

Servicebasierte Datenintegration

Präsentation zur Seminararbeit

Christoph Aßmann



UNIVERSITÄT LEIPZIG

Abteilung Datenbanken
am Institut für Informatik

Seminar Cloud Data Management
Servicebasierte Datenintegration

Aßmann, Christoph Leipzig, 26.01.2010

Folie 1

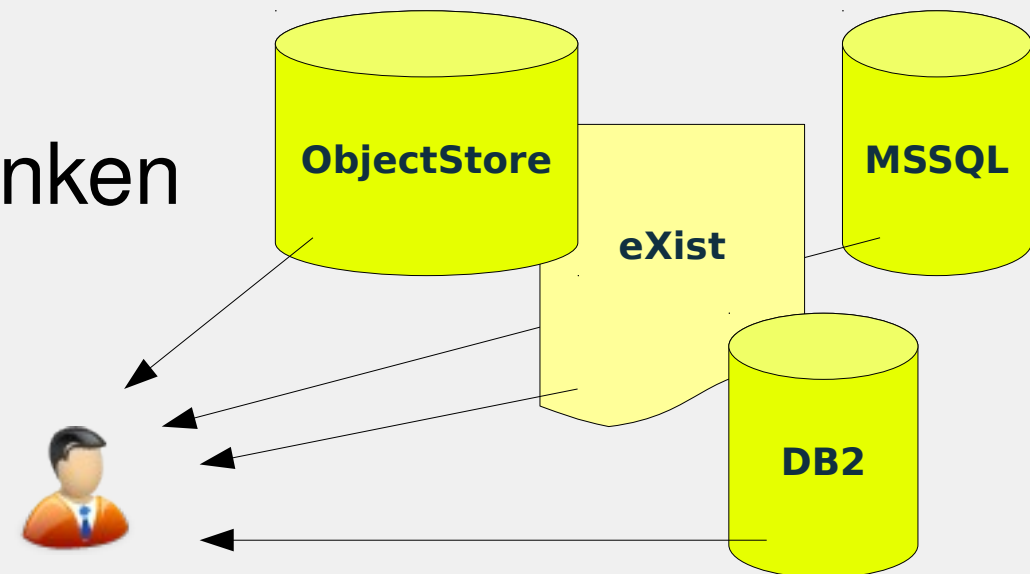


- Begriffe
- Motivation
- Abgrenzung Grid – Cloud
- OGSA: Architektur servicebasierter Grids
 - Standardisierung
 - Evolution / Erweiterungen
- Informatica Cloud Data Integration Solutions
- Zusammenfassung

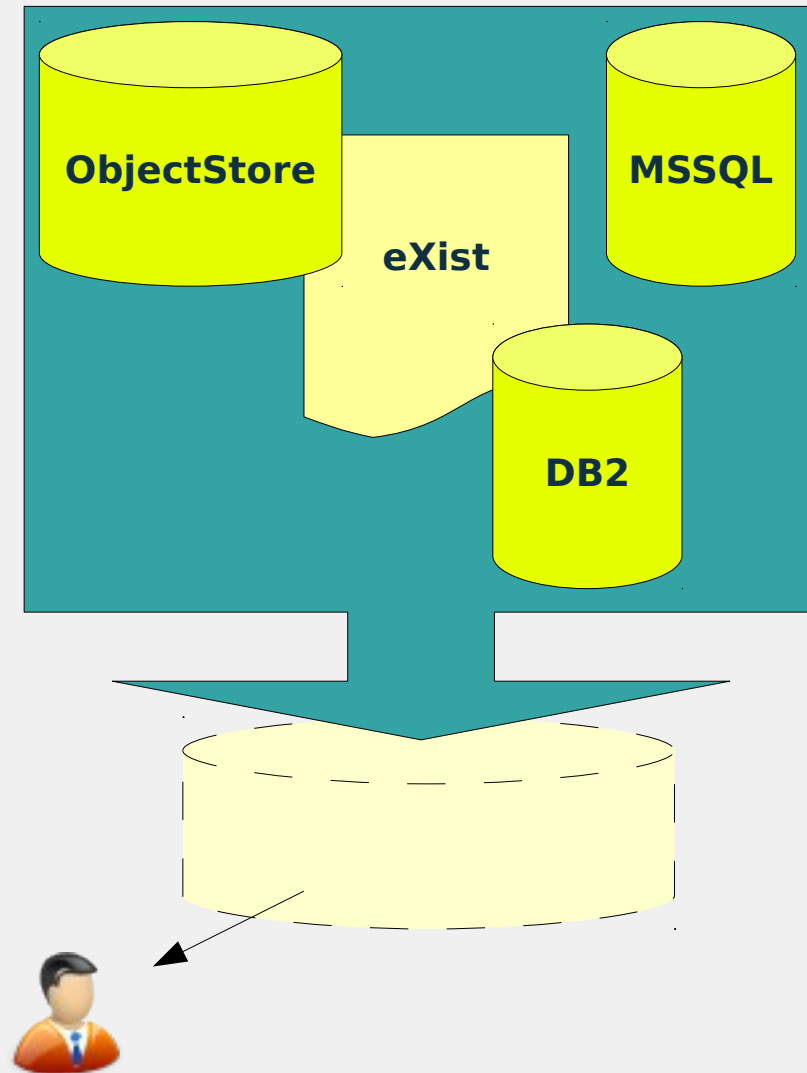


- Ziel: einheitlicher Zugriff auf mehrere DBs
 - Verteilte Datenbanken
 - Dezentral verwaltete Datenbanken

