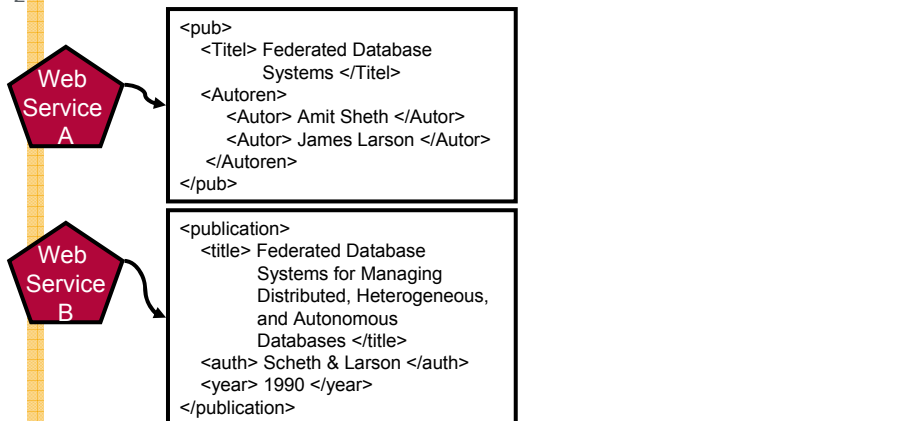


Informationsintegration Beispiel

17.4.2007
Felix Naumann

Informationsintegration

2



Integration

Identifikation

Fusion

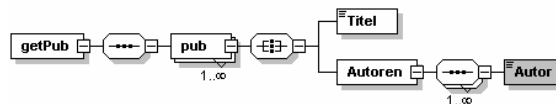
Optimierung

Visualisierung

Beispiel – Web Service A

3

- Standort: Trier
- Operation:
 - getPubByAuthor(firstName, lastName)
 - getPubByTitle(title)
- Output-Struktur:



Felix Naumann | Informationsintegration | SoSe 2007

Beispiel – Web Service A Output

4

```

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

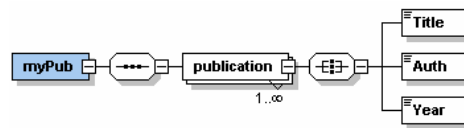
```

Felix Naumann | Informationsintegration | SoSe 2007

Beispiel – Web Service B

5

- Standort: Humboldt-Universität
- Operation: myPubs(Autor, Jahr)
- Struktur:



