



Symposium Forschungsdaten-Infrastruktur (FDI 2013)
22. Januar 2013, GeoForschungsZentrum Potsdam (GFZ)

**Gemeinsam organisiert von den DFG-Projekten
Radieschen, re3data.org, KomFor, EWIG und BoKeLa**

Symposium - Materialien

doi:10.2312/RADIESCHEN_002

Gefördert von:

DFG Deutsche
Forschungsgemeinschaft

Inhalt

Symposium-Bericht	4
Agenda.....	16
Vortragsprogramm (Materialien)	17
Keynote – Dr. Torsten Reimer (JISC).....	17
Session 1 - Projektpräsentationen.....	27
DFG-Projekt EWIG – Tim Hasler (ZIB)	27
DFG Projekt Radieschen – Dr. Jochen Klar (AIP).....	29
DFG-Projekt KomFor – Dr. Michael Diepenbroek (MARUM)	34
DFG-Projekt re3data.org – Frank Scholze (KIT)	40
DFG-Projekt BoKeLa – Dr. H.-J.Wallrabe-Adams	43
Session 2 – Data Curation Continuum – Teil 1	48
„Private Domäne ‚Die Sicht des Forschers‘“ – Prof. Frederik Tilman (GFZ)	49
„Gruppendomäne ‚VfUs‘“ – Dr. Peter Bartelheimer (Uni Göttingen).....	58
Session 3 – Data Curation Continuum – Teil 2	62
„Persistente Domäne ‚PID, DLZA, Zertifikate‘“ – Reiner Mauer (GESIS)	62
„Zugangsdomäne ‚Portale, Best Practices‘“ - Dr. Hans Pfeiffenberger (AWI)	72



Symposium Forschungsdaten-Infrastruktur (FDI 2013) 22. Januar 2013, GeoForschungsZentrum Potsdam (GFZ)

Gemeinsam organisiert von den DFG-Projekten
Radieschen, re3data.org, KomFor, EWIG und BoKeLa

Symposium-Bericht

Die Flut an digitalen Daten in Wissenschaft und Forschung wächst rasant. Die dauerhafte Speicherung dieser Daten für zukünftige Generationen von Forscherinnen und Forschern stellt das gesamte Wissenschaftssystem vor Herausforderungen. Jüngst hat die EU-Kommission Empfehlungen an ihre Mitgliedsstaaten zur dauerhaften Zugänglichkeit von Forschungsdaten verabschiedet. Doch noch sind viele Fragen ungelöst. So müssen Aspekte der Finanzierung, Organisation und Technologie der zu schaffenden Forschungsdaten-Infrastrukturen, sowie deren rechtliche und politische Rahmenbedingungen geklärt werden. Diese Themenfelder wurden im Rahmen eines gemeinsam organisierten Symposiums der DFG-Projekte Radieschen, re3data.org, KomFor, EWIG und BoKeLa diskutiert.

Im Hauptprogramm des Symposiums wurden, mit dem Datenlebenszyklus als Leitlinie, die jeweiligen Herausforderungen und neuesten Erkenntnisse durch Vorträge und daran anschließende Diskussionen erörtert. In begleitenden Workshops wurden zentrale Themenfelder von Expertinnen und Experten bearbeitet und deren Ergebnisse anschließend dem Publikum vorgestellt.

Das Symposium versammelte die Experten-Community zu Forschungsdaten aus dem deutschsprachigen Raum. Kern der Veranstaltung bildeten Vorträge zu verschiedenen Aspekten der Forschungsdaten-Infrastrukturen, sowie Präsentationen laufender Forschungsprojekte. Darüber hinaus wurden zentrale Herausforderungen der Weiterentwicklung von Forschungsdaten-Infrastrukturen in Workshops diskutiert. Desweiteren bot das Symposium eine Plattform zum Gedankenaustausch und zum Netzwerken.

Das Symposium fand am 22. Januar 2013 am Helmholtz-Zentrum Potsdam, Deutsches GeoForschungsZentrum GFZ statt. Eröffnet wurde die Veranstaltung durch den Vorstand des GeoForschungsZentrums **Prof. Dr. Dr. h.c. Reinhard F.J. Hüttl** und durch ein Geleitwort von **Dr. Stefan Winkler-Nees (DFG)**. Das Vortragsprogramm begann mit einer Keynote von Dr. Torsten Reimer vom Joint Information Systems Committee (JISC).



Abb. 1: Begrüßung durch den Vorstand des GFZ – Prof. Dr. Dr. h.c. Reinhard F.J. Hüttl

Keynote

Dr. Thorsten Reimer (JISC) beschrieb in seiner Keynote die Herausforderungen im Aufbau von Forschungsdaten-Infrastrukturen und erläuterte die Vorteile der dauerhaften Zugänglichkeit von Daten für die verschiedenen Stakeholder. Nach einer kurzen Vorstellung von JISC erläuterte Dr. Reimer Visionen für den Entwicklungsstatus von Forschungsdaten - Infrastrukturen mit 10 Jahresperspektive.

JISC verfolgt in Großbritannien einen 4-Phasenplan für die Entwicklung von Forschungsdaten-Infrastrukturen und dem Management von Forschungsdaten, an dem sich auch das zugehörige vierstufige Förderprogramm orientiert. Momentan befindet sich das Programm in der Konsolidierungsphase, in der es darum geht, gefundene Lösungen zu verfestigen und institutionelle Kapazitäten weiter auszubauen. In diesem Kontext beschrieb er zwölf damit verbundene Herausforderungen und verschiedene Ansätze, diesen Herausforderungen zu begegnen.



Abb. 2: Keynote-Sprecher Dr. Thorsten Reimer, JISC

Präsentation der DFG-Projekte Radieschen, re3data.org, KomFor, EWIG und BoKeLa

Frank Scholze, KIT, stellte das Projekt re3data.org – Registry of Research Data Repositories vor. Das Projekt hat das Ziel Daten-Repositoryn in einem web-basierten Verzeichnis zu erschließen und so eine Orientierung über bestehende Datensammlungen zu bieten. Zentrales Anliegen von re3data.org ist es, Wissenschaftlern eine Orientierung in der heterogenen Landschaft der Forschungsdaten-Repositoryn zu geben. Dabei wird sowohl die Rolle der Erhebenden, die z. B. durch Förderorganisationen oder wissenschaftliche Zeitschriften aufgefordert sind Daten zugänglich machen, als auch die Rolle der Nutzenden, z. B. Forschende, die Daten von Dritten nachnutzen möchten, fokussiert. Weiter soll infrastrukturellen Dienstleistern, wie Datenzentren, Rechenzentren und Bibliotheken, eine Übersicht über die Landschaft der Forschungsdaten-Repositoryn gegeben werden. Darüber hinaus kann der Dienst von Förderorganisationen genutzt werden: Fordern diese Mittelempfänger auf, gewonnene Forschungsdaten offen zugänglich zu machen, kann eine Recherche in re3data.org helfen, potenzielle Repositoryn zur Speicherung der Forschungsdaten zu finden. Weiterführende Informationen finden sich auf der Projekt-Website (<http://www.re3data.org>).



Abb. 3: Präsentation „re3data.org“, Frank Scholze, KIT

Tim Hasler, ZIB, präsentierte das DFG-Projekt EWIG - Entwicklung von Workflowkomponenten für die Langzeitarchivierung von Forschungsdaten in den Geowissenschaften. EWIG zielt auf Probleme speziell bei der Datenübergabe an Langzeitarchive. Dazu werden drei Themen beleuchtet. Zum einen werden auf institutioneller Ebene bei den Projektpartnern GFZ und Institut für Meteorologie der FU Berlin Policies entwickelt, die die Grundlage für eine genaue Übergabedefinition zwischen Datenlieferant und Archiv / Repository sind. Am Zuse-Institut Berlin liegt der Fokus auf der technischen Qualitätssicherung gelieferter Daten. Der dritte Schwerpunkt der Arbeiten soll das Thema Forschungsdatenmanagement bereits in der fachwissenschaftlichen Ausbildung etablieren. Dazu werden in EWIG Konzepte für Lehr- und Weiterbildungsveranstaltungen entworfen. <http://ewig.gfz-potsdam.de/>.

Dr. Jochen Klar, AIP, präsentierte das DFG-Projekt „Rahmenbedingungen einer disziplin-übergreifenden Forschungsdateninfrastruktur“, kurz Radieschen. Ziel des Projekts ist die Erstellung einer Roadmap mit Handlungsempfehlungen für eine disziplinübergreifende Infrastruktur für Forschungsdaten in Deutschland.

Sie identifiziert und behandelt die Anforderungen an generische Komponenten einer Infrastruktur und die Vernetzung mit disziplinspezifischen Komponenten. Die Analyse basiert auf einer Bestandsaufnahme mit bestehenden und neuen Projekten sowie Maßnahmen zum Community Building. Zentrale Dimensionen der Analyse sind Technik, Organisation und Kosten sowie die Untersuchung von Querschnittsthemen. Das Projektkonsortium ist ebenso disziplin- und organisationsübergreifend wie die Analyse. Weiterführende Informationen finden sich auf der Projekt-Website unter <http://www.forschungsdaten.org/>.

Dr. Michael Diepenbroek, MARUM, präsentierte das DFG-Projekt KomFor – Kompetenzzentrum für Forschungsdaten aus Erde und Umwelt (<http://www.komfor.net>). Das Kompetenzzentrum ist als Bindeglied zwischen Forschungseinrichtungen, Verlagen, Bibliotheken und einem bestehenden Archivnetzwerk für Daten aus Erd- und Umweltforschung geplant. Allgemeines Ziel ist die nachhaltige Verbesserung von Datenverfügbarkeit und –qualität. Konkret sollen nachhaltige und verlässliche Wege zur Publikation wissenschaftlicher Daten geschaffen werden die den Qualitätsstandards wissenschaftlichen Publizierens entsprechen. Das Kompetenzzentrum soll hierfür wissenschaftliche Projekte, Institute, Forschergruppen aber auch einzelne Wissenschaftler in allen Fragen des Datenmanagements begleiten – von der Planungsphase über die Datenerhebung, Qualitätssicherung, Registrierung und Langzeitarchivierung bis zur Publikation der Daten. Weiterführende Informationen zum Projekt auf <http://www.komfor.net/>.



Abb. 4: Präsentation ‚KomFor‘, Dr. Michael Diepenbroek, MARUM

Dr. Hans-Joachim Wallrabe-Adams, MARUM, erläuterte in seinem Vortrag die Ziele und aktuellen Ergebnisse des Projekts BoKeLa - Aufbau des Dateninformationssystems für das GESEP Kern- und Probenlager zur Erfassung und Verwaltung von Bohrkernen und Nachweis der Bestände in einem Internetportal (<http://www.gesep.de>). Das Projekt BoKeLa wurde vom German Scientific Earth Probing Consortium (GESEP e.V.) initiiert. Dem Konsortium gehören 15 geowissenschaftliche Institute mit einer breiten Expertise im Bereich Forschungsbohrungen an Land, im Meer und im Eis an. BoKeLa wird von drei Projektpartnern getragen: der Bundesanstalt für Geowissenschaften und Rohstoffe (BGR), dem Deutschen GeoForschungszentrum (GFZ) und dem Zentrum für Marine Umweltwissenschaften (MARUM). Ausgangspunkt der Entwicklung eines Bohrkern- und Probenverwaltungssystems ist das Drilling Information System (DIS) das bereits seit Jahren in diesem Bereich eingesetzt wird, aber für die hier geplante Nutzung umfangreicher Anpassungsarbeiten bedarf. Dieses System unterstützt das Datenmanagement bereits an der Bohrung, dann im Labor und dient später der Verwaltung im Kernlager. Alle Kerne und Proben werden mit einem eindeutigen Identifikator versehen, der International Geo Sample Number (IGSN). Damit wird es möglich sein Kerne und Proben dauerhaft identifizierbar und in der Fachliteratur zitierbar zu machen, sowie wissenschaftliche Daten direkt Proben zuzuordnen. Über ein Webportal werden die Kernlager-Informationen nach dem Open Access Prinzip zugänglich gemacht und mit Datenzentren (z.B. PANGAEA) und Publikationen verknüpft.

Vorträge

Die Vorträge des Hauptprogramms spiegelten die einzelnen Stationen im Datenlebenszyklus wieder. Abb. 5 zeigt das Domänenmodell nach Treloar mit integrierter dauerhafter Domäne.

Als ein Beispiel für die Private Domäne erläuterte **Prof. Frederik Tilman (GFZ)** den Datenzyklus in der Wissenschaft aus der Perspektive der Forscher anhand seines Forschungsgebiets Seismologie. Als Herausforderung beschrieb er das Auffinden von früher datierten Datensets, die Verwaltung der Daten, das Vermeiden von Überschreibungen, sowie die Nicht-Reproduzierbarkeit der Daten aufgrund fehlender oder veralteter Software. Eingesetzte Software-Lösungen sind zumeist ‚low-tech solutions‘.

