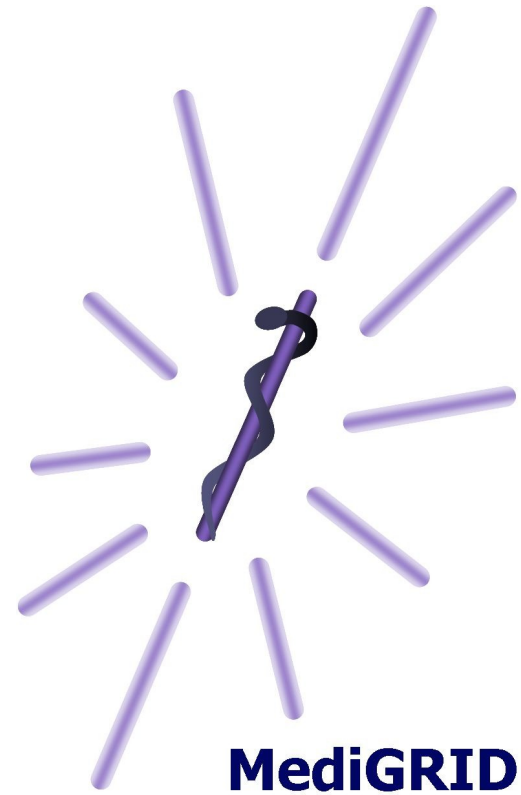


Verwaltung von Dateien und Datenbanken im Grid mit SRB und OGSA-DAI



Kathrin Peter, Zuse-Institut Berlin

***Datenmanagement Workshop
Berlin, 18-19.05.2006***

- ❑ MediGRID Projekt
- ❑ Anwendungen in MediGRID
- ❑ Architektur und Datenmanagement
- ❑ Verwendung von SRB und OGSA-DAI
- ❑ Anforderungen und nächste Schritte

- ❑ Anwender aus Bioinformatik, medizinischer Bildverarbeitung und klinischer Forschung.
- ❑ Ziel: Arbeitsplattform für Biomedizinische Forschung am Beispiel ausgewählter Anwendungen
- ❑ Momentan keine gemeinsame Datennutzung wegen Heterogenität der Anwendungen und der Daten.
- ❑ Perspektivisch Vernetzung medizinischer und biologischer Daten.
- ❑ Aufgaben:
 - Globale Datenspeicherung auf verteilten heterogenen Ressourcen
 - Bereitstellung von Speicherplatz
 - Bereitstellung und transparenter Zugriff auf Datenbanken
 - Bereitstellung eines Metadatenmanagementsystems
 - Unterstützung der Anwender bei der Organisation und Speicherung großer Datenmengen

- ❑ Patientenakte in computerlesbarer Form.
- ❑ Sammlung und Dokumentation aller Daten zur Gesundheitsversorgung eines Patienten, z.B. Befunde, Behandlungen.
- ❑ Die gewaltige Datenflut benötigt geeignete Mechanismen zur Verwaltung der Daten.
- ❑ Sortier- und Suchkriterien zur Navigation und zur Unterstützung einer multiplen Nutzung der Daten.
- ❑ Erfüllung strenger Sicherheitsanforderungen zur Langzeitarchivierung, Zugriffskontrolle und kryptografischen Sicherung der Kommunikationswege und Speicherung selbst.
- ❑ Vorteile:
 - Schnelle, umfassende Verfügbarkeit aller relevanten Informationen.
 - Effiziente Verwaltung.
- ❑ Nachteile:
 - Technikabhängigkeit, Zuverlässigkeit, Sicherheit, Missbrauch
 - Entscheidungsfindung auf Untersuchungsergebnissen anderer Ärzte.

