenität – Beispiel

---

```
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

Felix Naumann  
Information Integration  
Winter 2019/20

## Schnittstellen Heterogenität

---

- Schnittstellen von Informationssystemen sind im wesentlichen deren Anfragensprache:
  - HTML Formular,
  - „Google“-Sprache (+, - , ...),
  - SQL,
  - XQuery,
  - etc.

**Jetzt hier!**

## Schnittstellen Heterogenität

---

- 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)
  - Oder auch joins
- Gebundene und freie Variablen [RSU95,LC00,YLGU99]
  - später
- Andere Einschränkungen
  - Joins über maximal 3 Relationen
  - z.B. Prädikate nur über eine Auswahl von Werten

# Schnittstellen Heterogenität - Beispiel

The screenshot displays the Netscape 7.0 email client interface. The main window title is "VLDB-Teilnahme - Posteingang für naumann@informatik.hu-berlin.de - Netscape 7.0". The search bar at the top contains the text "integration". An orange arrow labeled "Suche" points to this search bar. Below the search bar, a search dialog box titled "Nachrichten durchsuchen" is open. It shows search criteria for "Lokale Ordner" with the option "Untergeordnete Ordner durchsuchen" checked. The search criteria are set to "Mindestens einer der folgenden Faktoren muss zutreffen". Three criteria are listed: "Betreff" "enthält" "integration", "Betreff" "enthält nicht" "Humboldt", and "Betreff" "enthält" (empty). A dropdown menu is open for the third criterion, showing options: "enthält" (checked), "enthält nicht", "gleich", "ungleich", "beginnt mit", and "endet mit". An orange arrow labeled "Konjunktion/Disjunktion" points to the dropdown menus for the search criteria. Another orange arrow labeled "gleich/ungleich" points to the open dropdown menu.

Felix Naumann  
Information Integration  
Winter 2019/20

# Schnittstellen-Heterogenität – Beispiel

amazon.de

WUNSCHZETTEL | MEIN KONTO | HILFE

HOME FELIX' SHOP **BÜCHER** ENGLISH BOOKS ELEKTRONIK & FOTO **NEU** KÜCHE & HAUSHALT MUSIK DVD VIDEO SOFTWARE COMPUTER & VIDEOSPIELE

ERWEITERTE SUCHE | STÖBERN | BESTSELLER | NEUHEITEN | FACHBÜCHER | ZEITSCHRIFTEN | PREIS-HITS | GEBRAUCHT

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

[Suche in Zeitschriften](#) | [Suchtipps](#)

### 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 "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:

Erscheinungsdatum:  (z.B. 1999)

Suche in:  Gebundene Bücher  E-Books

Gebundene Variablen

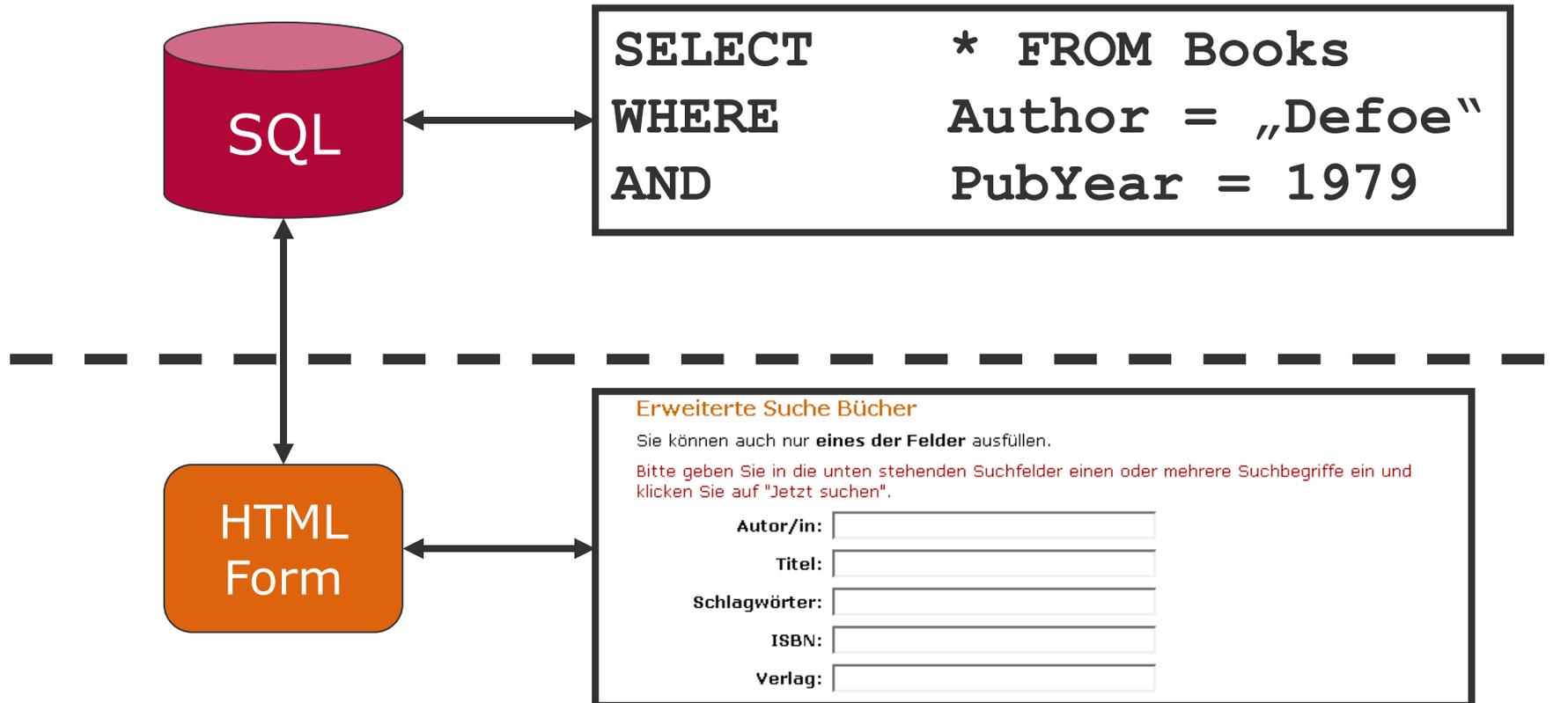
Prädikat nur mit  
Auswahl von Werten

## Schnittstellen Heterogenität

---

- 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

# Mächtige globale Anfragesprache



Felix Naumann  
Information Integration  
Winter 2019/20

# Mächtige globale Anfragesprache

```
SELECT * FROM Books  
WHERE Author = „Defoe“  
AND PubYear = 1979
```

```
Daniel Defoe, Robinson Crusoe, 1979
```



```
PubYear = 1979
```

```
Daniel Defoe, Robinson Crusoe, 1986  
Daniel Defoe, Robinson Crusoe, 1979  
Daniel Defoe, Moll Flanders, 1933
```



## Erweiterte Suche Bücher

Sie können auch nur **eines der Felder** ausfüllen.

Bitte geben Sie in die unten stehenden Suchfelder einen oder mehrere Begriffe ein und klicken Sie auf "Jetzt suchen".

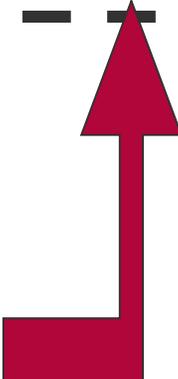
Autor/in:

Titel:

Schlagwörter:

ISBN:

Verlag:



# Mächtige globale Anfragesprache

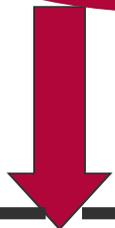
```
SELECT * FROM Books  
WHERE Author = „Defoe“  
AND PubYear > 1979
```

Daniel Defoe, Robinson Crusoe, 1986



PubYear > 1979

Daniel Defoe, Robinson Crusoe, 1986  
Daniel Defoe, Robinson Crusoe, 1979  
Daniel Defoe, Moll Flanders, 1933



Erweiterte Suche Bücher

Sie können auch nur eines der Felder ausfüllen.

Bitte geben Sie in die unten stehenden Suchfelder einen oder mehrere Begriffe ein und klicken Sie auf "Jetzt suchen".

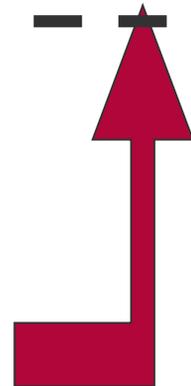
Autor/in:

Titel:

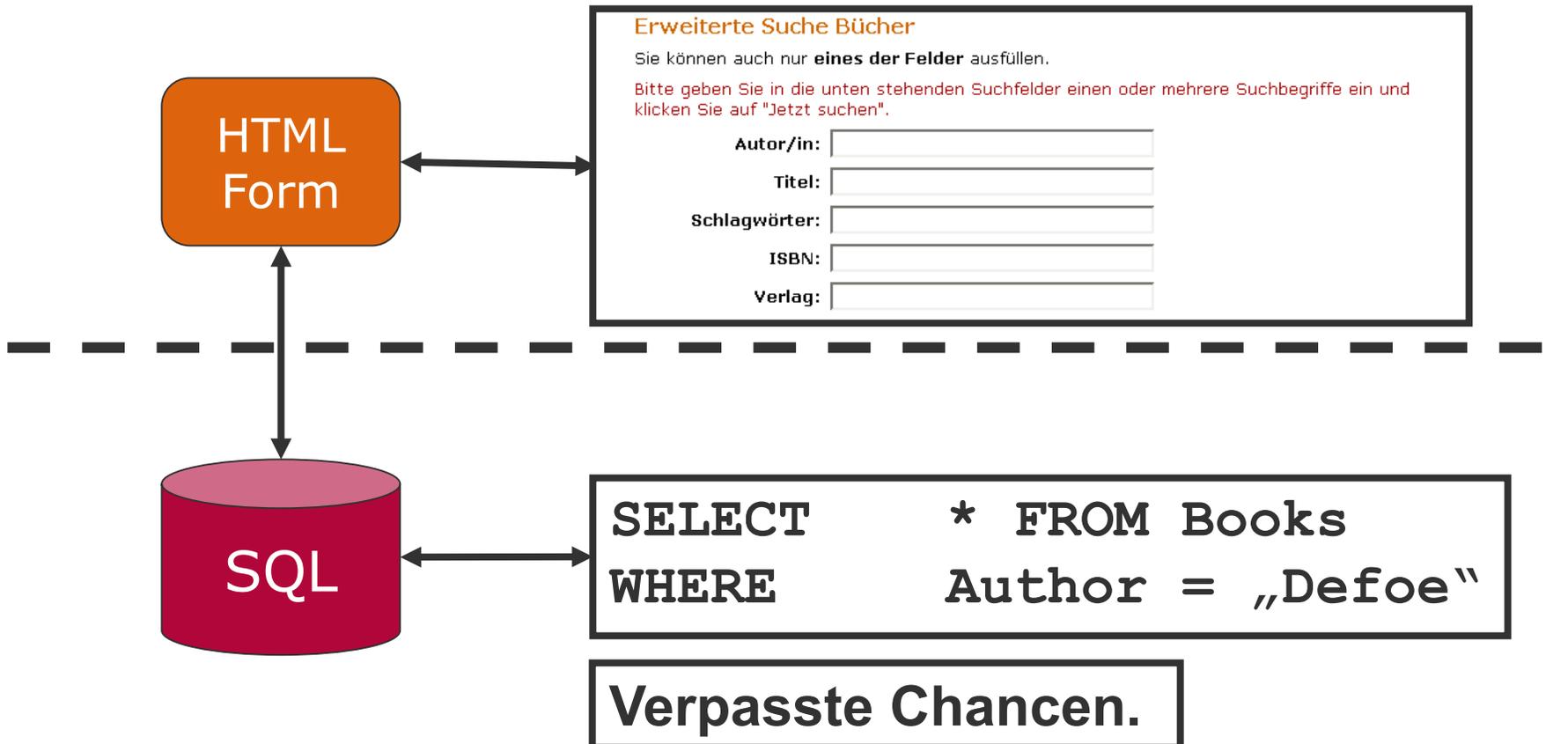
Schlagwörter:

ISBN:

Year:



# Mächtige lokale Anfragesprache



## Gebundene & Freie Variablen

---

- **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

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

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

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

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

# Gebundene & Freie Variablen – Beispiel & Ausblick

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...

