



Arbeitspaket 2: Blaupausen und Beratung

Aufbaukonzepte für neue Community-Grids, Aktivitäten und Ergebnisse der Fachberaterteams¹

Deliverable	2.2.2 Aufbaukonzepte Community-Grids
Autoren	Arbeitspaket 2: Blaupausen und Beratung
Editoren	I. Agapov, F. Schlünzen
Datum	10-04-2011
Dokument Version	1.2.0

A: Status des Dokuments

Deliverable 2.2.2, Version 1.2.0, Release

Dieses Dokument ist ein „Living Document“, insofern es Zusammenfassungen von Berichten enthält, die vom Fachberater-Team von WissGrid unternommen wurden. Es gibt daher auch nur ein Improvement des AP2, kein RFC des gesamten Projektes für eine neue Version des Dokumentes.

B: Bezug zum Projektplan

Basierend auf den Blaupausen und Materialien werden Konzepte für den Aufbau neuer Community-Grids entwickelt.

¹ This work is created by the WissGrid project. The project is funded by the German Federal Ministry of Education and Research (BMBF).

C: Abstract

Die Fachberater-Teams sollen den Aufbau neuer Community-Grids unterstützen, und neuen User-Communities den Zugang zu der Grid-Infrastruktur erleichtern. Die Aufbaukonzepte für die neuen Communities werden vorgestellt, sowie die Aktivitäten der Fachberater dokumentiert.

D: Änderungen

Version	Date	Name	Brief summary
0.1.0	14.01.2010	F.Schlünzen	Working Draft Creation
0.1.1	28.09.2010	I. Agapov	Überarbeitung
0.1.2	12.11.2010	T. Rathmann	Überarbeitung
0.1.3	22.11.2010	I. Agapov	Überarbeitung
1.1.0	10.04.2011	F. Viezens	Dokumentation Fachberater-Workshop 01.2011
1.1.1	10.04.2011	F. Schlünzen	Überarbeitung und Aktualisierung
1.2.0	23.05.2011	H. Enke	Publikation/Release

E:

Inhaltsverzeichnis

1	Overview	5
2	Aufbaukonzepte	6
2.1	Compute-Grid	7
2.2	Daten-Grid	7
2.3	Prozesse	8
2.4	Software-Entwicklung, Adaptierung von existierender Software	8
2.5	Testumgebungen	8
2.5.1	Die D-Grid-Referenzinstallation	9
2.5.2	dech VO	9
2.6	Deployment	9
2.7	Support	10
2.7.1	WissGrid Support	10
2.7.2	Community Support	10
3	Aktivitäten und Ergebnisse	11
3.1	Forschung mit Photonen	11
3.1.1	Virtuelle Organisationen	11
3.1.2	Authentifizierung	12
3.1.3	Prototypische Anwendungen	13
3.1.4	Anforderungen an eine virtuelle Datenmanagement-Umgebung	18
3.1.5	Datamanagement VRE	21
3.1.6	Ausblick	22
3.2	Sozialwissenschaften	24
3.3	Biostatistik und Epidemiologie	24
3.3.1	Beschreibung	24
3.3.2	Status der Beratungstätigkeiten	24
4	Workshop Virtuelle Forschungsumgebungen aufbauen - mit D-Grid	25

4.1	Dokumentation des Workshops	25
4.1.1	Communities und die Bildung Virtueller Organisationen	25
4.1.2	Instrument-, Compute-, Data-Sharing: Ressourcen gemeinsam nutzen statt selber aufbauen	26
4.1.3	Gemeinsame Standards, Verfahren und Arbeitsabläufe nutzen	27
4.1.4	Über rechtliche, finanzielle und organisatorische Rahmenbedingungen	27
4.2	Schlussbemerkungen und Tendenzen zum Workshop	28
4.3	Weitere geplante Aktionen	28
	Quellen	28

1 Overview

Große Forschungscommunities wie die Hochenergiephysik, welche die LHC Experimente betreibt, haben meistens Grid-ähnliche Lösungen für das verteilte Datenmanagement und Rechnen entwickelt, weil beides für die Datenanalyse essenziell ist.

In kleinen Forschungscommunities und Communities, in denen traditionell kein großer Bedarf an Rechen- und Speicherkapazität besteht, fehlen solche Erfahrungen in der Regel. Forschungslandschaft und Bedarf an IT-Ressourcen können sich aber schnell ändern.

Erstens, dank moderner Technologien wie z.B. bei Experimenten mit Synchrotronstrahlung, ist die Menge der entstehenden primären Forschungsdaten so groß geworden, dass sie nicht mehr lokal an den Mess-PC's oder Laptops verwaltet werden können, sondern in Datenzentren aufbewahrt werden müssen.

Zweitens, mehr und mehr Forschung wird im Rahmen von internationalen Kollaborationen durchgeführt. Die Forschungsdaten sollen weltweit für mehrere Nutzergruppen zugänglich sein. Dies lässt sich nur innerhalb von Grid-Infrastrukturen bzw. virtuellen Forschungsumgebungen vernünftig lösen.

WissGrid stellt Expertise für Aufbau von neuen Community-Grids bereit und bietet entsprechende Beratung an. In diesem Dokument werden generische Aufbaukonzepte präsentiert, die nicht Community-spezifisch sind und in einem breiten Spektrum von Umgebungen umgesetzt werden können.

Im Kapitel 2 werden typische Anwendungsszenarien und die Fragen erörtert, die vor dem Aufbau der Infrastruktur geklärt werden müssen. Verfahren zum Aufbau von Community-Grids werden diskutiert. Im Kapitel 3 werden Aktivitäten und Ergebnisse des Fachberaterteams anhand der Beispiele Photonenphysik, Klimafolgenforschung, Epidemiologie und Sozialwissenschaften vorgestellt.

2 Aufbaukonzepte

Der Grund, warum lokales Rechnen und Datenmanagement nicht genügen und eine Grid-Lösung erwünscht ist, muss klar verstanden werden. Es kann, z.B., Mangel an Speicherkapazität, Mangel an Rechenkapazität oder entstehende wissenschaftliche Kollaborationen sein. Die Ziele der Umsetzung von Grid-Infrastruktur oder Virtueller Forschungsumgebung müssen definiert werden. Das Ziel kann z.B. weltweiter Zugriff auf die Daten oder schnelleres Rechnen sein. Die folgende Liste kann ein Startpunkt für die Aufstellung eigener Ziele sein:

- **Erforderliche Rechenkapazität:**
Reichen die lokalen Rechenressourcen für die aktuellen Forschungsaufgaben? Wenn nicht, wie viele Ressourcen fehlen (in FLOPS [1], SPEC [2], etc.)?
Wie viel wird die Beschaffung der entsprechenden Ressourcen kosten?
Gibt es Partnerinstitute, die Rechenkapazität teilen können?
- **Software:**
Wird überwiegend Standardsoftware (wie R [3]) oder spezifische Software (wie LHC ATLAS Software) in der Community verwendet?
Wie aufwändig ist es, diese Software zu unterstützen?
Sind Lizenzbestimmungen zu beachten?
- **Datenmengen, Anforderungen an Speicher:**
Wie viele Daten sind zu speichern?
Wie hoch ist die gewünschte Zugriffsgeschwindigkeit?
Welche Latenz ist zu tolerieren?
- **Datenzugriff:** Sollen die Daten direkt von der Anwendungssoftware gelesen werden?
Soll auf die Daten über das Web, bzw. via FTP zugegriffen werden können?
Wie oft werden die Daten gelesen/geschrieben?
Ist der Zugriff konsequent oder random?
- **Datenmanagement:** Sind übliche Metadaten in der Community definiert?
Existieren schon Standards und Policies? Werden community-spezifische Datenformate verwendet?
- **Sicherheitsanforderungen:**
Sind die Anwendungen und Daten sicherheitskritisch?
Welche Schäden können Datenverluste und andere Sicherheitsvorfälle anrichten?
- **Organisatorische Rahmenbedingungen:**
Müssen bestimmte Regelungen (Land, DFG, EU Projekte) in Sachen Sicherheit, Nutzerverwaltung usw. eingehalten werden?
- **Personal und Finanzierung**

Je nachdem, ob verteiltes Datenmanagement oder effiziente Nutzung von Rechenkapazität im Mittelpunkt steht, wird von einem Daten-Grid oder einem Compute-Grid gesprochen.

2.1 Compute-Grid

Beim Compute-Grid befindet sich im Fokus die Möglichkeit, Rechenjobs zwischen verschiedenen Rechenzentren zu verteilen. Dies wird von Grid-Middlewares wie Globus oder gLite unterstützt. Der Hauptanwendungsfall ist die Kooperation mehrerer Rechenzentren zum Zwecke der Lastverteilung. Am meisten profitieren davon jedoch kleinere Forschungsgruppen, die selber keine vollständige Recheninfrastruktur aufbauen können.

Das Grid ist nur dann nützlich, wenn dort Community-spezifische Anwendungssoftware zur Verfügung steht. Allein Rechenkapazität zu teilen reicht nicht. Dann könnten sich kleinere Forschungsgruppen theoretisch auch Rechenzeit bei kommerziellen Anbietern kaufen.

Es ist wichtig, dass im Compute-Grid ein Rechenzentrum die Führungsrolle übernimmt. Ein Beispiel für ein gut geführtes Compute-Grid ist das LHC Grid [4]. Dort leitet das CERN als Tier-0 Zentrum Betrieb und Entwicklung. Große Tier-1 Zentren leisten dabei bedeutende Beiträge. Kleinere universitäre Forschungsgruppen profitieren stark vom Grid, weil sie dort physikalische Analysen ohne zu großen administrativen Aufwand laufen lassen können. Grid-Infrastrukturen, in denen kein Partner eine führende Rolle übernimmt und die Hauptverantwortung für den Betrieb trägt, scheinen viel schwieriger aufzubauen und zu betreiben zu sein.

Compute-Grids können auch eingesetzt werden, um spezifische Hardware, wie z.B. GPU's, zur Verfügung zu stellen. Solche Systeme sind aber noch wenig verbreitet.

Die Wirtschaftlichkeit des Compute-Grid wurde in letzter Zeit durch sinkende Hardware-Preise und entstehende Cloud-Computing- und Virtualisierungstechnologien begrenzt. Wikipedia zufolge [1] betrug der Preis pro GFLOP im Jahr 2000, als das Grid-Konzept etwa entstanden ist, um \$ 1000, im Jahr 2010 aber nur noch etwa \$ 0.10 - \$ 0.50. Es wird erwartet, dass die Anzahl der Cores pro Prozessor weiter steigen und der Bedarf an föderiertem Computing im Nicht-Hochleistungsbereich entsprechend zurückgehen wird.

2.2 Daten-Grid

Ein Daten-Grid ist eine Infrastruktur, in der eine oder mehrere Nutzer-Gruppen ihre Daten verteilt managen und auf die Daten zugreifen können. Dabei kann die Infrastruktur auch zur Datenverarbeitung genutzt werden, diese steht aber nicht im Vordergrund. Beispiele für solche Infrastrukturen sind:

- Klimadatenarchive
- Archive für medizinische Daten
- Archiv für Primärdaten aus Synchrotron-Strahlenquellen

Beim Aufbau solcher Infrastrukturen ist Folgendes zu beachten:

- Bei großen Datenmengen muss geprüft werden, ob Netzwerk-Bandbreite, Storage-Throughput usw. ausreichen. Beispielsweise können die Daten von Experimenten wie LHC [4] oder XFEL [5, 6] nur in einer spezifischen Grid-Topologie und Netzwerk-Infrastruktur gespeichert werden.

- Sollen die Daten in lokalen Rechenzentren repliziert werden? Bei solchen föderierten Lösungen ist es wichtig, dass die lokalen Sites eine ähnliche Infrastruktur besitzen (Metadaten-Datenbanken, Zugriffsprotokolle, AAI usw. unterstützen).
- Der Zugriff soll benutzerfreundlich sein, z.B. über ein Web-Portal.

