

Auffinden von und Zugriff auf Datenquellen

Seminar Informationsintegration und Informationsqualität

Dragan Sunjka

TU Kaiserslautern

30. Juni 2006

Gliederung

Autonome Datenquellen

- Autonomieklassen
- Folgen der Autonomie

Mediatorbasierte Systeme

- Mediation
- Wrapper

Hidden Web

- Einleitung
- Automatische Klassifikation von Hidden-Web-Quellen
- Beispiel

Data Management in Grids

- Einleitung
- OGSA-DAI
- DynaGrid

Zusammenfassung

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

Hidden Web
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

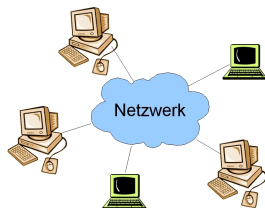
Data Management
in Grids
Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Autonomie

Verteilung führt zu Autonomie...

- ▶ Intra-Organisation: historisch
- ▶ Inter-Organisation: Internet



Autonomie

- ▶ Grad zu dem verschiedene DBMS unabhängig operieren

Autonomieklassen

- ▶ Entwurfsautonomie
- ▶ Kommunikationsautonomie
- ▶ Ausführungsautonomie

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Autonomieklassen

Entwurfsautonomie

- ▶ Datenmodell, Schema

Kommunikationsautonomie:

- ▶ Wahl mit *welchen* Systemen *wann* was kommuniziert wird, Anfragesprache

Ausführungsautonomie

- ▶ Wahl *wann* und *wie* Anfragen ausgeführt werden
- ▶ Wahl der Scheduling- und Optimierungs-Strategie

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Autonomie führt zu Heterogenität

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

Hidden Web
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids
Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Autonomie als Ursache für Heterogenität:

Autonome Systeme

- ▶ Gestaltungsfreiheit
 - ⇒ unterschiedliche Entscheidungen
 - ⇒ Heterogenität
 - ▶ technisch, logisch, semantisch

Gliederung

Autonome Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische Klassifikation von Hidden-Web-Quellen

Beispiel

Data Management in Grids

Einleitung

OGSA-DAI

DynaGrid

Zusammenfassung

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische
Klassifikation von
Hidden-Web-Quellen

Beispiel

Data Management
in Grids

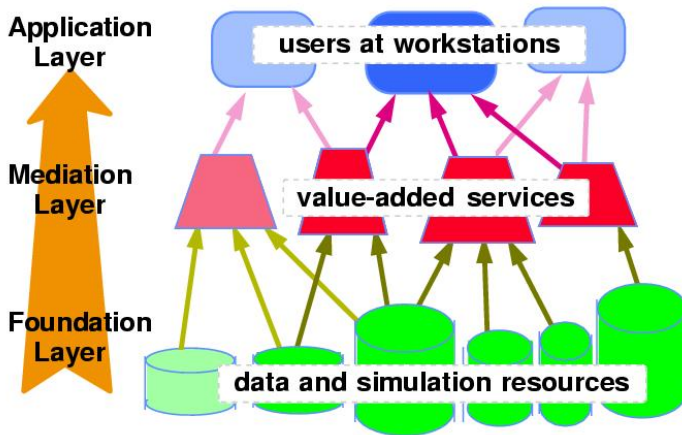
Einleitung

OGSA-DAI

DynaGrid

Zusammenfassung

Mediation



Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Mediation (2)

Funktionen der Mediation

- ▶ Suche und Auswahl von relevanten Datenquellen
- ▶ Transformation der Daten anhand von Metadaten
- ▶ Integration der transformierten Daten
- ▶ Zusammenfassung zur Präsentation

⇒ Transformation von Daten zu Informationen

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

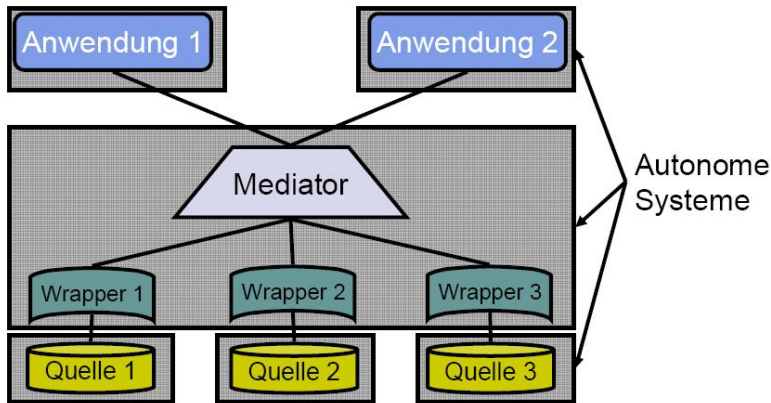
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Wrapper