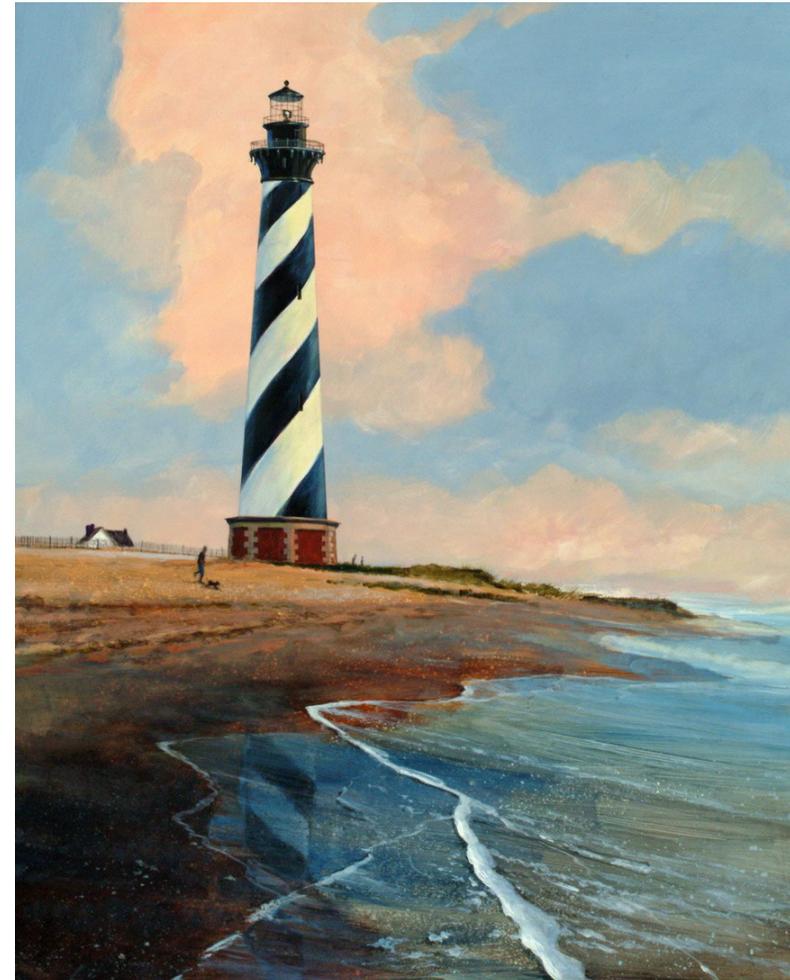
## Syntaktische Heterogenität - Zusammenfassung

---

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

# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. **Strukturelle Heterogenität**
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung



Felix Naumann  
Information Integration  
Winter 2019/20

## Strukturelle Heterogenität

---

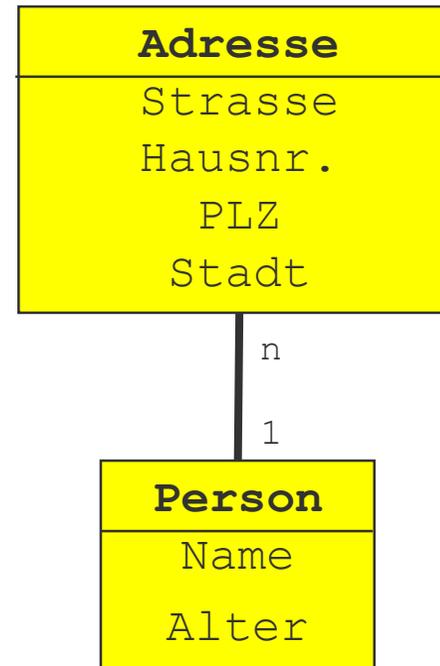
- Datenmodell-Heterogenität
  - Unterschiedliche Semantik
  - Unterschiedliche Struktur
  
- Schematische Heterogenität
  - Integritätsbedingungen, Schlüssel, Fremdschlüssel, etc.
  - Schema (Attribut vs. Relation etc.)
  - Struktur (Gruppierung in Tabellen)

# Datenmodell-Heterogenität

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

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

`Person ( Id, Name, Alter)`

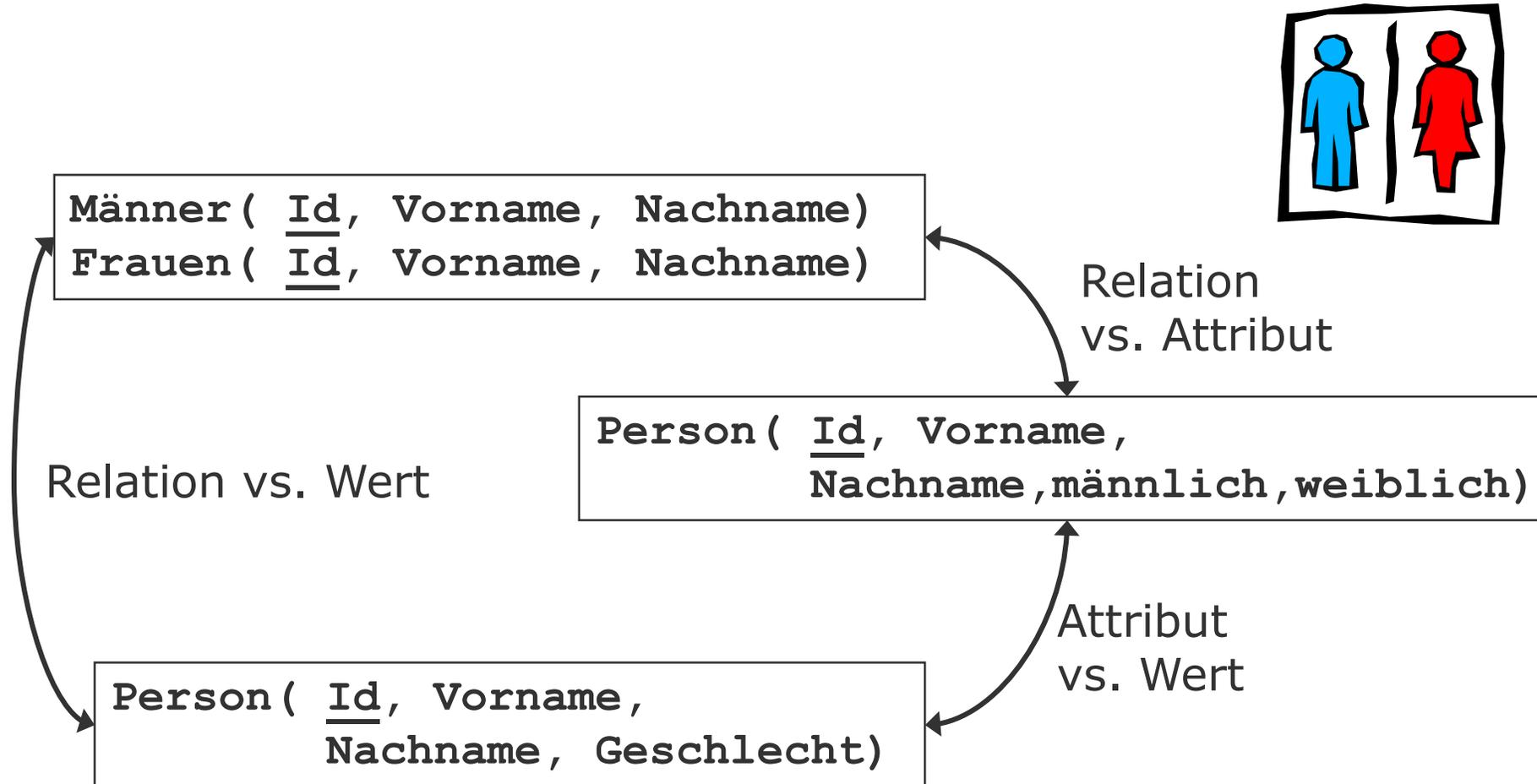


# Schematische Heterogenität

---

- 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!

# Schematische Heterogenität



## Schematische Heterogenität

---

- 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)

# Schematische Heterogenität – Beispiel

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)

## Schematische Heterogenität

---

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

# Schematische Heterogenität – Beispiel

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
...	...	...	...