2.3 Prozesse

In vielen Communities sind die IT-Prozesse nicht genügend formalisiert und standardisiert, um sie auf einer gemeinsamen Plattform wie einem Grid ausführen zu können. Dies muss geschehen, bevor von einer Grid-Lösung gesprochen werden kann. Solche Prozesse sind

- Nutzer-Management
- Software-Management
- Sicherheits-Policies

Die genannten Prozesse sind nicht nur von der Funktionalität her Voraussetzung für den Grid-Betrieb, sondern bergen auch viele Risiken – der Aufwand, manche IT-Prozesse zu betreiben, wird oft nicht rechtzeitig erkannt bzw. unterschätzt. Grid-enabled Workflow Engines wie GWES [7] können dabei helfen, Policies und Workflows technisch in Grid umzusetzen.

2.4 Software-Entwicklung, Adaptierung von existierender Software

Es kommt sehr selten vor, dass kein Bedarf an Software-Entwicklung besteht.

- Community-Grids können sich in existierenden Software-Stacks bedienen (z.B. iRODS [9], Fedora [10], Globus [11]). Es besteht aber fast immer Bedarf an Nutzer-Interface-Entwicklung (wie z.B. Web-Portale).
- Für spezifische Anwendungsfälle kann es einen großen Bedarf an Entwicklungen geben.

Bei geringem Bedarf an Entwicklung (z.B., wenn nur ein einfaches Web-Interface zum Daten-Archiv und einfache Workflows benötigt werden), können existierende Lösungen relativ mühelos adaptiert werden. Hier können neue Communities Hilfe beispielsweise bei WissGrid-Experten suchen. Die Entwicklung komplexer Software soll aber innerhalb der Communities geschehen.

2.5 Testumgebungen

Um die Funktionalität von Grid-Infrastrukturen besser zu verstehen und den Entwicklungs- und Betriebsaufwand schätzen zu können, brauchen neue Communities eventuell eine geeignete Testumgebung. WissGrid selbst hat keine Testumgebung, es können aber unter bestimmten Voraussetzungen die Referenzinstallation des D-Grid oder die dech VO des DECH-Verbundes für Tests genutzt werden. Eventuell kann die neue Community auch das Testgrid der betreuenden Community nutzen. Um Zugang zu solchen Testumgebungen zu bekommen, sollten sich neue Communities an das Fachberater-Team wenden.

2.5.1 Die D-Grid-Referenzinstallation

Das D-Grid ist ein Zusammenschluss einer Vielzahl von Sites. Die Referenzinstallation bildet ein Site mit allen unterstützten Software-Komponenten in Form eines Prototypen (D-Grid Standard-Site) ab.

Den Prototypen können Nutzer und Entwickler von D-Grid-Projekten zum Testen nutzen. Hierzu ist lediglich eine D-Grid-Registrierung erforderlich, weil dieselbe Gridmap-Datei wie im D-Grid verwendet wird. Die virtuellen Organisationen des D-Grid sind auch im Prototypen vorhanden. Außerdem dient die D-Grid-Referenzinstallation dem D-Grid zur Reproduktion von Fehlern in einem definierten Umfeld und als Testplattform für neue offiziell bereitgestellte Software-Versionen.

Die D-Grid-Referenzinstallation besteht nämlich nicht nur in dem Prototypen, sondern bietet auch Software an. Unter der Webadresse

<http://dgiref.d-grid.de>

werden Download und Dokumentation der unterstützten Software angeboten:

- Die drei Grid-Middlewares des D-Grid : Globus [11], gLite [12] und UNICORE [13]
- dCache [14] und OGSA-DAI [15] als Datengrid-Software
- Die Betriebssysteme Scientific Linux und OpenSUSE Linux
- Cluster-Software

Die Referenzinstallation wird auch künftig zur Verfügung stehen. Sie wird ab Ende 2010 nicht mehr wie zuvor zweimal jährlich, sondern nur noch bei Bedarf (Security-Fixes) erneuert. Bevor ein Release zur neuen Referenzinstallation erklärt wird, wird dieses auf einer von der Referenzinstallation getrennten Plattform getestet. Beta-Releases dienen nur noch zum Testen von Security-Fixes, so dass das stabile Release im Wesentlichen so bleibt, wie es ist.

2.5.2 dech VO

Die DECH-Region besteht aus den Staaten Deutschland und der Schweiz, die sich im Rahmen der internationalen Großprojekte EGEE und EGI zu einem Regionalverband zusammengeschlossen haben und eng zusammenarbeiten. Die dech VO ist ein Testbed für deutsche und schweizer ROC (Regional Operation Center) Mitglieder. Eine Community, die hier testen will, sollte international ausgerichtet sein. Die dech VO eignet sich auch, um erste Erfahrungen im Grid-Umfeld zu sammeln, und ist vergleichsweise gut dokumentiert mit einfachen [Tutorials](#) und [Support-Materialien](#).

Die schweizer Seite von DECH bietet darüber hinaus an, Federated IDs (Shibboleth) zu verwenden, um daraus Short-Lived User Certificates (SLCS) zu generieren. Diese Zertifikate haben eine Lebensdauer von 11 Tagen und ersparen die Generierung von persönlichen Grid-Zertifikaten und die damit verbundenen Probleme.

2.6 Deployment

Die produktive Umsetzung von Grid-Infrastrukturen soll außerhalb von WissGrid geschehen.

2.7 Support

Der Support durch das WissGrid-Fachberatererteam ist eingebettet in eine vernetzte Support-Struktur (für mehr Infos siehe auch WissGrid Deliverable 2.1.8 [16] „Blaupausen für den Aufbau einer technischen Support-Infrastruktur“).

2.7.1 WissGrid Support

WissGrid stellt den neuen Communities eine zentrale Anlaufstelle für Support-Anfragen bereit. Primär können Anfragen an die E-Mail-Adresse fachberater@wissgrid.de geschickt werden. Abonnenten dieser Mailing-Liste sind die Mitglieder des Fachberatererteams.

WissGrid ist als Support-Unit in den D-Grid User Support (DGUS) eingebunden und somit auch über das Helpdesk des D-Grid erreichbar. Tickets können im NGI-DE-Portal [17] geöffnet werden. Beim Aufruf des NGI-DE-Portals wird ein persönliches Zertifikat abgefragt. Wer noch kein Zertifikat hat, kann das Portal aber trotzdem nutzen. Für den Zugang zum NGI-DE-Portal ohne Zertifikat benötigt man eine Nutzerkennung und ein Passwort. Um diese zu bekommen, muss ein Webformular ausgefüllt werden.

Bei rein nationalen Community-Grids werden die Anfragen in der Regel innerhalb von WissGrid, DGUS oder den Community-eigenen Support-Strukturen behandelt. Einige Communities wie zum Beispiel die Photonenphysik sind hingegen grundsätzlich international ausgerichtet. Instrumente und Services der verschiedenen Einrichtungen wie zum Beispiel EMBL [18] werden von Wissenschaftlern weltweit genutzt. Ihre Support-Anfragen können daher nicht immer im (nationalen) Rahmen von WissGrid oder DGUS abgehandelt werden.

2.7.2 Community Support

Die Community-Grids sollen mittelfristig ihre eigenen Support-Strukturen aufbauen, mit Unterstützung von WissGrid. Der Community-spezifische Support ist essenziell, denn Applikations- oder Experiment-spezifische Probleme können nur innerhalb der Community selbst gelöst werden. Insbesondere wird auf diese Weise auch gewährleistet, dass Anfragen an entsprechende Bulletin-Boards und Ticketing-Systeme zielgerichtet weitergeleitet werden können.

3 Aktivitäten und Ergebnisse

3.1 Forschung mit Photonen

Die Photonenphysik-Communities haben bislang kaum Erfahrungen mit dem Grid. Grundsätzlich neigen diese Communities dazu, lokale Infrastrukturen zu nutzen, so unzugänglich diese auch sein mögen. Durch den Aufbau neuer Lichtquellen und Entwicklungen in der Detektor-Technologie werden die alten Probleme – wie Mangel an Speicherkapazität und ineffizientes Daten-Management – bereits jetzt dramatisch verschärft.

Bislang lag das Daten-Management weitgehend in der Hand der Anwender oder Experimentatoren. Nach erfolgtem Experiment werden die Daten entweder über das Internet in eines oder mehrere Heimatinstitute transferiert oder auf portable Festplatten kopiert und verschickt. Lokale Kopien werden nach einer nicht genau definierten Zeit gelöscht. Es existieren mit anderen Worten keine Strukturen um die langfristige Archivierung von Daten auch nur ansatzweise zu gewährleisten. Es gibt keine konsistente Erfassung von Metadaten und keine ausreichende Computing-Infrastruktur, um eine zeitnahe Analyse zu garantieren.

In den kommenden Monaten und Jahren werden die Datenraten dramatisch ansteigen und können ein Petabyte pro Woche überschreiten. Um mit der Entwicklung Schritt zu halten und eine langfristige Verfügbarkeit der Daten auch im Sinne einer Nachnutzung zu ermöglichen, haben das EMBL [18] und das European XFEL Consortium [5] begonnen, Möglichkeiten zu eruieren. Eine Grid-basierte Lösung bietet viele Vorzüge. Neben einer konsistenten und erprobten Daten-Managementstruktur erlaubt sie die nahtlose Integration von Computing-Infrastrukturen, die die Möglichkeiten für die beteiligten Wissenschaftler deutlich verbessert. Daher haben sowohl EMBL als auch der European XFEL starkes Interesse an der Beratung und Unterstützung durch WissGrid.

Im Folgenden sollen einige der grundsätzlichen Überlegungen skizziert und Aktivitäten der Fachberater dokumentiert werden.

3.1.1 Virtuelle Organisationen

EMBL hatte einige Anwender, die reine Bioinformatik-Anwendungen auf dem Grid realisieren wollten. Im Wesentlichen handelt es sich um Sequence-Alignment, 2D- und 3D-Strukturvorhersagen und Generierung von Struktur-Templates biologischer Makromoleküle. Innerhalb der biomed-VO existiert bereits eine Software-Umgebung, die praktisch alle Voraussetzungen mitbringt, um solche Anwendungen mit minimalem Aufwand im Grid zu realisieren. Daher haben wir in diesem Fall die biomed-VO, die von DESY auch unterstützt wird, als Einstiegspunkt empfohlen. Grid-Zertifikate wurden über DESY-RA erstellt. EMBL hat mittlerweile diese Applikationen im Grid in Produktion genommen.

Generell stellt sich aber die Frage, wie die VOs zu organisieren sind. Es würde sich zum Beispiel anbieten, VOs spezifisch für wissenschaftliche Disziplinen einzurichten. EMBL zum Beispiel bietet Experimente im Umfeld von Kleinwinkelstreuung (SAXS - Small Angle X-ray Scattering) und Protein-Kristallographie (MX - Macromolecular X-ray crystallography) an. Diese Experimente erfordern sehr unterschiedliche Analyse-Umgebungen. Andererseits finden sich zunehmend Anwendungen, die beide Disziplinen miteinander kombinieren, die also von einem einheitlichen Framework profitieren würden. European XFEL hingegen ist primär daran interessiert, das Datenmanagement im Grid zu

realisieren. In diesem Fall hätte eine Facility-spezifische VO ihre Vorteile.