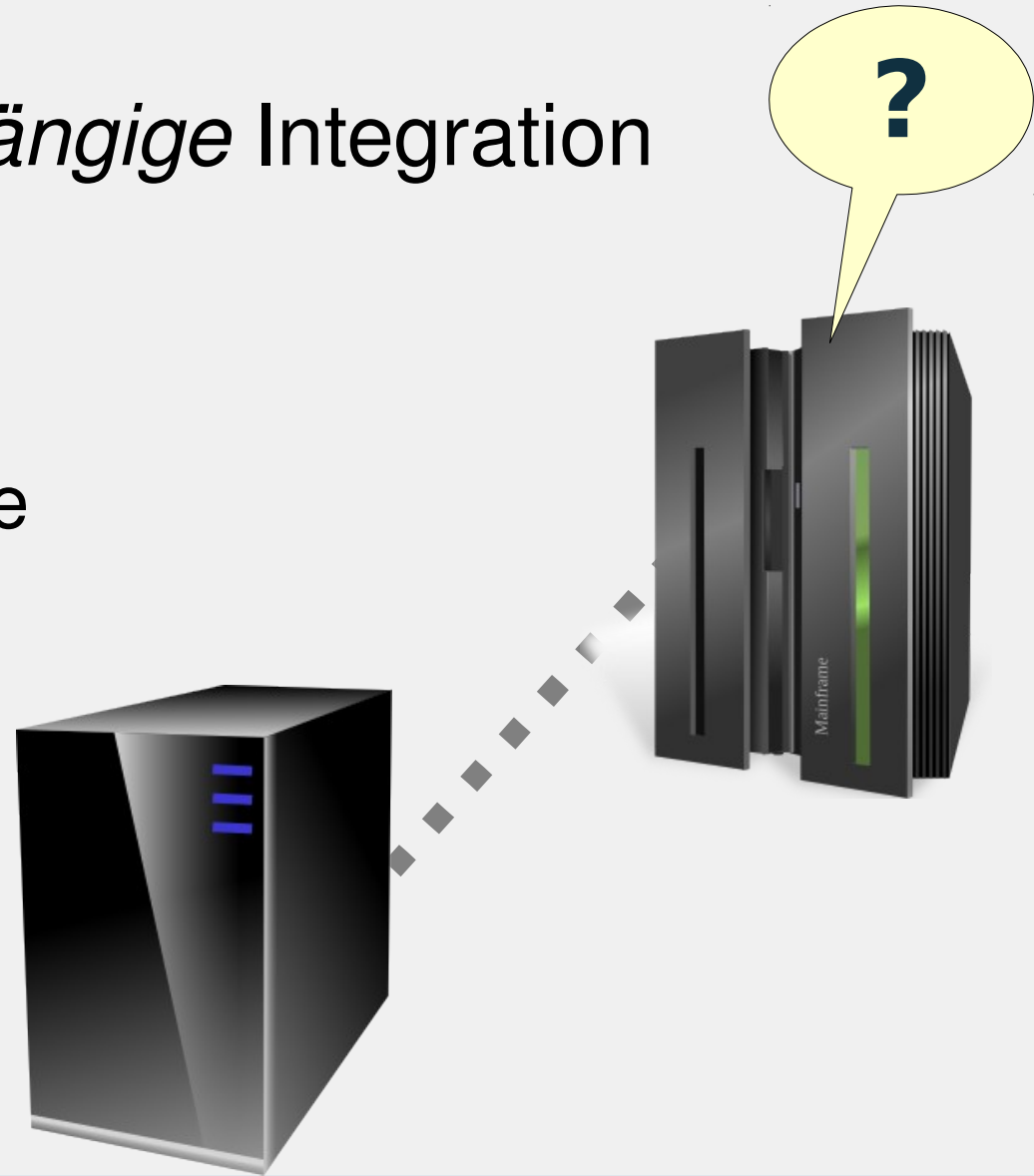
→ Problem:
Heterogene Datenbanken



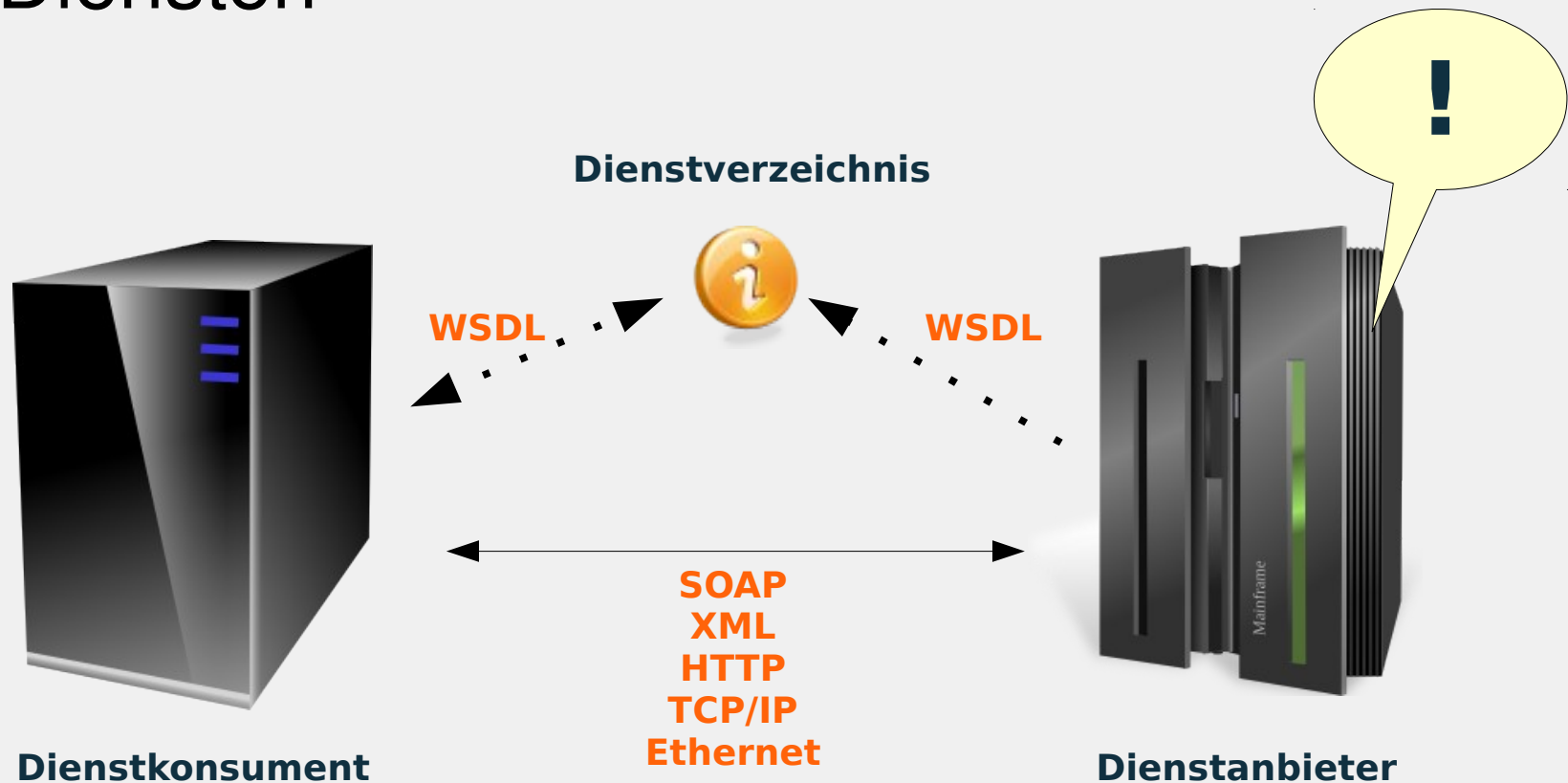
- Heterogenität: Ausprägungen
 - Zugriff
 - Syntax
 - Datenmodell
 - Schema
 - ...
- Data Cleaning
- Erstellung einheitlicher Sicht