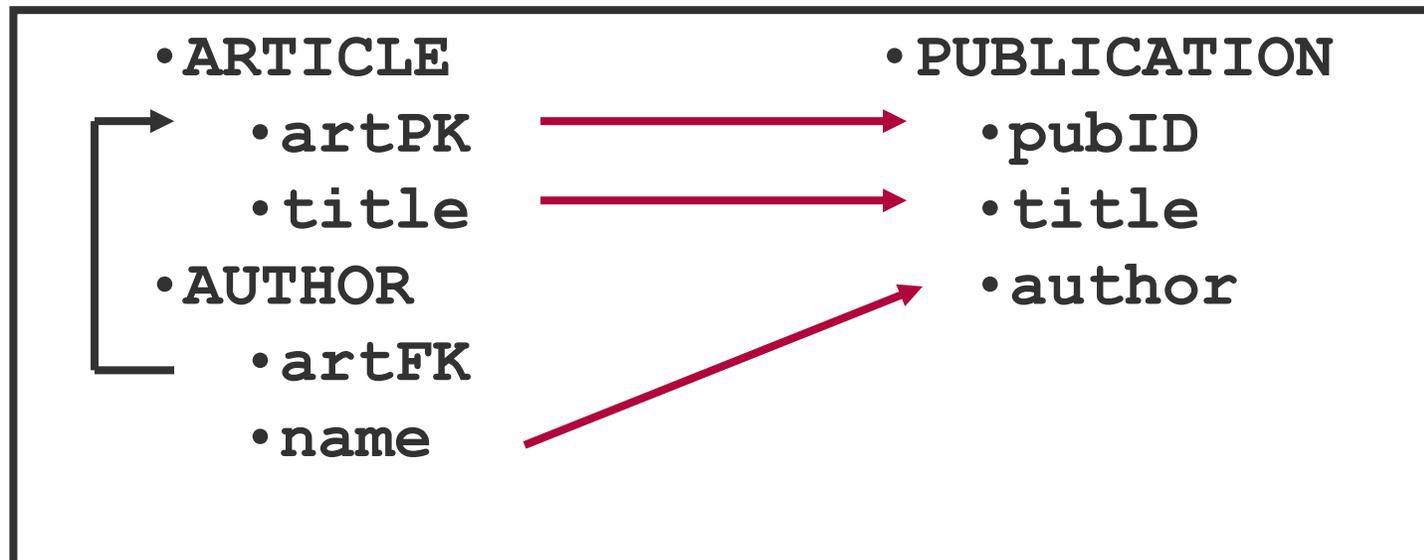
Synonym

Default Werte

Datentypkonflikt

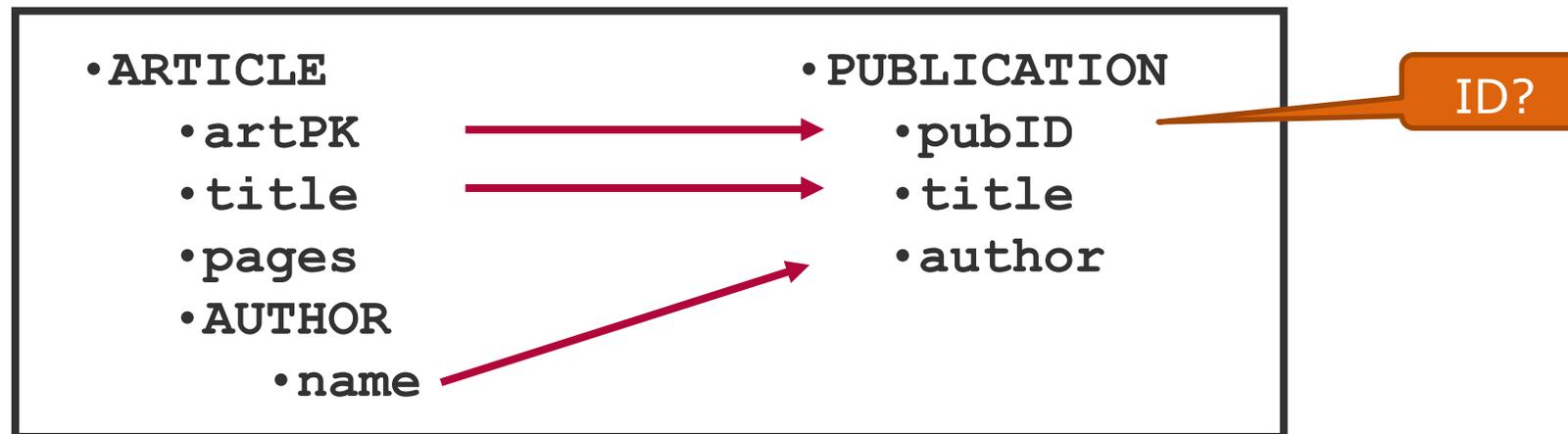
## Schematische Heterogenität - Beispiel

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



## Schematische Heterogenität - Beispiel

- Geschachtelt vs. Flach
  - 1:n Assoziationen werden unterschiedlich dargestellt
    - Als geschachtelte Elemente
    - Als Schlüssel-Fremdschlüssel Beziehung



Felix Naumann  
Information Integration  
Winter 2019/20

## Schematische Heterogenität - Lösungen

---

- Problem
  - Einheitlich auf beide Schemata zugreifen
    - Auf Schemaebene: Schema Mapping und Schema-Sprachen
    - Auf Datenebene: Virtuelle Integration
  - Beide Schemata 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.

## Schematische Heterogenität – Lösungen (Ausblick)

---

- SchemaSQL [LSS96]
  - Erweiterung von SQL
  - Daten und Metadaten werden gleich behandelt
  - Umstrukturierungen innerhalb der Anfrage
  - Dynamische Sicht-Definition
  - Horizontale Aggregation

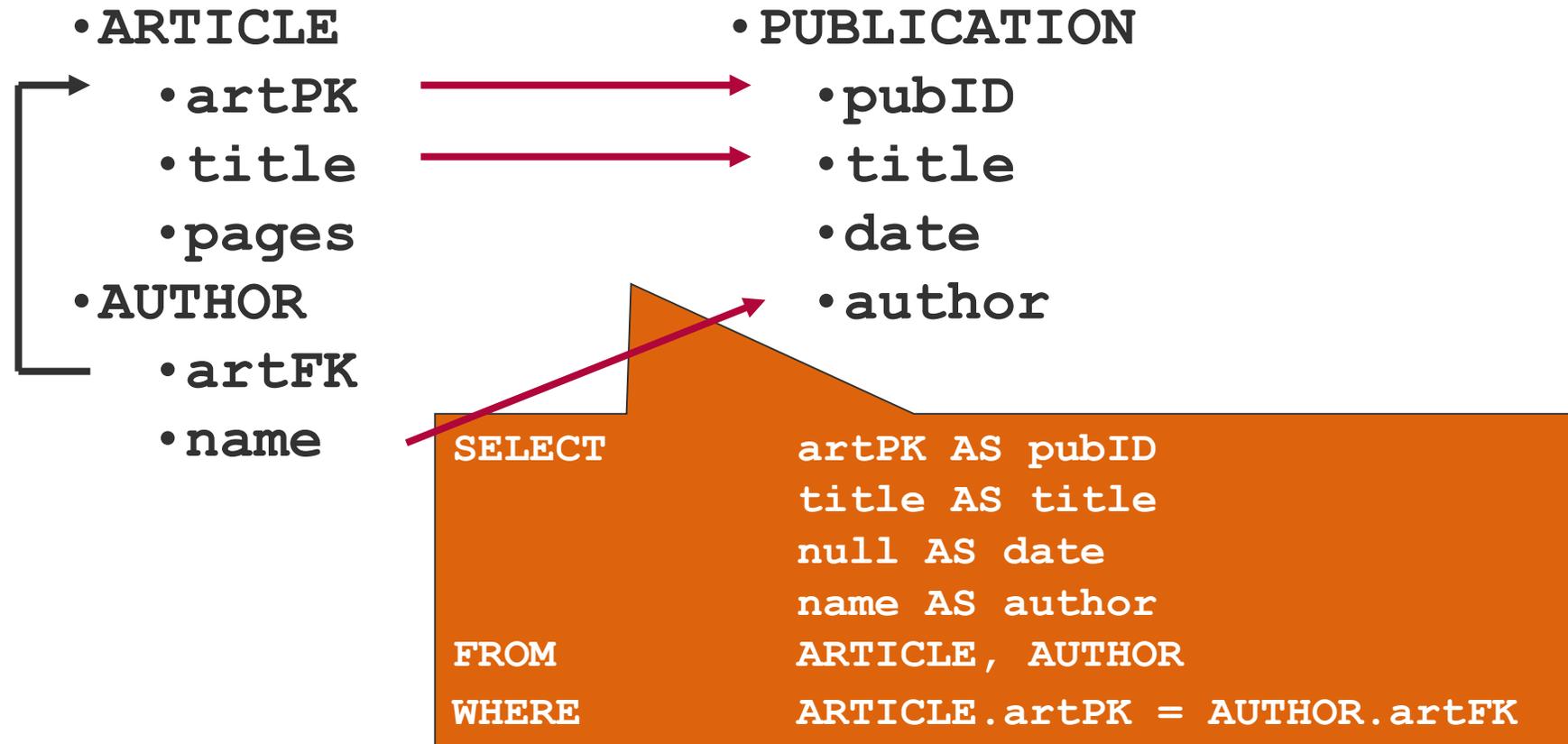
```
SELECT RelA
FROM uniA->RelA, uniA::RelA A, uniB::grundgehalt B
WHERE RelA = B.institut
AND A.Kategorie = „Student“
AND A.grundgehalt > B.Student
```



High-order Join

# Schematische Heterogenität – Lösungen (Ausblick)

## ■ Schema Mapping



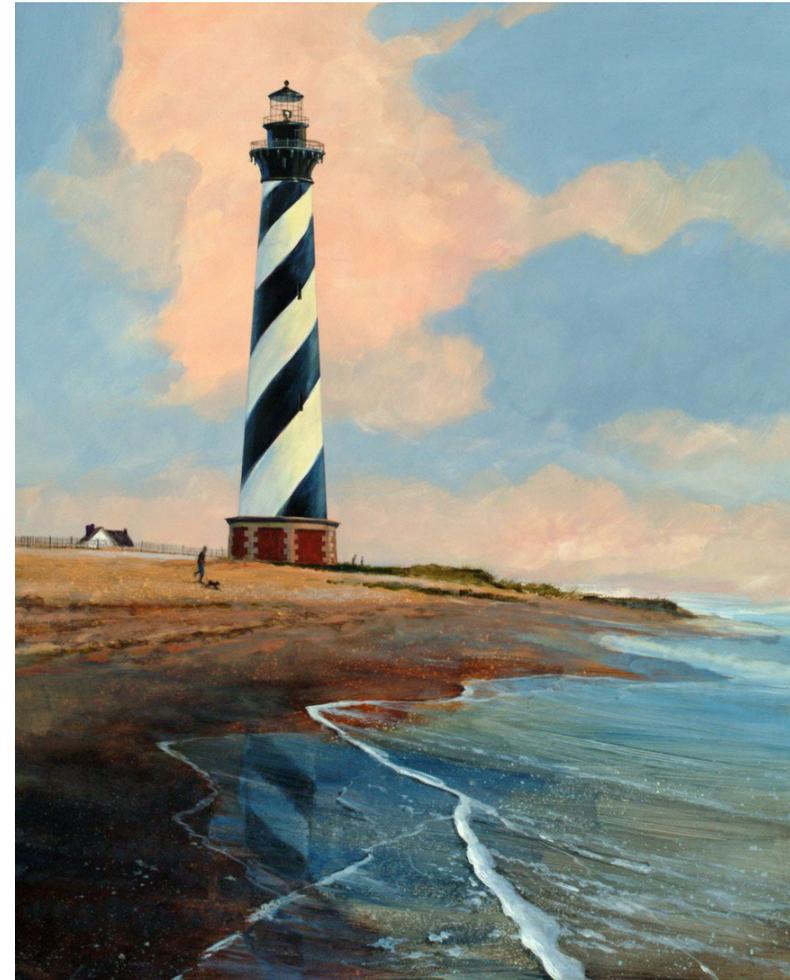
## Zusammenfassung

---

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

# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
- 5. Semantische Heterogenität**
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung



Felix Naumann  
Information Integration  
Winter 2019/20

- Fremdwörterduden “Semantik”:
  - Teilgebiet der Linguistik, das sich mit den Bedeutungen sprachlicher Zeichen und Zeichenfolgen befasst
  - Bedeutung, Inhalt eines Wortes, Satzes oder Textes
  
- *„Semantic heterogeneity is a fairly overloaded term. It basically refers to the differences among the databases that relate to meaning, interpretation, and intended use of data.”*  
[ÖV 2011]

# Semantik vs. Struktur

- Strukturelle Heterogenität
  - Betrifft Schemata
  - 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“

## Unterschiedliche Namen

---

- Die Probleme (Überblick)
  - 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