Felix Naumann | Informationsintegration | SoSe 2007

Beispiel – Web Service B Output

6

```

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

Felix Naumann | Informationsintegration | SoSe 2007

Integration von Web Services A & B

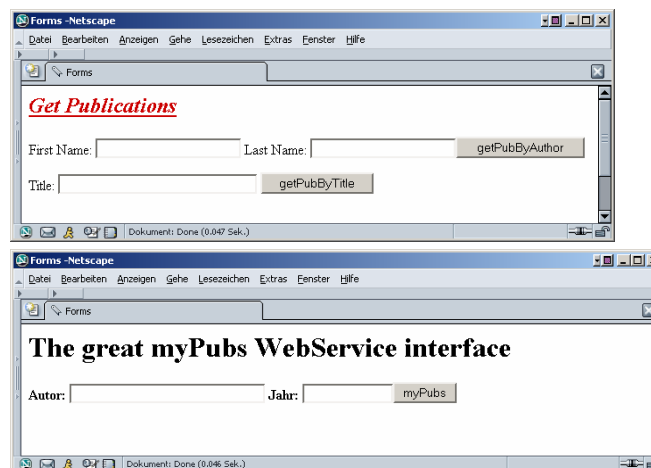
7

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

Felix Naumann | Informationsintegration | SoSe 2007

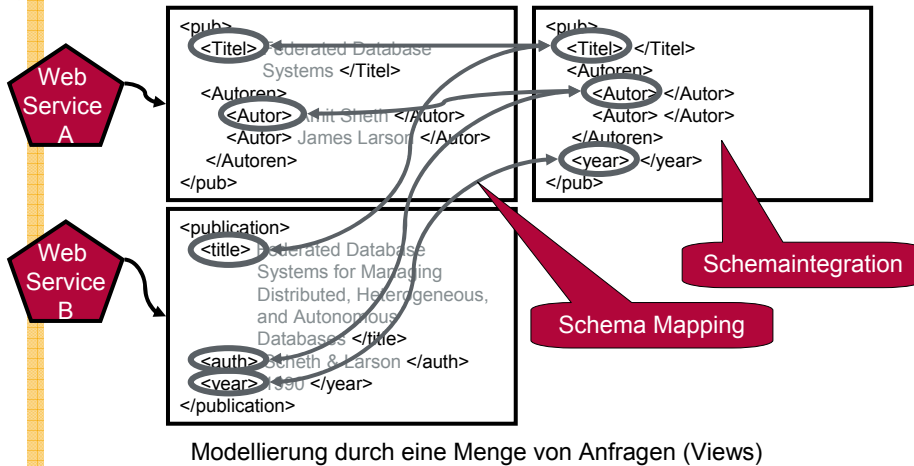
Nutzerschnittstellen

8



Felix Naumann | Informationsintegration | SoSe 2007

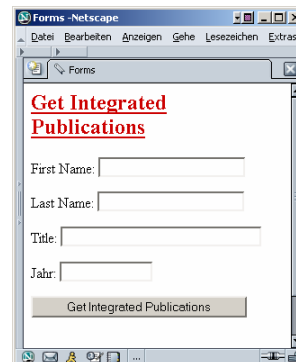
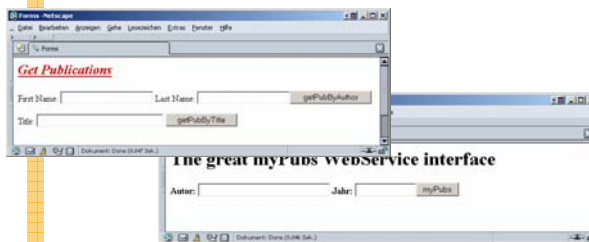
9



10

Integration der Anfrage durch Mediator:

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



Anfrageoptimierung

11

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

Zwei Resultate

12

```
<?xml version="1.0" encoding="UTF-8" ?>
<!-- edited with XMLSPY v2004 rel. 2 U (http://www.xmlspy.com) by Felix Naumann
(Universität zu Berlin) -->
<!-- Sample XML file generated by XMLSPY v2004 rel. 2 U (http://www.xmlspy.com)
- <getPub xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="C:\Dokumente und Einstellungen\Naumann\Eigene
  Dateien\Lehre\Winter03_04\VL InfoInt 1 WS03_04\getPub.xsd" ->
- <pub>
- <Title>Real-world Data is Dirty: The Merge/Purge Problem</Title>
- <Autorens
- <Autor>Mauricio Hernandez</Autor>
- <Autor>Salvatore Stolfo</Autor>
- </Autorens
- </pub>
- <pub>
- <Title>MAC: Merging Autonomous Content</Title>
- <Autorens
- <Autor>Felix Naumann</Autor>
- <Autor>Jens Bleilholder</Autor>
- </Autorens
- </pub>
- </getPub>
```

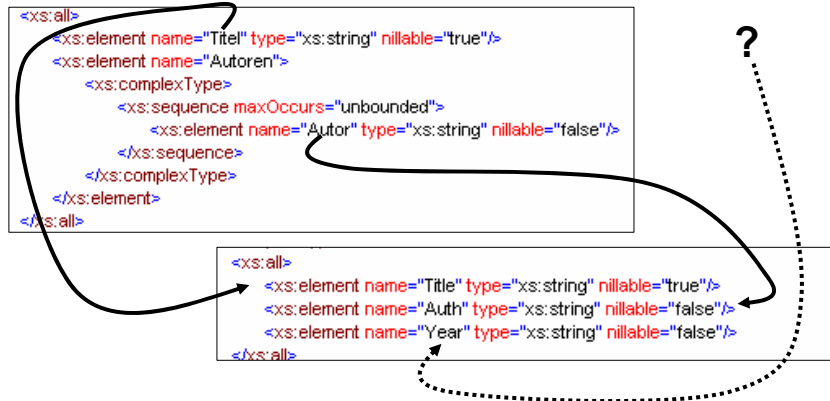
Web Service A