- Ziel: *technologieunabhängige* Integration heterogener Systeme
 - Transportprotokoll
 - Programmiersprache
 - Plattform

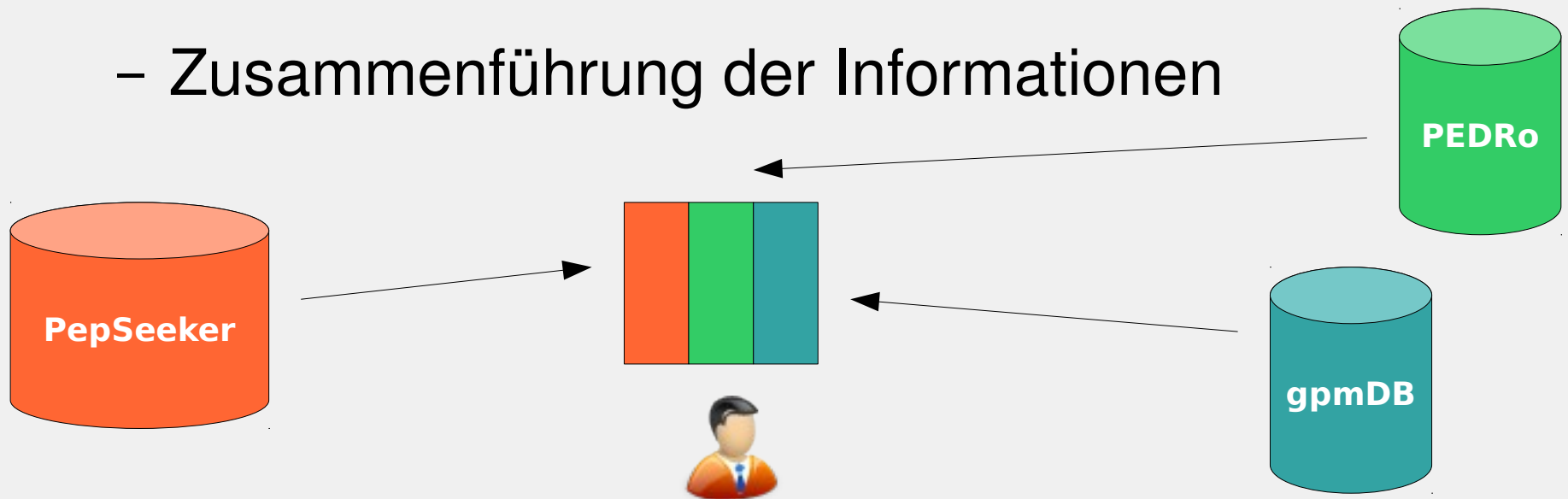


- Beschreibung, Veröffentlichung, Suche, Nutzung von Diensten



- ISPIDER Project

- Erfassung biochemischer Daten
- Identifizierung von Proteomen an verschiedenen Standorten
- Zusammenführung der Informationen

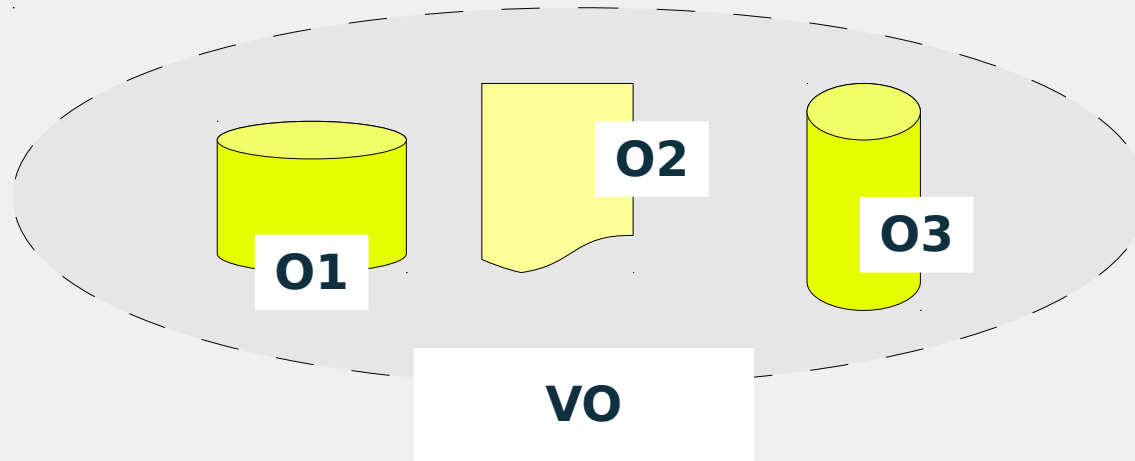


- Gemeinsamkeit: Dynamische Bereitstellung von Speicher und Rechenkapazität über Netzwerk
- Unterschiede:
 - Cloud
 - Zentraler Anbieter
 - Ökonomischer Ansatz
 - Grid
 - Virtuelle Organisation (VO) / dezentrale Admin.
 - Wissenschaftlicher Kontext