Um den Einstieg zu erleichtern, bieten wir verschiedene Test-Umgebungen an, die zum Teil spezifisch für die Photonophysik eingerichtet wurden:

- **DECH VO:** Ist ein Testbed für deutsche und schweizer ROC (Regional Operation Center) Mitglieder. Es eignet sich um erste Erfahrungen im Grid-Umfeld zu sammeln.
- **DESY VO:** Als lokale VO bietet DESY diese VO generell als Testumgebung an. Als Betreiber der VO kann DESY insbesondere in den ersten Phasen sehr zeitnahen Support bei Problemen anbieten, und die lokalen Computing- und Storage-Elemente sind unmittelbar ins lokale Netzwerk eingebunden. Dies hilft, so mancher anfänglicher Enttäuschung vorzubeugen. Insbesondere für die Photonophysik hat DESY ebenfalls eine AMGA-Instanz [19] mit Oracle als Backend implementiert. Dies dient als Testumgebung für FLASH und später XFEL, um die automatische Registrierung und Annotierung der experimentellen Daten zu testen. Angesichts der großen Datenmengen und Anzahl der Dateien, die pro Sekunde generiert werden, sind Performance-Tests hilfreich, um sich langfristig für eine Lösung zu entscheiden.
- **XFEL VO:** Für die Belange der European XFEL wurde eine spezifische VO eingerichtet. Bislang gab es wenig Aktivitäten in diesem Umfeld.
- **PS VO:** VOs definieren nicht nur in gewissem Maße den Einstiegspunkt, sondern bieten auch eine Sichtbarkeit nach außen. DESY wird daher eine generische VO für die Photonophysik einrichten und betreiben. Die Details müssen allerdings noch geklärt werden.

3.1.2 Authentifizierung

In der Regel erfolgt die Authentifizierung im Grid über persönliche Grid-Zertifikate. Authorisierungen können z.B. über GridMaps oder rollenbasiert vergeben werden. Alternative Authentifizierungs-Instanzen mit OpenID oder Shibboleth/SAML sind bislang nicht sehr verbreitet im Grid. Die Akzeptanz für Grid-Zertifikate ist allerdings in einigen Communities nicht sehr ausgeprägt. Im TextGrid wird eine Shibboleth/SAML basierte Authentifizierung bevorzugt. C3Grid vermeidet die Verwendung von Zertifikaten und setzt im Rahmen Earth System Grids (ESG) auf eine OpenID basierte Lösung [21]. In einigen Diskussionen mit vielen verschiedenen Vertretern der Photonophysik zeigen sich zwei konträre Standpunkte. Die einen vertreten die Auffassung, dass kein Weg um Zertifikate herumführt und dies generell eine sinnvolle Angelegenheit ist. Die anderen finden die Notwendigkeit Grid-Zertifikate zu verwenden komplett inakzeptabel.

Die Lichtquellen am DESY, FLASH, European XFEL, Petra3 und DORIS, werden rund 5000 Anwender pro Jahr zu bedienen haben, von denen rund 80 % jedes Jahr neu hinzukommen. Es wären also jedes Jahr rund 4000 Grid Zertifikate zu generieren, was mit erheblichem Aufwand verbunden wäre. Zur Zeit diskutierte Alternativen basieren auf Shibboleth und OpenID, die beide ihre Vor- und Nachteile mit sich bringen. OpenID ist einfach zu implementieren und cross-Site Trusts einfach zu realisieren. OpenID hat allerdings kaum Mechanismen, eine einmal aufgebaute Authentifizierungskette wieder abzubauen. Shibboleth bietet Federated IDs. Das ist aber in einem internationalen Umfeld nicht einfach zu realisieren, da es kaum Mechanismen gibt, Trusts über nationale Grenzen hinweg aufzubauen.

ESRF, PSI und das Helmholtzzentrum Berlin versuchen zur Zeit einen Shibboleth-basierten Prototyp zu implementieren, der die nötigen Voraussetzungen erfüllt. DESY verfolgt diese Bemühungen und ist durch ein pan-Europäisches Projekt namens [PaN-data](#) direkt involviert.

3.1.3 Prototypische Anwendungen

Wie erwähnt, gibt es eine Vielzahl sehr unterschiedlicher Anwendungen, die für die Photonophysik von Interesse sind. Nicht alle diese Anwendungen lassen sich sinnvollerweise im Grid realisieren. Eine der Aufgaben ist es also, Anwendungen zu identifizieren, die sich gut eignen (das sind nicht-grafische nicht-Echtzeit-Anwendungen, die relativ viel Rechenkapazität brauchen und keine Lizenzbeschränkungen haben). DESY hat mit EMBL und European XFEL einige Applikationen ausgewählt, die als prototypische Anwendungen implementiert werden sollen. Diese können dann wiederum als Demonstrations-Projekte und zur Dokumentation dienen. Im Rahmen der Gespräche mit Vertretern der Photonophysik haben sich, neben dem Datenmanagement, einige Anwendungen herauskristalliert, die im weiteren beschrieben werden sollen.

Generischer Prototyp Anhand der Materialien-Sammlung (D2.1.1 [20]) und den daraus extrahierten Architektur-Konzepten wurde ein generisches Konzept für die Photon Sciences erstellt (siehe Abb. 1). In Gesprächen mit den beteiligten Partner aus den Photon Science Communities hat sich gezeigt, dass diese generische Architektur als Test-Plattform den zur Zeit formulierten Anforderungen genügen kann. Erfahrungen mit dem Testbed in Kooperation mit ESRF und PSI haben gezeigt, dass nur die konkrete Implementierung einer Grid-Infrastruktur und der Anwendungen Aufschluss über das optimale Modell geben wird.

Gespräche mit EMBL legen nahe, dass mittelfristig eine Erweiterung der Architektur um einige Komponenten wünschenswert werden könnten (z.B. EDNA Workflow Engine). EMBL strebt an, die verschiedenen Instrumente und Prozesse in eine einheitliche, integrierte Umgebung einzubetten, so dass Informationen und Daten aus so verschiedenen Bereichen wie Nasslabor, Kristallisation oder Kleinwinkelstreuung verarbeitet werden können, wobei insbesondere die Remote-Steuerung von Instrumenten zunehmend an Bedeutung gewinnt. Dies wird für den Aufbau des Prototypen noch nicht unmittelbar von Interesse sein, allerdings sollte die Infrastruktur hinreichend modular aufgebaut sein, so dass eine spätere Integration physikalischer Instrumente und von Schnittstellen zu den wichtigsten Datenbanken problemlos möglich ist.

European XFEL Datamanagement Der European XFEL befindet sich zur Zeit in der Konstruktionsphase, die IT-Infrastruktur wurde jedoch noch nicht komplett definiert. Die groben Rahmenbedingungen wurden in einem Computing Technical Design Report (C-TDR) dokumentiert. Darin sind die Anforderungen an das Datenmanagement basierend auf Erfahrungen von FLASH und LCLS dargelegt. Die Computing-Anforderungen hängen stark von den avisierten Experimenten und den Analyse-Methoden ab. Diese lassen sich zur Zeit nicht abschließend abschätzen, so dass der Fokus auf dem Datenmanagement und der Archivierung liegt. Der European XFEL beabsichtigt, das Datenmanagement in eine Grid-Infrastruktur einzubetten und dCache als Storage-Backend zu nutzen. Authentifizierung und Authorisierung sollen zumindest anfänglich auf X.509 PKI basieren.

Da der XFEL noch im Aufbau begriffen ist, fallen zur Zeit keine experimentellen Daten an. Allerdings werden Daten bereits am LCLS und insbesondere in Simulationen des Strahlungsfeldes im XFEL erzeugt. Diese Simulationsdaten sind Grundlage für weitergehende Analysen, die den Bau des XFEL beeinflussen. Die Simulationsdaten sowie alle Metadaten sollen daher als Public Domain bereitgestellt werden. Sekundäre Simulationen müssen mit den Original-Daten verknüpft und annotiert werden.

Die Simulationsdaten liegen in Dateien von 50–250 GB Größe pro Datei vor. Das Dateiformat basiert

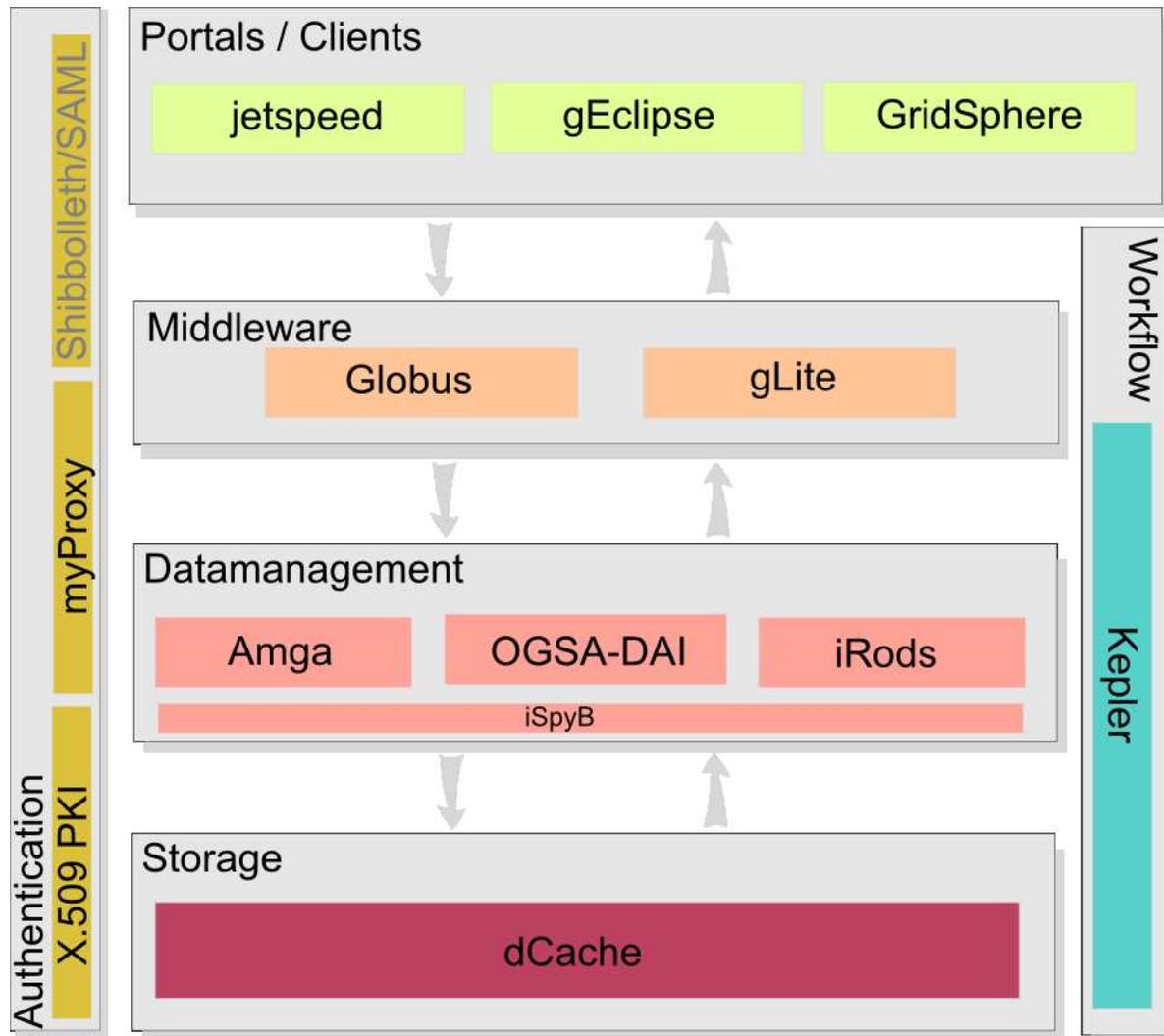


Abbildung 1: Generische Architektur-Konzept für Photon Sciences

auf Fortran Stream I/O, die dazugehörigen Metadaten sind in den Header-Records der Dateien selbst abgelegt. Die Dateigrößen und Formate erschweren die Handhabung erheblich. Daher sollen die Daten in HDF5 [22] umgewandelt werden, die Metadaten in einer Metadatenbank abgelegt werden. Die Daten sollen über ein Portal lokalisierbar und zugreifbar gemacht werden.

Da XFEL erwägt, das Data Acquisition (DAQ) System auf dem Glassfish Application Server aufzubauen. Java-basierte Portal-Lösungen (wie z.B. Liferay [25]) erlauben eine einfache Integration von Datenmanagement und DAQ. Ebenso werden zunächst keine Workflow-Frameworks benötigt. Daher lässt sich das Architektur-Konzept in diesem Falle auf einige wenige Komponenten reduzieren, die sich aber vollständig aus dem generischen Photon Science Konzept ableiten lassen (Abb. 2).