- ▶ Vermittlung zwischen Mediator und Datenquelle
- ▶ jeweils spezialisiert auf eine Ausprägung autonomer, heterogener Datenquellen

Auffinden von und Zugriff auf Datenquellen

Dragan Sunjka

Autonome Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische Klassifikation von Hidden-Web-Quellen
Beispiel

Data Management in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Wrapper (2)

Vorteile des Wrappings

- ▶ überwinden Heterogenitäten
- ▶ Wiederverwendbarkeit
- ▶ Unabhängigkeit der Datenquellen

Nachteile des Wrappings

- ▶ i.A. schlechtere Leistung
- ▶ Aktualität der Wrapper notwendig

Auffinden von und Zugriff auf Datenquellen

Dragan Sunjka

Autonome Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische Klassifikation von Hidden-Web-Quellen
Beispiel

Data Management in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

XML Wrapper in IBM DB2 II

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

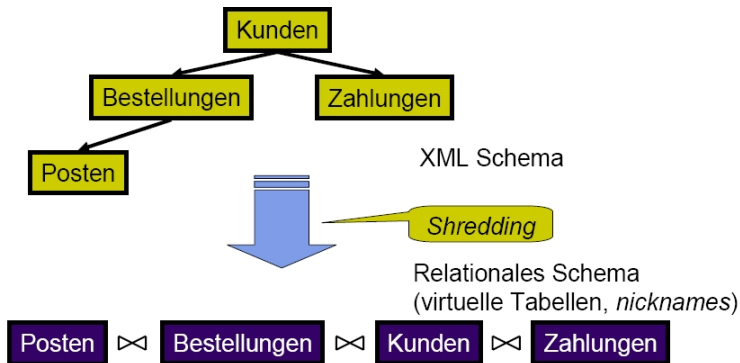
Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung



XML Wrapper in IBM DB2 II (2)

```
CREATE NICKNAME kunden_NN(  
name          VARCHAR(48) OPTIONS (XPATH './name/text()'),  
adresse       VARCHAR(48) OPTIONS (XPATH './address/text()'),  
kunden_NN_ID VARCHAR(48) OPTIONS (PRIMARY_KEY 'YES'))  
FOR SERVER xml_server  
  OPTIONS (XPATH '//customer', FILE_PATH 'customers.xml');
```

```
CREATE NICKNAME order_NN(  
amount        DOUBLE      OPTIONS (XPATH './amount/text()'),  
date          VARCHAR(48) OPTIONS (XPATH './date/text()'),  
order_NN_ID   VARCHAR(48) OPTIONS (PRIMARY_KEY 'YES'),  
customer_NN_FID VARCHAR(48) OPTIONS (FOREIGN_KEY 'CUSTOMER_NN'))  
FOR SERVER xml_server OPTIONS (XPATH './order');
```

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Gliederung

Autonome Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische Klassifikation von Hidden-Web-Quellen

Beispiel

Data Management in Grids

Einleitung

OGSA-DAI

DynaGrid

Zusammenfassung

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische
Klassifikation von
Hidden-Web-Quellen

Beispiel

Data Management
in Grids

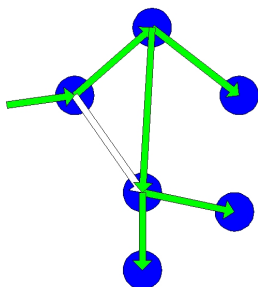
Einleitung

OGSA-DAI

DynaGrid

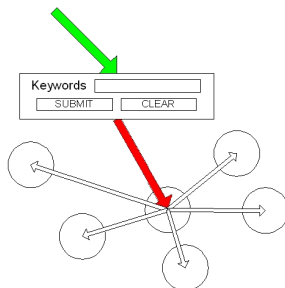
Zusammenfassung

Surface Web vs. Hidden Web



Surface Web

- ▶ Linkstruktur
- ▶ zum *Crawlen* geeignet



Hidden Web

- ▶ keine Linkstruktur
- ▶ Dokumente versteckt in DBMS

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Statistiken zum Hidden Web

- ▶ 550 mal größer als das Surface Web
- ▶ 7500 Terabyte im Hidden Web
- ▶ ca. 84% sind auf Textdokumente spezialisiert
- ▶ ca. 95% des Hidden Web ist öffentlich verfügbar
- ▶ am schnellsten wachsende Kategorie neuer Informationen im Internet