Anwendungsklassen

Bioinformatik

- Sequenzanalyse
- Genomdatenanalyse
- SequCorr
- RNAi Pipeline

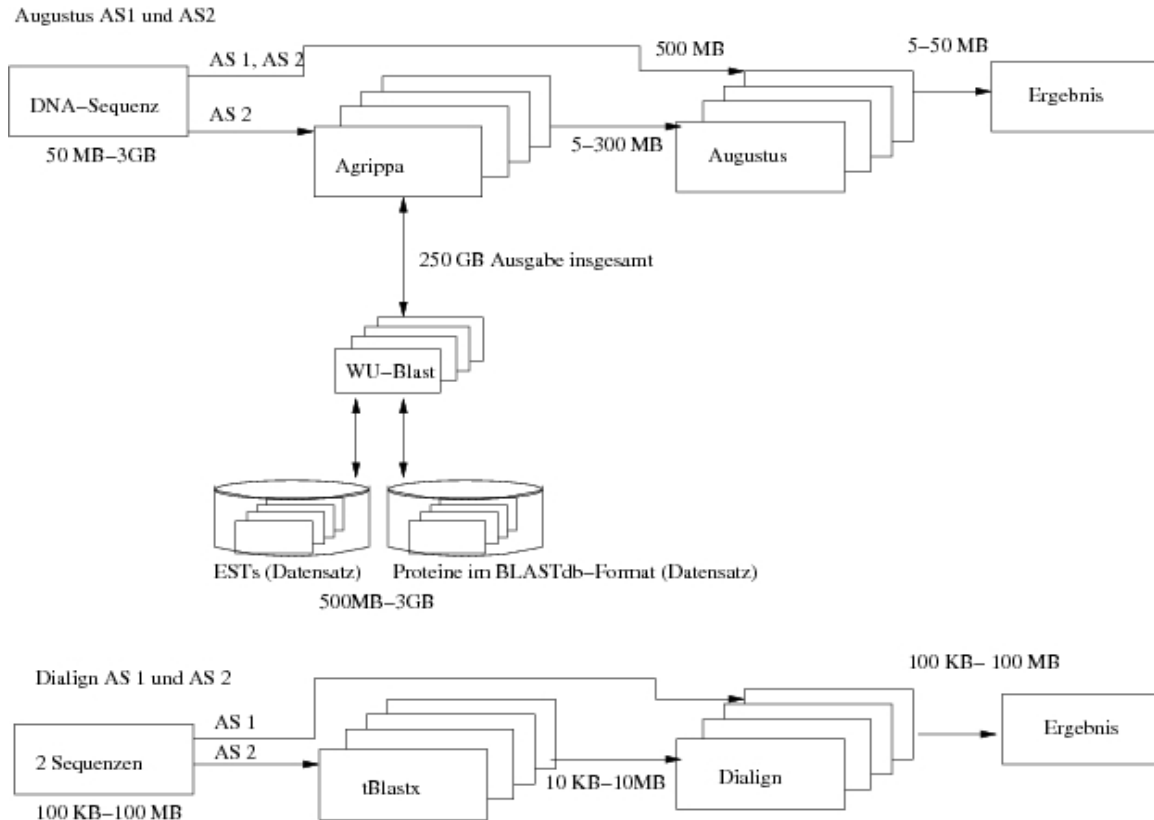
Medizinische Bildverarbeitung

- Med. Bildverarbeitung:
Prostatabiopsie
- Virtuelle Gefäßchirurgie
- Funktionelle Hirnbilddaten

- ❑ Datentransfer:
 - Eingabedatensatz (Upload)
 - Berechnungen (DB-Zugriffe, temporäre Zwischenergebnisse, Nutzer-Interaktion zur Auswahl und Einschränkung von Daten)
 - Ausgabedatensatz (Download)

- ❑ keine permanente Datenhaltung erforderlich
- ❑ Nutzer ist für Ein- und Ausgabedaten selbst verantwortlich