EMBL Daten- und Compute-Grid Im Rahmen des EU-Projektes ESRFUp [26] wurde untersucht, inwieweit das Grid sich für Photon Science spezifische Anwendungen nutzen ließe. Diese Evaluierung fiel durchweg negativ aus, was im wesentlichen auf vier Gründe zurückzuführen ist:

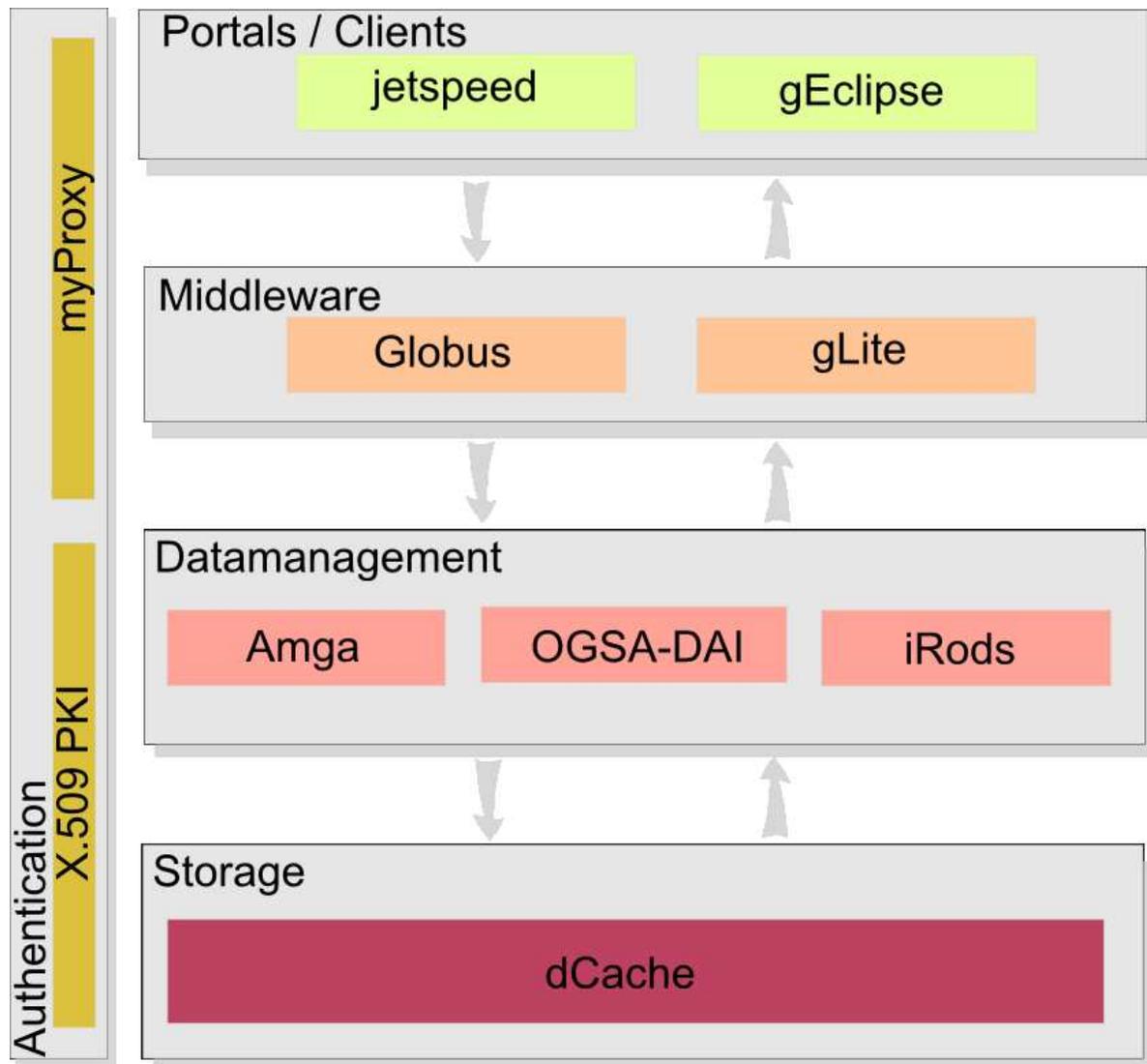


Abbildung 2: Reduziertes Architektur-Konzept für European XFEL

- **Datenmanagement:** Datentransfer-Tests wurden auf den sequentiellen Transfer vieler sehr kleiner Dateien beschränkt. Es ist nicht überraschend, dass Latenzen sich in diesem Fall negativ bemerkbar machen.
- **Computing:** natürlich profitiert nicht jede Anwendung automatisch von der Migration ins Grid. Eine Evaluierung der Grid-Tauglichkeit von Anwendungen wurde aber unterlassen.
- **Grid-Architektur:** die ESRFUp Grid-Infrastruktur hat sich ausschließlich auf EGEE und gLite als Middleware beschränkt, ohne deren Konzepte oder mögliche Alternativen auf ihre Eignung für das spezifische Anwendungsfeld zu prüfen.
- **Community:** ESRFUp hatte die individuellen Wissenschaftler, die die Experimente ausführen, als Zielgruppe gewählt. Aufgrund der Heterogenität dieser Communities und deren bisweilen recht geringen IT-Grundkenntnisse war diese keine gute Wahl, da die angebotene Grid-Infrastruktur die Anwender schlicht überfordert und frustriert zurückgelassen hat.

Die Materialiensammlung und die daraus konsensierten Blueprints und Konzepte erlauben diese Fehler zu vermeiden. Konsequenterweise wurden für die Photon Science Anwendungen und insbesondere für das EMBL ganz andere Konzepte vorgeschlagen:

- **Community:** wähle die Beamline- und Facility-Betreiber als Zielgruppe, die eine angepasste Grid-Infrastruktur als Service anbieten. Dies reduziert die Komplexität erheblich und vereinfacht die Definition der benötigten Standards.
- **Grid-Architektur:** verwende Architektur-Konzepte, die auf die Bedürfnisse der Anwender zugeschnitten ist.
- **Computing:** eine sorgfältige Auswahl der Anwendungen mindert die Frustration der Service-Anbieter wie der Anwender.
- **Datenmanagement:** die Verwendung von standardisierten Container-Formaten und die Integration von existierenden Metadaten-Engines hilft Throughput zu erhöhen, und erlaubt die Definition von standardisierten Workflows. De facto Standards für Datenformate in den Photon Sciences sind HDF5 basiert (wie HDF5 oder Nexus) oder CIF (wie imageCIF) basiert.
- **Kristallisation:** EMBL betreibt eine Kristallisations-Anlage, die prinzipiell als Service allen Anwendern weltweit zur Verfügung steht. Diese Kristallisations-Anlage ist in der Lage, tausende von Kristallisations-Ansätzen parallel zu verarbeiten. Insgesamt werden einige Millionen dieser Versuche parallel durchgeführt. Von jedem Versuch werden täglich typischerweise drei fotografische Aufnahmen gemacht, die die Grundlage für die Beurteilung der Versuche sind. Täglich sind also einige Millionen dieser Aufnahmen von weltweit verteilten Gruppen zu analysieren.

EMBL möchte zwei Aspekte dieser Analyse im Grid-Umfeld realisieren. Zum einen ist dies die automatische Zustellung der Bilddaten zu den einzelnen Gruppen, zum anderen ist die automatisierte Analyse der Daten ein Anliegen.

Die Zustellung der Bilddaten ist technisch unproblematisch. Für die Analyse existieren zur Zeit noch keine entsprechende Algorithmen. Dies hat den Vorteil, dass diese Algorithmen von vornherein mit Blick auf die Porabilität und Grid-Eignung implementiert werden können.



Abbildung 3: Kristallisations-Roboter am EMBL

- **AutoRickshaw:** EMBL-Hamburg hat eine Software-Suite zur automatischen Kristallstrukturbestimmung entwickelt. Diese Suite namens [Auto-Rickshaw](#) wird als Service für die Anwender angeboten und kann remote über einen Web-Server gestartet werden [28]. Die Analyse läuft zur Zeit auf einem lokalen 64-core Cluster. Die Verfügbarkeit des Clusters und die Performance sind nicht mehr ausreichend. Eine Erweiterung des Clusters wäre eine Option, andererseits ist diese Anwendung sehr gut ins Grid portierbar und soll daher prototypisch implementiert werden.

Die Analyse basiert auf einem iterativen Prozess, der unter anderem viele Struktur-Templates parallel zur Strukturbestimmung an die experimentellen Daten anpasst. Dieser Prozess lässt

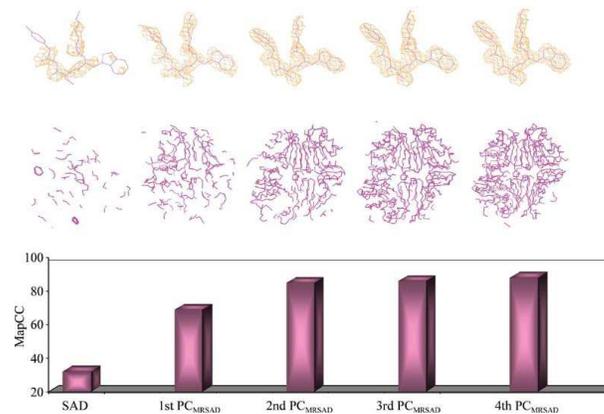


Abbildung 4: Illustration der iterativen Strukturbestimmung in Auto-Rickshaw [29]

sich relativ einfach in eine Vielzahl individueller Einzelprozesse zerlegen, die sich in einen wohl definierten Workflow abbilden lassen.

Diese Klasse von Anwendungen erfordert die Einbindung einer Vielzahl von Datenbanken und Informationsservices (z.B. PDB [31], NDB [32], EMDB [33], EBI [34] oder NCBI [35]). Für einige dieser Datenbanken wie z.B. PDB gibt es eine Java/WSDL-API, die eine Integration in ein Portal wie Liferay vereinfacht, für andere, wie NDB, müsste die Möglichkeiten einer Integration evaluiert werden.

3.1.4 Anforderungen an eine virtuelle Datenmanagement-Umgebung

Die Photon Science Communities sind wie mehrfach erwähnt heterogen, divers und volatil, was die Erstellung eines Anforderungsprofils kompliziert. Letztlich mussten die spezifischen Anforderungen und Rahmenbedingungen für jedes Instrument individuell ermittelt werden, und zumindest partiell in eine generische Datenmanagement-Lösung implementiert werden. Glücklicherweise sind wesentliche Aspekte Instrument- oder Facility-übergreifend. Diese Aspekte, das daraus resultierende generische Layout einer DM-Umgebung sowie die spezifische Implementierung für am DESY beheimatete Einrichtungen (Eur. XFEL, EMBL, CFEL und HASYLAB) sollen kurz aktualisiert dargestellt werden.

Policies Die gängige Praxis in den Photon Science Communities im Umgang mit wissenschaftlichen Daten überlässt alle Rechte und Pflichten den Anwendern, die diese Daten erzeugt haben. Daten werden auf Wechselmedien repliziert und anschließend gelöscht. Dieses Verfahren versetzt Daten in einen Zustand, die eine Referenzierung oder Nachnutzung ausschließen. Ein Datenmanagement-System ist unter diesen Bedingungen natürlich ebenso sinnlos wie überflüssig. Diese Praxis widerspricht zudem ganz grundsätzlich den Richtlinien von DFG [40] und OECD [39] sowie der Berliner Erklärung [38], die unter anderem von der Helmholtz-Gemeinschaft unterzeichnet wurde.