Herausforderungen

- ▶ Auffinden von relevanten Hidden-Web-Quellen
→ Klassifikation
- ▶ Zugriff auf Hidden-Web-Quellen
→ Anfragesprache lernen

Klassifikation: Zuordnung zu Kategorien in einer Hierarchie

- ▶ Manuell
 - ▶ Yahoo!, InvisibleWeb, SearchEngineGuide
- ▶ Automatisch

Zwei Arten von Klassifikation

Coverage (Abdeckung)-basierte Klassifikation

- ▶ #docs über das Thema

Specificity (Spezifität)-basierte Klassifikation

- ▶ #docs/|DB|

Classifier Learning

- ▶ Input: Menge von bereits klassifizierten Dokumenten
- ▶ Output: Menge von Klassifikationsregeln
 - ▶ IF **linux** THEN **Computers**
 - ▶ IF **ibm** AND **intel** THEN **Computers**
 - ▶ IF **jordan** AND **bulls** THEN **Sports**
 - ▶ IF **diabetes** THEN **Health**

Classifier Learning

- ▶ Input: Menge von bereits klassifizierten Dokumenten
- ▶ Output: Menge von Klassifikationsregeln
 - ▶ IF **linux** THEN **Computers** → *+linux*
 - ▶ IF **ibm** AND **intel** THEN **Computers** → *+ibm +intel*
 - ▶ IF **jordan** AND **bulls** THEN **Sports** → *+jordan +bulls*
 - ▶ IF **diabetes** THEN **Health** → *+diabetes*

Query Probing

Classifier Learning

- ▶ Input: Menge von bereits klassifizierten Dokumenten
- ▶ Output: Menge von Klassifikationsregeln
 - ▶ IF **linux** THEN **Computers** → *+linux*
 - ▶ IF **ibm** AND **intel** THEN **Computers** → *+ibm +intel*
 - ▶ IF **jordan** AND **bulls** THEN **Sports** → *+jordan +bulls*
 - ▶ IF **diabetes** THEN **Health** → *+diabetes*

Query Probing

- ▶ Schicke Anfragen an die Datenbank
- ▶ Hole (parse) Anzahl der Ergebnisse

Klassifikation einer Datenbank

1. Schicke Anfragen für Top-Level-Kategorien
2. Hole **Anzahl der Treffer für jede Anfrage**
3. Berechne geschätzte Spezifität und Abdeckung für jede Kategorie
4. Bewege die Datenbank in die sich qualifizierenden Kategorien (mit Spezifität $\geq T_s$, Abdeckung $\geq T_c$)
5. Wiederhole für alle sich qualifizierenden Subkategorien
6. Gib alle Kategorien zurück, die sich qualifiziert haben

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

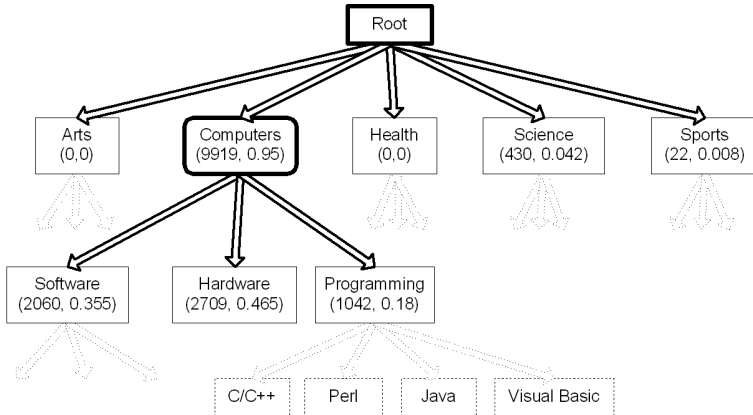
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Beispiel: ACM Digital Library ($T_c=100$, $T_s=0.5$)



Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Gliederung

Autonome Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische Klassifikation von Hidden-Web-Quellen