Bioinformatik: Sequenzanalyse

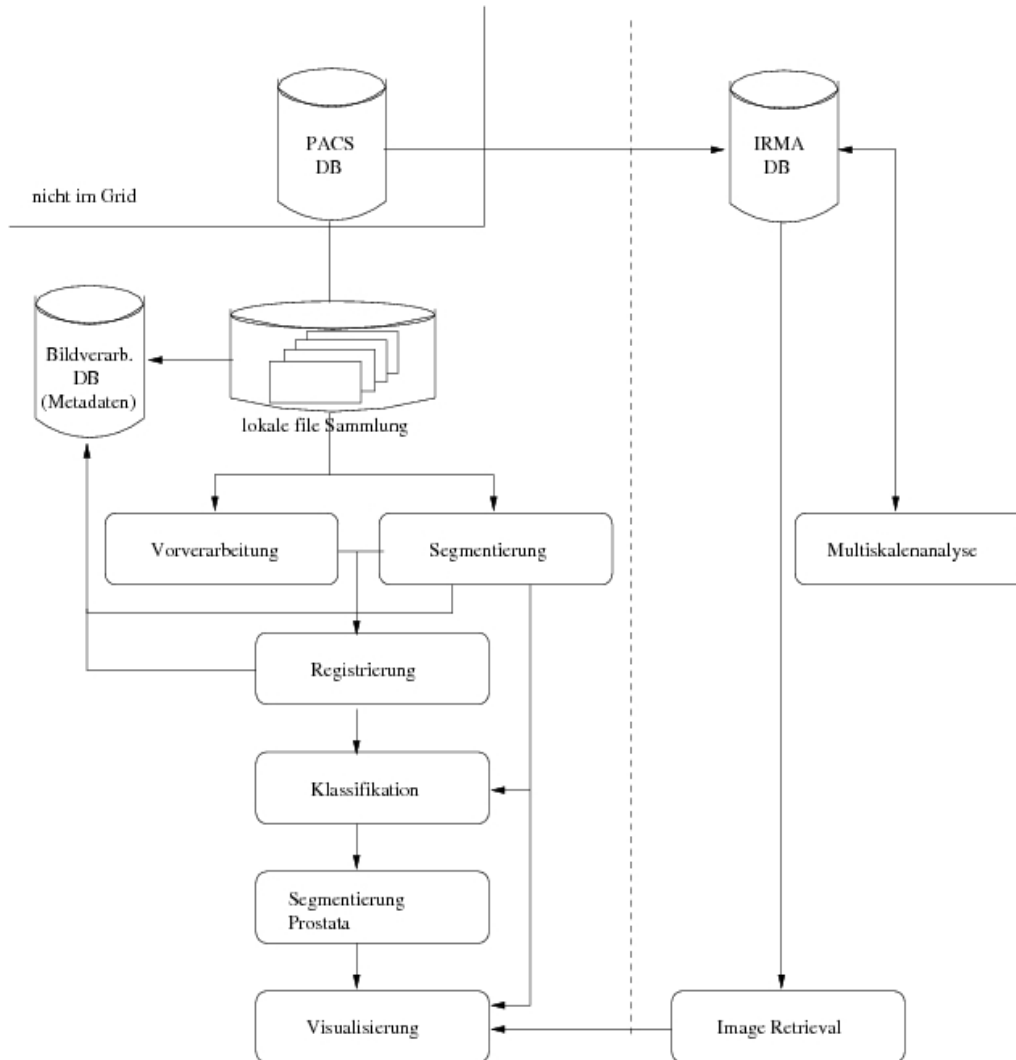


- ❑ Upload DNA-Sequenzen
- ❑ große Zwischenergebnisse, großer Hauptspeicherbedarf
- ❑ temporäre Bereitstellung der Ergebnisse zum Download