Unabdingbare Voraussetzung für die Implementierung einer DM-Infrastruktur ist die klare Definition verlässlicher Regeln, wie mit wissenschaftlichen Daten verfahren werden soll. Die einseitige Einführung einer solchen Daten-Policy kann jedoch leicht als gravierender Nachteil für eine Forschungseinrichtung verstanden werden, wenn konkurrierende Labore von der Einführung einer Daten-Policy absehen. Daher haben sich die meisten der europäischen Grossforschungseinrichtungen, die Photonen- oder Neutronenquellen betreiben, zusammengetan um eine gemeinsame Daten-Policy zu entwickeln und zu implementieren. Die daraus resultierende Daten-Policy (siehe [PaNdata Policy](#)) versucht zum einen den *best-practice* Empfehlungen von DFG und OECD Genüge zu tun, zum anderen für die Anwender akzeptable Rahmenbedingungen zu schaffen, wie in Abbildung 5 skizziert.

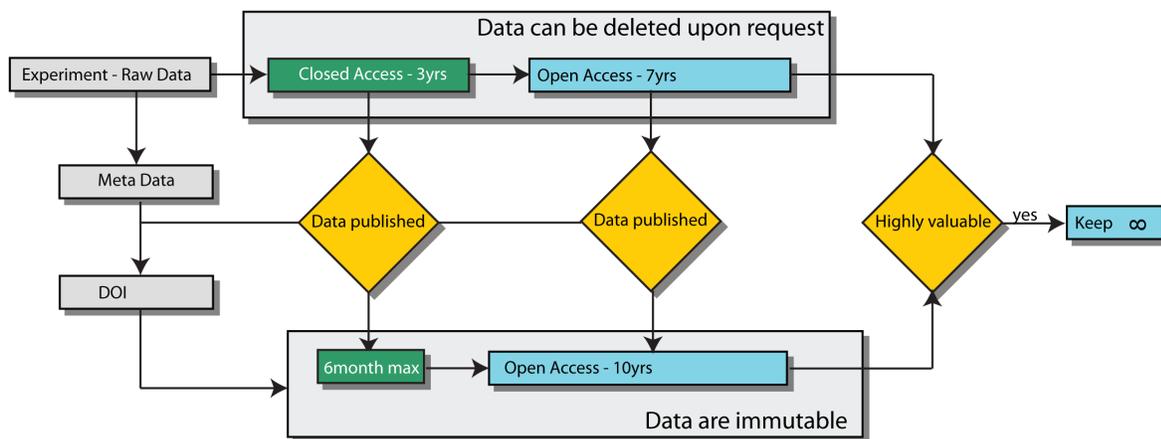


Abbildung 5: Vorschlag für den *Lifecycle* wissenschaftlicher Daten

Demnach werden alle experimentellen Daten zunächst einmal für mindestens 10 Jahre archiviert, wobei durchaus Ausnahmen zulässig sind, da Experimente nicht immer erfolgreich durchgeführt

werden können und Daten produzieren, die weder von wissenschaftlichem Wert noch nutzbar sind. Diese Daten können auf Anforderung jederzeit gelöscht werden.

Alle archivierten Daten werden spätestens 3 Jahre nach ihrer Erzeugung freigegeben (*Open Access*). Daten, die in einer Publikation referenziert werden, werden spätestens 6 Monate nach der Publikation freigegeben und verbleiben für mindestens 10 Jahre nach der Publikation im Archiv.

Bislang hat keine Einrichtung diese Daten-Policy verabschiedet. Dies ist zum einen den damit verbundenen Kosten geschuldet, die von keiner Seite finanziert werden. Zum anderen fehlt an den meisten Einrichtungen eine geeignete Infrastruktur, um die Daten und assoziierte Information nachhaltig zu archivieren.

Datenformate An den verschiedenen Synchrotrons und FEL-Quellen sind eine Vielzahl unterschiedlicher Detektoren bei unterschiedlichsten Experimenten im Einsatz. Die Detektoren schreiben in der Regel Dateien in proprietären Formaten, und die Anzahl der verwendeten Formate und Kompressions-Algorithmen ist unüberschaubar. Zudem können individuelle Datensätze aus hunderttausenden kleiner Dateien bestehen. Dies hat gravierende Nachteile sowohl für den Experimentator wie für den Archivar.

Um diesen Nachteil zumindest partiell zu mildern, gibt es seit vielen Jahren Bemühungen die verwendeten Datenformate zu standardisieren. In vielen Diskussionen mit Nutzern wie Betreibern lassen sich einige wesentliche Merkmale eines Standard-Datenformates benennen:

- **Performance:** Das Datenformat bzw. die Schnittstelle dazu muss gewisse Mindestanforderungen erfüllen. Aktuelle Challenges kommen unter anderem aus der Kristallographie, die 100 Bilder oder mehr als ein GByte pro Sekunde generieren. Andere Techniken wie Xray Photon Correlation Spectroscopy (XPCS) definieren ähnliche Anforderungen. Experimente an Freielektronen-Lasern wie LCLS oder in Zukunft der Eur.XFEL werden ein bis zwei Größenordnungen darüber liegen. Insofern muss ein Standard-Datenformat Potential für deutliche Beschleunigung mitbringen.
- **Support:** Da die Daten für einen langen Zeitraum archiviert werden sollen, ist ein langfristiger Support und die (kompatible) Weiterentwicklung essentiell.
- **Archivierbarkeit:** Zur Zeit generieren die Experimente eine extrem große Zahl ($\gg 10^9$ /Jahr) sehr kleiner Dateien. Für die Analyse der Daten ist das ungünstig, da die Performanz der Dateisysteme signifikant darunter leidet. Für die Archivierung der Daten ist dies eine Katastrophe, da Tape-Backends weniger durch den Datenfluss als vielmehr durch die Zahl der Dateien limitiert werden. Der Unterschied, eine 100 MB- oder eine 1 kB-Datei von Band zu extrahieren, ist vernachlässigbar. Zudem ist die Schreibdichte für viele kleine Dateien sehr viel geringer, was die Kosten für die Archivierung erhöht. Daher muss ein Dateiformat zum einen die individuellen Dateien zu Datensätzen aggregieren können. Zum anderen muss das Dateiformat hinreichend flexibel sein, um sehr große Datensätze in verdauliche Container separieren können ohne die Integrität eines Datensatzes zu verletzen.
- **Konsistenz:** Die Berechnung von Checksums ist an sich ein triviales Problem. Für sehr große Datensätze oder Dateien werden übliche Checksums zunehmend unzureichend und langsam. Datenformate sollten daher über Mechanismen verfügen, Checksums für einzelne Dateneinträge in einem Container zu unterstützen.

- **Portabilität:** Anwender aus den Photon Science Communities verwenden ein breites Spektrum verschiedener Unix- und Windows-Flavors. Nativer Support für diese Betriebssysteme sollte gegeben sein.
- **Schnittstellen:** Die Unterstützung der gängigen Programmiersprachen sollte gegeben sein. Native Bindings für C/C++, Fortran, Java und Python sollten nach Möglichkeit vorhanden sein, oder sich zumindest problemlos realisieren lassen.

Empfohlen wurde HDF5, das die meisten der Anforderungen erfüllt.

Authentifizierung und Authorisierung Wie erwähnt, erwarten die Photon Science User einen unkomplizierten Zugang zu Informationen und Ressourcen, mit einem zweistufigem Authorisierungs-Mechanismus, der validierte (z.B. Präsentation eines Ausweises) wie nicht-validierte (email-Adresse als ID) Identitäten unterstützt. Zertifikate haben trotz intensiver Diskussionen nur eine sehr geringe Akzeptanz. Zur Zeit wird die Authentifizierung in einem User-Management-System abgebildet, das auch entsprechende Rollen-Modelle für Anwender wie Support-Staff oder Proposal-Reviewer bereit stellt. In Europa sind nur zwei verschiedene Software-Systeme im Einsatz: das DUO-basierte System vom PSI, das auch am DESY verwendet (und entwickelt) wird, sowie ein System namens SMIS, das unter anderem beim EMBL sowie dem ESRF Verwendung findet. Aktuell macht es daher am meisten Sinn, Schnittstellen zu diesen Systemen zur Authentifizierung und Authorisierung bereitzustellen.

Letztlich ist dies aber im europäischen Rahmen unbefriedigend, da viele Gruppen verschiedene Facilities verwenden und ein Facility-übergreifende Nutzung von Ressourcen und/oder Daten nicht ohne weiteres möglich ist. Verschiedene Projekte wie [EuroFEL](#), [PaNdata](#) oder CRISP arbeiten an einer pan-Europäischen Lösung, die Shibboleth-basiert den Anforderungen der Communities auf der einen Seite und den Sicherheits-Anforderungen der Facilities andererseits gerecht werden soll. Diese Projekte befinden sich allerdings noch in einem sehr frühen Entwicklungsstadium.

Eine DM-Umgebung sollte daher eine möglichst flexible Anbindung an verschiedene Authentifizierungsinstanzen erlauben, und (primär) die User-Management-Systeme anbinden, aber auch Zertifikate und Shibboleth unterstützen können. Im Falle von DESY-Einrichtungen kommen Kerberos-Credentials hinzu.

Datenmigration und Data Access Daten, die an den verschiedenen Instrumenten erzeugt werden, durchwandern in der Regel verschiedene Dateisysteme und Storages. Von einem lokalen Cache werden die Daten auf einen Fileserver kopiert, wo die Analyse parallel zum Experiment auf dedizierten Knoten durchgeführt werden kann. Nach Beendigung des Experiments wird der Zugriff auf die Daten auf dem Fileserver gesperrt, um Störungen anschließender Experimente zu vermeiden. Die Daten müssen also von dem Fileserver auf entsprechende Storage-Elemente repliziert werden, um eine Offline-Analyse sowie den Transfer der Daten ans Heimat-Institut des Experimentators zu ermöglichen. Bei diesem Prozess sollen nach Möglichkeit Metadaten eingesammelt und in eine Datenbank eingespeist werden.

Dieser Prozess der Datenmigration wird nicht von dem Anwender durchgeführt, sondern entweder automatisiert oder unter Kontrolle der Beamline-Wissenschaftler durchgeführt. Dadurch wird der Kreis der Ingestoren und potentiellen Fehlerquellen eingeschränkt.

Metadaten Es gibt bislang noch keine Metadaten-Engine, die halbwegs automatisiert Informationen aus verschiedenen Quellen sammelt und in einer Datenbank aggregieren könnte. Tatsächlich überwiegt bislang die Neigung, Metadaten auf das absolute Minimum zu reduzieren, was die Analyse der Daten, die Einbettung in Workflows und die Nachnutzung massiv erschweren wird.

Metadaten, die in ein Antragsverfahren eingeflossen sind, können direkt über die User-Management-Systeme mit den experimentellen Daten verknüpft werden. Allerdings sind diese Daten nicht ausreichend, einzelne Experimente zu charakterisieren.

An diesem Punkt besteht weiterhin Informations- und Beratungsbedarf.

3.1.5 Datamanagement VRE

Basierend auf den minimalen Anforderungen der User wie der Beamline-Wissenschaftler wurde ein generisches Layout für eine DM-Umgebung entwickelt (siehe Abb. 6). Dies bildet die typischen Prozeduren in der Photon Science Umgebung ab, ohne konkret spezifische Implementierungen, Dienste oder Infrastrukturen vorauszusetzen.