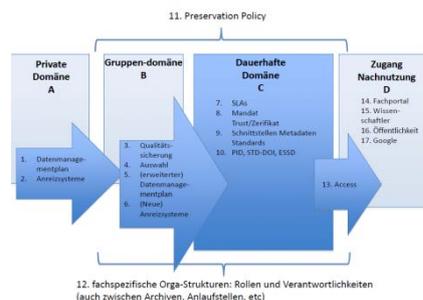


Abb. 5: Domänen-Modell
(abgeleitet von Treloar, A. &
Harboe-Ree, 2008)

Eine Herausforderung sei es, laut Prof. Tilman, die Wissenschaftler dazu zu bewegen, komplexere, aber besser geeignete Lösungen zum informellen Datenaustausch zu verwenden. Prof. Tilman zeigte persönliche, soziale und technologische Hinderungsgründe auf, die für viele Wissenschaftler dem freien Datenaustausch entgegenstehen. Als mögliche Anreize, um dieses Problem zu beheben, nannte er die gesteigerte Anerkennung von Datenpublikationen als wissenschaftliche Leistungen, die Einführung von Datenpublikationen als Voraussetzung für Publikationen generell, sowie mehrere Anreize, die auf das persönliche Umfeld der Forscher abzielen. Ein wichtiger Punkt sei auch die Nachvollziehbarkeit des zeitlichen Ablaufs sowohl für die Erstellung der Daten als auch für den Einsatz der Software. Er schloss seinen Vortrag mit dem Wunsch nach einfach verfügbaren und leicht auffindbaren Datensets, der Forderung nach wiederverwendbaren und rekonstruierbaren Datenprodukten und der Hoffnung auf bessere Anerkennung für die offene Weitergabe und das Publizieren von Daten.

Dr. Peter Bartelheimer (Universität Göttingen) betrachtete in seinem Vortrag die Datenbereitstellung und Virtuelle Forschungsumgebungen (VFU) als zwei Komponenten der IT-Forschungsinfrastruktur. Anhand des Projekts „Virtuelle Forschungsumgebung für die sozio-ökonomische Berichterstattung“ beschrieb er Funktionen der VFU in typischen Workflows der Datennutzung. Das Projekt befindet sich derzeit in der Einführungs- und Entwicklungsphase (2012 – 2014). Anwendungsfall für den Einsatz der VFU ist die Erstellung des Dritten Berichts zur sozio-ökonomischen Entwicklung in Deutschland. Die VFU wird projektbegleitend weiterentwickelt und operativ im Verbund der Projektpartner genutzt. Das Ziel ist die Bereitstellung einer nachhaltigen VFU-Plattform.

Dr. Bartelheimer befasste sich mit der Frage des Transfers der Daten in die VFU, sowie der Abläufe hinter der Datenschnittstelle. Besonderes Augenmerk galt hierbei dem Workflow in der VFU, sowie dem Transfer von Syntax- und Metadateien. Als entscheidend für die Akzeptanz einer VFU beschrieb Dr. Bartelheimer zum einen die Kontrollanforderungen an die datenhaltende Einrichtung, zum anderen die intuitive individuelle Anwendung für den Nutzer. Die VFU sollte eine einheitliche Arbeitsumgebung bieten und es erlauben, schrittweise mehr Daten zu integrieren. Die Forschungsdatenzentren sollten bei der Entwicklung von VFUs und der zugehörigen Nutzungsstudien in die technische und soziale Entwicklung integriert werden.

Reiner Mauer bot in seinem Vortrag einen Blick auf die Persistente Domäne und damit der dauerhaften Bereitstellung von Forschungsdaten. Der Vortrag begann mit einem Überblick über Forschungsdaten in den Sozialwissenschaften und den Aufgaben von GESIS als Infrastruktureinrichtung für die Sozialwissenschaften.

Die Langzeitarchivierung von Forschungsdaten ist ein explizites Ziel von GESIS und als solches in der Satzung verankert. Reiner Mauer positionierte die Aufgaben von GESIS im Kontext des Curation Continuum und beschrieb den Arbeitsablauf bei der Archivierung, der Datenaufbereitung und der Datenbereitstellung. Die Langzeitarchivierung dient der Substanzerhaltung, sowie dem Erhalt von Nutzbarkeit und Interpretierbarkeit. Er betonte, dass das Vertrauen der Nutzer in die dauerhaften Domäne von großer Bedeutung sei – Stichwort Zertifizierung. Herr Mauer erläuterte die Definition und die Aufgaben persistenter Identifikatoren und stellte einige Fragen zu deren Implementierung zur Diskussion.

Als Herausforderungen beschrieb Herr Mauer die immer stärker werdende Bedeutung strukturierender und kooperativer Dienste. Die Bedeutung komplexer Forschungsdesigns und neuer Datenformen wachse, so dass die LZA sich vor der Herausforderung sieht, mit immer komplexeren Forschungsdaten umzugehen. Desweiteren nimmt die Verlinkung bzw. Integration von Daten aus unterschiedlichen Datenquellen und Disziplinen zu. Dies resultiert in höhere Anforderungen an Daten und Metadaten. Rechtliche und Lizenzfragen gewinnen an Bedeutung. Bei der Qualität sind insbesondere Aktualisierung, Versionierung, und autoritative Versionen relevant.

Dr. Hans Pfeiffenberger behandelte in seinem Vortrag Portale und Best Practices und bot damit ein Beispiel für Datennutzung im Kontext der Öffentlichen Domäne. Nach einem blick auf Data Sharing in der Wissenschaft und dem geäußerten Wunsch, über einen singulären Zugangspunkt für möglichst alle wissenschaftlichen Datensätze zu verfügen, ging Dr. Pfeiffenberger näher auf die oftmals sehr unterschiedliche Herkunft der Daten und die Diversität der Datenformen ein. In mehreren Beispielen zeigte er die vielfältigen und verstreut existierenden Datenquellen am Beispiel von Forschungsschiffen und Bojen mit Meßfunktion. Weiterhin stellte er ein Modell zur Klassifizierung und Zuordnung der Daten zu Expeditionen, Forschern und Art der Daten vor. Wichtig sei, alle notwendigen Kontext- und Abstammungsinformationen mit den Daten zu erfassen.

Besonders hervor hob Dr. Pfeiffenberger die Notwendigkeit zur Reproduzierbarkeit der Daten. Er beschrieb den Weg von Rohdaten zu qualitätsgesicherten Primärdaten und daraus abgeleiteten Datenprodukten, und betonte die Notwendigkeit Datenakquise-Praktiken zu harmonisieren.

Workshop-Programm

Die Palette der angebotenen Workshops spiegelte die Bandbreite der von Forschungsdaten-Infrastrukturen angesprochenen Themen wieder. Das Spektrum der Themen reichte von Qualifizierung und Policies über eher anwendungsorientierte Themen wie Software-Lösungen und Datenmanagement bis hin zu rechtlichen Rahmenbedingungen. Die Workshops organisierten sich in drei Themenblöcken „Realisierung“, „Daten“ und „Rahmenbedingungen“. Die Workshops waren organisiert als einstündige Diskussionsrunden mit dem Ziel, offensichtliche Lücken (konzeptionell, organisatorisch, technisch) zu identifizieren, vorbildliche Lösungen aufzuzeigen und eine Vision für das Jahr 2020 (Horizon 2020) zu entwerfen.



Abb. 6: Workshop Impression

Workshop WS 1 – Qualifizierung

Der Workshop wurde eröffnet durch ein Impulsreferat von **Maxi Kindling (Humboldt-Universität zu Berlin)**. In ihrem Impulsreferat erläuterte sie die Integration von Datenmanagement als Thema in fachwissenschaftlichen Studiengängen und gab einen Einblick in die Aus- und Weiterbildungsaktivitäten am Institut für Bibliotheks- und Informationswissenschaft (IBI). Die Workshop-Teilnehmer diskutierten anschließend über die Schaffung eigener Studiengänge für das Datenmanagement, sowie die Integration des Datenmanagements in bibliotheks- und informationswissenschaftliche Studiengänge unter Berücksichtigung des Praxisbezugs. Die Teilnehmer kamen überein, dass die Breite des Qualifizierungsangebots insgesamt erweitert werden müsse. Eine Qualifizierung solle auf den Ebenen des Studiums, des weiterführenden Studiums, sowie als Weiterbildung außerhalb der Universitäten und Fachhochschulen angeboten werden. Der Workshop wurde moderiert von **Prof. Dr. Peter Schirmbacher (Humboldt-Universität zu Berlin)**.

Workshop WS 2A – Datenmanagement

Der Workshop begann mit einem Impulsreferat vom Moderator **Dr. Jens Nieschulze (Universität Göttingen)**. Schwerpunkte der anschließenden Diskussionen waren vor allem Hemmnisse und Anreize für ein gutes Datenmanagement sowie die Nachnutzung der Daten.

Als Hemmnisse wurden unter anderem der erhöhte Zeitaufwand bei der Pflege der Daten für Archivierung und Nachnutzung genannt, die Ungewissheit ob die Daten überhaupt nachgenutzt werden und sich der Aufwand deshalb lohnt, die Form der Daten, welche oftmals nur in analogem Format vorlägen (besonders die Metadaten in Form eines Laborbuches), sowie Geheimhaltungswünsche und -pflichten. Als Anreize wurden die Richtlinien der Forschungsförderer sowie ein „Impaktfaktor“ durch die Publikation von Daten genannt. Allerdings wird ein hoher Nutzen für den Mehraufwand erwartet. Auch erfolgreiche Beispiele werden als ein wichtiger Anreiz gesehen. Lösungen zum Datenmanagement sollten sich nahtlos in die Arbeitsumgebungen der Wissenschaftler integrieren und nachhaltig sein, da sie immer weiterentwickelt werden müssen.

Für die Zukunft wurde der Wunsch nach Datenzentren geäußert, die eng mit bestimmten Communities zusammenarbeiten, miteinander vernetzt sind und eine Weiterbildung im Datenmanagement anbieten. Die Universitäten müssen mehr Routine entwickeln, um mit solchen Datenzentren zusammenzuarbeiten - derzeit werden Datenzentren eher als Zusatzlast empfunden.

Workshop WS 2B – Datenmanagement

Wie kann das Datenmanagement in der privaten bzw. ruppen-Domäne (nach Treloar) besser unterstützt werden? Wie können hier Angebote über punktuelle Lösungen hinaus fruchtbar gemacht und besser koordiniert werden? Welche Aspekte gilt es hier über den der technischen Infrastruktur hinaus zu beachten? Diese Fragen standen im Mittelpunkt des von **Dr. Harry Enke (Leibniz-Institut für Astrophysik, AIP)** moderierten Workshops „Datenmanagement“. Kontrovers diskutiert wurde der Ansatz, Meß- und Beobachtungsdaten bereits bei der Entstehung zu speichern, ohne einen korrigierenden oder erläuternden Eingriff durch den Wissenschaftler abzuwarten. Der Moderator betonte, dass Forschungsdatenmanagement die Begleitung des gesamten Prozesses wissenschaftlichen Arbeitens umfassen müsse und in Datenmanagementplänen festgehalten werden müsse. Auch das Formulieren solcher Pläne könne einen kulturellen Wandel in den Disziplinen beschleunigen, die bislang dem Management der Daten noch nicht die nötige Bedeutung zubilligten.

Workshop WS 3 – Policies

Diskutiert wurde die Vielschichtigkeit von Policies - einerseits als Policies praktischer Natur, welche Handlungsempfehlungen darstellen, andererseits als Policies politischer Natur, welche eine Zielsetzung definieren. Betrachtung in diesem Workshop fanden Policies politischer Natur. In vielen Erklärungen wird davon ausgegangen, dass der Wille zur Wiederverwendung von Forschungsdaten vorhanden ist. Ganz eindeutig ist dies jedoch nicht geklärt, da sich die Gemeinschaft der potentiellen Nachnutzer in zwei Gruppen teilt - den eher observations-orientiert arbeitenden Disziplinen, bei denen die Nachnutzung eine besondere Rolle spielt und den eher experimentell arbeitenden Wissenschaften, bei denen die Beschreibung des Verfahrens eine große Bedeutung hat, um Ergebnisse nachvollziehen zu können. Diskutiert wurde die Frage, ob es legitim sei, eine Veröffentlichung zu fordern und ob Daten in größerem Maß fachübergreifend nachgenutzt werden, da es keine wissenschaftlichen Untersuchungen zu dem Problem gäbe. Eine Nachnutzung von Daten bedeutet immer auch deren Aufbereitung und die Kuratierung der zugehörigen Metadaten - ein teurer Prozess, der zur Differenzierung der Relevanz der Datensätze führt. Der Workshop wurde moderiert von **Dr. Christoph Bruch (Helmholtz Open Access Koordinationsbüro)**.

Workshop WS 4 – Software Lösungen

Der von **Dr. Robert Hauser (FIZ-Karlsruhe)** moderierte Workshop thematisierte Software-Lösungen für Forschungsdateninfrastrukturen (FDI). Bei einer Vorstellung der Teilnehmer zeigte sich, dass viele Institutionen sich gerade in der Evaluationsphase für Software für Forschungsdateninfrastrukturen befinden. In der weiteren Diskussion bestätigte sich die Annahme, dass es keine generelle einzelne Softwarelösung gibt, welche die Anforderungen der unterschiedlichen Fachrichtungen erfüllen kann.

Erkenntnisse des Projektes re3data.org zu existierende Forschungsdatenrepositorien zeigen, dass eine Vielzahl von Lösungen im Einsatz sind, meist Eigenentwicklungen oder Weiterentwicklungen existierender Software. Oft wird der Aufwand von Eigenentwicklungen geringer angesehen als die Integration und Anpassung existierender Lösungen. Wichtige Punkte bei der Entwicklung eigener Software-Lösungen sollten die Zusammenarbeit von Wissenschaftlern und Infrastrukturbetreibern, die Abbildung des gesamten Forschungsprozesses und wissenschaftlichen Workflows in der Software, sowie die Nachhaltigkeit sein. Desweiteren wurde die allgemeine Situation der Entwicklung von Softwarelösungen in Deutschland besprochen. Die Teilnehmer schlugen den Aufbau eines nationalen und zentral gesteuerten Netzwerks von Forschungsdatenzentren vor. Es wurde festgestellt, dass die derzeitigen Akteure und Initiativen (DINI, Wissenschaftsorganisationen, etc.) sich zwar in Netzwerken sammeln, jedoch nicht koordiniert agieren.