Abgrenzung Grid - Cloud

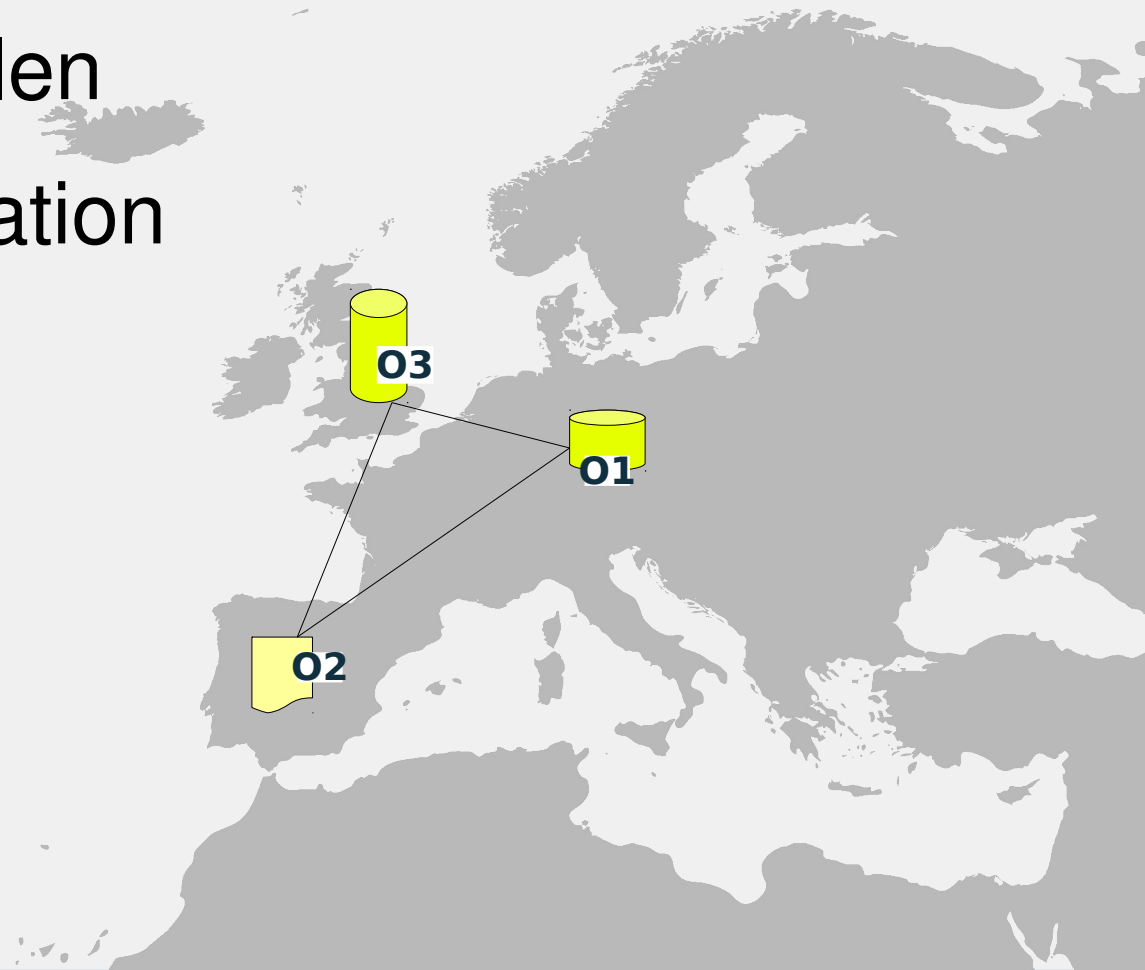
Grid



Cloud



- Hoher Grad an Heterogenität
 - Autonome Datenquellen
 - Dezentrale Administration
- Standardisierung



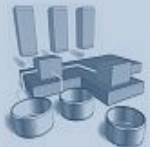
- Standardisierungsgremium:
Open Grid Forum (OGF)
- Diverse Arbeitsgruppen, u.a.:
Database Access and Integration Services (DAIS-WG)
- Architektur:
Open Grid Services Architecture (OGSA)



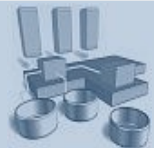
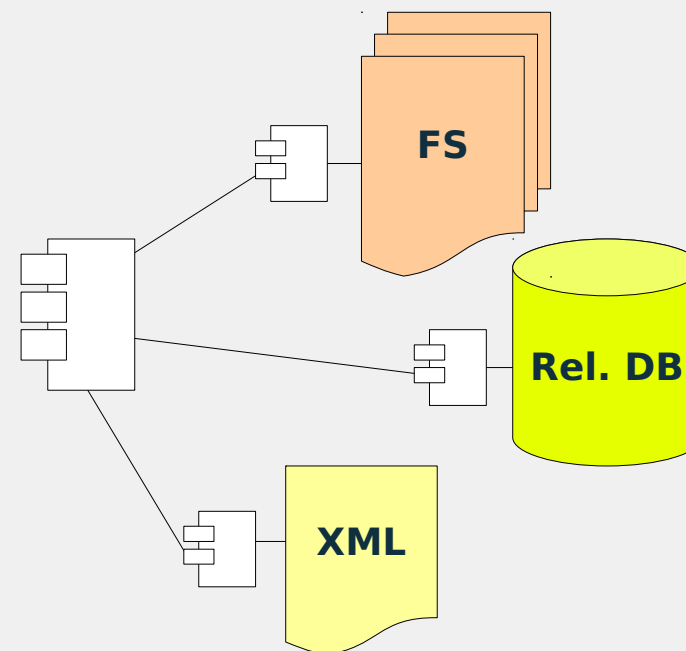
- Repräsentation von Ressourcen durch Dienste
→ Grid Services
- Standardisierung von Schnittstellen
 - Auffinden von Grid Services
 - Erzeugen / Beenden von Grid-Dienstinstanzen
 - Nachrichtenaustausch