Anwendungsklasse: Medizinische Bildverarbeitung

- ❑ Datentransfer:
 - Eingabedatensatz: Bilddaten + Metadaten
 - Verarbeitung (Produktion einer großen Menge neuer Daten, Interaktion zur Parameterauswahl und zur Wiederholung von Berechnungen usw.)
 - Ergebnis (weitere Interaktion zur Ergebnisauswahl und zum Postprocessing)
- ❑ Unterstützung durch Datenmanagement bei permanenter Datenspeicherung, Datensuche über Metadaten
- ❑ Zentrale Ergebnisbereitstellung für Lernzwecke

Medizinische Bildverarbeitung: Prostatabiopsie



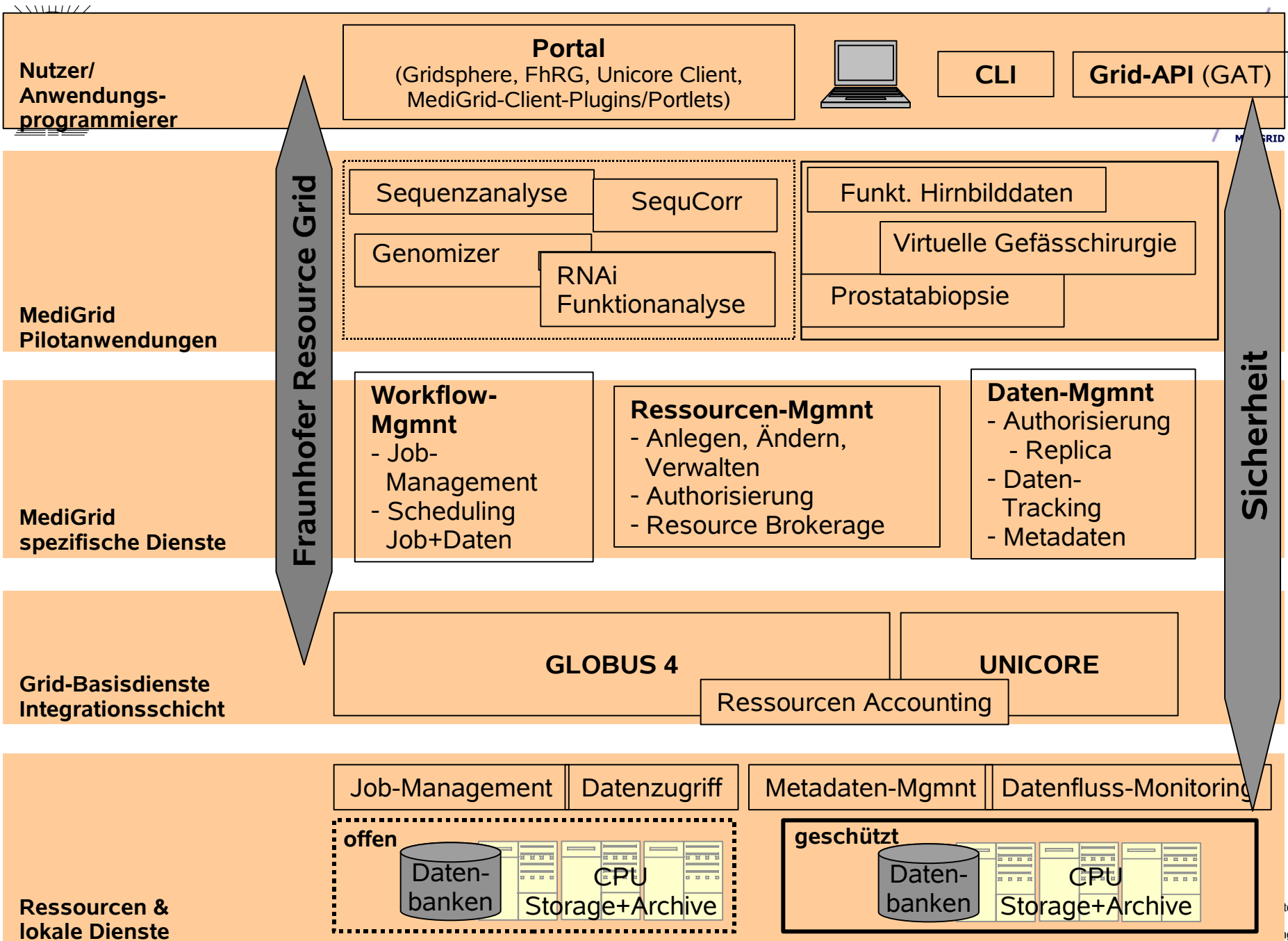
- ❑ Datei Sammlung (flat-files) mit Bilddaten
- ❑ verschiedene Bearbeitungsschritte
- ❑ Speicherung von Zwischenergebnissen; Ergebnisspeicherung
- ❑ Lerndatenbank IRMA