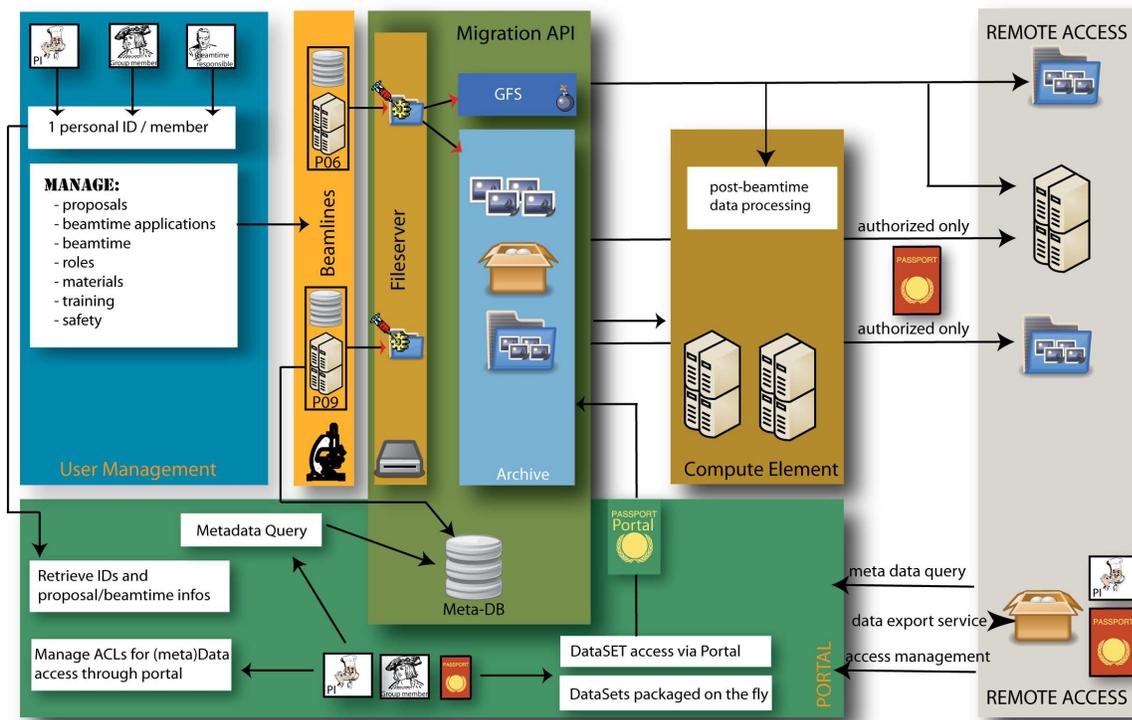


Abbildung 6: Konzept für *Lifecycle*, eine generische Photon Science DM-Umgebung

Basierend auf dem generischen Layout und der Erstellung typischer Use Cases, wurde eine spezifische Version für die am DESY ansässigen Einrichtungen entwickelt. Abbildung 7 illustriert die spezifischen Anforderungen und Anbindungen an existierende Infrastrukturen. Die Implementierung, die sich zur Zeit in einer Prototyp-Phase befindet, wird es allen Anwendern erlauben, ihre Daten remote zu managen. Obwohl dieses Konzept DESY-spezifische Infrastrukturen wie AFS oder dCache einbindet, erlaubt der modulare Aufbau eine Anbindung an beliebige Filesysteme oder SRMs. Die

Anbindung an iRods zum Beispiel wurde erfolgreich getestet. Ebenso gibt es kaum Beschränkungen für Authentifizierungs- oder Datentransfer-Protokolle. Aktuell unterstützt die prototypische Implementierung Grid-Zertifikate, DOOR/DUO-IDs oder Kerberos-Credentials. Eine Anbindung an Shibboleth wurde bislang nicht implementiert, wird aber möglich sein.

Der Daten-Transfer wird primär über HTTP erfolgen, da die Anwender in der Regel Probleme mit komplexeren Anwendungen haben. Prinzipiell könnten aber Protokolle wie Grid-FTP, WebDAV oder DCCP realisiert werden. Einige Charakteristika die zur Zeit implementiert sind:

- **Default Policies:** per Default werden automatisch die Policies wie oben dargelegt für jedes Objekt, das in das Archiv eingespeist wird, angewendet.
- **Aggregation:** Daten können automatisch aggregiert werden. Die Wandlung ganzer Datensätze in HDF5 ist zur Zeit in Arbeit.
- **Metadaten:** wie gesagt noch ein echtes Problem. Zur Zeit werden Metadaten über einfache Textdateien mit Key-Value-Paaren in die Datenbank eingespeist. Mit der Wandlung nach HDF5 lassen sich sehr viel mehr Metadaten automatisch erfassen.
- **Staging:** Daten können auf beliebige URLs restauriert werden.
- **ACLs:** Access Control kann durch die Anwender beliebig konfiguriert werden. Defaults (Zugriffsrechte für alle Teilnehmer an einem Experiment sowie den Principal Investigator) sind implementiert.
- **Shopping:** Daten können individuell in einem Shopping-Cart für den Transfer zusammengestellt werden.
- **History:** Versionierung und History werden unterstützt.
- **External:** es können Daten von externen Einrichtungen remote und/oder lokal eingespeist werden.

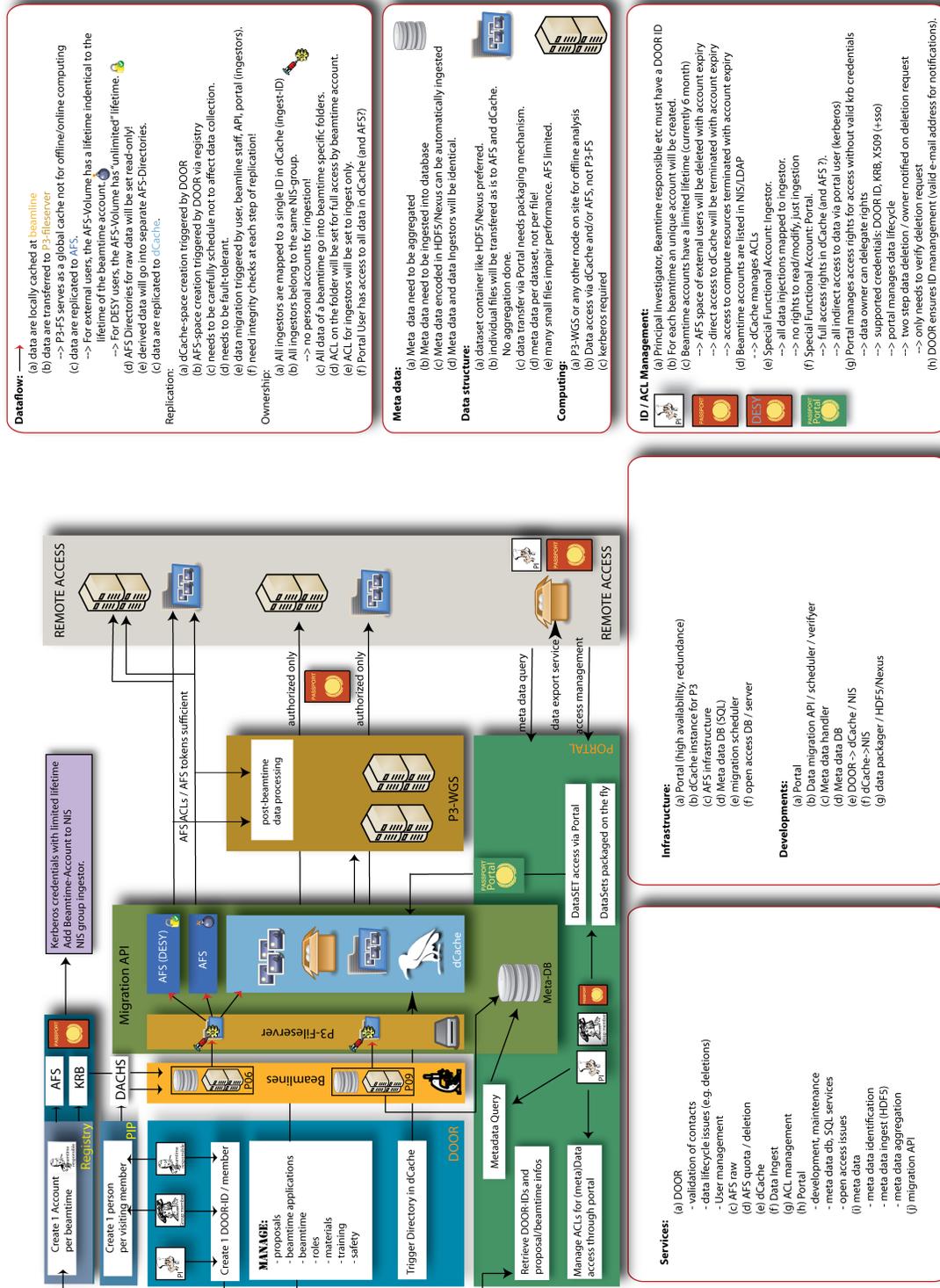
Was hingegen bislang weitgehend fehlt, sind Betriebs- und Support-Modelle.

3.1.6 Ausblick

Diverse Punkte sind noch offen oder bedürfen weiterer Diskussionen, wie zum Beispiel

- Betriebsmodell, Supportmodell
- Metadaten Handling (Erfassung, Struktur etc)
- Format-Konvertierung nach HDF5
- Ingest-Management

Andererseits ergibt sich die Möglichkeit, die DM-Umgebung für einen breiteren Anwenderkreis nutzbar zu machen. Zum einen sind dies die Partner des PNI-HDRI Projektes, die nach einer gemeinsamen oder zumindest kompatiblen DM-Plattform suchen. Zum anderen nutzt auch die Hochenergie-Physik (namentlich ILC) dieselben Beschleuniger (FLASH) für Experimente und Entwicklung, so dass eine Einspeisung der Daten in dieselbe (oder eine unabhängige Instanz derselben) DM-Umgebung erwogen wird. Dies war einer der Gründe für die Unterstützung der ungeliebten Grid- (und DFN-) Zertifikate.



Dataflow:

- (a) data are locally cached at beamline
- (b) data are transferred to P3-Fileserver
- (c) data are replicated to AFS
- (d) For external users, the AFS-Volume has a lifetime identical to the lifetime of the beamtime account
- (e) For DESY users, the AFS-Volume has "unlimited" lifetime
- (f) AFS Directories for raw data will be set read-only
- (g) derived data will go into separate AFS-Directories
- (h) data are replicated to dCache

Replication:

- (a) dCache-space creation triggered by DOOR
- (b) AFS-space creation triggered by DOOR via registry
- (c) needs to be carefully schedule not to affect data collection
- (d) needs to be fault-tolerant
- (e) data migration triggered by user, beamline staff, API, portal (ingestors)
- (f) need integrity checks at each step of replication!

Ownership:

- (a) All ingestors are mapped to a single ID in dCache (ingest-ID)
- (b) All ingestors belong to the same NIS-group
- (c) All data of a beamtime go into beamtime specific folders
- (d) ACL on the folder will be set for full access by beamtime account
- (e) ACL for ingestors will be set to ingest only
- (f) Portal User has access to all data in dCache (and AFS?)

Meta data:

- (a) Meta data need to be aggregated
- (b) Meta data need to be ingested into database
- (c) Meta data encoded in HDF5/Nexus can be automatically ingested
- (d) Meta data and data Ingestors will be identical

Data structure:

- (a) dataset container like HDF5/Nexus preferred
- (b) individual files will be transferred as is to AFS and dCache
- (c) No aggregation done
- (d) data transfer via Portal needs packaging mechanism
- (e) meta data per dataset, not per file!
- (f) many small files impair performance. AFS limited.