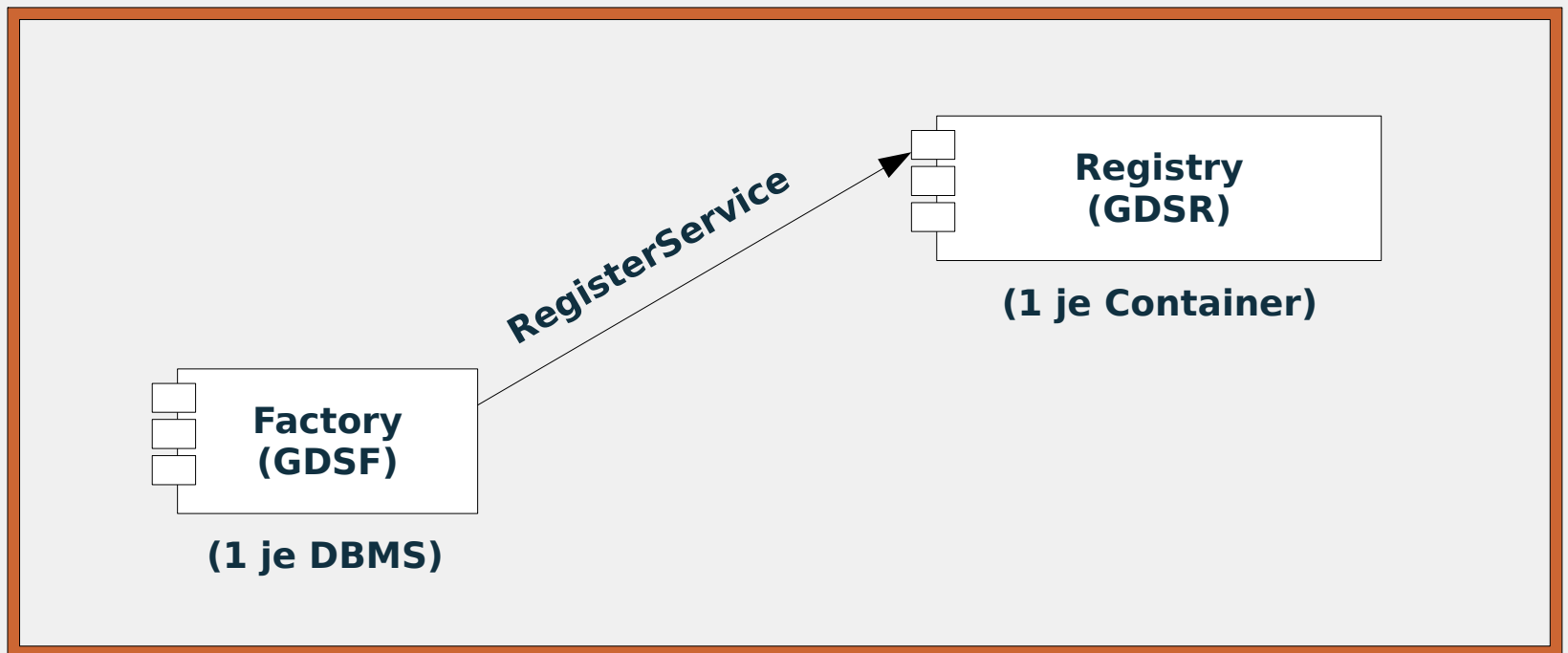
- OGSA-DAI:
dienstbasierter Zugriff auf Datenbanken
- OGSA-DQP:
Koordinierung des Zugriffs auf mehrere
OGSA-DAI-Ressourcen



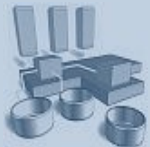
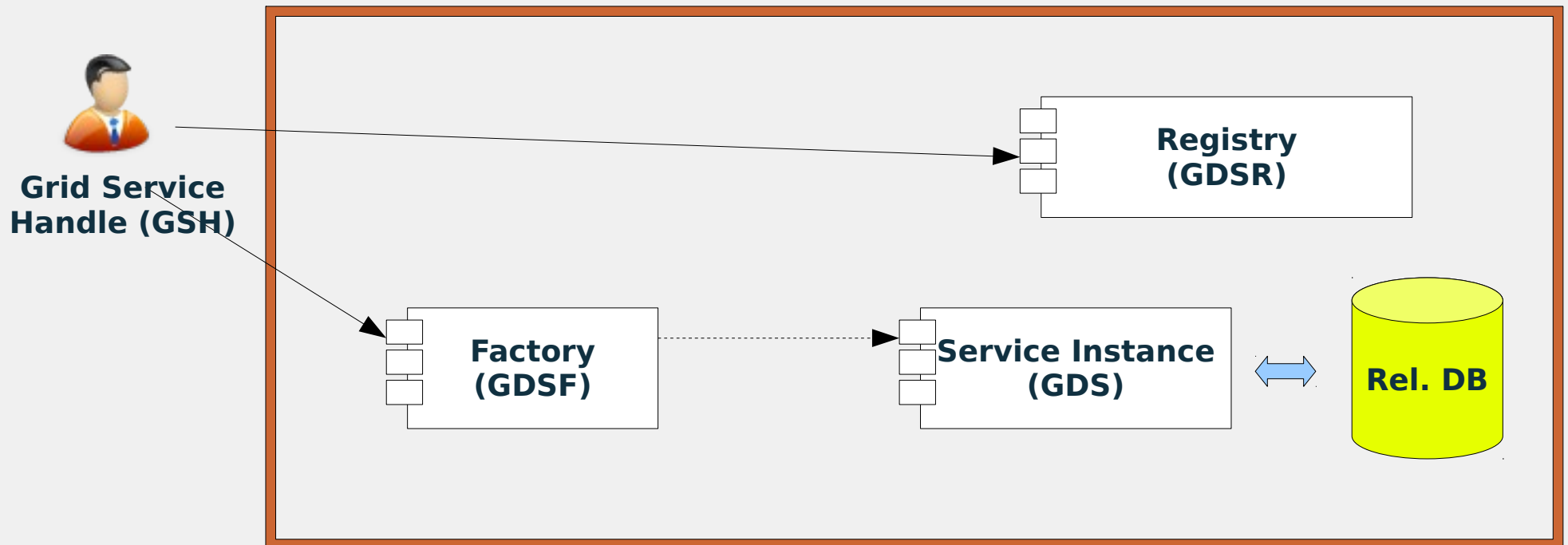
- Data Access and Integration (OGSA-DAI)
- Java-basierte Middleware zur Integration heterogener Datenquellen auf Basis der OGSA