## Konzept

---

- 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

## Konzept am Beispiel „Wie viele Mitarbeiter hat IBM?“

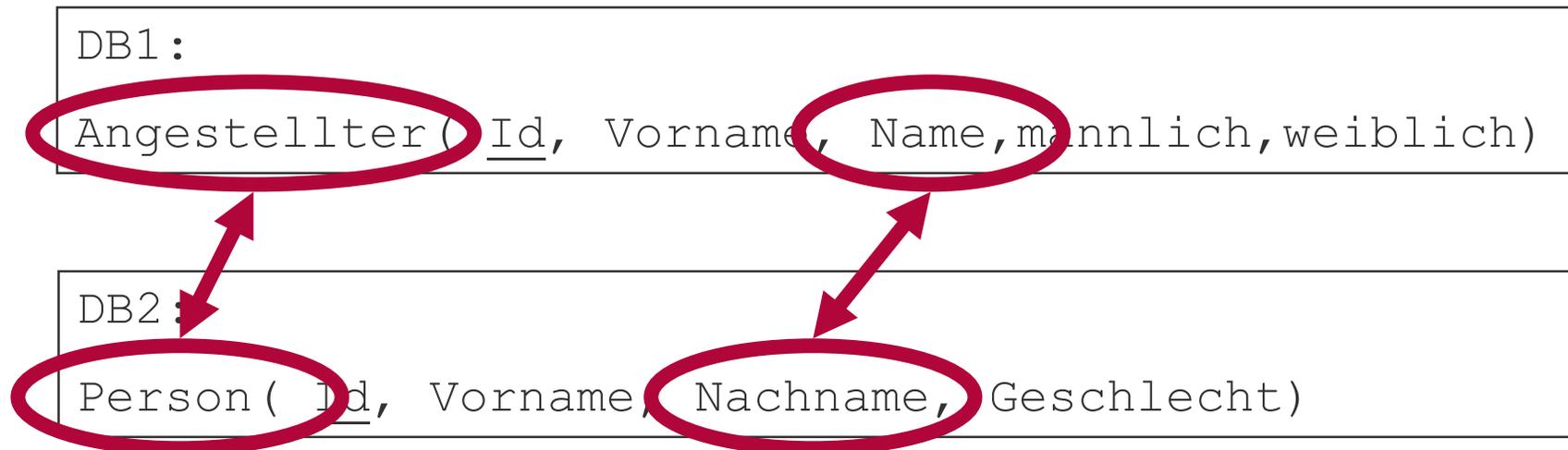
- Definition Mitarbeiter:
  - temporäre MA
  - Diplomanden
  - Berater
  - Studentische Mitarbeiter
  - Stellen oder Köpfe?
- Definition IBM
  - Welche Region? Welcher Geschäftsbereich?
  - Informix?
  - PWC?
- Welcher Zeitpunkt?
- Definition der Zählung:
  - Doppelte Zählung bei mehreren Anstellungen?
- Beispiel 2: Wie viel Hardware haben wir an das HPI verkauft?



1983 IBM East Fishkill employees celebrate the introduction of the IBM 4381 processor. Poughkeepsie Journal photos

## Synonyme

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



# Homonyme

- Gleiche Worte verschiedener Bedeutung
- Andere Domäne
- Andere Bedeutung

Sekr., Sachbearbeiter,  
Bereichsleiter, etc.

DB1:

Angestellter( Id, Vorname, Name, m, w, Titel)

DB2:

Person( Id, Vorname, Nachname, Titel, Geschlecht)

Herr, Frau, Dr., Dipl.Ing., etc.

## Andere -nym Wörter

---

- **Synonym**
  - Verschiedene Wörter, gleiche Semantik
- **Homonym**
  - Gleiche Wörter, verschiedene Semantik
- **Antonym**
  - Verschiedene Wörter, gegenteilige Semantik
- **Auto-Antonym:**
  - Gleiche Wörter, gegenteilige Semantik
  - Transparenz, aufheben, umfahren, Quantensprung, Kriegsgegner, ...
  - Overlook, comprise, dust, left, fast, wind up...
- **Heteronym**
  - Gleiche Schreibung, verschiedene Aussprache, verschiedene Semantik
    - Read, lead, bass, close, ...
- **Autonym (selbstbeschreibend, Wort = Semantik, „Substantiv“)**
- **Pseudonym u.v.a.m.**

- [http://www.fun-with-words.com/nym\\_words.html](http://www.fun-with-words.com/nym_words.html)

## Einheiten

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

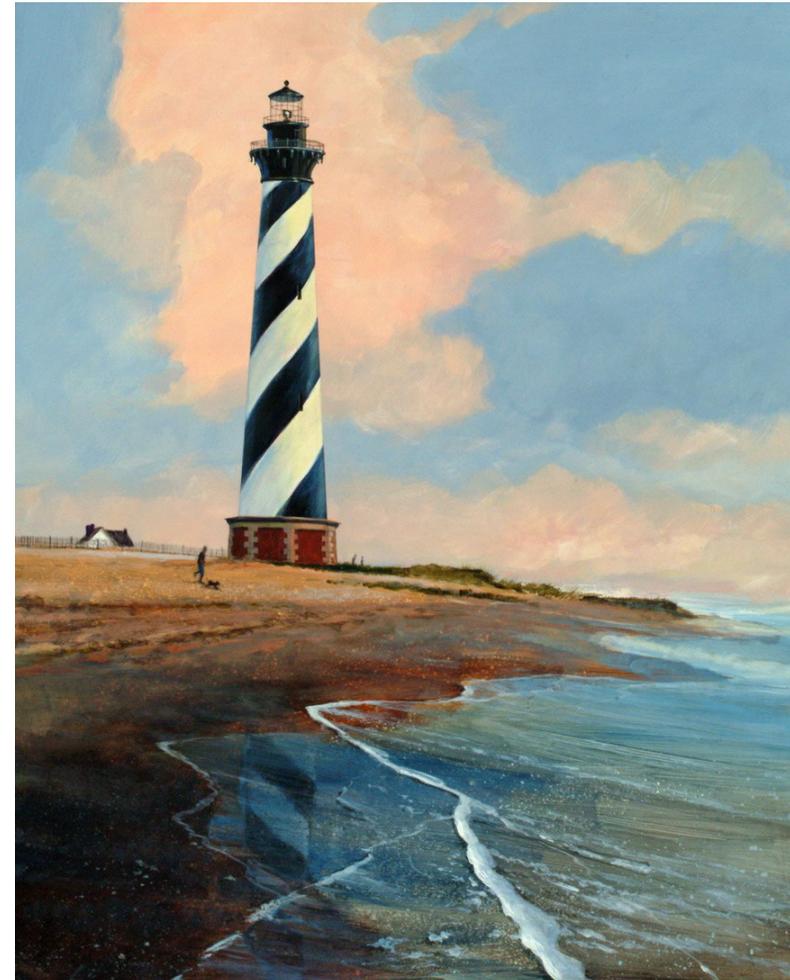
DB1 (USA) : USD  
Angestellter( Id, Vorname, Name, gehalt)

DB2 (Spanien) :  
Person( Id, Vorname, Nachname, gehalt) EUR



# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - **Identität**
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung

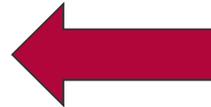


Felix Naumann  
Information Integration  
Winter 2019/20

# Identität

---

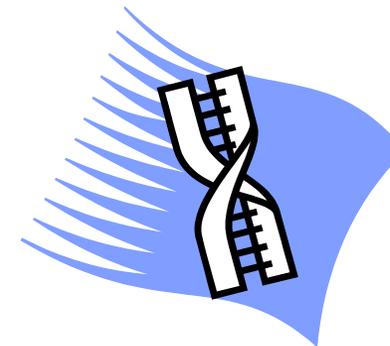
- Drei zentrale Fragen
  - Was ist ein Objekt?