Bestehende Initiativen in Deutschland sollten gebündelt werden, um durch gemeinsame Veranstaltungen mit Vorträgen und Experten-Workshops den Austausch zu fördern.

Workshop WS 5A – Datenpublikationen

Der Workshop begann mit einer Präsentation von Moderator **Dr. Michael Lautenschlager (Deutsches Klimarechenzentrum)**, welche den Lebenszyklus der Daten in der Langzeitarchivierung, sowie Daten als Teil virtueller Forschungsumgebungen beschrieb. Laut Moderator gehe der Trend zur Anonymisierung von Datenproduzent und Nutzer. Ziel sei ein interdisziplinäres, verteiltes, selbstbeschreibendes, qualitätsgeprüftes Langzeitdatenarchiv, um Daten ohne Kontakt zum Datenautor nachnutzen zu können. Nachfolgend wurde vor allem über die Thematik der Qualitätssicherung von Daten, sowie der Zertifizierung von Archiven und Prozessen diskutiert.

Workshop WS 5B – Datenpublikationen

Diskutiert wurden die vielen Gesichter der Datenpublikationen, das Thema Qualitätssicherung und Visionen für die Zukunft der Datenpublikationen im Jahr 2020. Gefordert wurden Lektorate für Forschungsdaten, da sich die Bewertung und Qualitätssicherung bislang recht schwierig darstellt. Kritisiert wurden bislang eher unbefriedigende, recht technisch ausgelegte Frameworks für Qualitätskriterien. Desweiteren müsse die Nutzbarkeit und Art der Metadaten bedacht werden. Das Publizieren von Daten wurde als kompliziert bewertet. Moderiert wurde der Workshop von **Roland Bertelmann (Deutsches GeoForschungsZentrum GFZ)**.

Workshop WS 6 – Kosten

Fragestellungen dieses Workshops waren die Kostenermittlung für die Langzeitarchivierung von Forschungsdaten und wie Kosten für das Datenmanagement für einen Projektantrag kalkuliert werden sollten. Desweiteren wurde über potentielle Kostenträger, Preismodelle und noch fehlende Richtlinien zum Datenmanagement für die Förderung von Forschungsaufträgen diskutiert. Moderiert wurde der Workshop von **Dr. Henk Harmsen (DANS)**.

Workshop WS 7 – Persistente Identifier

Der Workshop, moderiert von **Dr. Jens Klump (GFZ)**, wurde durch ein Impulsreferat von **Dr. Janna Neumann (TIB Hannover)** eröffnet. Die Teilnehmer widmeten sich zunächst der Terminologie zur Beschreibung von Identifikatoren. Die Definition des Begriffs Persistente Identifier (PID) ist ausschlaggebend für die Wahl der Anwendungsfelder und die Auswahl des darunter liegenden technischen Dienstes. Die nachfolgende Diskussion gliederte sich in drei Hauptthemen: Die Wahl des Anbieters eines PID-Systems, die Anwendungsfälle von Systemen im Kontext von Forschungsdaten, und die Probleme der Interoperabilität zwischen Systemen und Anbietern.

Eine Klärung der Begriffe war notwendig, um die Frage zu diskutieren, wie noch nicht veröffentlichte Objekte eindeutig referenziert werden, welches System dafür geeignet und ob es überhaupt notwendig ist, für verschiedene Phasen des Lebenszyklus von Forschungsdaten unterschiedliche PID-Systeme zu verwenden. Hinter dieser Frage stehen bisher ungeklärte Anforderungen an Qualität, Vertrauenswürdigkeit und Kontrollierbarkeit der Systeme. Aber auch bei praktischen Fragen der Identität, Granularität und Versionierung digitaler Forschungsobjekte wäre eine Vereinheitlichung der Begriffe hilfreich und würde zur konzeptionellen Klarheit beitragen.

Ein aktuelles Beispiel für unterschiedliche Ansätze für den Betrieb eines PID-Systems ist die Identifikation von Autoren. Hier stehen heute mehrere Systeme zur Auswahl, die teilweise miteinander interoperabel sind. Der Hauptunterschied zwischen den verschiedenen Systemen zur Identifikation von Autoren liegt jedoch in ihrer Organisation und in der Rolle der Verlage und Bibliotheken.

Workshop WS 8A – Daten in Virtuellen Forschungsumgebungen

In diesem Workshop, moderiert von **Dr. Heike Neuroth (Niedersächsischen Staats- und Universitätsbibliothek Göttingen)**, wurden zunächst Definitionen von VFUs besprochen und anschließend Probleme zum Thema VFU aus einer vorgestellten Liste zu priorisiert und Lösungswege besprochen. Im Detail diskutiert wurden der Ablauf des Dialogs zwischen (Fach-)Wissenschaftler und VFU-Betreiber, der Gegensatz von „Einfache Nutzung vs. Alleskönner“, sowie Möglichkeiten zur Sicherstellung der Nachhaltigkeit. Für die Zukunft von VFUs wurde die Wunschvorstellung geäußert, alle Arbeitsabläufe der Forschung so zu integrieren, dass der Wechsel zwischen normaler Arbeitsumgebung und VFU für den Nutzer praktisch nicht merkbar ist. Dafür müsste allerdings die Benutzerfreundlichkeit der VFUs erheblich verbessert werden.

Workshop WS 8B – Daten in Virtuellen Forschungsumgebungen

Zu Beginn des Workshops wurde gemeinsam versucht, Begriffsbestimmungen für Forschungsdaten und Virtuelle Forschungsumgebungen (VFU) zu finden. Die Teilnehmerinnen und Teilnehmer waren sich einig, dass für die Kategorisierung von Forschungsdaten vor allem praktische Gesichtspunkte den Ausschlag geben und dass deren Form, Format und Umfang je nach Wissenschaftsdisziplin sehr unterschiedlich sein können.

Virtuelle Forschungsumgebungen stellen komplexe Systeme dar, die disziplinübergreifend oder disziplinspezifisch verschiedene Tätigkeiten des wissenschaftlichen Arbeitsprozesses abbilden. Dabei sind sie meist kein abgeschlossenes Produkt, sondern ein Framework oder Baukastensystem, das zur Unterstützung der Gesamtheit oder eines Teils des Wissenschaftsprozesses dient. In der weiteren Diskussion wurde die Abbildung der Vorlieben und Bedürfnisse der Nutzer als eine besonders wichtige Anforderung an virtuelle Forschungsumgebungen betont. Bei der institutsübergreifenden Nutzung virtueller Forschungsumgebungen ist es wichtig, wer Zugriff auf die enthaltenen Forschungsdaten haben soll, wer dies festlegt und wie der Zugang organisatorisch und technisch geregelt wird. Wenn Daten öffentlich zugänglich gemacht werden sollen, können Sicherheitsaspekte eine Rolle spielen. Virtuelle Forschungsumgebungen sollten als Dienst flexibel und anpassungsfähig sein um eine möglichst einfache Zusammenarbeit bei der Nutzung der Inhalte zu erreichen. Das Vertrauen der Nutzerinnen und Nutzer ist eine wichtige Voraussetzung für die Annahme virtueller Forschungsumgebungen als Arbeitsmittel und damit auch für die Nachhaltigkeit der entstehenden Strukturen. Bei der Finanzierung müssen gerechte Lösungen für die Umlage der Kosten gefunden werden.

Zusammenfassend wurde festgestellt, dass virtuelle Forschungsumgebungen auch in Zukunft als Abbild des bisherigen Wissenschaftsprozesses ihre Berechtigung haben werden, wenn der Mensch als Nutzer im Mittelpunkt steht. Der Workshop wurde moderiert von **Prof. Dr. Peter Schirnbacher (Humboldt-Universität zu Berlin)**.

Workshop WS 9 – Rechtliche Rahmenbedingungen

Der Umgang mit digitalen Forschungsdaten wird häufig von rechtlichen Fragestellungen berührt. Vor diesem Hintergrund widmete sich der Workshop „Rechtliche Rahmenbedingungen“ der juristischen Dimension der Forschungsdaten-Infrastruktur. Moderiert wurde der Workshop von **John H. Weitzmann (Creative**

Commons). Im Kern der Diskussion standen Fragen wie „Wem gehören (meine) Daten?“ und angrenzende Themen wie die Haftung für Folgeschäden bei „falschen Daten“. Der spannende Workshop zeigte, dass noch viele Fragen ungeklärt sind. Aus Sicht der Wissenschaft scheint der Wunsch nach rechtsicheren Verfahren im Umgang mit Forschungsdaten höchste Priorität zu haben. Zum Abschluss der Diskussion wurde noch ein Blick auf den gesetzgeberischen Handlungsbedarf geworfen. Hier wurde u.a. festgestellt, dass bei den bisherigen Diskussionen über eine weitergehende Wissenschaftsschranke der Themenkomplex Forschungsdaten noch wenig Berücksichtigung fand.



Abb. 7: Die Teilnehmer des FDI 2013-Symposiums

Fazit

In seinem Schlusswort fasste **Dr. Jens Klump (GFZ)** die Ergebnisse des Symposiums zusammen. Mit etwa 140 Teilnehmern, 10 Vorträgen und 12 Workshops war es den einladenden Projekten gelungen, eine große Zahl an Experten zu einem intensiven und anregenden Symposium zusammenzubringen. In den Workshops und am Rande der Veranstaltung konnten die Teilnehmer neue Kontakte knüpfen und so die Vernetzung der Akteure verbessert werden. Die große Zahl der Teilnehmer zeigt, dass das Thema Forschungsdaten als wichtiges Thema wahrgenommen wird. In den Diskussionen war zu beobachten, dass die Entwicklung von Forschungsdaten-Infrastrukturen zwar nach wie vor heterogen verläuft, und auch viele Fragen auch weiterhin unbeantwortet bleiben. Die Diskussion zeigte jedoch auch, dass das Thema Forschungsdaten-Infrastrukturen eine konzeptionelle Reife erreicht hat, die mit Konzepten in anderen europäischen Staaten, den USA oder Australien vergleichbar ist.

Insbesondere in der technischen und der konzeptionellen Entwicklung ist eine gewisse Konvergenz zu beobachten. Daneben zeigen sich in anderen Bereichen konzeptionelle Unschärfen, bei denen noch Forschungsbedarf besteht, so z.B. bei den Themen „Qualität“, „Vertrauen“ und „Kosten- und Preismodelle“. Lücken bestehen auch noch bei Datenmanagement-Werkzeugen und deren Integration in die Arbeitsabläufe der Wissenschaftler. Die Teilnehmer bemängelten, dass viele Akteure im Datenmanagement nicht weit genug bekannt seien und auch die Vernetzung innerhalb der Community optimiert werden könnte. Auch sei der Bedarf an Angeboten in den Bereichen Qualifizierung und Beratung besonders ausgeprägt.

Auch bei dieser Veranstaltung wurde deutlich, dass eine Verbesserung des Umgangs mit Forschungsdaten nicht nur eine Frage der unterstützenden technischen Infrastrukturen ist, sondern auch einen kulturellen

Wandel in der Wissenschaft erfordert. Der kulturelle Wandel im Hinblick auf Forschungsdaten bewegt sich in Richtung eines offeneren Umgangs mit diesem Teil der wissenschaftlichen Überlieferung, wie es im Bericht „*Science as an Open Enterprise*“ der Royal Society (2012) vorgeschlagen wird. Der Grad der Offenheit wird bestimmt durch die Spannung zwischen Vertrauen in die „Peers“ und Kontrolle über das eigene „Werk“. Neben den sozialen Normen muss auch noch der rechtliche Rahmen für Forschungsdaten weiterentwickelt werden.

Für eine Verbesserung der Situation ist es notwendig, die Erstellung von Daten, Software und Infrastrukturen neben den bisher üblichen Literaturveröffentlichungen als Beitrag im wissenschaftlichen Wertesystem zu verankern. Dies bedingt, dass der Nutzen einer Forschungsdateninfrastruktur für Forscher offensichtlich sein muss. Datenpolicies und die Berücksichtigung dieser Leistungen in den institutionellen Bewertungssystemen können diesen Wandel unterstützen.

Positiv bewerteten die Teilnehmer die Möglichkeit, neue und vorbildliche Lösungen für den Umgang mit Forschungsdaten kennenzulernen. Einige Projekte und vielversprechende Ergebnisse wurden vorgestellt. Durch Evaluation der Ergebnisse und Erfahrungsaustausch kann hier eine Entwicklung angestoßen werden, welche über das Stadium des Experimentierens hinausgeht. Wichtig ist, dass der Austausch zwischen den Akteuren weitergeht und auch Praxisvermittlung einschließt, um den Kreis der Akteure zu erweitern.

Sowohl die Organisatoren als auch die Teilnehmer des Symposiums bewerteten die Veranstaltung als äußerst hilfreich für den Austausch von Ideen, der Generierung neuer Impulse und für die Vernetzung der Akteure untereinander. Über Folgeveranstaltungen wird bereits nachgedacht.