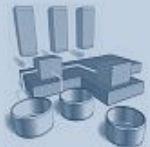
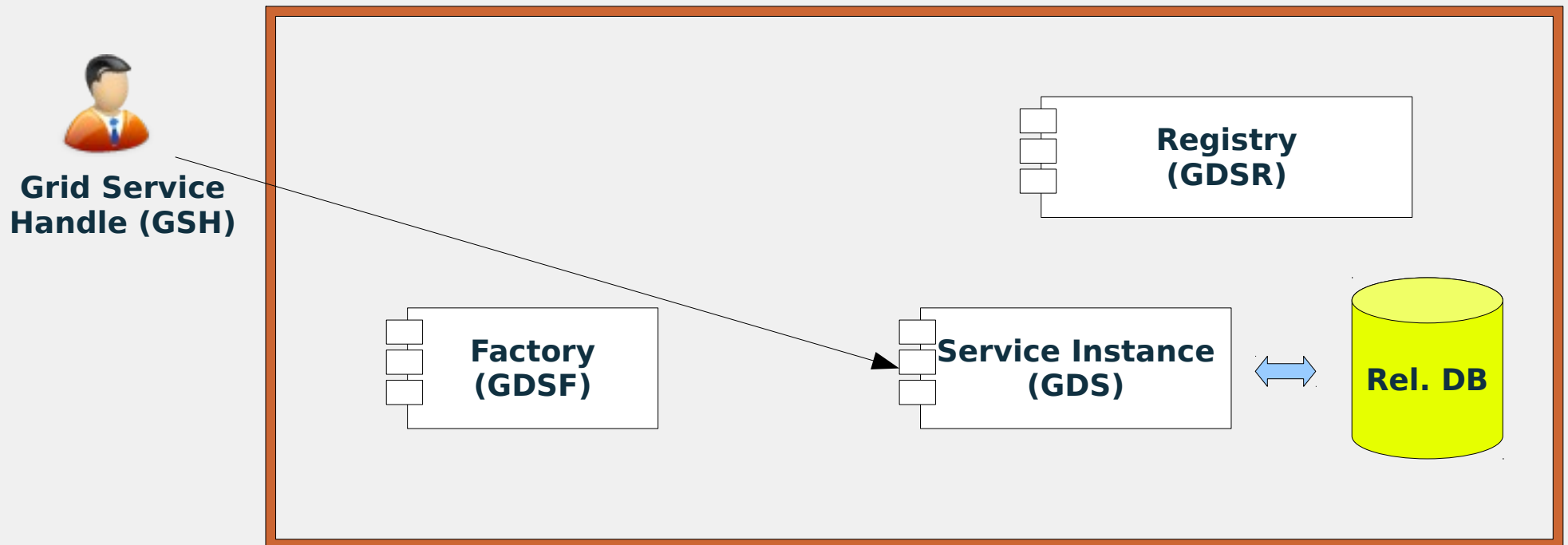
- Ablauf einer Anfrage gegen das Grid
(a) Container-Start



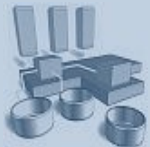
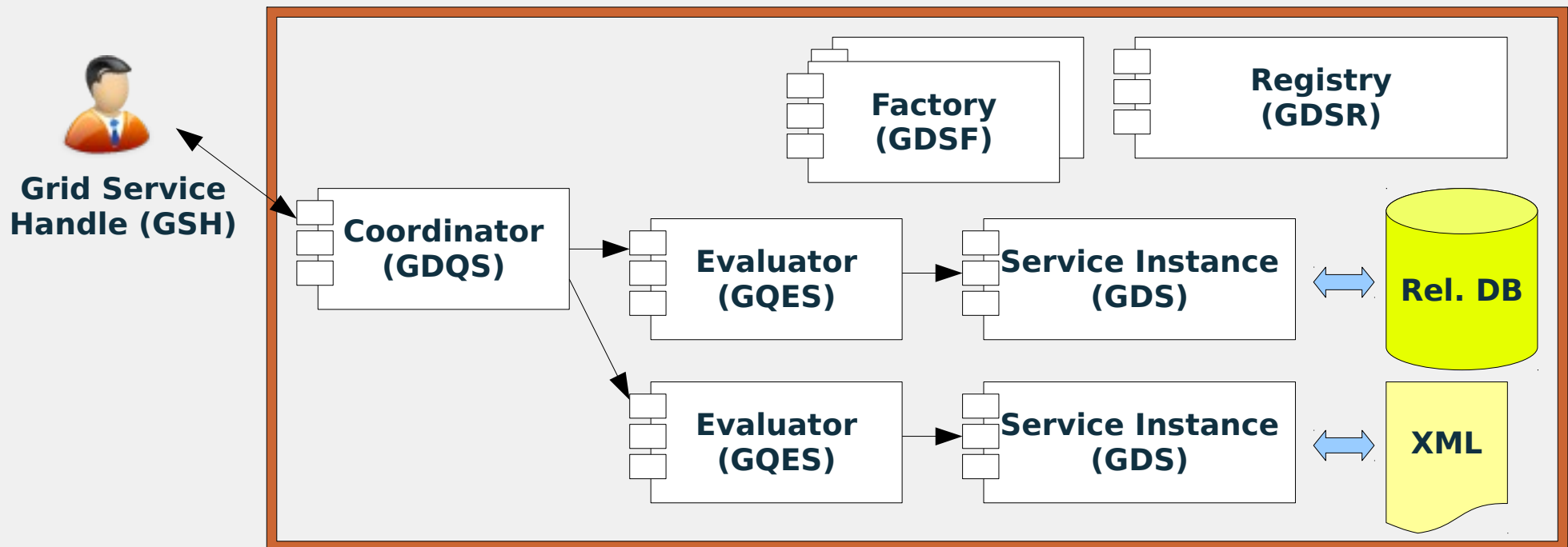
- Ablauf einer Anfrage gegen das Grid
(b) Dienstlokalisierung



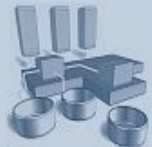
- Ablauf einer Anfrage gegen das Grid
(c) Durchführen einer Anfrage