    - XML: Über mehrere Schachtelungsebenen hinweg
    - Relationales Modell: Über mehrere Relationen hinweg
  - Repräsentiert Objekt A die gleiche Entität wie Objekt B?
  - Wie finde ich effizient gleiche Repräsentationen?
- Namen des Problems
  - Duplikaterkennung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



## Typische Anwendungen

---

- 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.



# Duplikaterkennung

---

- 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
  - Eine Tabelle und ein Anfragedatensatz (matching)
  - ...

# Duplikaterkennung

---

- Praktisches Problem 1
  - Wie entscheide ich, ob zwei Tupel das gleiche Objekt repräsentieren?
  - Ähnlichkeitsmaße und Klassifikation
    - Edit-Distance
    - N-grams
    - IDs
    - Wahrscheinlichkeitstheoretische Ansätze
    - Maschinelles Lernen
    - Augenschein

# Duplikaterkennung

- Praktisches Problem 2

- 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)`
- `FROM R as C1, R as C2`
- `WHERE M(C1,C2)`

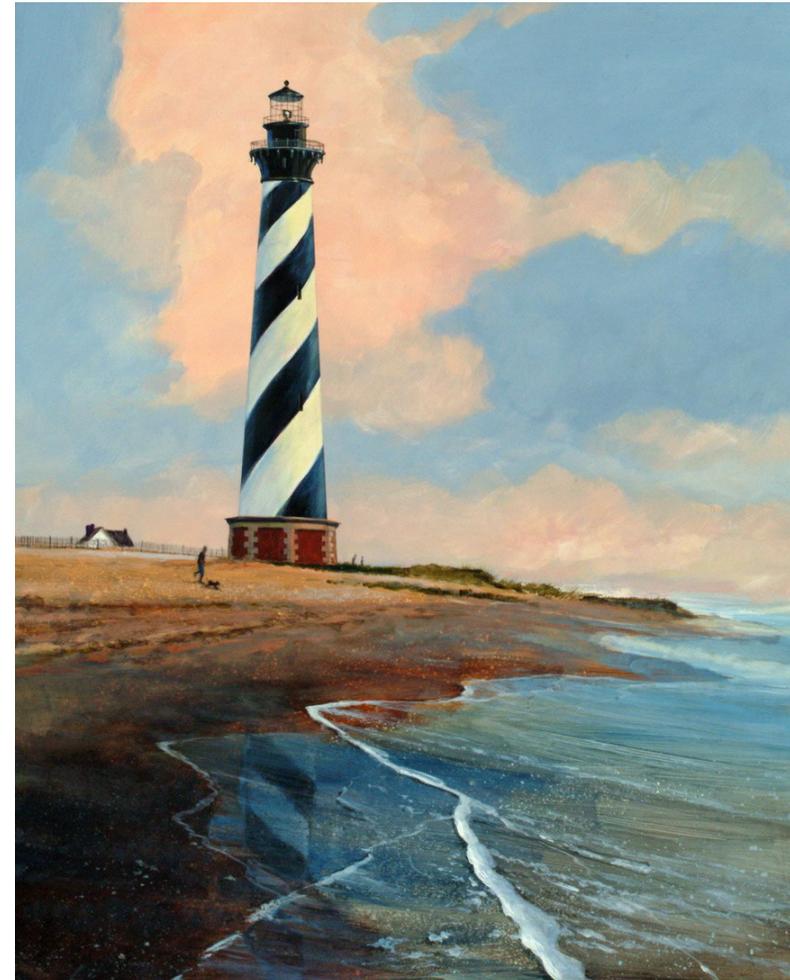


- Schwieriger als normaler Join

- Ähnlichkeitsmaß ist nicht nur Gleichheit

# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - **Datenkonflikte**
6. Gebundene und Freie Variablen
  - Adornments
  - Anfrageplanung



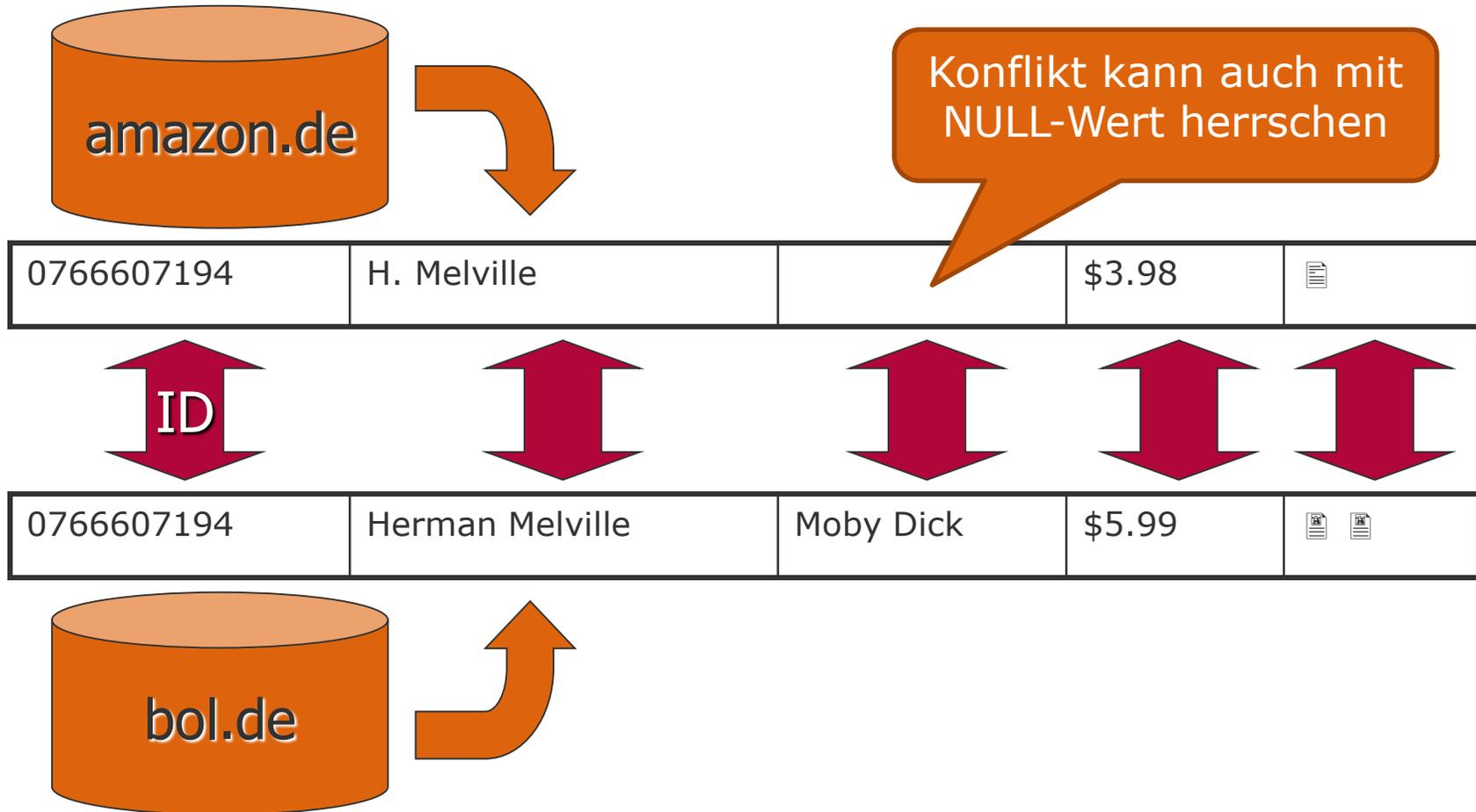
Felix Naumann  
Information Integration  
Winter 2019/20

## Datenkonflikte

---

- 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.

# Datenkonflikte - Beispiel



## Datenkonflikte – Entstehung

---

Auch innerhalb eines Informationssystems

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

## Datenkonflikte – Entstehung

---

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 (OCR)
    - U vs. V
    - 0 vs. O
    - 1 vs. l vs. I
    - ...

## Datenkonflikte – Behebung

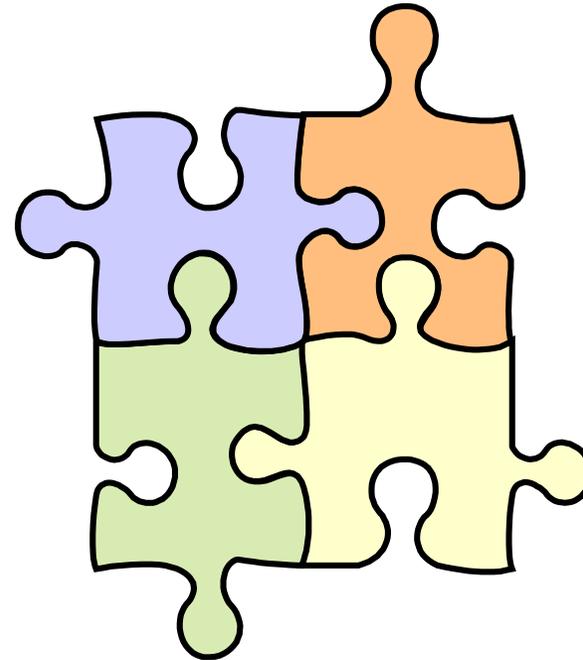
---

- Referenztabelle für exakte Wertabbildung
  - Z.B. Städte, Länder, Produktnamen, Codes...
  - Wiederum unter Verwendung von Ähnlichkeitsmaßen
- Standardisieren und transformieren
- Nutzung von Hintergrundwissen (Metadaten)
  - bzgl. von Konventionen (landestypische Schreibweisen)
  - Ontologien zur Behandlung von Zusammenhängen
  - Thesauri, Wörterbücher zur Behandlung von Homonymen, Synonymen, ...