Beispiel

Data Management in Grids

Einleitung

OGSA-DAI

DynaGrid

Zusammenfassung

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen

Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation

Wrapper

Hidden Web

Einleitung

Automatische
Klassifikation von
Hidden-Web-Quellen

Beispiel

Data Management
in Grids

Einleitung

OGSA-DAI

DynaGrid

Zusammenfassung

Begriffsklärung: Grid

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Grid-Paradigma: Virtualisierung von Ressourcen

- ▶ Ursprüngliche Idee: Rechenleistung aus der Steckdose
- ▶ Ziel: High Performance Super-Computing
- ▶ CPU-Rechenleistung und Arbeitsspeicher als einzige Ressourcen im Grid

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Grids

- ▶ anfangs: Datentransfer nur dateibasiert, Input/Output-Dateien
- ▶ später: Speicherplatz als Ressource
- ▶ heute: (semi-)strukturierte Daten als Ressourcen

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

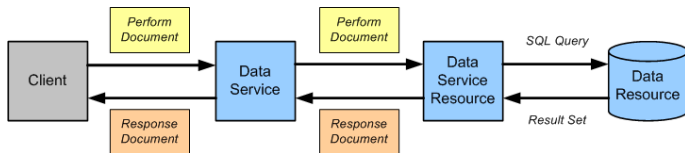
Open Grid Services Architecture - Data Access and Integration

- ▶ ein in Java entwickeltes Framework für den Zugriff auf Datenbanken
- ▶ Ziel: einheitliche Schnittstelle zum Zugriff auf heterogene Datenbanken
- ▶ versucht nicht, die Kluft zwischen SQL und XML zu überbrücken
⇒ Überwindung technischer Heterogenität

Anfragen an Quellen

- ▶ erfolgen mittels XML-basierten Perform-Dokumenten
 - ▶ Basisaktivitäten (Anfragen, Transformationen, Übertragung)

Interaktion von Komponenten



Herausforderungen

- ▶ Common Data Model?
- ▶ SQL vs. XML
- ▶ XML WebRowSet ist flexibel aber ineffizient

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Dynamische Integration in Grid-Umgebungen



Dynamische Integration heterogener,
autonomer, verteilter Datenquellen

Merkmale

- ▶ Zum Zeitpunkt der Anfrage gibt es kein globales Schema
- ▶ automatische Integration von neuen Quellen
- ▶ OGSA-DAI als Infrastruktur

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

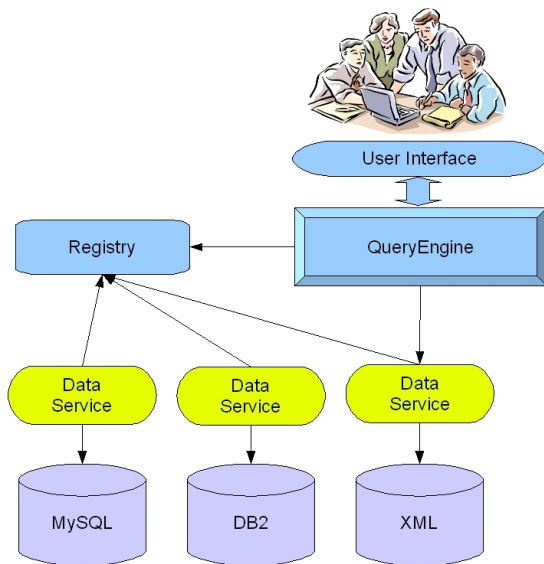
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

DynaGrid Architektur



Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen
Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme
Mediation
Wrapper

Hidden Web
Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids
Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Zusammenfassung

- ▶ viele Varianten klassischer Mediator/Wrapper-Architektur
 - ▶ Wrapping zur Überbrückung der Heterogenität
- ▶ Hidden Web wächst schnell
 - ▶ Wie findet man solche Eintrittspunkte?
 - ▶ Semantic Web?
- ▶ zahlreiche autonome Datenquellen und schnell wechselnde Benutzeranforderungen
 - ▶ dynamische Discovery und flexible Schema-Matching-Verfahren

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung

Vielen Dank für die Aufmerksamkeit

Fragen?

Auffinden von und
Zugriff auf
Datenquellen

Dragan Sunjka

Autonome
Datenquellen

Autonomieklassen
Folgen der Autonomie

Mediatorbasierte
Systeme

Mediation
Wrapper

Hidden Web

Einleitung
Automatische
Klassifikation von
Hidden-Web-Quellen
Beispiel

Data Management
in Grids

Einleitung
OGSA-DAI
DynaGrid

Zusammenfassung