Computing:

- (a) P3-WGS or any other node on site for offline analysis
- (b) Data access via dCache and/or AFS, not P3-FS
- (c) Kerberos required

ID / ACL Management:

- (a) Principal Investigator, Beamtime responsible etc must have a DOOR ID
- (b) For each beamtime an unique account will be created
- (c) Beamtime accounts have a limited lifetime (currently 6 month)
- (d) AFS space of external users will be deleted with account expiry
- (e) direct access to dCache will be terminated with account expiry
- (f) access to compute resources terminated with account expiry
- (g) Beamtime accounts are listed in NIS/LDAP
- (h) dCache manages ACLs
- (i) Special Functional Account: Ingestor
- (j) no rights to read/modify, just ingestion
- (k) all data injections mapped to ingestor
- (l) Special Functional Account: Portal
- (m) full access rights in dCache (and AFS ?)
- (n) Portal manages access rights via portal user (kerberos)
- (o) indirect access to data via portal user (kerberos)
- (p) data owner can delegate rights
- (q) supported credentials: DOOR-ID, KRB, X509 (+sso)
- (r) portal manages data lifecycle
- (s) two step data deletion / owner notified on deletion request
- (t) may need ID verification
- (u) DOOR ensures ID management (valid e-mail address for notifications)

Infrastructure:

- (a) Portal (high availability, redundance)
- (b) dCache instance for P3
- (c) AFS infrastructure
- (d) Meta data DB (SQL)
- (e) migration scheduler
- (f) open access DB / server

Developments:

- (a) Portal
- (b) Data migration API / scheduler / verifier
- (c) Meta data handler
- (d) Meta data DB
- (e) DOOR -> NIS
- (f) dCache->NIS
- (g) data packager / HDF5/Nexus

Services:

- (a) DOOR
 - validation of contacts
 - data lifecycle issues (eg, deletions)
 - User management
- (b) AFS raw
- (c) AFS quota / deletion
- (d) dCache
- (e) Portal
- (f) ACL management
- (g) Portal
- (h) Portal
 - development, maintenance
 - meta data db, SQL services
 - open access issues
- (i) meta data
 - meta data identification
 - meta data ingest (HDF5)
 - meta data aggregation
- (j) migration API

Abbildung 7: Konzept für Lifecycle eine spezifische DM-Umgebung

3.2 Sozialwissenschaften

Nach einer vorangegangenen allgemeinen Information über die Grid-Umgebung und WissGrid-Aktivitäten wurde im Rahmen der WissGrid-AP3-Begutachtung der Anwendungsfall der Sozialwissenschaften besprochen und vorgestellt. Im Anschluss wurden im Rahmen eines Workshops die fachwissenschaftlichen Anforderungen an eine virtuelle Arbeitsumgebung für SOEB erarbeitet. Weiterhin wurde WissGrid im Februar 2010 damit beauftragt, eine Expertise für (VirtAug) zur Grid-Nutzung zu erstellen. Diese liegt inzwischen vor.

3.3 Biostatistik und Epidemiologie

3.3.1 Beschreibung

Die Universitätsmedizin Göttingen hat neben der Biostatistik zudem auch die Genetische Epidemiologie in Fragen zu Grid-Technologien beraten. In beiden Fällen ist eine konkrete Eingrenzung der Community nicht möglich gewesen, so dass nicht die gesamte Community der Biostatistik bzw. Genetische Epidemiologie beraten wurde, sondern eine Untergruppe.

Die Biostatistik wird in den Beratungstätigkeiten von der Gruppe von Prof. Reißbarth aus der Medizinischen Statistik der Universitätsmedizin Göttingen repräsentiert. In der Biostatistik werden die im Rahmen klinischer Studien erhobenen Daten ausgewertet, um z.B. Korrelationen zwischen Phänotyp und genetischen Merkmalen zu erforschen. Die Basis dafür sind anonymisierte bzw. pseudonymisierte Daten zu Genetik, Biomaterial, Bildmaterial und Krankheitsbildern, auf denen z.B. Mustererkennungsalgorithmen und statistische Analysen im Grid durchgeführt werden. Ein typisches Szenario ist hierbei der Erhalt bestimmter in einer Studiendatenbank vorgehaltenen Daten und die anschließende Auswertung mithilfe von statistischen Entwicklungsumgebungen wie R oder SAS.

Die Leiterin der Genetischen Epidemiologie der Universitätsmedizin Göttingen, Prof. Bickeböller, ist zugleich 1. Vizepräsidentin der GMDS (Deutsche Gesellschaft für Medizinische Informatik, Biometrie und Epidemiologie). Zudem ist Prof. Bickeböller Mitglied im GMDS-Fachausschuss für Epidemiologie. Somit ist die beratende Untergruppe der Genetischen Epidemiologie eine repräsentative Untermenge einer Gesamtcommunity. In der Epidemiologie werden Ursachen und Folgen sowie Verbreitung von Krankheitsbildern in Populationen untersucht. Dabei werden vornehmlich statistische Methoden angewendet. Daher besteht eine große methodische Nähe zur Biostatistik.

3.3.2 Status der Beratungstätigkeiten

Sowohl in der Biostatistik als auch in der Genetischen Epidemiologie bestehen rechen- und datenintensive Aufgaben. Da keine ausreichende IT-Infrastruktur vorhanden ist, um diese Aufgaben mit lokal verfügbaren Systemen optimal auszuführen, besteht ein Bedarf an zusätzlicher Rechenleistung. Grid-Technologie wird hier als eine mögliche Lösung erachtet. Für ein zukünftiges Projekt der Genetischen Epidemiologie wurde daher ein Support Letter über das Angebot von Grid-Ressourcen sowie entsprechender Schulungsmaßnahmen erstellt.

Weiterhin wurden die Mitarbeiter und Mitarbeiterinnen der Biostatistik sowie der genetischen Epidemiologie durch gezielte Gespräche und Workshops auf die Nutzung von Grid-Technologien vor-

bereitet. Hierbei kam es zu Beratungsgesprächen durch Mitglieder des WissGrid-Fachberaterenteams, bei denen die bestehende Situation analysiert und das Potenzial von Grid-Technologie zur Optimierung herausgestellt worden ist. Zudem haben die Biostatistik sowie die Genetische Epidemiologie Workshops, wie der AP2-Workshop „Virtuelle Forschungsumgebungen aufbauen - mit D-Grid“ [44] oder einem speziell für die Biostatistik ausgerichteten Workshop „Grid and Cloud Computing for Computational (Bio-) Statistics“ [45] besucht, um einen Einstieg in das Thema Grid Computing zu erhalten.

Als erster Schritt auf dem Weg zur Nutzung der Grid-Infrastruktur, wurden typische Anwendungen aus der Biostatistik sowie der Genetischen Epidemiologie untersucht. Neben der Statistikanwendung R, die in beiden Abteilungen ein Standardwerkzeug bilden, wurde die Aufgabe der Imputation von genetischen Daten betrachtet. Exemplarisch wurden diese Berechnungen auf den Einsatz eines lokalen D-Grid-Rechenclusters angepasst. Für das Beispiel der Imputation konnten dadurch erhebliche Berechnungszeiteinsparungen erzielt werden [46]. Im nächsten Schritt soll nun eine Überführung dieser typischen Anwendungen in das Grid-Umfeld erfolgen. Die Mitarbeiter der Biostatistik sowie genetischen Epidemiologie wurden daher bereits bei dem Bezug eines Grid-Zertifikats beraten. Als nächste Maßnahme ist es erforderlich, notwendige Schritte zur Nutzung der D-Grid-Infrastruktur abzustimmen. Hierzu sind Beratungstätigkeiten bei der Verteilung der Anwendungen auf Grid-Knoten, der Erstellung von Grid-Workflows zur verteilten Abarbeitung voneinander abhängiger Einzelschritte einer Aufgabe sowie die Vereinfachung der Nutzung durch die Einbettung der Grid-Anwendungen in eine Benutzeroberfläche geplant. Dabei werden die neuen Communities in die MediGRID-VO integriert.

4 Workshop Virtuelle Forschungsumgebungen aufbauen - mit D-Grid

4.1 Dokumentation des Workshops

Am 19.01.2011 wurde der erste Fachberater-Workshop zum Thema *Virtuelle Forschungsumgebungen aufbauen - mit D-Grid* an der SUB Göttingen veranstaltet. Die Zielsetzung des Fachberater-Workshops war es, Gespräche mit neuen Community-Grids zu initiieren und die vorhandene Infrastruktur bzw. einhergehende Möglichkeiten aus D-Grid vorzustellen. Weitere Informationen sowie alle Beiträge zu dem Fachberater-Workshop finden sich auf den [WissGrid-Seiten](#).

4.1.1 Communities und die Bildung Virtueller Organisationen

Eröffnet wurde der Fachberater-Workshop durch Harry Enke (AIP) mit seinem Vortrag zu Communities und der Bildung von Virtuellen Organisationen (VO). Der Begriff der Virtuellen Forschungsumgebungen wurde definiert und die dazugehörige Infrastruktur vorgestellt. Der Einsatz einer solchen Infrastruktur erfolgt momentan nur in der akademischen Forschung.

- Über Communities mit rechenintensiven Aufgaben
Nach dem Beitrag zu den vorhandenen Formen und Strukturen berichtete Illya Agapov (DESY) über Communities mit rechenintensiven Aufgabenstellungen in einer Best-Practice-Anwendung (LHC). Ergänzend bemerkte Harry Enke, dass sich daten- und rechenintensive

Anwendungen nicht genau abgrenzen lassen. Die Übergänge sind fließend. Die zentrale Frage ist, wie verteilt man die Daten?

- Über Communities mit datenlastigen Aufgaben
In einem zweigeteilten Vortrag wurde dann über Communities mit datenlastigen Aufgaben berichtet. Benjamin Löhnhardt (UMG) stellte Praxisbeispiele aus dem Bereich Bildverarbeitung und Statistik vor. Eines der Beispiele war die Imputation nicht-genotypischer Marker auf Gen-Chips. Dabei gab es eine Frage zur Vergleichbarkeit solcher Chips, die von Melanie Sohns (Genetische Epidemiologie) mit nein beantwortet wurde. Jens Ludwig (SUB) stellte die TextGrid-Basisinfrastruktur mit 275 TB Speichervolumen ins Verhältnis zur Menge der digitalisierten Bücher in Europa mit 100 PB. Herausgestellt wurde die Bedeutung des Rechtemanagements. In Bezug auf Nutzerakzeptanz setzt man auf einfache Tools, Schulung, Workshops, Benutzerhandbücher und gemeinsame Repositorien. Heike Neuroth (SUB) erklärte auf Anfrage, dass es nur lesenden Zugriff gibt. Es werden Kopien angelegt, um weiter daran arbeiten zu können mit dem Hinweis auf den jeweiligen Veröffentlichungszeitpunkt.

4.1.2 Instrument-, Compute-, Data-Sharing: Ressourcen gemeinsam nutzen statt selber aufbauen

Frank Schlünzen (DESY) zeigte aktuelle Beispiele zum Thema Ressourcen-Sharing. Dabei ging es im Kern um das Anliegen der Communities, Ressourcen gemeinsam zu nutzen statt selber lokale Systeme aufzubauen. Essentiell dafür ist eine Definition von sogenannten KO-Kriterien (z.B. Usability). Ein weiterer wichtiger Punkt sind Richtlinien für die Veröffentlichung von Daten. Diese können über Förderungsbedingungen erzwungen werden, aber ist das die richtige Strategie? Harry Enke bemerkte den darin begründeten kulturellen Wandel in den Communities, es fehle aber noch die Unterstützung bei der Durchsetzung solcher Deklarationen. Heike Neuroth fügte hinzu, dass Wissenschaftler und Gutachter noch nicht sensibilisiert seien für das Thema. Bernadette Fritzsch fragte nach den Methoden für das Datenmanagement. Es sollten keine Sanktionen verhängt werden, es müssten eher Standardisierungen/Richtlinien aus den Communities kommen, z.B. die Zitierfähigkeit der produzierten Daten als betriebsartbedingter Benefit, denn die Daten müssen referenzierbar sein (gute wissenschaftliche Praxis).

- Gemeinsame Nutzung von Instrumenten
Im Anschluss berichtete Jens Klump (GFZ Potsdam) über die gemeinsame Nutzung von Massenspektrometern. Dabei wurden die Fernsteuerung der Instrumente und die Bereitstellung der Ergebnisse vorgestellt. Es wurde auf die DOI für Experimentaldaten und die IGSM als Probenreferenz verwiesen; bei der TIB Braunschweig sind bereits DOI für Geodaten hinterlegt. Tanja Schmidt (SOFI) stellte den Use Case der Sozialforschung vor. Gegenstand seien Outputs von in SPSS/ R-generierten Datensätzen; dabei sind datenschutzrelevante Bedingungen unbedingt einzuhalten (Mikrozensusdatensätze). Die Syntax (R/SPSS) sei die eigentliche wissenschaftliche Leistung. Sharing werde von der Community benötigt, aber die Umsetzung sei bisher ungeklärt. Gabriele Dröge (BGBM) stellte die DNA-Bank-Netzwerke vor. Es wird ein Wrapper benutzt und der Zugriff erfolgt über Webservices.
- Gemeinsame Nutzung von Daten
Neela Enke (BGBM) berichtete über eine Umfrage in der Biodiversitäts-Community. Die Befragten sahen primär den erhöhten Aufwand im Zusammenhang mit einem Data-Sharing,

obwohl dies effektiv eher weniger zutrifft. Die Problemstellung in der Community ist das Mergen und die Visualisierung der Daten auf übergreifender Ebene. Zustimmend zu den Vortragend wurde die Integration des Datenmanagements in den Workflow vor der Publikation als notwendig angesehen, ebenso die benötigte Änderung in der Publikationspolitik.