## Relationale Objektintegration

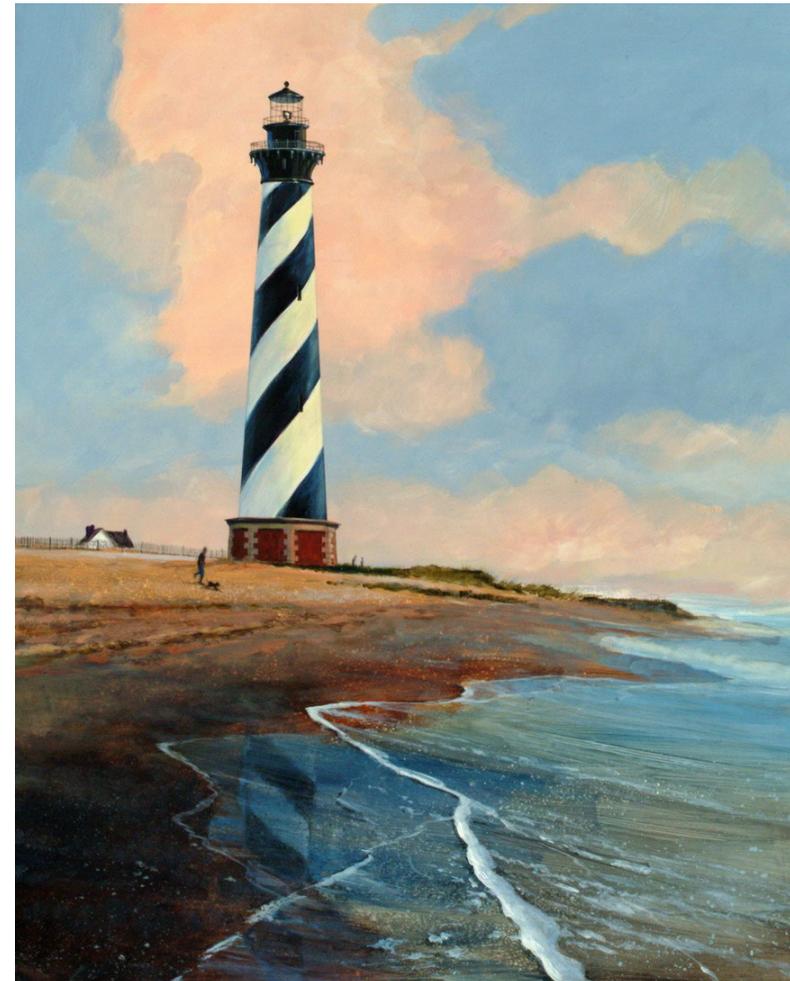
---

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



# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
- 6. Gebundene und Freie Variablen**
  - Adornments
  - Anfrageplanung



Felix Naumann  
Information Integration  
Winter 2019/20

## Gebundene & Freie Variablen

**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

Bitte korrigieren oder ergänzen Sie die ROT markierten Eingabefelder

---

ENERGIESPAREN: Bestellung des Sonderheftes  
Preis: € 5.80 (zzgl. € 2,50 Versandkosten)

Anzahl

---

**Persönliche Daten**  
Felder mit einem \* sind Pflichtfelder und müssen ausgefüllt werden.

Anrede:\*  Alter:

Vorname:\*  Nachname:\*

Land:\*

Straße und Hausnummer:\*  Telefonvorwahl / Telefonnummer:\*

PLZ:\*  Ort:\*  Email:\*

## Gebundene und Freie Variablen – Adornments

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

### Computing Capabilities of Mediators

Ramana Yerneni, Chen Li, Hector Garcia-Molina, Jeffrey Ullman  
Department of Computer Science  
Stanford University  
{yerneni,chenli,hector,ullman}@cs.stanford.edu

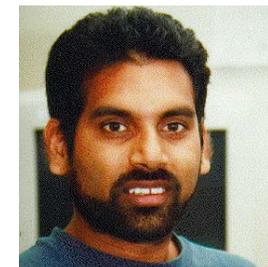
#### Abstract

Existing data-integration systems based on the mediation architecture employ a variety of mechanisms to describe the query-processing capabilities of sources. However, these systems do not compute the capabilities of the mediators based on the capabilities of the sources they integrate. In this paper, we propose a framework to capture a rich variety of query-processing capabilities of data sources and mediators. We present algorithms to compute the set of supported queries of a mediator, based on the capability limitations of its sources. Our algorithms take into consideration a variety of query-processing techniques employed by mediators to enhance the set of supported queries.

for other mediators. Furthermore, users have a difficult time understanding the set of supported mediator queries in these systems. Consequently, users must endure a frustrating trial-and-error approach, submitting queries that are rejected until finally hitting upon a query that is answered by the mediator.

In this paper, we present algorithms to precompute mediator capabilities automatically, so that users and other mediators know which queries are supported. In addition, we extend the types of source limitations that can be handled by existing mediation systems. For example, we handle attributes that can only be queried with values from a fixed menu of constants.

The World Wide Web is a prime example of a context where we need to handle many different



Felix Naumann  
Information Integration  
Winter 2019/20

## Gebundene und Freie Variablen – Adornments

Z gebunden

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):

$R_1(X, Y, Z)$

$R_2(X, Y, Z)$

$R_3(X, Y, Z)$

$R_4(Z, U)$

$R_5(U, V, W)$

## Gebundene und Freie Variablen – Adornments

---

- Anfragefähigkeiten der Quellen modelliert als *templates*
  - Wie ein WWW Formular
  - Templates bestehen aus einem *adornment* für jedes Attribut
- Anhänge (adornments = Verzierungen/Ausschmückung) 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.

## Adornments - Beispiele

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[ $z_1, z_2$ ]

# Adornments – Anfragebearbeitung

- Anfragebearbeitung
  - $R1(X,Y,Z)$ : bff, ffb
  - $R2(X,Y,Z)$ : fbf
  - Sei  $M = R1 \cup R2$  eine integrierte Sicht
    - Anfragen gegen diese Sicht (Schema: X,Y,Z)
  - 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?

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

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

# Adornments – Verknüpfung durch UNION

3 Sichten und deren Adornments:

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

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

$R_3(X,Y,Z)$ : ffc[s<sub>1</sub>], c[s<sub>2</sub>]ff

$R_1 \cup R_2$ :

bff  $\cup$  fb = bbf

ffb  $\cup$  fb = fbb

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

bbf  $\cup$  ffc[s<sub>1</sub>] = bbc[s<sub>1</sub>] usw.

$\cup$	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	o[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	u
<b>o[s<sub>1</sub>]</b>	o[s <sub>1</sub> ]	o[s <sub>1</sub> $\cap$ s <sub>3</sub> ]	c[s <sub>1</sub> ]	c[s <sub>1</sub> $\cap$ s <sub>4</sub> ]	u
<b>b</b>	b	c[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	-
<b>c[s<sub>2</sub>]</b>	c[s <sub>2</sub> ]	c[s <sub>2</sub> $\cap$ s <sub>3</sub> ]	c[s <sub>2</sub> ]	c[s <sub>2</sub> $\cap$ s <sub>4</sub> ]	-
<b>u</b>	u	u	-	-	u

## Adornments – Verknüpfung durch Join ( $\bowtie$ )

---

- Unterschied zu UNION
  - Nicht jedes Attribut der integrierten Sicht ist auch Attribut jeder beteiligten Quelle.
  - Beispiel:  $R1(X,Y,Z)$  und  $R4(Z,U)$
  - Sicht:  $M(X,Y,Z,U) = R1(X,Y,Z) \bowtie R4(Z,U)$
  
- Berechnung des Templates der Sicht
  1. Adornments der nicht-Join-Attribute werden kopiert.
  2. Adornments der Join-Attribute werden gemäß der UNION Tabelle vereint.

## Adornments – Selektion und Projektion

---

### ■ 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

- Problem
  - UNION-Matrix zu restriktiv

	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	o[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	u
<b>o[s<sub>1</sub>]</b>	o[s <sub>1</sub> ]	o[s <sub>1</sub> ∩ s <sub>3</sub> ]	c[s <sub>1</sub> ]	c[s <sub>1</sub> ∩ s <sub>4</sub> ]	u
<b>b</b>	b	c[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	-
<b>c[s<sub>2</sub>]</b>	c[s <sub>2</sub> ]	c[s <sub>2</sub> ∩ s <sub>3</sub> ]	c[s <sub>2</sub> ]	c[s <sub>2</sub> ∩ s <sub>4</sub> ]	-
<b>u</b>	u	u	-	-	u

- Idee: Erhöhung der Menge beantwortbarer Anfragen
  - durch *post-processing*
  - durch *passing bindings*

## Adornments und Postprocessing

$R_1(X,Y,Z)$ : 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?