Danksagungen

Wir möchten uns herzlich bei unseren Vortragenden, Workshop-Moderatoren und den fast 140 Teilnehmern bedanken, die mit ihren abwechslungsreichen und interessanten Beiträgen unsere Veranstaltung bereichert und neue Impulse für die Weiterentwicklung von Forschungsdaten-Infrastrukturen gegeben haben. Insbesondere danken wir den Protokollanten der Workshops für ihre Mitschriften:

Dr. Torsten Rathmann, Dr. Daniela Koudela, Dr. Beate Rusch, Damian Ulbricht, Paul Vierkant, Vivien Hollad, Dr. Peter Löwe, Britta Dreyer, Markus Schnalke, Dr. Wolfgang Peters-Kottig, Dr. Paul Schulze-Motel und Heinz Pampel.

Lizenz



Alle Texte dieser Veröffentlichung, ausgenommen Zitate, sind unter einem Creative Commons Namensnennung 3.0 Deutschland Lizenzvertrag lizenziert: <http://creativecommons.org/licenses/by/3.0/de>

Zitationsvorschlag

Schäfer, Leonie; Pampel, Heinz; Klump, Jens; Häsler, Tim (2013), Bericht Symposium „Forschungsdaten-Infrastrukturen (FDI 2013)“, Workshop Report, Helmholtz-Zentrum Potsdam Deutsches GeoForschungs-Zentrum, Potsdam, Germany. [online] Available from: http://dx.doi.org/10.2312/RADIESCHEN_003

Agenda

Agenda Symposium Forschungsdaten-Infrastruktur (FDI 2013) - Dienstag, 22.01.2013			
	Hörsaal		
9:30 - 10:00	Registration & Kaffee		
10:00 - 10:30	Eröffnung durch den Vorstand Prof. Dr. Dr. h.c. Reinhard F. J. Hüttl des GFZ - Geleitwort Dr. Winkler-Nees (DFG) - Vorstellung des Tagesablaufs		
10:30 - 11:00	Keynote "Policy" Dr. Torsten Reimer (JISC)		
	Hörsaal	Tagungsraum 1	Tagungsraum 2
11:00-12:00	Session 1 Projektpräsentationen "EWIG" Tim Hasler "Radieschen" Jochen Klar(AIP) "KOMFOR" Dr. Michael Diepenbroek "re3data.org" Frank Scholze (KIT) und "Bokela" Dr. H.-J. Wallrabe-Adams	Workshop 1 "Qualifizierung" Prof. Michael Seadle (HUB)	Workshop 2 "Datenmanagement inkl. Datenpublikationen" ⇒ WS 2A: Dr. Jens Nieschulze (Uni Göttingen) ⇒ WS 2B: Dr. Harry Enke (AIP)*
12:00 - 12:15	Kaffee		
12:15 - 13:15	Session 2 Vortrag 1: "Private Domäne 'Die Sicht des Forschers'" Prof. Frederik Tilmann (GFZ) Vortrag 2: "Gruppendomäne 'VFUs'" Dr. Peter Bartelheimer (Uni Göttingen)	Workshop 4 "Software Lösungen" Dr. Robert Hauser (FIZ)	Workshop 5 "Daten-Publikationen + QS" ⇒ WS 5A: Dr. Michael Lautenschlager (DKRZ) ⇒ WS 5B: Roland Bertelmann (GFZ)*
13:15 - 14:15	Mittagessen - Kantine		
14:15 - 15:15	Session 3 Persistente Domäne Vortrag 1: "PID, DLZA, Zertifikate" R.Mauer, GESIS Vortrag 2: "Zugangsdomäne Portale, Best Practices" Dr. Hans Pfeiffenberger (AWI)	Workshop 7 "Persistente Identifier" Dr. Jens Klump (GFZ) / Dr. Janna Neumann (TIB)	Workshop 8 "Daten in VFUs" ⇒ WS 8A: Dr. Heike Neuroth (SUB) ⇒ WS 8B: Prof. Peter Schirmbacher*
15:15 - 15:30	Kaffee		
	Hörsaal		
15:30-16:30	Zusammenfassende Berichte aus den Workshops		
16:30 - 17:00	Diskussion und Zusammenfassung - Dr. Jens Klump (GFZ)		

* Raum A 17/10.27 Treffpunkt am Empfang 11:00/12:10/14:10 Uhr

Vortragsprogramm (Materialien)

Keynote – Dr. Torsten Reimer (JISC)

JISC 22/01/2013
Symposium Forschungsdaten-Infrastrukturen, Potsdam

Dr Torsten Reimer
Programme Manager, Digital Infrastructure, Jisc

Folie 1

JISC Benefits – Jisc’s “Keeping Research Data Safe”

Dimension 1	Direct Benefits	Indirect Benefits (Costs Avoided)
	<ul style="list-style-type: none"> -New research opportunities -Scholarly communication/access to data -Re-purposing and re-use of data -Increasing research productivity -Stimulating new networks/collaborations -Knowledge transfer to industry - Increasing skills base of researchers/students/staff -Increasing productivity/economic growth -Verification of research/research integrity -Fulfilling mandate(s) 	<ul style="list-style-type: none"> -No re-creation of data -No loss of future research opportunities -Lower future preservation costs -Re-purposing data for new audiences -Re-purposing methodologies -Use by new audiences -Protecting returns on earlier investments
Dimension 2	Near-Term Benefits	Long-Term Benefits
	<ul style="list-style-type: none"> -Value to current researcher & students -No data lost from Post Doc turnover -Short-term re-use of well curated data -Secure storage for data intensive research -Availability of data underpinning journal articles 	<ul style="list-style-type: none"> -Secures value to future researchers & students -Adds value over time as collection grows and develops critical mass -Planned management from an early stage in the research lifecycle is ultimately more cost-effective than late intervention (providing proper selection of what to keep is done)
Dimension 3	Private Benefits	Public Benefits
	<ul style="list-style-type: none"> -Benefits to sponsor/funder of research/archive -Benefits to researcher -Fulfill grant obligations -Increased visibility/citation -Commercialising research 	<ul style="list-style-type: none"> -Input for future research -Motivating new research -Catalysing new companies and high skills employment

http://www.beagrie.com/kRDS_Factsheet_0910.pdf

Folie 2

- We drive innovation in UK education and research, and have been doing so for more than 15 years.
- We are a registered charity and work on behalf of UK higher education, further education and skills to champion the use of digital technologies.
- Jisc offerings and services include:
 - Janet, the UK's high speed network for education and research
 - Jisc Collections, providing digital content for education and research
 - (Advisory) services on the use of digital technology
 - An innovation programme to help our customers make the most of new technologies



Folie 3



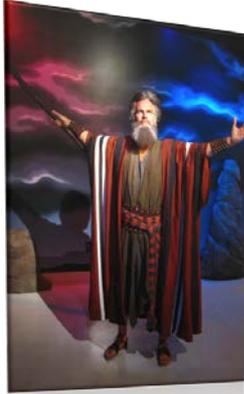
<http://www.flickr.com/photos/lorenjavier/5686207364/> CC BY ND 2.0

Folie 4

- "In ten years, it is expected that..."
 - Leonardo Candela (2011): *Virtual Research Environments*. In: *Technological & Organisational Aspects of a Global Research Data Infrastructure - A View from Experts*.
- http://www.grdi2020.eu/Pages/SelectedDocument.aspx?id_documento=eb0e8fea-c496-45b7-a0c5-831b90fe0045



Folie 5



In ten years, it is expected that the trend characterizing science and scientific collaborations discussed above continues, thus becoming the “default” approach for scientific investigations as well as for any societal collaboration-based activity.

Virtual Research Environments will be **integrated into standard practices** and tools used by communities of practice.



The creation and management of Virtual Research Environments **will be a very straightforward process** that relies on specific services – VRE Management Services – **built atop a “global virtual infrastructure”**.

The VRE definition phase will guide an authorised actor of an application domain in **characterizing the expected VRE** service in very abstract terms.

The VRE **deployment phase will be almost automatic**.

The VRE **monitoring / maintenance phase will require little direct human control**.

The resulting Virtual Research Environment will be **very flexible and customizable**.

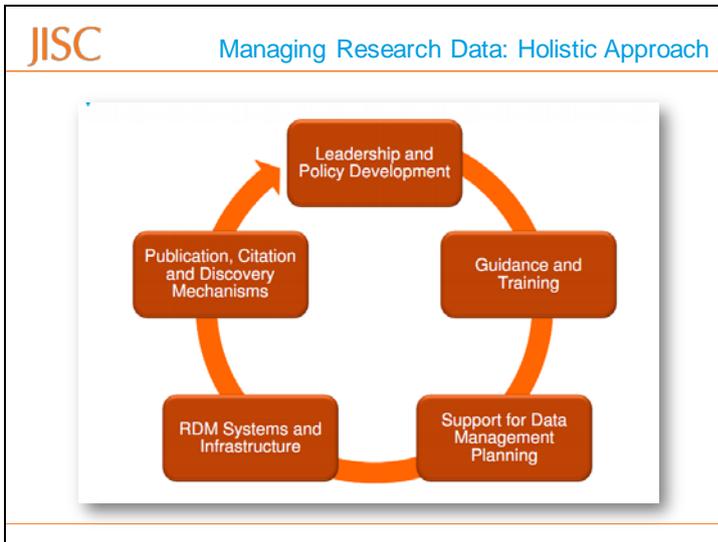


Thus Virtual Research Environments creation and management **will become a societal and organisational process rather than a technological one**.

Fundamentally, the most important point to have emerged from our study is that VREs need to be conceptualised as **community building projects rather than technology projects**. [...] By far the most important challenge faced by VREs is **sustainability**.

Carusi / Reimer: VRE Landscape Study
<http://www.jisc.ac.uk/media/documents/publications/vrelandscapeport.pdf>





Folie 9

JISC Managing Research Data at Jisc

Understanding the problem (pre-2007-2009)

Prototyping solutions (2009-11)

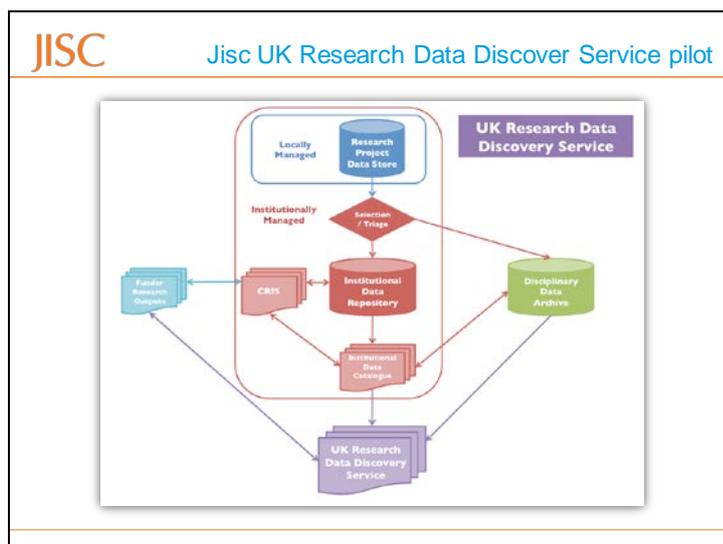
Hardening solutions and building institutional capacity (2011-13)

Developing elements of national infrastructure (2013+)

- MRD Phase 1, 2009-2011
 - Infrastructure, Planning, Support, Citation and Publishing, Training
- MRD Phase 2, 2011-2013
 - institutional projects infrastructures and policies;
 - disciplinary best practice, implementation of data management plans
 - customised data management planning tools for institutional use
 - innovative data publication
 - disciplinary training materials

slide 10

Folie 10



Folie 11

● Full Coverage ● Partial Coverage ○ No Coverage

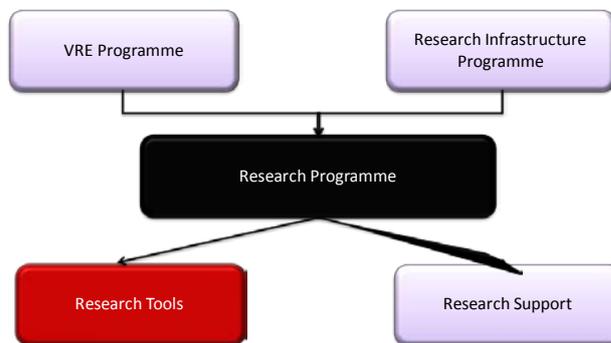
Research Funders	Policy Coverage			Policy Stipulations				Support Provided			
	Published outputs	Data	Time limits	Data plan	Access/sharing	Long-term curation	Monitoring	Guidance	Repository	Data centre	Costs
AHRC	●	●	●	●	●	●	○	●	○	○	○
BBSRC	●	●	●	●	●	●	●	●	●	●	●
CRUK	●	●	●	●	●	●	●	○	●	○	○
EPSRC	●	●	●	○	●	●	●	○	○	○	○
ESRC	●	●	●	●	●	●	●	●	●	●	○
MRC	●	●	●	●	●	●	○	○	●	○	○
NERC	●	●	●	●	●	●	●	●	●	●	○
STFC	●	●	●	●	●	●	●	○	●	○	○
Wellcome Trust	●	●	●	●	●	●	●	●	●	○	○