4.1.3 Gemeinsame Standards, Verfahren und Arbeitsabläufe nutzen

Bernadette Fritsch (AWI) eröffnete die nächste Session mit Grundüberlegungen beim Aufbau von VREs, z. B. das vorhandene Standards, Verfahren und Workflows in den einzelnen Communities gemeinsam genutzt werden sollen.

- Ein Blick auf die Grid-Architektur
Christian Grimme (TU Dortmund) gab einen generalisierten Überblick über die aktuelle Grid-Architektur, wobei z.B. auf die Verwendung von homogenen und heterogenen Ressourcen eingegangen wurde. Heike Neuroth bemerkte zusätzlich, dass die Daten auch zur Ressourcenschicht zu zählen seien.
- VRE-Management und Sicherheit
Aufbauend darauf stellte Harry Enke die zentralen Thesen zum Thema Management und Sicherheit im Kontext Virtueller Forschungsumgebungen vor. Wichtig für den Anwender seien benutzbare Arbeitsumgebungen, aber auch Sicherheitsaspekte, wozu z.B. eine Klassifizierung der Provider oder Leistungsgarantien zählen.
- Informations- und andere Dienste
Torsten Rathmann (DKRZ) stellte die wichtigsten Informationsdienste aus C3Grid und AstroGrid vor. Im C3Grid gehört zu den Aufgaben des Metadatensuchdienstes nicht nur die Suche in den Metadaten, sondern auch das Harvesten von Metadaten, d.h. das Einsammeln aus lokalen Metadatenkatalogen. Der Datenmanagementdienst hingegen ist für alle Datenmanagementaktivitäten im C3Grid zuständig, wie z.B. das Holen und Zwischenspeichern von Daten. Eine Besonderheit unter den Informationsdiensten ist das Datenstrom-Management von AstroGrid. Bei einem Datenstrom kann es sich z.B. um Messdaten eines Teleskopes handeln. Mehrere Teleskope können mit Hilfe dieses Dienstes ihre Beobachtungen aufeinander abstimmen. Wenn ein interessantes Ereignis erkannt wird, z.B. eine Gravitationslinsen-Anomalie, kann dieses gemeinsam beobachtet werden.

4.1.4 Über rechtliche, finanzielle und organisatorische Rahmenbedingungen

Oliver Schmid (Uni Trier) begleitete die letzte Session über rechtliche, finanzielle und organisatorische Rahmenbedingungen von VREs. Das Problem des Datenschutzes und der Datensicherheit in VREs ist erkannt und wird bearbeitet, z.B. bei der Verarbeitung von Patientendaten bei der biomedizinischen Community.

- Einige rechtliche Aspekte
Die Definition einer Rechtsform bei einer VRE ist ein vielschichtiges Thema, vor allem bei internationaler Kooperation muss das Tatortprinzip berücksichtigt werden, da auf internationaler Ebene sehr unterschiedliche Datenschutzrichtlinien existieren. Was den Punkt der

Daten betrifft, ist die Situation innerhalb der Uni geklärt. Wie ist das im Grid? Urheber- und Patentrecht sind im Grid eher geklärt, bei den Daten ist die Situation hingegen ungeklärt. Das *Sense of Ownership* ist im Gridumfeld eher ein soziales als ein rechtliches Problem, wird aber zum IT-Problemfall. Im Anschluss wurde über Potentiale und Möglichkeiten nationaler und europäischer Förderung berichtet.

- **Nationale Förderung**
Frank Dickmann stellte fest, dass die Förderlandschaft bzw. Förderpolitik an die Anforderungen durch künftige VREs angepasst werden müsse. Heike Neuroth bemerkte, dass die momentan genutzte IT-Infrastruktur nicht in den Programmpauschalen der Förderer (BMBF, DFG) enthalten sei. Jens Klump fügte hinzu, dass ein entsprechendes Angebot von Rechenleistung monetären Austausch voraussetze. Bernadette Fritsch stellte fest, dass dazu rechtliche Grundlagen geschaffen werden müssten, wie z.B. die Rechnungslegung von IT-Providern auszuweisen hat.
- **Europäische Förderung**
Heike Neuroth machte Ausführungen zu den europäischen Initiativen ESFRI für Infrastrukturen und ERIC als Gesellschaftsform mit Mehrwertsteuerbefreiung im europäischen Raum, wobei fünf Jahre Betrieb und Laufzeit garantiert sein müssen.

4.2 Schlussbemerkungen und Tendenzen zum Workshop

Harry Enke (AIP) resümierte, dass Virtuelle Forschungsumgebungen auch primär soziale Aspekte adressieren. Bezüglich monetärer Ausgleichsströme sei noch offen, wie in einer solchen Umgebung mit der Mehrwertsteuer umgegangen werden soll. Ebenso sei momentan nur die Nutzung von Rechenleistung im D-Grid möglich, nicht möglich seien z.B. 10 Jahre Datenarchivierung. Zwischen den beiden Feldern Infrastruktur und Basisleistung gebe es keine Kommunikation, was als Problem erkannt und herausgestellt wurde. Das Spektrum einer VRE müsse breiter vorgestellt und das Bedürfnis zur Vertiefung der Themen berücksichtigt werden in der zukünftigen Arbeit. Auf die Frage nach einer Minimalgröße für Grids antwortete Bernadette Fritsch (AWI), dass die Größe nicht fix sei und von der wissenschaftlichen Fragestellung abhängen. Das Schlussfazit drehte sich um die Frage, ob sich die Nutzung einer Grid-Infrastruktur für eine VRE lohnt. Das hängt vor allem davon ab, wie durch das Grid die Prozesse in der Community optimiert werden können und wie sich dies gegenüber der Anschaffung von lokalen Ressourcen darstellt.

4.3 Weitere geplante Aktionen

Im Anschluss wurden mögliche flankierende Veranstaltungen/Aktivitäten in die allgemeine Diskussion gegeben. Darunter war ein allgemeiner Workshop für vertiefte Beratung der neuen Communities, ein Workshop zur Einführung in das Grid, ein Workshop über vorhandene Sicherheitsmechanismen, ein Workshop für Anwender mit praktischen Übungen (mit kleiner Teilnehmerzahl für intensivere Beratungsgespräche) und die Sammlung der Themen und Bedarfe über die Fachberater-Liste, um die Inhalte für die Workshops einzugrenzen.

Quellen

- [1] <http://en.wikipedia.org/wiki/FLOPS>
- [2] Standard Performance Evaluation Corporation <http://www.spec.org/>
- [3] <http://www.r-project.org/>
- [4] Worldwide LHC Computing Grid <http://lcg.web.cern.ch/lcg/>
- [5] The European XFEL
<http://www.xfel.eu/>
- [6] The European XFEL Computing Technical Design Report (TDR)
- [7] GWES Grid Workflow <http://www.gridworkflow.org/snips/gridworkflow/space/GWES>
- [8] Kepler Workflow
<https://kepler-project.org>
- [9] iRods
<https://www.irods.org/>
- [10] Fedora Commons
<http://fedora-commons.org/>
- [11] Globus Middleware
<http://www.globus.org/toolkit/>
- [12] gLite Middleware
<http://glite.web.cern.ch/glite/>, <http://en.wikipedia.org/wiki/GLite>
- [13] UNICORE Middleware
<http://unicore.eu>,
- [14] dCache
<http://www.dcache.org/>
- [15] OGSA-DAI
<http://www.ogsadai.org.uk/>
- [16] Blaupausen für den Aufbau einer technischen Support-Infrastruktur
[Support-Infrastruktur.pdf](#)
- [17] NGI-DE-Portal
<https://helpdesk.ngi-de.eu/>
- [18] EMBL: European Molecular Biology Laboratory, Outstation Hamburg
<http://www.embl-hamburg.de/>
- [19] AMGA
amga.web.cern.ch/amga/
- [20] Community-Grids - Überblick und Report
[WP2 Deliverable 2.1.1](#)

- [21] F. Siebenlist, R. Ananthakrishnan, D. E. Bernholdt, L. Cinquini, I. T. Foster, D. E. Middleton, N. Miller, D. N. Williams: Earth System Grid Authentication Infrastructure: Integrating Local Authentication, OpenID and PKI, *The 2009 TeraGrid Conference*, www.teragrid.org/tg09/files/tg09_submission_79.pdf
- [22] HDF5 Data Model
<http://www.hdfgroup.org/HDF5/>
- [23] GridSphere
<http://www.gridisphere.org>
- [24] Jetspeed
<http://portals.apache.org/jetspeed-2/>
- [25] Liferay portal
<http://www.liferay.com/>
- [26] ESRF-Up WP11, EGEE feasibility study
<http://www.esrf.eu/Infrastructure/Computing/Grid>
- [27] GridShib: Bridging SAML/Shibboleth and X.509 PKI for campus and grid interoperability
<http://gridshib.globus.org/>
- [28] Auto-Rickshaw: The EMBL-HH Automated Crystal Structure Determination Platform
<http://www.embl-hamburg.de/Auto-Rickshaw/>
- [29] Panjikar S, Parthasarathy V, Lamzin VS, Weiss MS, Tucker PA. Auto-Rickshaw: an automated crystal structure determination platform as an efficient tool for the validation of an X-ray diffraction experiment. *Acta Crystallographica D*61:449-457, 2005.
- [30] Information System for Protein CrystallographY Beamlines (ISPYB)
<http://www.esrf.eu/UsersAndScience/Experiments/MX/Software/isyb>
- [31] PDB: The Protein DataBase
<http://www.pdb.org/>
- [32] NDB: The Nucleic acid DataBase
<http://ndbserver.rutgers.edu/>
- [33] The Electron Microscopy DataBase
www.ebi.ac.uk/pdbe/emdb/
- [34] EBI: European Bioinformatics Institute
<http://www.ebi.ac.uk/>
- [35] NCBI: The National Center for Biotechnology Information
<http://www.ncbi.nlm.nih.gov/>
- [36] GlassFish: open source application server implementing Java EE
<https://glassfish.dev.java.net/>
- [37] PaNdata: Photon and Neutron Data Infrastructure
http://www.pan-data.net/Main_Page

- [38] Berliner Erklärung
<http://oa.mpg.de/lang/de/berlin-prozess/berliner-erklarung/>
- [39] OECD Principles and Guidelines for Access to Research Data from Public Funding
<http://www.oecd.org/dataoecd/9/61/38500813.pdf>
- [40] Empfehlungen zur gesicherten Aufbewahrung und Bereitstellung digitaler Forschungsprimärdaten: [DFG Empfehlungen](#)
- [41] Climate Service Center (CSC),
<http://www.climate-service-center.de/>
- [42] Klimabüro für Polargebiete und Meeresspiegelanstieg,
<http://www.klimabuero-polarmeer.de/>
- [43] Institut für Umweltphysik an der Universität Bremen (IUP),
<http://www.iup.physik.uni-bremen.de/>
- [44] Fachberater Workshop: *Virtuelle Forschungsumgebungen aufbauen - mit D-Grid*,
<http://www.wissgrid.de/workgroups/ap2/workshop-2011-01.html>,
- [45] Workshop *Grid and Cloud Computing for Computational (Bio-)Statistics*,
<http://rostlab.org/cms/cloudstat2011/>
- [46] Löhnhardt B, Quade M, Skrowny D, Sohns M, Bickeböller H, Sax U:
Hochleistungsrechencluster zur Unterstützung der biomedizinischen Forschung.
In: 55. GMDS-Jahrestagung 2010, Mannheim, 2010, pp. 447-449.