Anwendungsklassen

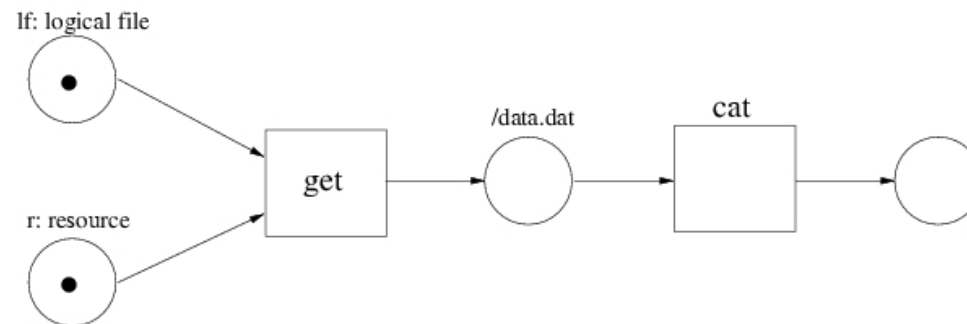
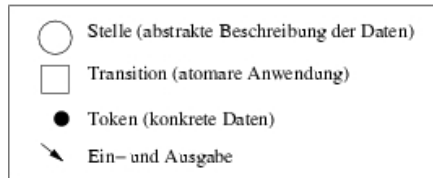
- ❑ 250 GB Plattenplatz für Zwischenergebnisse (Sequenzanalyse)
- ❑ bis 100 GB Plattenplatz für Datenbank, 100 GB Plattenplatz für Dateien (Genomdatenanalyse, RNAi-Pipeline)
- ❑ 300 MB pro Datensatz + 3*300 MB zu archivierende Daten pro Berechnung (funktionelle Hirnbilddaten)
- ❑ Art der Datenoperationen:
 - Datenupload: von nicht Grid-Rechner auf Grid-Rechner
 - Datentransfer zwischen Grid-Rechnern: Übertragung von Zwischenergebnissen und temporären Daten
 - Langfristige Datenspeicherung mit Ablage der Daten im Grid-Storage-System. Erzeugung von Metadaten u.a. für Suche nach Daten.
 - Datendownload: Ergebnisdownload von Grid-Rechner auf nicht Grid-Rechner

MediGRID-Architektur

- ❑ Basis-Technologie: Globus Toolkit 4
- ❑ Übernahme und Weiterentwicklung von Komponenten des Fraunhofer Resource Grid (FhRG)
- ❑ Hauptzugang über Portal, optional Nutzung von Diensten über Commandline Interfaces
- ❑ Ressourcenmetadaten mit MDS4, Ganglia Cluster Monitoring System und Ressourcenverwaltungssystem aus dem FhRG
- ❑ Datenmanagement zur Einführung einer Datenabstraktionsschicht. Datenressourcen: Datenbanken und Filesysteme