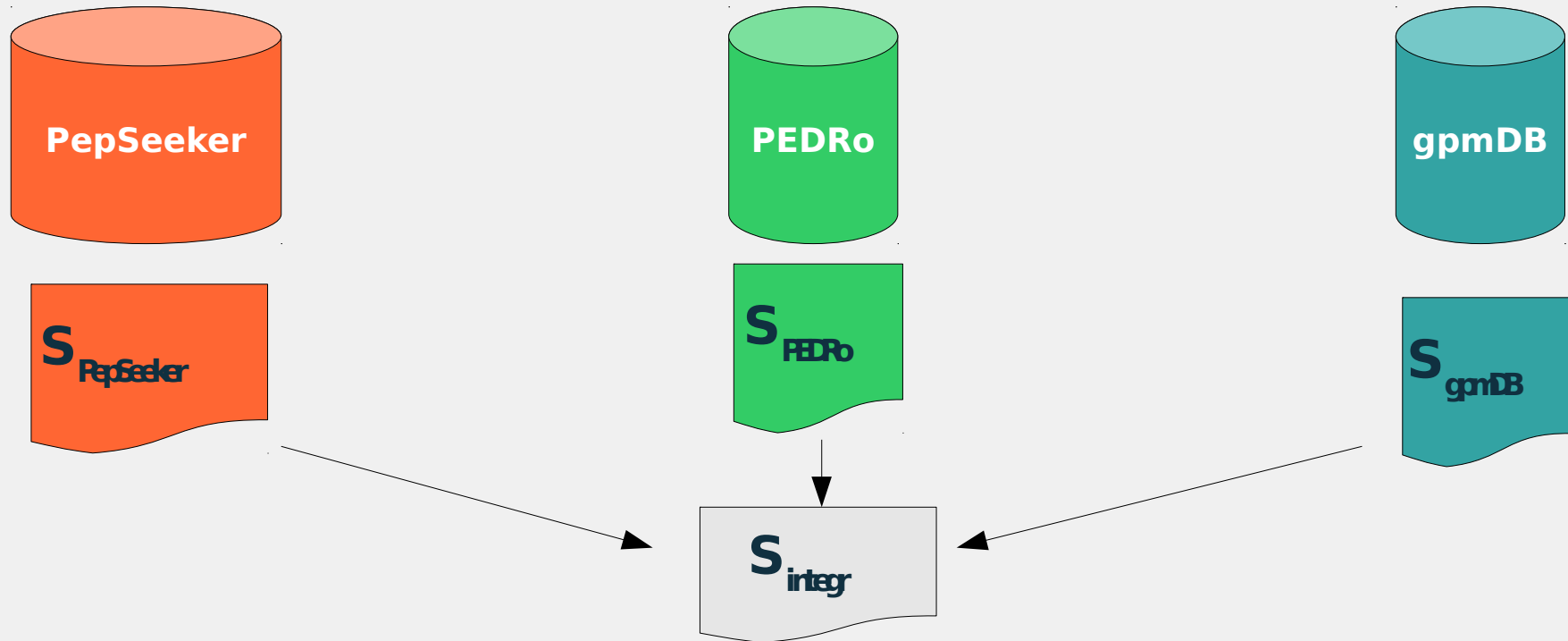
- Ablauf einer Anfrage gegen das Grid
(c) Durchführen einer verteilten Anfrage via DQP



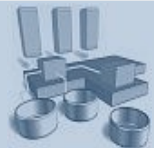
- Dezentrale Administration
 - Schemaevolution
 - einmalige Schemaintegration unzureichend
- Grid Data Integration System (GDIS)
 - Mapping-Katalog
 - Änderung / Hinzufügen einer Ressource: Update
- Ausführliche Beschreibung s. Ausarbeitung



ISPIDER: virtuelle Datenintegration



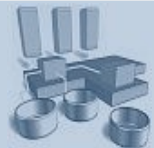
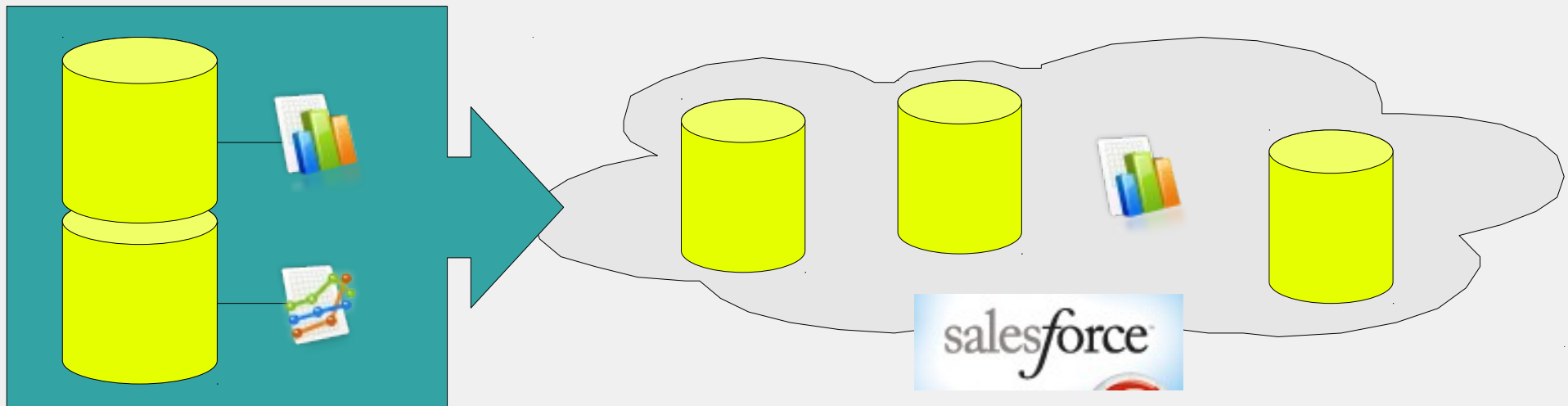
LSID	src_id	acc_nr	col_xy
URN:LSID:ispider.man.ac.uk:pepseeker.protein:1	1	ENSP00000339074	...
URN:LSID:ispider.man.ac.uk:pepseeker.protein:2	2	ENSP00000339074	...
URN:LSID:ispider.man.ac.uk:pedro.protein:1	1	ENSP00000339074	...
URN:LSID:ispider.man.ac.uk:gpmdb.protein:1	1	ENSP00000339074	...