<http://www.dcc.ac.uk/resources/policy-and-legal/overview-funders-data-policies>

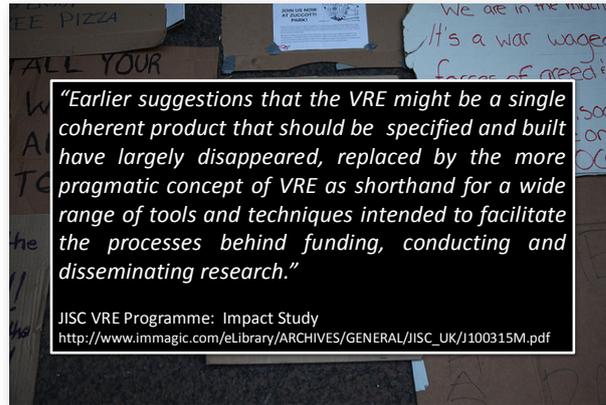
Folie 12

VRE1	VRE2	VRE3	VRE RI
2004-2007	2007-2009	2009-2011	2010
15 Projects	4 Demonstrators	10 Projects	14 Projects
Technology Focused	User and Research Practice focused	Broadening Use, across institutions and disciplines	Rapid Innovation throughout Research Lifecycle
Experimental	Developmental	Embedding	Rapid Innovation
Diverse design and developmental approached	Unified design and development models	Diverse design – community and challenge driven	'Scratching itches', solution driven
Standalone solutions	Integrated pilots	Tools, frameworks and interoperability	Technical Solutions, Business Community Engagement
Collaboration			
Large and small scale research			
Single and Multi-disciplinary Research			

Folie 13



Folie 14



<http://www.flickr.com/photos/jimkiernan/6168587606/>

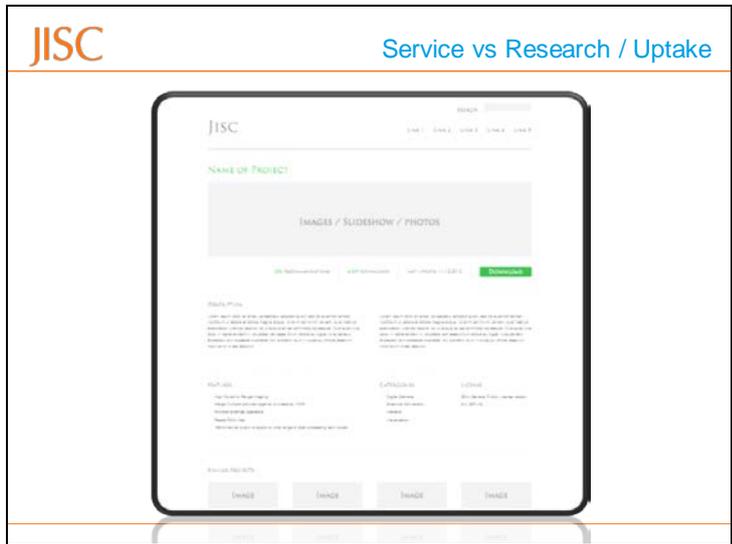
Folie 15

VREs still in operation?	Red
Uptake of code beyond the project?	Yellow
Interoperable with other developments?	Orange
New collaborations?	Yellow
New research facilitated?	Yellow
Sustainable technology/code?	Orange
Sustainable business model?	Red
Skills/capacity gained?	Green
Lessons learned?	Green

Folie 16

Business Models	Service vs Research	Cost	Staff Roles, Careers
Support	Skills	Uptake	Institutional Policies
Interoperability	Law	Students	Usability

Folie 17



Folie 18



Folie 19

JISC Law

- Jisc study into *The Value and Benefits of Text Mining*
- “the market intervention of copyright, originally intended to protect creative producers, is itself becoming a barrier to new creative production and may be inhibiting new knowledge discovery and innovation”

<http://www.jisc.ac.uk/publications/reports/2012/value-and-benefits-of-text-mining.aspx>

Folie 20

JISC Institutional Policies

<http://www.flickr.com/photos/rpenalozan/5367289745/> CC BY NC SA 2.0

Folie 24

JISC Usability

Usability Service Enhancements to Digimap

Usability testing should be an important part of the development of any user interface. Ensuring that the interface is intuitive and easy to use is critical for its success. However, running usability sessions with real users often strikes fear into project teams. They assume that it will be a costly and time-consuming process and will confuse as much as it clarifies the design process. This article aims to demonstrate how easy it is to set up an effective usability lab on a shoestring budget.

Usability lab on a shoestring budget

Photo 4 on January 27, 2012 by Sarah Ebbels

Search by date
January 2012

W	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

<http://used.blogs.edina.ac.uk/2012/01/27/usability-lab-on-a-shoe-string-budget/>

Folie 25

JISC Business Models

JISC ITT: Increasing the capacity of the JISC community to assess the economic costs and benefits of digital infrastructure innovation

JISC invites tenders for work to improve the capacity of the sector, especially those closely associated with JISC-funded innovation work in the area of Digital Infrastructure, to assess the economic costs and benefits of innovation work in these areas, to improve programmatic and project design, management and evaluation.

The purpose of the work is to enable JISC and those with whom it works to be better able to capture and use data that provides evidence of the economic impact of JISC-funded innovation work on digital infrastructure. One part of this will be to develop shared understanding of this evidence between the JISC Executive and those working within JISC innovation programmes.

Summary

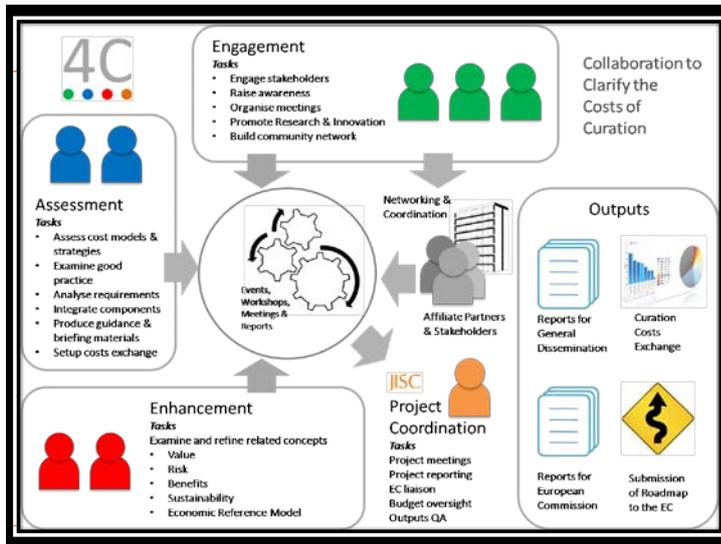
Submission Deadline
16 December 2011 12:00

Topic
Data & Text Mining

Evaluation
Research & Innovation

Strategic Themes
Information Environment

Folie 26



Folie 27



Folie 28

Session 1 - Projektpräsentationen

DFG-Projekt EWIG – Tim Hasler (ZIB)




EWIG
„Entwicklung von Workflowkomponenten für die
Langzeitarchivierung von Forschungsdaten im Bereich
Erd- und Umweltwissenschaften“

Symposium Forschungsdaten-Infrastrukturen

21. Januar 2013

Tim Hasler

Gefördert durch die 

Folie 1




EWIG Symposium Forschungsdaten-Infrastrukturen

Drei Komponenten:

- Entwicklung von institutionellen **Policies**
- ‚**Toolbaukasten**‘ zur technischen Qualitätssicherung (im Ingest) als Webservice
- Entwicklung von **Lehrveranstaltungen**

Gefördert durch die 

Folie 2

EWIG Symposium Forschungsdaten-Infrastrukturen

GFZ
Helmholtz-Zentrum
POTSDAM

ZITB
Freie Universität
Berlin

Policies – Tools – Lehre



Definition von Workflows entlang der Treloar Domänen durch zwei Projektpartner mit unterschiedlich organisiertem Datenmanagement

Schwerpunkt Qualitätssicherung und definierte Erfassung von Metadaten

Begleitung der Policyentwicklung in Lehrveranstaltungen

Gefördert durch die **DFG**

Folie 3

EWIG Symposium Forschungsdaten-Infrastrukturen

GFZ
Helmholtz-Zentrum
POTSDAM

ZITB
Freie Universität
Berlin

Policies – **Tools** – Lehre



„Vermenschlichung“ vorhandener Tool im Pre-Ingest – Augenmaß – Anspruchsvoll gegenüber sich selbst aber milde gegenüber dem was von außen kommt

Zugangserleichterung über ein Bewertungssystem (5-star rating)

Angebot, verschiedene Tools als Webservice zu testen

Gefördert durch die **DFG**

Folie 4

EWIG Symposium Forschungsdaten-Infrastrukturen

GFZ
Helmholtz-Zentrum
POTSDAM

ZITB
Freie Universität
Berlin

Policies – Tools – **Lehre**



Veranstaltung im Bachelorstudiengang Meteorologie im Modul Statistik WS 12/13 mit dem Schwerpunkt „Datenmanagement“

Fachübergreifende Veranstaltungen um auch die Postgraduates zu erreichen

Veranstaltung Datenmanagement auf den PhD days im GFZ

Verwendung „künstlich gealterter“ DIPs als „worst case scenario“

Gefördert durch die **DFG**

Folie 5

Vielen Dank ...

ewig.gfz-potsdam.de

Tim Hasler
Konrad-Zuse-Zentrum
für Informationstechnik Berlin (ZIB)
hasler@zib.de

DFG Projekt Radieschen – Dr. Jochen Klar (AIP)

Rahmenbedingungen einer disziplinübergreifenden Forschungsdateninfrastruktur

(Radieschen)

Jochen Klar
22.1.2013
FDI 2013



Hintergrund

- Es besteht Konsens über die herausragende Bedeutung des Umgang mit Forschungsdaten
- Diverse Initiativen resultierten in einer Vielzahl von verschiedensten Projekten
- Anforderungen zum Teil projekt- oder disziplinspezifisch, aber auch disziplinübergreifend



Folie 2

Projektziel

- Bestandsaufnahme und Analyse der bestehenden und geplanten Projekte im Bereich der Forschungsdateninfrastruktur
- Erstellung einer Roadmap mit Handlungsempfehlungen für eine disziplinübergreifende Forschungsdateninfrastruktur in Deutschland.



Folie 3

Förderung & Projektpartner

- 2 Jahre Laufzeit, bis April 2013
- Förderung durch DFG (LIS)
- 6 Partner aus unterschiedlichen Disziplinen und Wissenschaftsorganisationen



Folie 4

1. **Bestandsaufnahme**
 - Konzeption und Durchführung von 9 Experteninterviews
2. **Technik**
 - Anforderungen und Entwicklungsbedarf
 - Dokumentation vorbildlicher, disziplinspezifischer Lösungen
3. **Organisation**
 - Akteure und Strukturen in den einzelnen Disziplinen
 - Einbettung des Forschungsdatenmanagements
4. **Kosten**
 - Kosten des Forschungsdatenmanagements
 - Ableitung generische Faktoren und Risikoabschätzungen
5. **Community**
 - Website www.forschungsdaten.org
 - Workshop 17. 04. 2012, FDI 2013
6. **Synthese**
 - Zusammenfassung der Ergebnisse
 - Entwicklung von Handlungsempfehlungen



Folie 5

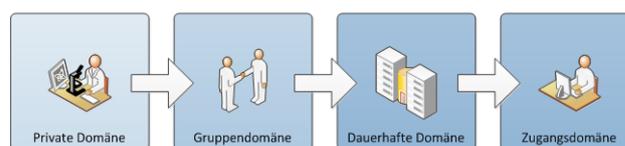
Bestandsaufnahme

- Interviews mit verschiedenen Akteuren
Bereich Forschungsdateninfrastruktur.
- Fokus auf:
 - Projekte der Ausschreibung
„Informationsinfrastrukturen für
Forschungsdaten“
 - INF Projekte in den SFB und TRR der DFG
- Literaturrecherche



Folie 6

Domänenmodell



- **Private Domäne**
 - Datenerzeugung und Datenverarbeitung in einer Hand.
- **Gruppendomäne**
 - Eine Kollaboration arbeitet auf einem gemeinsamen Satz von Daten.
- **Dauerhafte Domäne**
 - Langzeitarchivierung ohne direkten Bezug zum erzeugenden Akteur.
- **Zugangsdomäne**
 - Veröffentlichung der Daten für eine Gruppe oder die Öffentlichkeit.



Folie 7

Technik

- Höchst unterschiedliche Anforderungen, selbst innerhalb von Disziplinen. Workflows sind stark Disziplinabhängig.
- Manuelle Arbeitsschritte (z.B. zur Qualitätskontrollen) noch verbreitet.
- Der Großteil der verwendeten Software sind Eigenentwicklungen, vor allem im Bereich Datenanalyse.
- Disziplinübergreifender Bedarf ergibt sich hauptsächlich bei Fremdsoftware (z.B. im Bereich der Virtualisierung).

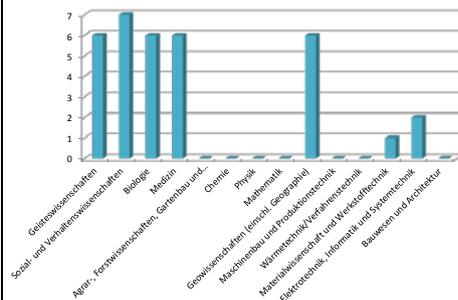


Folie 8

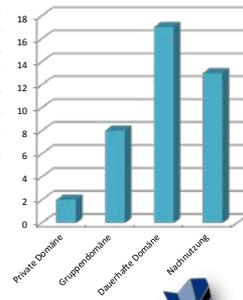
Organisation

DFG Projekte „Informationsinfrastrukturen für Forschungsdaten“

Verteilung nach Fachgebieten:



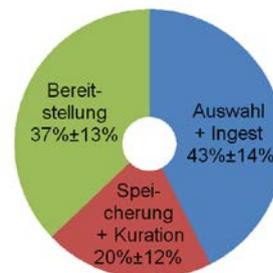
Verteilung nach Domänen:



Folie 9

Kosten

- Belastbare Zahlen sind schwer zu erhalten und zu vergleichen.
- Kosten von Archivierung
 - Personalaufwand für die einzelnen Schritte der Langzeitarchivierung
- Risiken der Nicht-Archivierung
- Risiken der Archivierung



Folie 10

Synthese

- Überblick über den aktuellen Status in der Entwicklung der Forschungsdateninfrastruktur
- Identifikation von Defiziten und Entwicklungsbedarf
- Entwicklung von Zukunftsszenarien
- Erarbeitung von Handlungsempfehlungen



Folie 11

www.forschungsdaten.org

**Informationsportal
Forschungsdaten**

Über uns
Forschungsdatenzentrum
– Eine Einleitung
Zusammenfassung für die
Forschungscommunity
Aktuelle Meldungen

Willkommen!