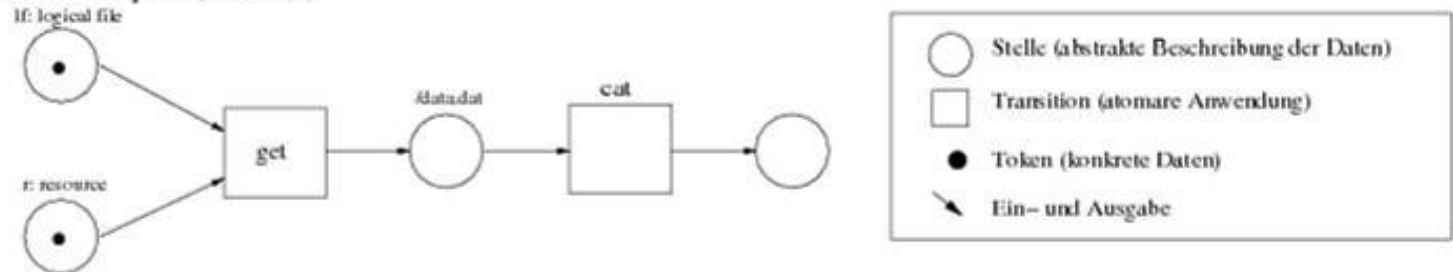
- ❑ Sicherheit: Nur autorisierte Benutzer dürfen Zugriff auf bestimmte Daten erhalten.
Portal: Anmeldung mit Passwort und Erzeugung kurzlebiger Credentials mit Hilfe des auf dem MyProxy-Server abgelegten MyProxy.
-> SRB GSI Authentifizierung Mapping, OGSA-DAI Mapping
- ❑ Metadatenmanagement: Verwaltung von Metadaten über Dateien, Verzeichnisse, Benutzer usw., benutzereigenen Metadaten und Replikainformationen.
-> SRB MCAT
- ❑ Datenmanagement bietet verschiedene Services, vorzugsweise als Webservices.
- ❑ Modellierung von Datentransfers und komplexeren Datenmanagementfunktionalitäten im Workflow mit Hilfe von Petri-Netzen.



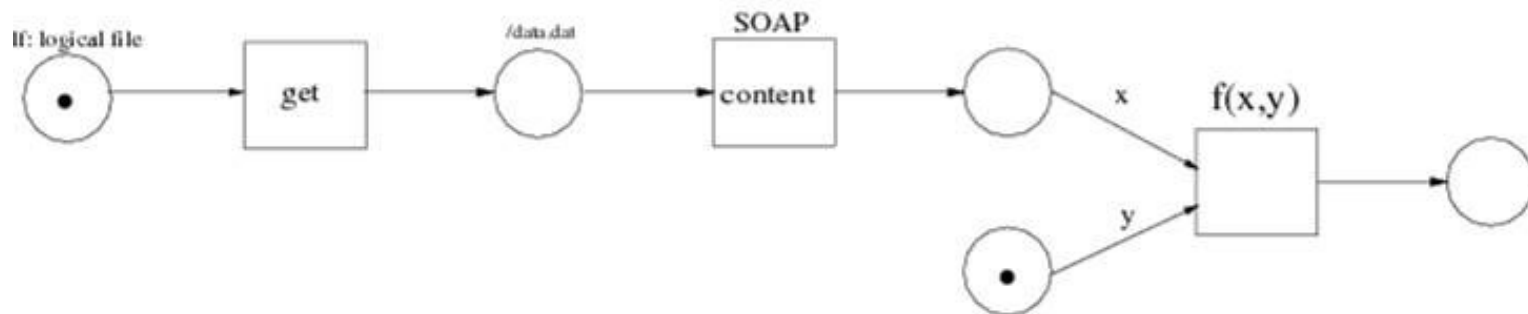
Beispiel für einen Workflow modelliert durch ein Petri-Netz:

- Transport einer Datei lf aus dem Grid-Storage-System zur Ressource r.
- Auf Ressource r wird Job cat ausgeführt.