Ja

Nicht ohne  
weiteres

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 (post-processing)

$o[s] = f$  : falls Bindung nicht in  $s$ , weglassen und später Filtern

Zusammen:  $R_1 \cup R_2 = bff$

# Adornments – Verknüpfung durch UNION

Vorher:

	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	o[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	u
<b>o[s<sub>1</sub>]</b>	o[s <sub>1</sub> ]	o[s <sub>1</sub> ∩ s <sub>3</sub> ]	c[s <sub>1</sub> ]	c[s <sub>1</sub> ∩ s <sub>4</sub> ]	u
<b>b</b>	b	c[s <sub>3</sub> ]	b	c[s <sub>4</sub> ]	-
<b>c[s<sub>2</sub>]</b>	c[s <sub>2</sub> ]	c[s <sub>2</sub> ∩ s <sub>3</sub> ]	c[s <sub>2</sub> ]	c[s <sub>2</sub> ∩ s <sub>4</sub> ]	-
<b>u</b>	u	u	-	-	u

Nachher:

	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	f	b	c[s <sub>4</sub> ]	f
<b>o[s<sub>1</sub>]</b>	f	f	b	c[s <sub>4</sub> ]	f
<b>b</b>	b	b	b	c[s <sub>4</sub> ]	b
<b>c[s<sub>2</sub>]</b>	c[s <sub>2</sub> ]	c[s <sub>2</sub> ]	c[s <sub>2</sub> ]	c[s <sub>2</sub> ∩ s <sub>4</sub> ]	c[s <sub>2</sub> ]
<b>u</b>	f	f	b	c[s <sub>4</sub> ]	f

## Adornments und Passing Bindings

JOIN über templates ohne  
*passing bindings*

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

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*:

$R_1 \bowtie R_5 = \text{fbff}$

# Adornments and Passing Bindings

Vorher:

	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	$o[s_3]$	b	$c[s_4]$	u
<b>o[s<sub>1</sub>]</b>	$o[s_1]$	$o[s_1 \cap s_3]$	$c[s_1]$	$c[s_1 \cap s_4]$	u
<b>b</b>	b	$c[s_3]$	b	$c[s_4]$	-
<b>c[s<sub>2</sub>]</b>	$c[s_2]$	$c[s_2 \cap s_3]$	$c[s_2]$	$c[s_2 \cap s_4]$	-
<b>u</b>	u	u	-	-	u

Zweite Quelle

Nachher:

	<b>f</b>	<b>o[s<sub>3</sub>]</b>	<b>b</b>	<b>c[s<sub>4</sub>]</b>	<b>u</b>
<b>f</b>	f	f	f	$c[s_4]$	f
<b>o[s<sub>1</sub>]</b>	f	f	f	$c[s_4]$	f
<b>b</b>	b	b	b	$c[s_4]$	b
<b>c[s<sub>2</sub>]</b>	$c[s_2]$	$c[s_2]$	$c[s_2]$	$c[s_2 \cap s_4]$	$c[s_2]$
<b>u</b>	f	f	f	$c[s_4]$	f

Erste Quelle

## Adornments – Selektion mit Post-processing

---

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

Variante 2:  $R_1(X,Y,Z), X < x_1$ : 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  
post-processing

# Adornments – Selektion mit Post-processing

Vorher      Nachher

Base View Adornment	Sel. Attribute Adornment
f	f
o[s <sub>1</sub> ]	f
b	f or b
c[s <sub>1</sub> ]	f or c[s <sub>1</sub> ]
u	f

## Viele Templates

- Problem: Quellen exportieren oft mehrere templates
  - Beispiel: Amazon(Autor, Titel, Schlagwort, ISBN, Verlag)
  - bffff, fbfff, ffbff, fffbf, ffffb
  - Beispiel: Verlage(Verlag, Ort)
  - bf, fb
  - Sicht im IIS: Amazon  $\bowtie$ Verlag Verlage
  - Templates der Sicht aus jeder Kombination:
    - bffffff, fbffff, ffbfff, fffbff, fffbf
    - bffffb, fbfffb, ffbffb, fffbfb, fffbb
    - + fffffb (ffffb  $\bowtie$  fb mit passing binding)
- Beobachtung:
  - Einige templates sind redundant: Jede Anfrage, die dieses template erlaubt wird auch von einem anderen template erlaubt.

**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 "Jetzt suchen".

Autor/in:

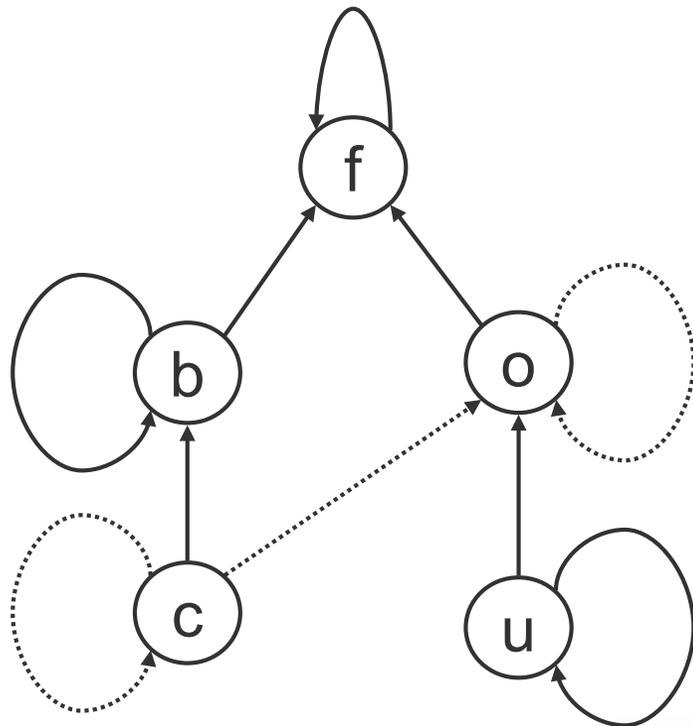
Titel:

Schlagwörter:

ISBN:

Verlag:

# Redundanz in Templates



—————> Weniger restriktiv  
 .....> Weniger restriktiv falls Auswahllisten Teilmengen sind } transitiv

- bffffff, fbffff, ffbfff, fffbff, fffbf
- bffffb, fbfffb, ffbffb, fffbfb, ffffbb

**Subsumption Test:** A template  $T$  is subsumed by another template  $T'$  if for every attribute  $X$  the adornment of  $X$  in  $T$  is at least as restrictive as the adornment of  $X$  in  $T'$  (based on Figure 1).

## Adornments – Fallbeispiel

---