Auf dieser Seite finden Sie einen Überblick über das Thema Forschungsdaten und die wichtigsten Informationsquellen für das Datenmanagement in Forschungsinstituten.

Die Informationsquellen sind sowohl in allen Themenbereichen repräsentativ, aber werden fortwährend ergänzt und aktualisiert. Das Informationsportal Forschungsdaten wird aktuell vom Projekt Radischen Datenbedingungen einer disziplinübergreifenden Forschungsdateninfrastruktur (FDI) betreut, das hier auch weiter Aktionen versteht.

Aktuelle Meldungen

Radischen Experten-Workshop am 17. April 2012
wester/Dagstuhl Winter School 2011

Büchse Veranstaltungen

17. April:
Radischen Experten-Workshop am 17. April 2012



Folie 12



KomFor
Centre of Competence for Research Data
in the Earth & Environmental Sciences

Michael Diepenbroek
PANGAEA® / MARUM, University Bremen

Potsdam 2013 www.komfor.net

Folie 1

Partner

- World Data Center for Remote Sensing of the Atmosphere (WDC-RSAT)
- World Data Center on Climate (WDC-C)
- Helmholtz-Zentrum Potsdam Deutsches GeoForschungsZentrum (GFZ)
- PANGAEA® - Data Publisher for Earth & Environmental Science
- Technische Informationsbibliothek (TIB)



Potsdam 2013 www.komfor.net

Folie 2

Ziele

- Aufbau eines Kompetenzzentrum
 - Bindeglied zwischen wissenschaftlichen Einrichtungen, Verlagen, Bibliotheken & Datenarchiven
 - Vergleichbar ANDS oder DCC
- Publikation wissenschaftlicher Daten



Potsdam 2013 www.komfor.net

Folie 3

Arbeitspakete

Potsdam 2013

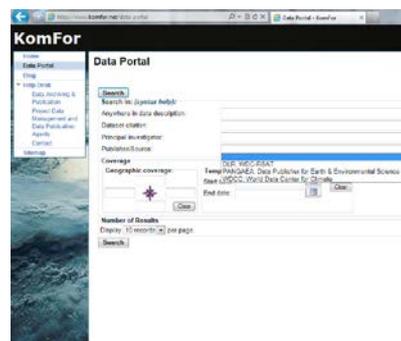
www.komfor.net



Folie 4

Mandantenfähige Service Plattform

- Helpdesk, Redaktionssysteme und Wissensbasis
- Datenportal



Potsdam 2013

Folie 5

Datenpublikation

- Organisatorische Voraussetzungen und Zusammenhänge (Workflows)
- Modelle für ein Peer-Review wissenschaftlicher Daten

Potsdam 2013

www.komfor.net



Folie 6

Katalog- und Registrierungsdienste

- Datenkatalogdienste (cross-linking)
- Bibliometrische Dienste
- Zertifizierung und Akkreditierung neuer Datenarchive / Publikationsagenten

Potsdam 2013

www.komfor.net



Folie 7

Geschäftsmodelle

- Kostenmodelle für Langzeitarchivierung & Datenpublikation
 - Open Access
- Service Modell
 - Projektdatenmanagement

Potsdam 2013

www.komfor.net



Folie 8

Data publishing - building blocks

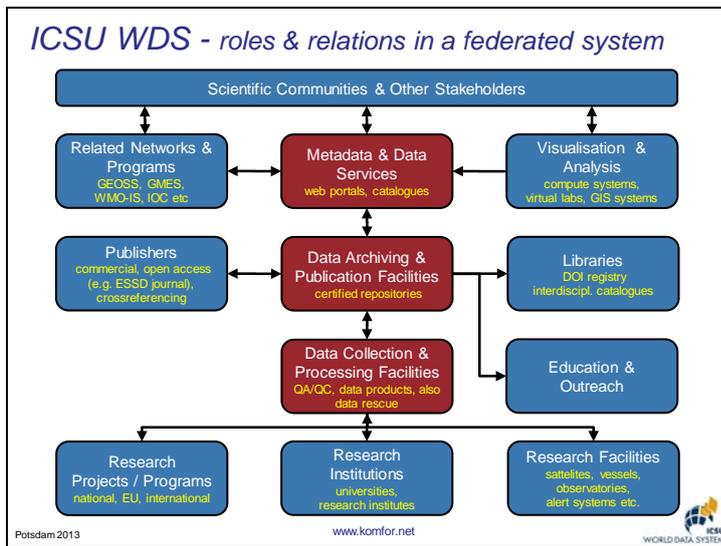
- DataCite -> DOI Registry for scientific data \Rightarrow 
- Collaborations with science publishers (Elsevier, Viley, Thompson Reuters etc.)
 - ✓ with science journals
 - ✓ linking & crossreferencing data & articles
- Thomson Reuters Data Citation Index (2012) \Rightarrow 
- ORCID registry for researchers (2012) \Rightarrow
- ICSU World Data System (WDS) for long-term data stewardship and publication

Potsdam 2013

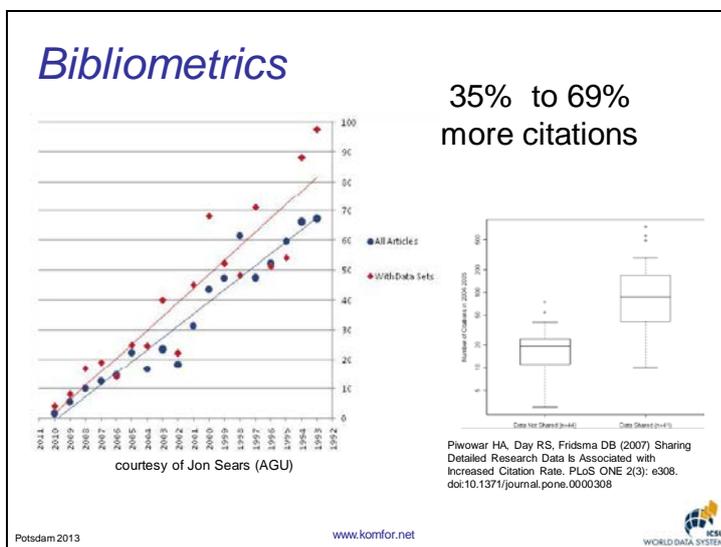
www.komfor.net



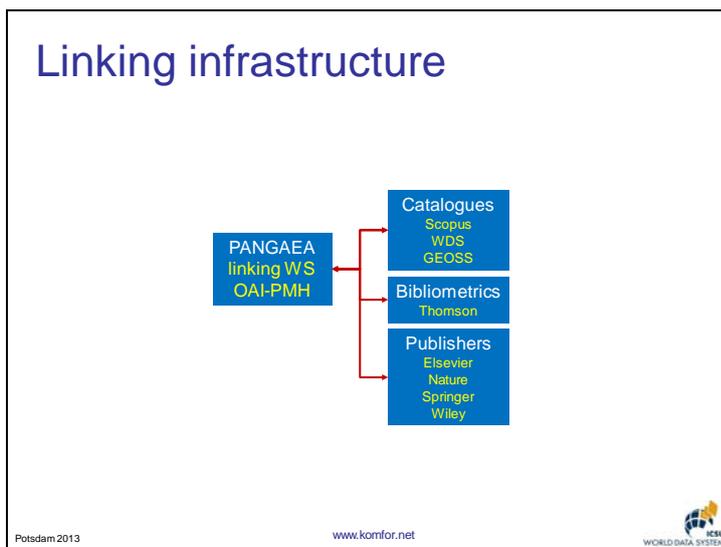
Folie 9



Folie 10



Folie 11



Folie 12

Data Publishing – Cross-referencing

PANGAEA®
Data Publisher for Earth & Environmental Science

Citation: Mohtadi, M et al. (2010): Surface sediment samples from several fore-arc basins west and southwest of the Indonesian Archipelago: analyzed by planktonic foraminifera, stable oxygen and carbon isotopic signals and opal and CaCO₃ contents in bulk sediment. doi:10.1554/PANGAEA.733340.

Supplement to: Mohtadi, Mahyar; Max, Lars; Hebbeln, Dierk; Baumgart, Anne; Krück, Nils; Jennerjahn, Tim C (2007): Modern environmental conditions recorded in surface sediment samples off W and SW Indonesia: Planktonic foraminifera and biogenic compounds analyses. *Marine Micropaleontology*, 65(1-2), 96-112. doi:10.1016/j.microp.2007.06.004

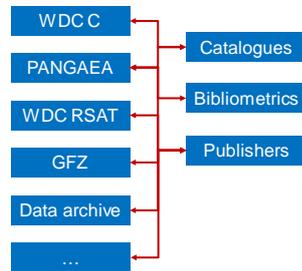
Abstract: A total of 69 surface sediment samples from several fore-arc basins located west and southwest of the Indonesian Archipelago was analyzed with respect to the faunal composition of planktonic foraminifera, the stable oxygen and carbon isotopic signal of a surface-dwelling (Globobulimina ruber) and a thermocline-dwelling (Neobulimina dutertrei) species, and the opal and CaCO₃ contents in bulk sediment. Our results show that the distribution pattern of opal in surface sediments corresponds well to the upwelling-induced chlorophyll concentration in the upper water column and thus, represents a reliable proxy for marine productivity in the coastal upwelling area off S and SW Indonesia. Present day oceanography and marine productivity are also reflected in the faunal to subsurface and upwelling assemblages of planktonic foraminifera in the surface sediments, which in part differ from previous studies in the region probably due to different coring methods and distribution effects. The average stable oxygen isotope values (δ¹⁸O) of B. ruber in surface sediments vary between 2.9 per mil and 3.2 per mil from basin to basin and correspond to the oceanographic settings during the SE monsoon (July-October) off west Sumatra, whereas off southern Indonesia, they reflect the NW monsoon (December-March) or annual average conditions. The δ¹³C values of B. dutertrei show a stronger interbasinal variation between 1.5 per mil and 2.2 per mil and correspond to the upper thermocline hydrology in July-October. In addition, the difference between the small carbon isotopic values (δ¹³C) of G. ruber and B. dutertrei (Δδ¹³C) appears to be an appropriate productivity recorder only in the non-upwelling areas off west Sumatra. Consequently, joint interpretation of the isotopic values of these species is obstructive for different fore-arc basins W and SW of Indonesia and should be considered in paleoceanographic studies.

Project(s): Center for Marine Environmental Sciences (MARUM)

Coverage: Median Latitude: 2.448391 ° Median Longitude: 103.924024 ° South-bound Latitude: -9.012150 ° West-bound Longitude: 85.331100 ° North-bound Latitude: 3.874500 ° East-bound Longitude: 121.902200

Event(s):
Geob109024

Linking infrastructure



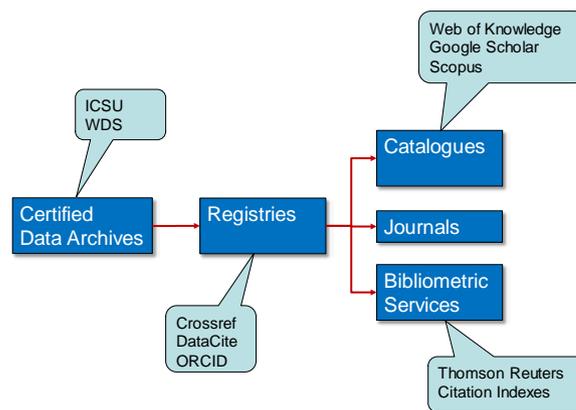
Potsdam 2013

www.komfor.net



Folie 16

ICSI WDS perspective



Potsdam 2013

www.komfor.net



Folie 17



re3data.org
Registry of Research Data Repositories

Frank Scholze, Hans-Jürgen Goebelbecker, KIT-Bibliothek
Heinz Pampel, GFZ Potsdam

Symposium „Forschungsdaten-Infrastrukturen“, Potsdam, 22.01.2013



Folie 1

Forschungsdaten-Repositoryen

- EC: ICT infrastructures for e-science
„The landscape of data repositories across Europe is fairly heterogeneous, but there is a solid basis to develop a coherent strategy to overcome the fragmentation and enable research communities to better manage, use, share and preserve data.“



European Commission. (2009). ICT infrastructures for e-science. Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions. COM(2009) 108 final. Retrieved from <http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=COM:2009:0108:FIN:PDF>

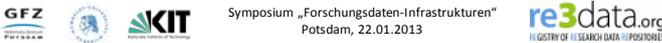


Folie 2

Ziele

- Aufbau eines “Registry of Research Data Repositories”
- Weltweites, disziplinübergreifendes, web-basiertes Verzeichnis von Forschungsdaten-Repositoryen
- Struktur- und Inhaltsanalyse der Angebote
- Orientierung für Wissenschaftler, Förderer, Verlage, Infrastruktur
- Vernetzung und Standardisierung von Forschungsdaten-Repositoryen
- Förderung der “culture of sharing”

Symposium „Forschungsdaten-Infrastrukturen“
Potsdam, 22.01.2013



Folie 3

Aktueller Stand

- Erste Version von **re3data.org** seit Dezember 2012 online.
 - ca. 120 voll erschlossene Forschungsdaten-Repositoryen
 - ca. 210 Forschungsdaten-Repositoryen in Bearbeitung
-> Basisdatensatz sichtbar
 - Die einfache Suche kann durch diverse Filtermöglichkeiten verfeinert werden
- Die Dateneingabe und -haltung erfolgen im XML-Format.
 - Dazu wurde ein eigenes XML-Schema entwickelt.