Datei wird auf Ressource r kopiert (/data.dat)



Dateiinhalt wird in SOAP Nachricht transportiert (x)

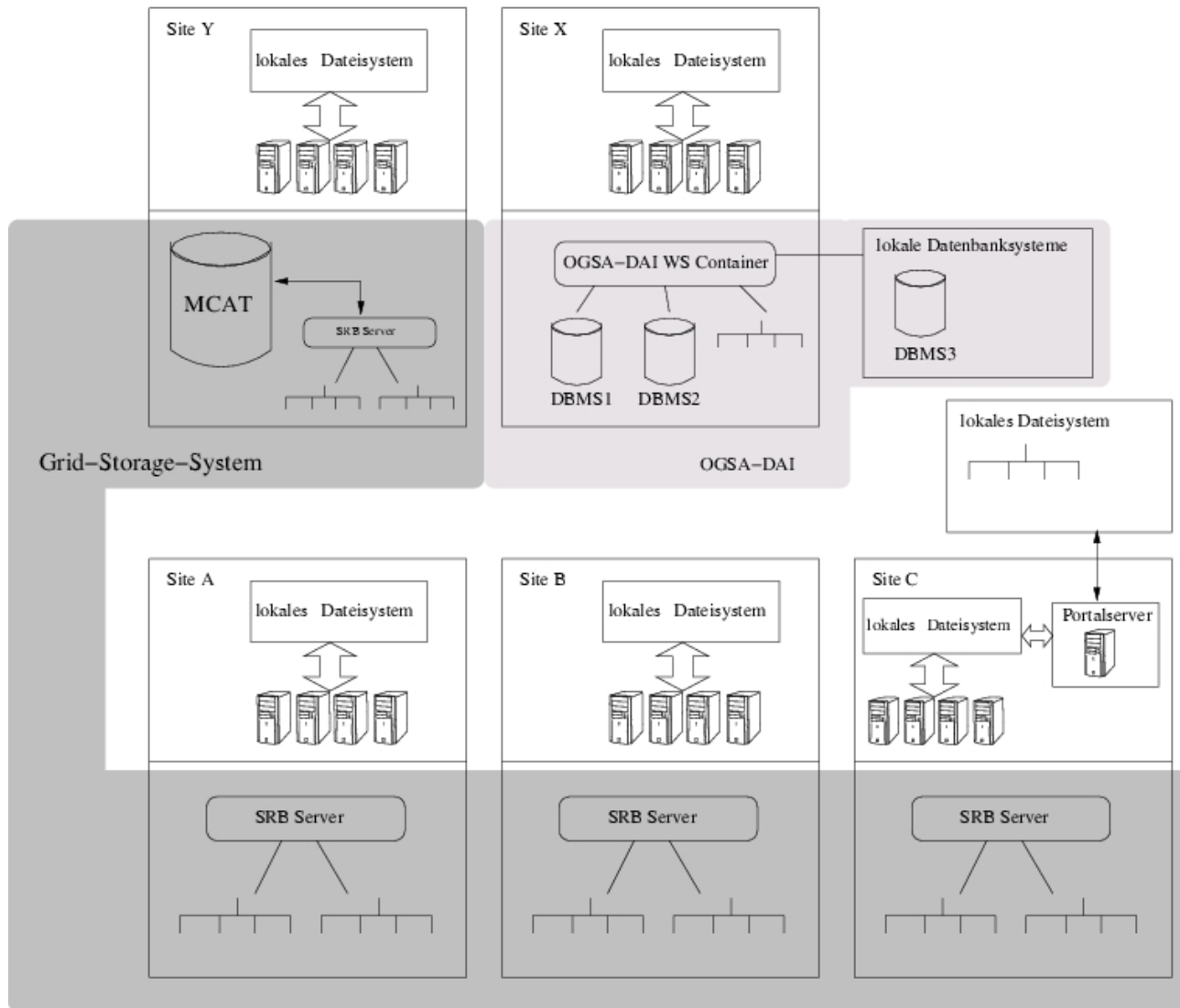


Beispiel für einen Workflow modelliert durch ein Petri-Netz:

- Der Inhalt von Datei lf wird in einer SOAP Nachricht an den Webservice $f(x,y)$ übertragen

- Datenmanagementservices um neue Funktionalitäten erweiterbar
- Datenmanagementservices innerhalb des Grid:
 - Datentransfer mit gridFTP zwischen Grid-Rechnern
 - Datentransfer zwischen Grid-Rechnern und Grid-Storage-System
 - Zugriff auf Datenbanken über WebServices

SRB und OGSA-DAI in MediGRID



- ❑ Wofür soll der SRB genutzt werden?
 - Grid-Storage-System mit globalem Namensraum. Virtualisierung der Speicherung durch Abstraktion vom physikalischen Ort.
 - Langfristige sichere Datenspeicherung, Möglichkeit zur Replizierung über mehrere Speicherressourcen.
 - Metadatenverwaltung mit MCAT (verteilt oder zentral)
 - Zugriffsrechteverwaltung
 - SRB-Clients: SCommands, JARGON, SRBAdmin
- ❑ Noch offen:
 - Eine MediGRID-Zone mit einem MCAT und pro Site ein SRB-Server.
 - Jede Site bildet eine SRB-Zone mit je einem MCAT. Verbindung der Zonen durch Zonenföderation.
 - Stresstest: Dateigrößen von mehreren Gigabyte?
- ❑ Interessante Zusatztools:
 - GridFTP Schnittstelle zum SRB
 - Matrix: Workflow aus SRB-Diensten mit SOAP/WSDL Schnittstelle (in Planung)

- ❑ Wofür soll OGSA-DAI benutzt werden?
 - Bereitstellung von Datenbanken im Grid.
 - Generischer Zugriff unabhängig vom Ort der Datenbank.
 - 1. Realisierung: GeneOntology-Datenbank und Thesaurus im MediGRID-Portal
- ❑ OGSA-DAI WSRF verfügt über Webservice Schnittstellen.
- ❑ OGSA-DAI Webservice Container kann unter Globus betrieben werden.
- ❑ Modellierung des Datenbankzugriffs als Webservice im Workflow.
- ❑ OGSA-DAI-Clients: JavaClient, Administration per GUI oder CLI
- ❑ Noch offen:
 - Funktionalität des Zugriffs auf Dateisysteme mit OGSA-DAI.
Einbindung von Dateiressourcen und Verwaltung der Zugriffsrechte.
 - Neues Release: OGSA-DAI WSRF 2.2 bietet Zugriffsverwaltung für Ressourcen?

Anforderungen

- ❑ SRB:
 - Datenaustausch zwischen Grid-Rechnern und SRB

- ❑ OGSA-DAI:
 - Dateizugriff nur auf strukturierte Daten, Metadaten zu Dateien?
 - Zugriffskontrolle für Dateien: Mapping auf Unix-Account unter Benutzung des Globus-Mapfiles?
 - Wie erfolgt der Datenbankzugriff wenn mehrere OGSA-DAI WS Container im Grid betrieben werden?

- ❑ Möglichkeit zum Upload/Download von lokalen Daten ohne Portal und ohne Installationen von Client-Interfaces auf lokalem Rechner?

Nächste Schritte

- ❑ Speicherung von Daten im Grid-Storage-System:
 - Webservice für SRB put und SRB get bereitstellen.
 - GSI-Authentifizierung, Delegation der Rechte im Portal und Workflow.

- ❑ Aufbau des Grid-Storage-Systems:
 - Installationen von SRB-Servern bei den Ressourcenprovidern.
 - Integration und Bereitstellung von weiteren Datenbanken mit OGSA-DAI.

- ❑ Planung weiterer Datenmanagementfunktionalitäten:
 - Bereitstellung als Webservice.
 - Modellierung im Workflow.

Danke.