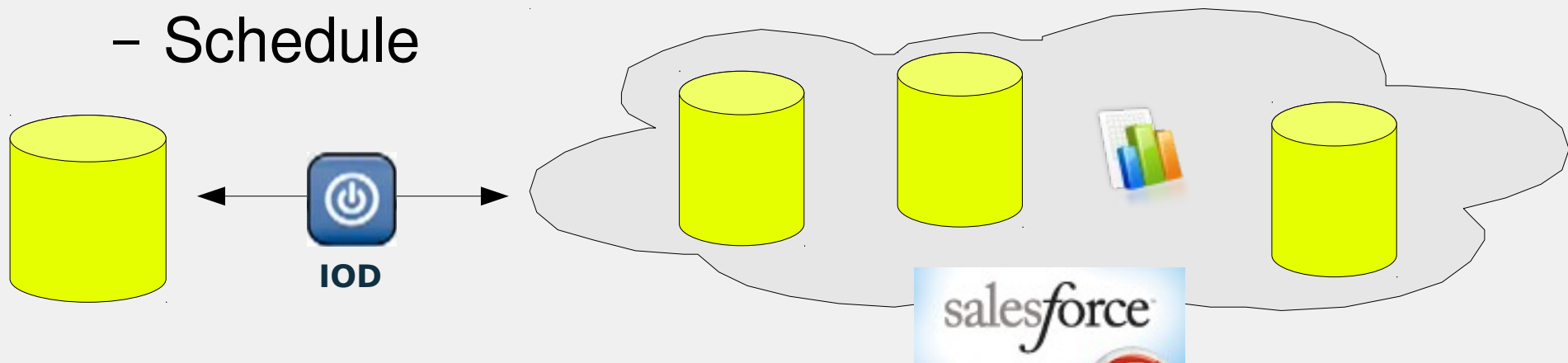
- Nutzung von OGSA-DAI und OGSA-DQP
- Manuelles Erstellen des globalen Schemas
 - Basis: S_{FEDo}
 - Beschreibung der Korrespondenzen $S_{\text{local}} \rightarrow S_{\text{global}}$
 - IQL, Transformation Pathways
- Speicherung in *Schemas & Transformations Repository*
 - Schemaevolution möglich



- SaaS Anwendung: Salesforce CRM
- Cloud Plattform: Force.com
- Problematik: Migration lokaler Datenbestände aus Legacy Software nach Salesforce CRM



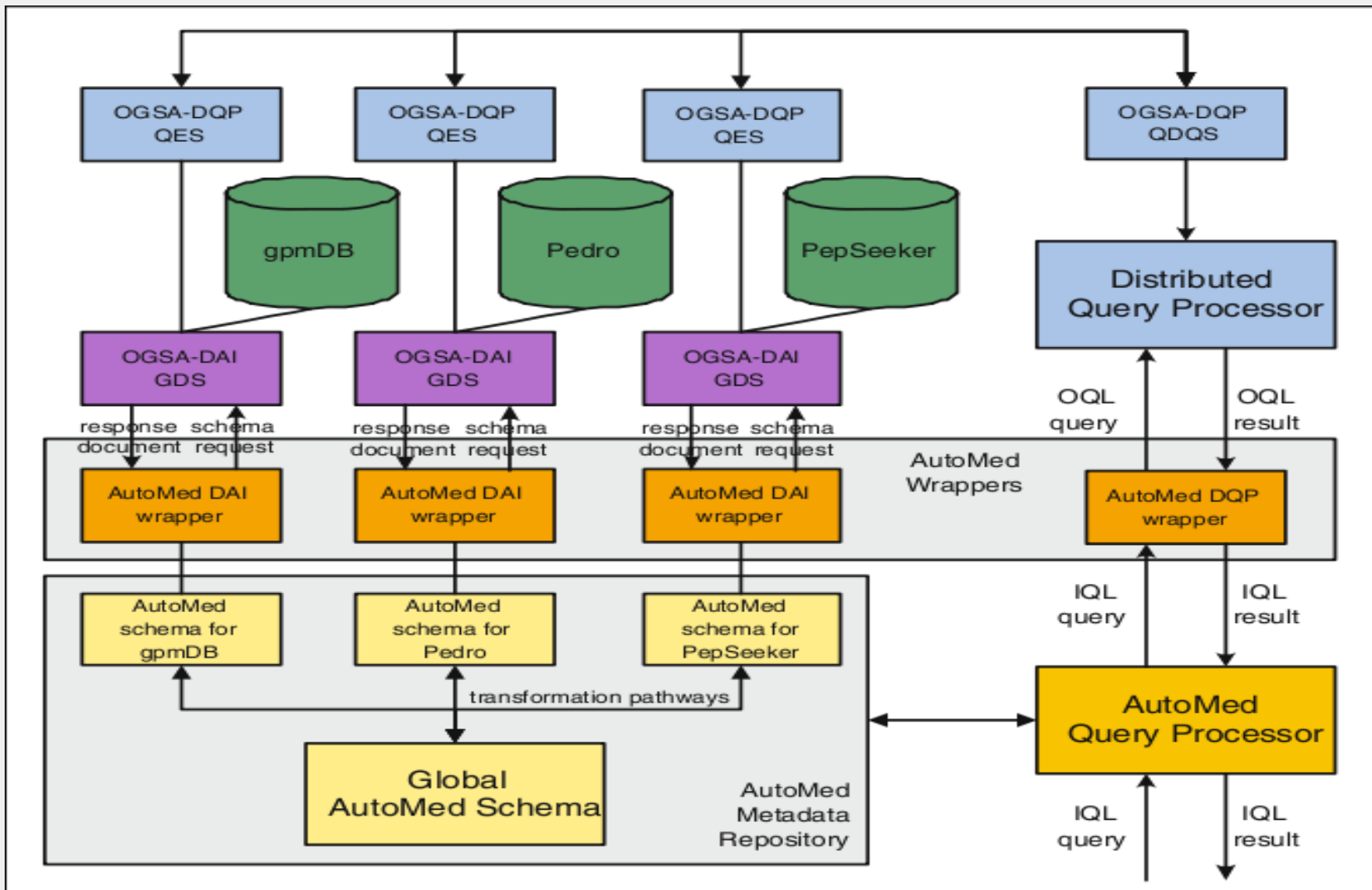
- Informatica On Demand (IOD):
Webbasierte Konfiguration von Verbindungen
 - Source Connection, Target Connection
 - Filter Definition
 - Field Mappings, Transformation
 - Schedule



- Integrationsproblematik in Grids und Clouds vorhanden
- Lösung in Grids: Standardisierung
 - Virtuelle Integration
 - Dienstbasierte Ansätze
 - OGF
- Lösung in Clouds: Eigenentwicklungen
 - Physische Integration
 - Migration in die Cloud
 - Proprietäre Ansätze



ISPIDER Architecture



[ZFB+ 06]



- [ZFB+ 06] Lucas Zamboulis, Hao Fan, Khalid Belhajjame, Jennifer A. Siepen, Andrew Jones, Nigel J. Martin, Alexandra Poulouvassilis, Simon J. Hubbard, Suzanne M. Embury und Norman W. Paton. Data Access and Integration in the ISPIDER Proteomics Grid. In Ulf Leser, Felix Naumann und Barbara A. Eckman, Hrsg., DILS, Jgg. 4075 of Lecture Notes in Computer Science, Seiten 3–18. Springer, 2006.