Vokabular



- Vierkant, P. et al. (2012).
Vocabulary for the Registration
and Description of Research
Data Repositories. Version 2.0.
doi:10.2312/re3.002

Erschließung

- Erschließung durch Autopsie (review)
 - 37 Hauptkriterien
 - 22 weitere Unterkriterien
 - 25 kontrollierte Vokabulare
- Icons als einfache Außendarstellung
 - Mehrwert für Repositoryen-Betreiber und wissenschaftliche Nutzer



Screenshot: Ergebnisliste

The screenshot shows the search results for 'mineralogy' on the re3data.org website. It includes a search bar, filter options for subjects, content types, and countries, and a list of search results such as 'Crystallography Open Database', 'Inorganic Crystal Structure Database', and 'PANGAEA'.



Symposium „Forschungsdaten-Infrastrukturen“
Potsdam, 22.01.2013



Folie 7

Screenshot: Vollanzeige

The screenshot shows the detailed view of the PANGAEA repository. It lists various metadata fields such as Name of repository, Repository URL, Subjects, Description, Content type(s), Key word(s), Repository type, Policy URL, Research data repository language(s), Name of the responsible institution, URL of the responsible institution, Contact, Additional name of the institution, Country, Type(s) of responsibility, Type of responsible institution, Data and/or service provider, Type of access to research data repository, Type of access to data, Data license name, Data license URL, Type of data upload, and Data valid restriction type.



Symposium „Forschungsdaten-Infrastrukturen“
Potsdam, 22.01.2013



Folie 8

Nächste Schritte

- Entwicklung der Webformulare
- Aufbau eines Redaktions-Workflows
- Einrichtung von Import- und Export-Schnittstellen (XML, Dublin Core, Statistiken)
- Öffentlichkeitsarbeit und Vernetzung
 - Dialog mit diversen Akteuren
 - Daten werden nachnutzbar sein
 - DINI-Wiki zum Thema im Aufbau
 - Mailingliste
 - Twitter



Symposium „Forschungsdaten-Infrastrukturen“
Potsdam, 22.01.2013



Folie 9



re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

Home Search Support FAQ About Editorial Contact Register

Search for repositories (date only)

info@re3data.org
http://re3data.org

 With the exception of all photos and graphics, this slide is licensed under the "Creative Commons Attribution 3.0 Germany (CC BY 3.0)" Licence.



GFZ Helmholtz-Zentrum POTSDAM
UNIVERSITY OF GIESSEN
KIT Karlsruhe Institute of Technology
DINI
DFG Deutsche Forschungsgemeinschaft

Folie 10

DFG-Projekt BoKeLa – Dr. H.-J.Wallrabe-Adams



Deutsches Forschungsbohrkonsortium GESEP e.V.

**BoKeLa:
Aufbau des Dateninformations-
systems für das GESEP Kern- und
Probenlager**

 German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 1



„Aufbau des Dateninformationssystems für das GESEP Kern- und Probenlager“

DFG-Vorhaben im Rahmen LIS* Förderprogramm:
„Informationsmanagement“

Projektstart August 2011 für die Dauer von drei Jahren
initiiert durch GESEP

*Wissenschaftliche Literaturversorgungs- und Informationssysteme (LIS)



German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 2



Was ist GESEP?

Kooperation 15 geowissenschaftlicher Institute, die

- im Bereich Forschungsbohrungen (Meer, Land, Eis) stark engagiert sind
- breite Expertise und Infrastruktur vernetzen
- gemeinsam neue Verfahren und Modelle entwickeln



German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 3



Drei Projektpartner – drei Aufgabenfelder

BGR (Bundesanstalt für Geowissenschaften und Rohstoffe):

Aufbau und Betrieb eines Bohrkernlagers am Standort Berlin-Spandau für kontinentale Kerne ohne Kühlung (Festgesteinskerne)

GFZ (Deutsches GeoForschungsZentrum):

Entwicklung eines Bohrkern- und Probenverwaltungssystems

MARUM (Zentrum für Marine Umweltwissenschaften, Bremen):

Betrieb eines Bohrkernlagers in Bremen für Kerne, die der Kühlung bedürfen (+4°C) (See-Sedimente)

Entwicklung eines Web-Portals mit entsprechenden Web-Diensten für die Integration der beiden Standorte Berlin-Spandau und Bremen sowie weiterer wissenschaftlicher Bohrkernlager im Umfeld von GESEP



German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 4

Kernarchiv für kontinentale Kerne

2009 Aufruf zur Einrichtung des Kernlagers
Einreichen eines gemeinsamen Konzepts von MARUM und BGR

2010 Start des Kernlagers
erste „Sampling Party“ am MARUM durch das Lake Van-Projekt, Türkei

09.2012 Eröffnung des nationalen Bohrkermlagers am BGR-Standort in Berlin-Spandau



Nationales Bohrkermlager

German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 5

GFZ: Entwicklung eines Bohrkern- und Probenverwaltungssystems

➔ **Drilling Information System (DIS)**

- zur Unterstützung des Datenmanagements an der Bohrung, im Labor und später zur Verwaltung im Kernlager
- dient der Dokumentation und Erfassung von
 - **Basis Daten** (Expedition, Lokation,...)
 - **Messungen und Berichten** (lithologische Beschreibung,...)
 - **Probenverwaltung**
- zur Bereitstellung einer
 - **gemeinsamen Referenz für alle beteiligten Wissenschaftler des Projektes**
- so früh wie möglich einsetzen – bereits an der Bohrung
- auch um nicht-synchronisierte und nicht-authorisierte Datensätze zu vermeiden



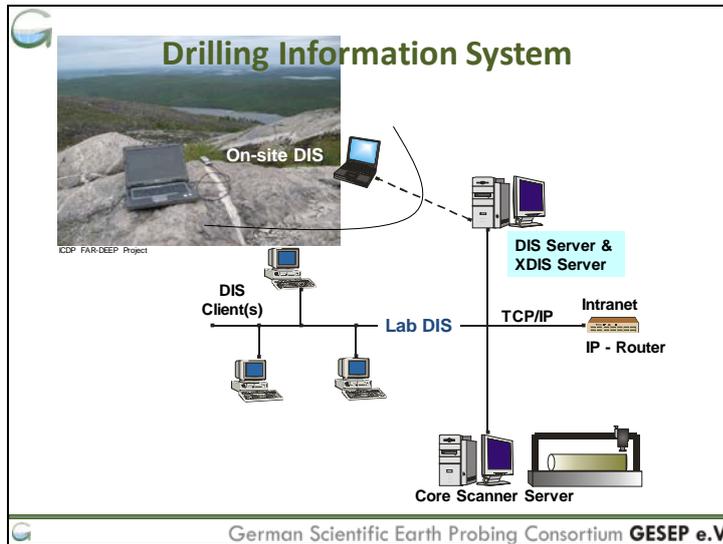

German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 6



German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 7



Folie 8

Bei der Eingabe der Daten in das DIS erhalten alle Kerne und Proben einen eindeutigen Identifikator, die sogenannte International Geo Sample Number (IGSN):

- Proben bzw. Daten dauerhaft identifizierbar zu machen
- Proben und Daten in Fachliteratur zitierbar zu machen

analog den bereits bekannteren DOIs.

IGSN

Science Reports

**Scientific Drilling Into the San Andreas Fault Zone
—An Overview of SAFOD's First Five Years**

by Mark Zoback, Stephen Hickman, William Ellsworth,
and the SAFOD Science Team

doi:10.2204/iodp.sd.11.02.2011

Abstract

Detailed planning of a research experiment focused on drilling, sampling, and downhole measurements directly within the San Andreas Fault Zone began with an intern...

The San Andreas Fault Observatory at Depth (SAFOD)

Folie 9

MARUM: Web-Portal und Web-Dienste

- für die Integration der beiden Standorte Bremen und Spandau sowie weiterer wissenschaftlicher Kernlager im Umfeld von GESEP in ein virtuelles GESEP Bohrkernlager
- zum Nachweis von verfügbarem Kern- und Probenmaterial
- zur Bereitstellung dieser Informationen auch an andere Web-Portale nach dem Open Access Prinzip
- als Einstieg für den Wissenschaftler/ Nutzer zur
 - Planung von Bohrprojekten
 - Durchführung von Core Opening- und Sampling-Parties
 - Einlagerung von Kernen und Proben
 - Anforderung von Proben
 - Durchführung von Schulungen und Workshops
- Eingliederung in übergeordnete Portale, um Kerne, Publikationen und Daten miteinander zu verknüpfen (z.B. PANGAEA)

German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 10

PANGAEA®
Publishing Network for Geoscientific & Environmental Data

You are not logged in (LOG IN)

Advanced Search

Search

Search terms: Geographic coverage:

Ar PANGAEA®
Data Publisher for Earth & Environmental Science

Er Data Description

Ci Citation: Zolitschka, B et al. (2012): Lake level reconstructions at Laguna Potrok Aike (southern Patagonia, Argentina) for the last 40 ka (doi: 10.1564/PANGAEA.801041). Supplement to Zolitschka, Bernd; Anselmetti, Flavio S; Ariztegui, Daniel; Corbella, Hugo; Francus, Pierre; Lüchge, Andreas; Maidana, Nora; Ohlendorf, Christian; Schabitz, Frank; Westgaard, Stefan: Environment and climate of the last 51 000 years - new insights from the Potrok Aike maar lake sediment archive-drilling project (PASADO). Quaternary Science Reviews, accepted

Abstract: Four seismic surveys and a stratigraphic record from southernmost Patagonia (Argentina) based on 51 AMS-14C dates obtained in the framework of ICDP expedition 5022 "Potrok Aike Maar Lake Sediment Archive Drilling Project" (PASADO) provide a database to compare the 106 m composite profile from the lake centre with piston cores from the littoral and outcrops in the catchment area. Based on event correlation using distinct volcanic ash layers with unique geochemical composition and optically stimulated luminescence (OSL) dates on kilobars, sediment records are firmly linked. This approach allows to match the sediment record with water levels during the past ca. 40 ka providing evidence for lake level variations. Reconstructed lake levels were 20 m higher than today during the last Glacial until the early Holocene. With the migration of the Southern Hemisphere Westerlies over this site the lake level dropped ca. 15 m for a period of few millennia. Thereupon the water balance was more positive again causing a stepwise rise of the lake level until the maximum was reached during the Little Ice Age with a subsequent lowering since the 20th century. We suggest that the mid- to late-Holocene lake level variation is caused by intensity changes of the Southern Hemispheric Westerlies.

Project(s): Potrok Aike Maar Lake Sediment Archive Drilling Project (PASADO)

Ci Coverage: Latitude: -51.563100 ° Longitude: -70.379400 ° Elevation: 116.0 m ° Location: Patagonia, Province of Santa Cruz, Argentina ° Device: Hydraulic piston corer °

Er License: CC BY Creative Commons Attribution 3.0 Unported

Search

Folie 11





Ziel

- Verbesserung der allgemeinen Verfügbarkeit und Qualität der Daten
- projektunabhängige, zentrale Datenbasis, die Primärdaten, Metadaten und die dazugehörigen Proben und Publikationen zusammenführt



 German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 12



ENDE



 German Scientific Earth Probing Consortium **GESEP e.V.**

Folie 13

Session 2 – Data Curation Continuum – Teil 1

Die Vorträge des Hauptprogramms spiegelten die einzelnen Stationen im Datenlebenszyklus wieder. Abb. 8 zeigt das Domänenmodell nach Treloar¹ mit integrierter dauerhafter Domäne.