### ■ 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

# Adornments - Fallbeispiel

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'

Felix Naumann  
Information Integration  
Winter 2019/20

# Adornments - Fallbeispiel

IIS

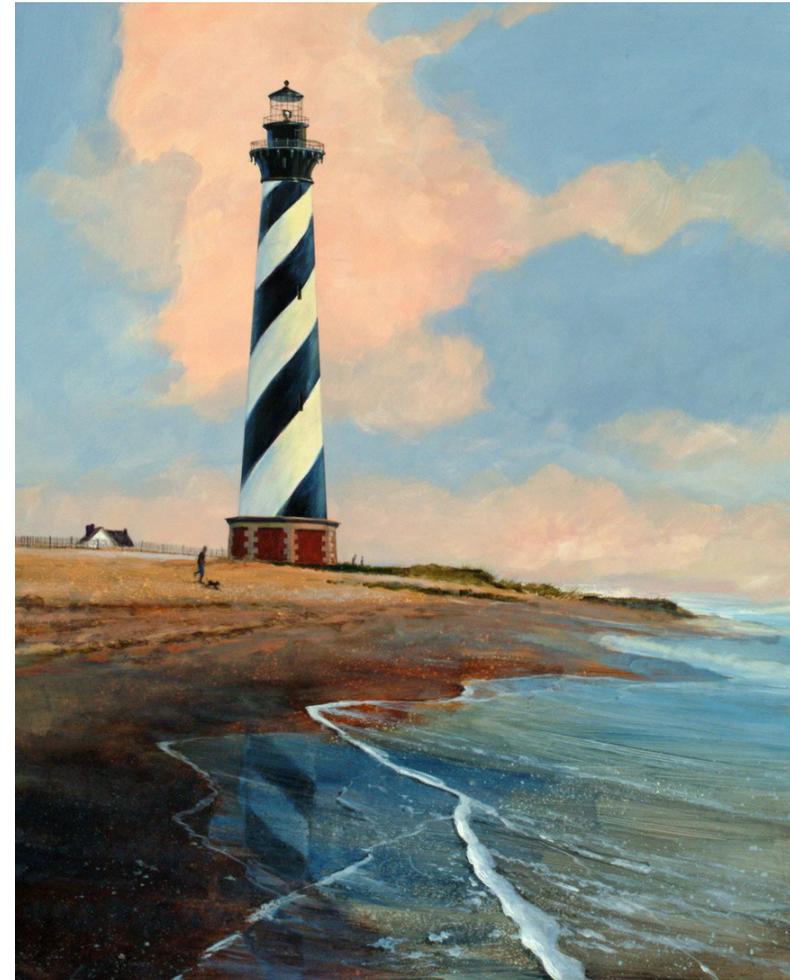
author	title	format	subject	KW	ISBN	pub	date	price	ship	age
f	f	f	u'	u'	b	f	f	f	f	u'
f	f	f	u'	b'	f	f	f	f	f	u'
b	f	f	o'	u'	f	f	f	f	f	u'
f	b	f	o'	u'	f	f	f	f	f	u'
b	f	f	u'	f'	f	b	f	f	f	u'
f	b	f	u'	f'	f	b	f	f	f	u'
f	b	f	u'	f'	f	f	b	f	f	u'
b	f	f	u'	f'	f	f	b	f	f	u'

## 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)

# Überblick

1. Verteilung
2. Autonomie
3. Syntaktische Heterogenität
4. Strukturelle Heterogenität
5. Semantische Heterogenität
  - Namenskonflikte
  - Identität
  - Datenkonflikte
6. Gebundene und Freie Variablen
  - Adornments
  - **Anfrageplanung**



Felix Naumann  
Information Integration  
Winter 2019/20

## Gebundene & Freie Variablen – Beispiel

---

$v_1(\text{Song}, \text{CD})$

$\langle \text{Friends}, \text{Life} \rangle$

$\langle \text{Friends}, \text{Love} \rangle$

$v_2(\text{CD}, \text{Artist}, \text{Price})$

$\langle \text{Love}, \text{Lucy}, \$15 \rangle$

$\langle \text{Story}, \text{Snoopy}, \$14 \rangle$

$v_3(\text{CD}, \text{Artist}, \text{Price})$

$\langle \text{Story}, \text{Lucy}, \$13 \rangle$

$\langle \text{Love}, \text{Snoopy}, \$10 \rangle$

$\langle \text{Life}, \text{Charlie}, \$8 \rangle$

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

## Gebundene & Freie Variablen – Beispiel

$v_1(\text{Song}, \text{CD})$

$\langle \text{Friends}, \text{Life} \rangle$

$\langle \text{Friends}, \text{Love} \rangle$

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

$v_2(\text{CD}, \text{Artist}, \text{Price})$

$\langle \text{Love}, \text{Lucy}, \$15 \rangle$

$\langle \text{Story}, \text{Snoopy}, \$14 \rangle$

$v_3(\text{CD}, \text{Artist}, \text{Price})$

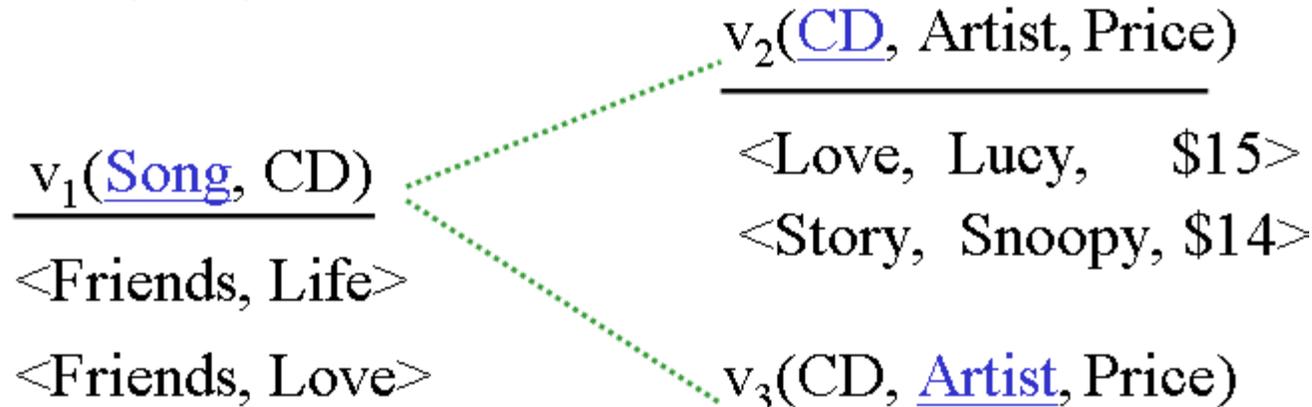
$\langle \text{Story}, \text{Lucy}, \$13 \rangle$

$\langle \text{Love}, \text{Snoopy}, \$10 \rangle$

$\langle \text{Life}, \text{Charlie}, \$8 \rangle$

# Gebundene & Freie Variablen – Beispiel

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



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

# Gebundene & Freie Variablen – Beispiel

$v_1(\text{Song}, \text{CD})$  

<Friends, Life>

<Friends, Love>

$v_2(\text{CD}, \text{Artist}, \text{Price})$

<Love, Lucy, \$15>

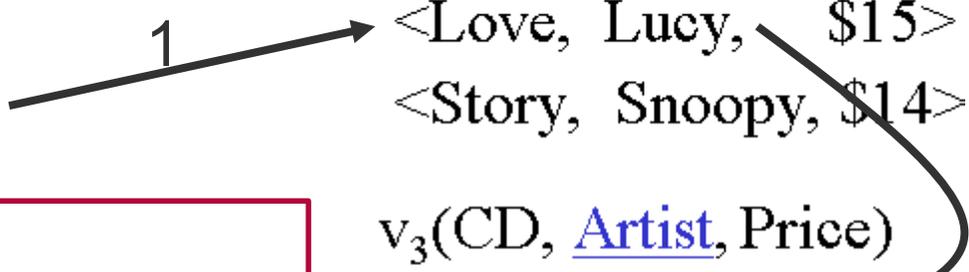
<Story, Snoopy, \$14>

$v_3(\text{CD}, \text{Artist}, \text{Price})$

<Story, Lucy, \$13>

<Love, Snoopy, \$10>

<Life, Charlie, \$8>



Idee 1:

$(v_1 \bowtie_{\text{CD}} v_2) \bowtie_{\text{Artist}} v_3$

Ergibt zusätzlich:

<Story, Lucy, \$13>

^---PECH

# Gebundene & Freie Variablen – Beispiel

$v_1(\text{Song}, \text{CD})$

<Friends, Life>

<Friends, Love>

$v_2(\text{CD}, \text{Artist}, \text{Price})$

<Love, Lucy, \$15>

<Story, Snoopy, \$14>

$v_3(\text{CD}, \text{Artist}, \text{Price})$

<Story, Lucy, \$13>

<Love, Snoopy, \$10>

<Life, Charlie, \$8>

Idee 2:

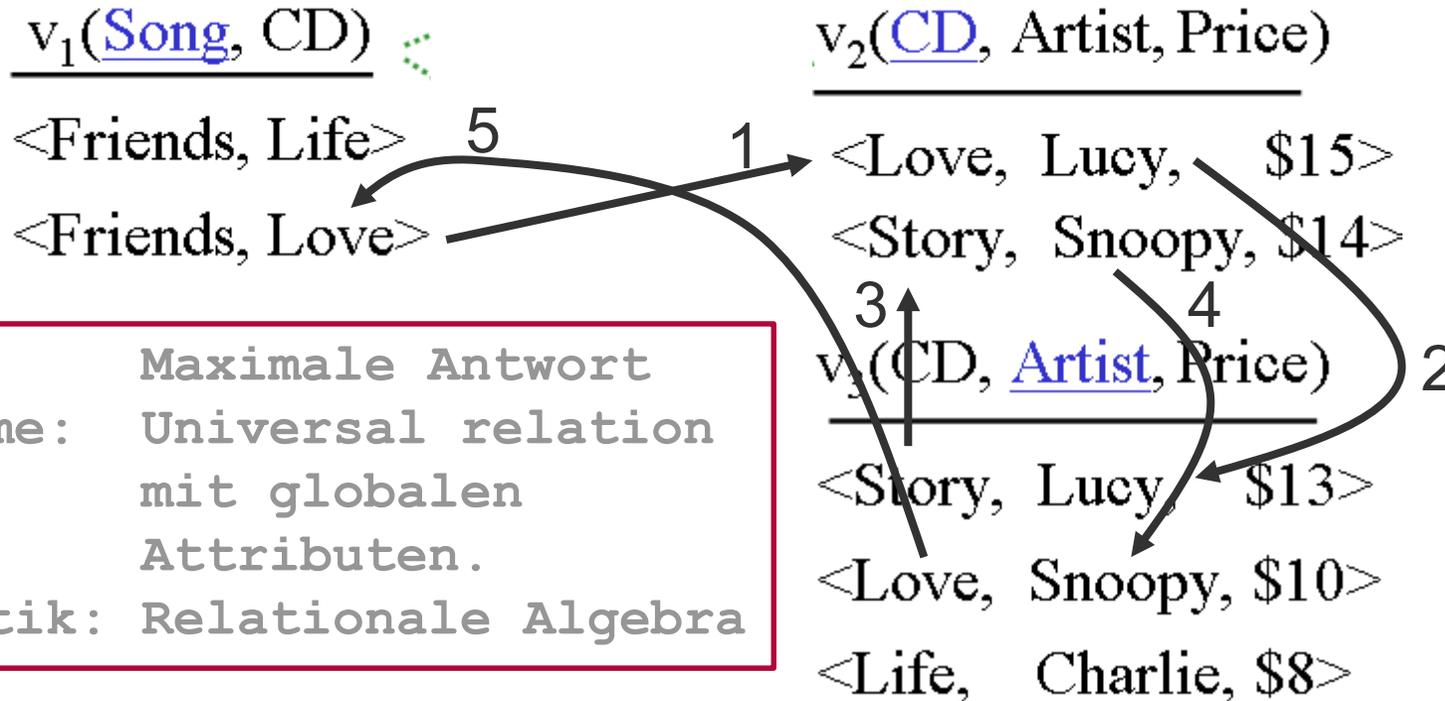
$((v_1 \bowtie v_2) \bowtie v_3) \bowtie v_2 \bowtie v_3 \bowtie v_1$

Ergibt zusätzlich:

<Love, Snoopy, \$10>

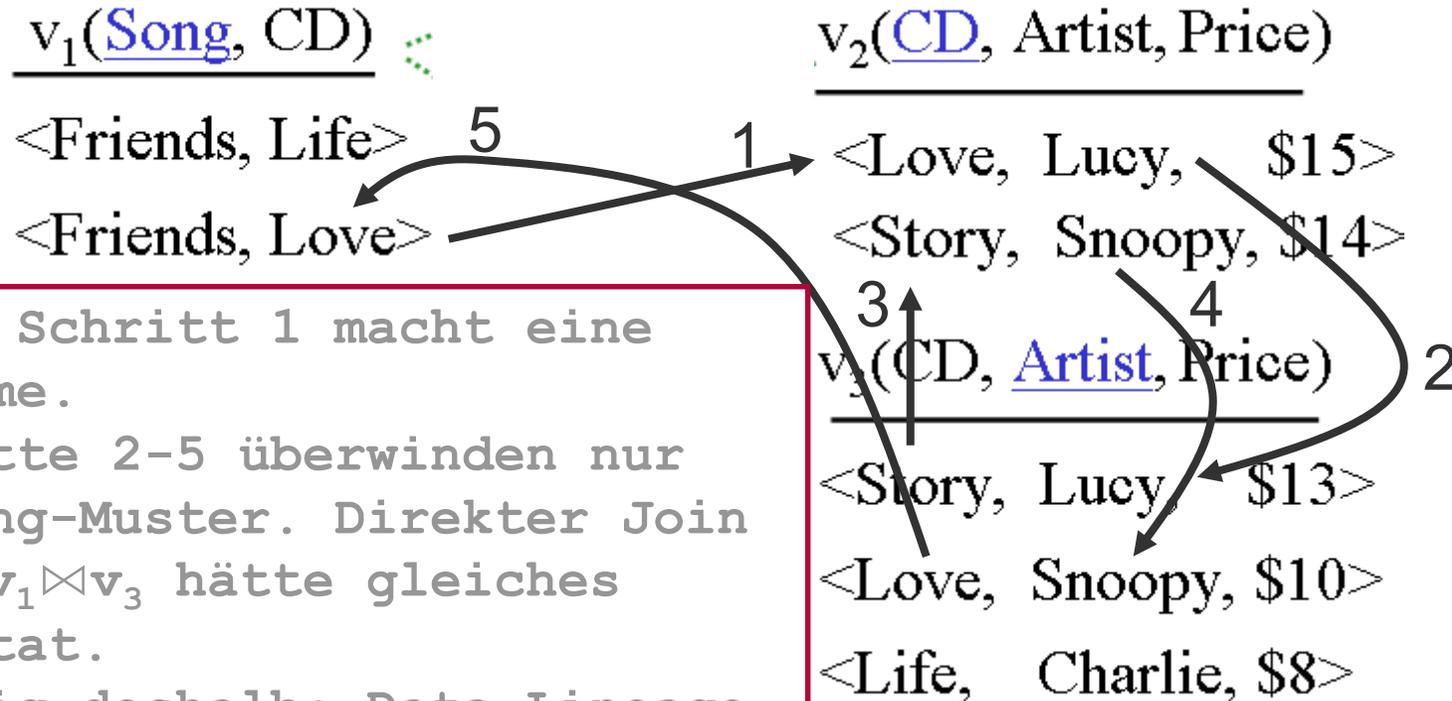
^--- Glück, weil auch in  $v_1$ !

Gebundene & Freie Variablen – Beispiel: Semantik



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

# Gebundene & Freie Variablen – Beispiel: Semantik



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

## Literatur

---

### Wichtigste Literatur

- [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](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