```
<?xml version="1.0" encoding="UTF-8" ?>
<!-- Sample XML file generated by XMLSPY v2004 rel. 2 U (http://www.xmlspy.com)
- <myPub xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="C:\Dokumente und Einstellungen\Naumann\Eigene
  Dateien\Lehre\Winter03_04\VL InfoInt 1 WS03_04\myPubs.xsd" ->
- <publication>
- <Title>Merging Autonomous Content</Title>
- <Auth>Naumann</Auth>
- <Year>2003</Year>
- </publication>
- <publication>
- <Title>Object Matching for Information Integration</Title>
- <Auth>Doan</Auth>
- <Year>1999</Year>
- </publication>
- </myPub>
```

Web Service B

Schema Matching

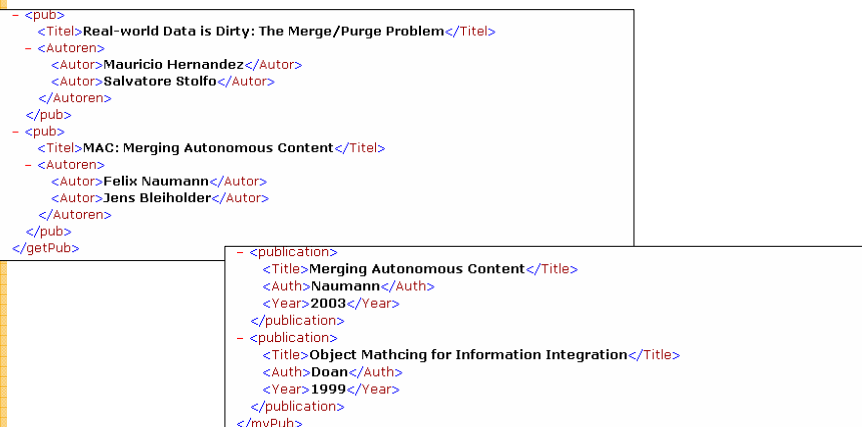
13



Felix Naumann | Informationsintegration | SoSe 2007

Objektidentifikation

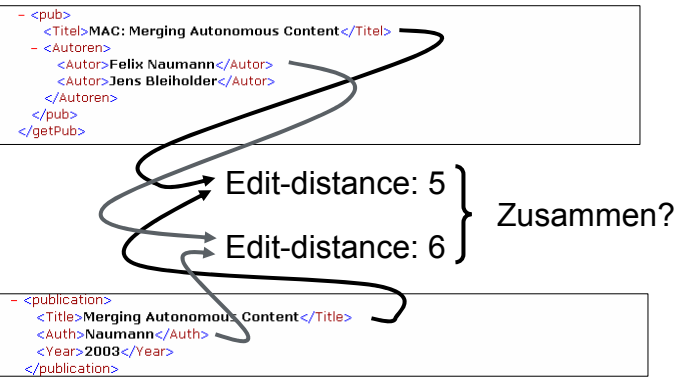
14



Felix Naumann | Informationsintegration | SoSe 2007

Objektidentifikation

15



Felix Naumann | Informationsintegration | SoSe 2007

Stand der Dinge

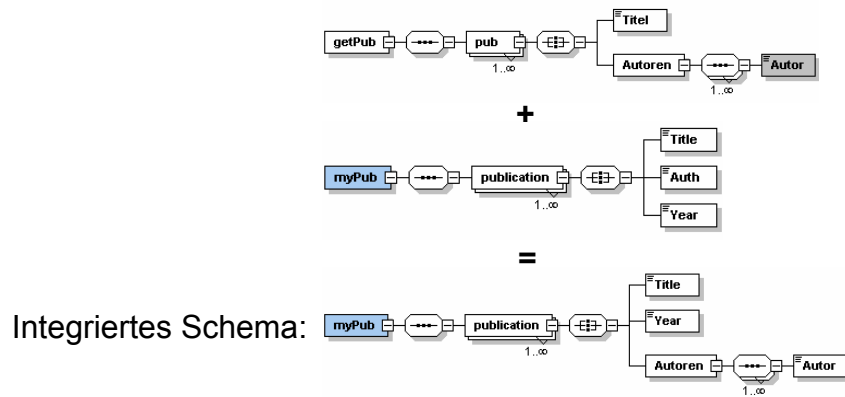
16

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

Felix Naumann | Informationsintegration | SoSe 2007

Angestrebtes Integrationsergebnis

17



Angestrebtes Integrationsergebnis

18

Integrierte Daten:

```

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

```

Integrierte Daten – was ist passiert?

19



Implementierung

20

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

Anzeige beim Nutzer

21

Konflikt-
lösung

Visualisierung der

- Datenherkunft
- Qualität
- veränderten Daten
- Operationen

```

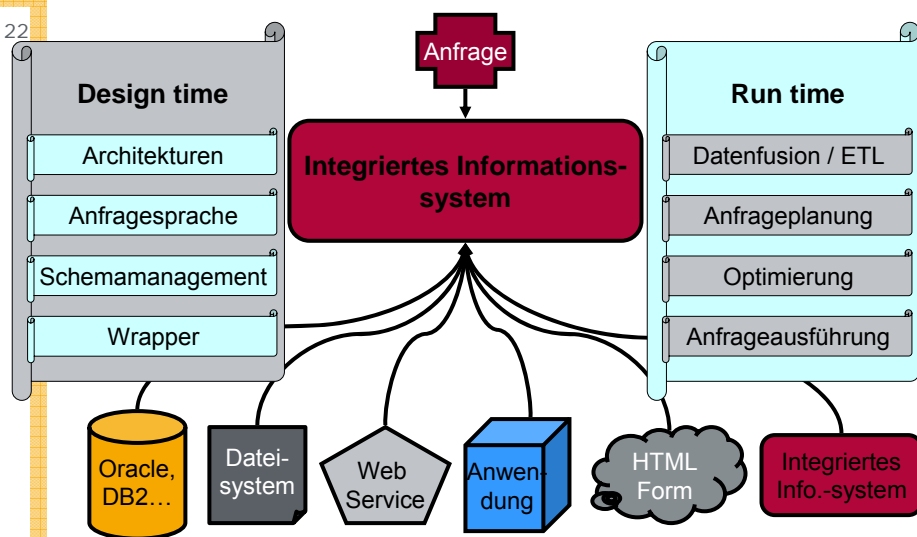
- <publication>
  <Title>Real-world Data is Dirty: The Merge/Purge Problem</Title>
  <Year>1999</Year>
- <Autoren>
  <Autor>Mauricio Hernandez</Autor>
  <Autor>Salvatore Stolfo</Autor>
</Autoren>
</publication>
publication
<Title>Merging Autonomous Content</Title>
<Year>2003</Year>
- <Autoren>
  <Autor>Felix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

Vorher:
„Naumann“

Felix Naumann | Informationsintegration | SoSe 2007

Integrierte Informationssysteme

22



Felix Naumann | Informationsintegration | SoSe 2007

- Einführung in die Informationsintegration
- Szenarien der Informationsintegration
- Verteilung und Aufnahme
- Heterogenität
- Materialisierte und virtuelle Integration
- 5-Schichten-Architektur
- Mediator/Wrapper-Architektur / PDMS
- Schema Mapping
- Schema Matching
- SchemaSQL
- Global-as-View und Lokal-as-View Modellierung
- Global-as-View-Anfragebearbeitung

- Containment & Local-as-View-Anfragebearbeitung
- Bucket Algorithmus
- Verteilte Anfragebearbeitung
- Dynamische Programmierung in verteilten Datenbanken
- Top-N Anfragen
- Duplikaterkennung
- Datenfusion - Union & Co.
- DWH, ETL & Data Lineage
- Informationsqualität
- Hidden Web
- Semantic Web

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

The end.