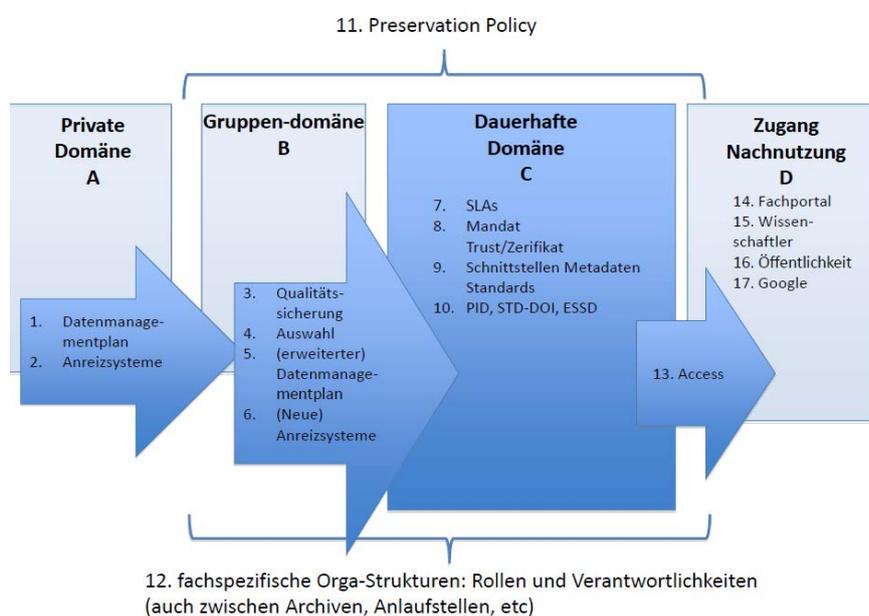


Abb. 8: Domänen-Modell (abgeleitet von Treloar, A. & Harboe-Ree, 2008)

¹ Treloar, A. & Harboe-Ree, C., 2008. Data management and the curation continuum. How the Monash experience is informing repository relationships. VALA2008 14th Biennial Conference. Melbourne, Australia 5.-7. Feb. 2008

Earth and Planetary Science Letters

Earth and Planetary Science Letters

After shock seismicity of the 27 February 2010 Mw 8.8 Maule earthquake rupture zone

Dierich Lange^{1,*}, Frederik Tilmann^{1,†}, Sergio E. Barrientos², Eduardo Contreras-Reyes³, Pascal Metzke^{4,†}, Marcos Moreno⁵, Ben Holt⁶, Hans Aguirre⁷, Pascal Bernard⁸, Jean-Pierre Villette⁹, Susan Beck¹⁰

¹University of Bremen, Germany
²Geological Survey of Chile, Santiago, Chile
³Universidad de Chile, Santiago, Chile
⁴University of Liverpool, UK
⁵University of Chile, Santiago, Chile
⁶University of Texas at Austin, USA
⁷University of Chile, Santiago, Chile
⁸University of Chile, Santiago, Chile
⁹University of Chile, Santiago, Chile
¹⁰University of Texas at Austin, USA

On 27 February 2010 the Mw 8.8 Maule earthquake in Central Chile ruptured a segment of the Nazca plate subducting beneath the South American plate. We present the aftershock distribution and first rupture geometry for an extensional fault on the subducting Nazca plate, and its kinematics. We present the aftershock distribution and first rupture geometry for an extensional fault on the subducting Nazca plate, and its kinematics. We present the aftershock distribution and first rupture geometry for an extensional fault on the subducting Nazca plate, and its kinematics.

© 2011 Elsevier B.V. All rights reserved.

Gives rise to derived results by data reduction and combination with prior results:

Folie 7

```

2010 315 0042 50.0 L -34.267 -72.124 24.8 XXX 16 0.4 4.25XXX 4.0XXXX 1
GAP=191 0.71 2.8 7.4 4.1 -0.1312E+02 0.1339E+02 -0.2639E+01E
ACTION:SPS 11-10-14 20:00 OP:CD STATUS: ID:20100315004122 1
STAT SP IPHASE 0 SEISM SOURCE CODA AMPLTY PVEL AKIMO VELO AIN AN TRGS W SIS CAS?
Q020 EE P OMO 0042 54.68 152 -0.2610 10.6 112
Q016 EE P OMO 0042 56.81 125 0.2610 26.2 58
Q017 EE P OMO 0042 57.99 113 -0.0710 37.2 146
Q012 EE P OMO 0042 58.10 113 -0.0310 37.9 88
Q018 EE P JMO 0043 02.25 100 0.42 5 64.0 50
Q013 EE P OMO 0043 05.20 78 -0.0410 87.4 77
Q014 EE P OMO 0043 05.25 78 -0.0410 84.1 171
Q009 EE P JMO 0043 10.12 73 -0.71 185 34
Q004 EE P AMO 0043 20.47 73 0.11 0 191 45
Q003 EE P JMO 0043 22.90 73 0.21 5 208 18
Q010 EE P AMO 0043 25.31 73 -0.21 7 228 186
Q011 EE P JMO 0043 25.48 73 -0.0410 228 186
Q005 EE P AMO 0043 27.22 73 1.36 0 229 29
Q001 EE P AMO 0043 29.77 73 -0.50 0 260 26

2010 315 0158 12.7 L -35.120 -72.675 37.9 XXX 8 1.0 4.25XXX 4.0XXXX 1
GAP=211 3.15 16.1 103.3 90.8 -0.3868E+03 0.8751E+04 -0.4130E+02E
  
```

Folie 8

```

2010 315 0042 50.0 L -34.267 -72.124 24.8 XXX 16 0.4 4.25XXX 4.0XXXX 1
GAP=191 0.71 2.8 7.4 4.1 -0.1312E+02 0.1339E+02 -0.2639E+01E
ACTION:SPS 11-10-14 20:00 OP:CD STATUS: ID:20100315004122 1
STAT SP IPHASE 0 SEISM SOURCE CODA AMPLTY PVEL AKIMO VELO AIN AN TRGS W SIS CAS?
Q020 EE P OMO 0042 54.68 152 -0.2610 10.6 112
Q016 EE P OMO 0042 56.81 125 0.2610 26.2 58
Q017 EE P OMO 0042 57.99 113 -0.0710 37.2 146
Q012 EE P OMO 0042 58.10 113 -0.0310 37.9 88
Q018 EE P JMO 0043 02.25 100 0.42 5 64.0 50
Q013 EE P OMO 0043 05.20 78 -0.0410 87.4 77
Q014 EE P OMO 0043 05.25 78 -0.0410 84.1 171
Q009 EE P JMO 0043 10.12 73 -0.71 185 34
Q004 EE P AMO 0043 20.47 73 0.11 0 191 45
Q003 EE P JMO 0043 22.90 73 0.21 5 208 18
Q010 EE P AMO 0043 25.31 73 -0.21 7 228 186
Q011 EE P JMO 0043 25.48 73 -0.0410 228 186
Q005 EE P AMO 0043 27.22 73 1.36 0 229 29
Q001 EE P AMO 0043 29.77 73 -0.50 0 260 26

2010 315 0158 12.7 L -35.120 -72.675 37.9 XXX 8 1.0 4.25XXX 4.0XXXX 1
GAP=211 3.15 16.1 103.3 90.8 -0.3868E+03 0.8751E+04 -0.4130E+02E
  
```

Can use the same arrival time data set for tomographic study:

- Velocity model
- Improved locations

→ Second publication?

Lange et al., in preparation

Folie 9

An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures

Jon Claerbout
(Stanford geophysicist)



*An article about **observational** science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures as well as reference to the underlying data, which must be open-access.*

paraphrasing Jon Claerbout
(Stanford geophysicist)



*An article about **observational** science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures as well as reference to the underlying data, which must be open-access.*

paraphrasing Jon Claerbout
(Stanford geophysicist)



For any given observation-based publication, only time can tell whether the ideas and hypothesis inspired by the data or the observations themselves will have the more lasting impact and will be the reason for which the paper is cited

Scientific value chain for Lange et al (2012)

Process	Product	Archival/distribution
Field work	Raw data, often proprietary format	USB drive
Conversion and metadata generation (geometry, instrument calibration)	Standardised waveform data	Seismological data archive (z.B. GEOFON) – possibly time-limited access restrictions
Data analysis 1: Picking of travel times, Earthquake location using 1D-model	Bulletin (Picks + Lokationen)	Some institutional repositories, e.g. ISC Temporary exp, project internal: local, email, dropbox
Data analysis 2: Tomographic inversion	3D-Geschwindigkeitsmodell, verbesserte Lokationen	Local, email, dropbox
Seismotectonic interpretation; higher order relationships	New scientific knowledge (Cartoons, integrative plots)	ISI publication (Commercial publisher or scientific agency; (p)reprint servers)

Private data organisation (computer-based analysis)

- Lab book: still widely used but varying in quality – more notebook than record of reference
 - Unstructured file-based notes: `README` and `doit.csh` files
 - In geoscience little danger of disputes of priority; proof of exact timing not usually required
 - **Challenges:** finding old data; avoiding overwriting successful runs with new parameter sets; inability to reproduce prior results due to missing obsolete software
- Personal 'low-tech' solutions, mostly based on UNIX tools:
- `unison/rsync/rsnapshot`
 - `locate + grep`
 - `subversion (svn)`
- Also popular
- Dropbox

Private data organisation (computer-based analysis)

Evaluation of existing practice

- Existing standard tools generally up to the job; challenge is in convincing people to use them (`svn/Trac` system in my section only used by 2-3 people)
- Necessity to move from `doit.csh` type note-taking to make based build system

Informal data sharing

Why are email and similar informal approaches so popular?

- Immediate use; no need to register, no limitation on data type (except size → Dropbox)
- Full control over data release in each instance; can send 'latest version'
- Email allows informality in instructions and imposed restrictions.

Disadvantages:

- **Maintenance** does not scale well: each request generates work; hard to keep up-to-date with improvements to derived data
- **Danger of data loss** if data generator leaves science (common for PhD students, postdocs, but also senior faculty retire etc)

In spite of disadvantages and availability of superior tools, real challenge to get adoption in the wider community (PanMetaDocs)

Vision I: static data and software post-publication

- Data **available** in form suitable for further processing at all stages of the value chain with suitable meta-information (geometry etc)
- Results **reproducible** using either standard open source software, or software supplied; parameter sets controlling software behaviour fully specified in digestible form
- Data **discoverable** by being fully tagged and geo-referenced

Vision I: static data and software post-publication

- Data **available** in form suitable for further processing at all stages of the value chain with suitable meta-information (geometry etc)
- Results **reproducible** using either standard open source software, or software supplied; parameter sets controlling software behaviour fully specified in digestible form
- ~~Data **discoverable** by being fully tagged and geo-referenced~~

Barriers to data sharing

Technological

- Lack of storage
- Proprietary data formats
- Interactive software

Personal

- Time and effort necessary for organising and documenting
- Reservation for future steps of scientific value chain
- Feeling of ownership (fear of data being stolen)
- Hiding irregularities / selective data presentation
- Software: retain edge / enforce collaboration by not releasing of software

• Social

- Lack of attribution, impact hard to track
- Data production given low value compared to publication record
- Danger to reputation due to data misuse by others

Overcoming barriers 'Carrots and sticks'

• Improving attribution and recognition

- DOIs are the technical pre-requisite; great for consistent datasets early in the value chain, but they are not on the same footing as publications (e.g. don't count for *h* factor)
- Formats for data publications (but in practice used rarely)
- Journals editors and reviewers must pay more attention to proper data attribution (author guidelines; targeted questions to reviewers)

Overcoming barriers 'Carrots and sticks'

• Improving attribution and recognition

- DOIs are the technical pre-requisite; great for consistent datasets early in the value chain, but they are not on the same footing as publications (e.g. don't count for *h* factor)
- Formats for data publications (but in practice used rarely)
- Journals editors and reviewers must pay more attention to proper data attribution (author guidelines; targeted questions to reviewers)

Vision II: attribution

Use of data and data products **fully referenced**; data production **rewarded** on equal footing with paper production

Overcoming barriers 'Carrots and sticks'

- Supplementary data as pre-requisite of publication (journals/funding bodies)
 - Require authors to supply underlying data and software for all figures, moving toward *Reproducible Research*
 - Respect exclusive use for intermediate data: force archival but allow time-limited restrictions similar to primary data to replace distribution-on-demand

NB: Many journal guidelines are already requiring data sharing on demand, but this policy is hard to police, and often fails for reasons outlined earlier

Enn, Vol. 84, No. 36, 9 September 2003

Complete PostScript: An Archival and Exchange Format for the Sciences?
—PAUL WENZEL, School of Ocean and Earth Science and Technology, University of Hawaii at Manoa, Honolulu

Comp. in Sci. & Engineering, 2000
Making scientific computations reproducible
Matthias Schwab, Martin Karrenbach, Jan Claerhout

GUEST EDITORS
INTRODUCTION

COMPUTING IN SCIENCE & ENGINEERING
Reproducible Research

GFZ
GEOLOGISCHES FORSCHUNGSZENTRUM
POTSDAM

HELMHOLTZ
ASSOCIATION

Folie 22

Modest steps

- Aim for reproducible research in the private sphere
- Better enforcement of existing journal policies
- End practice of converting vector graphics to jpg's for publication (in post-production, but also by users in composing WORD documents)
- Encourage sharing and reproducible research by awarding openness grades

GFZ
GEOLOGISCHES FORSCHUNGSZENTRUM
POTSDAM

HELMHOLTZ
ASSOCIATION

Folie 23

Gold:

- Underlying primary data open and data usage described precisely (time frame, scope)
- All necessary software open and supplied or accessible from repository with full versioning history
- Secondary data products supplied
- Complete build sequence specified ('Reproducible research'); (by necessity) interactive steps clearly described, and results of all interactive steps supplied

Silver:

- As Gold, but also allow proprietary software (if versioned and usage documented) and reasonable time-restrictions on distribution of data to be in place

Bronze:

- Only require (1) and (3) from Gold standard, limited time-restrictions allowed

To encourage adoption, can make this voluntary, and at a later stage make „Bronze“ or „Silver“ requirement of publication.

Possible exceptions:

- Collaborations with industry: allow longer time restrictions, or selective obfuscation of data (e.g. geographic location)
- Can define standard proprietary software which is allowed even for Gold standard, e.g. Excel, matlab (BUT: versioning problem)

GFZ
GEOLOGISCHES FORSCHUNGSZENTRUM
POTSDAM

HELMHOLTZ
ASSOCIATION

Folie 24

What about dynamic data sets/products?

Time-dependency through

- 1) **Growth**: natural growth of most continuous datasets (e.g. continuous seismological data and earthquake catalogues)
- 2) **Corrections**: Fixing of obvious mistakes (e.g. polarity, timing for seismological data; discovery of bugs in algorithms)
- 3) **Improvements**: improved results due to: (data products)
 - Growth of and corrections to primary data
 - Improved methodology (e.g. location in 3D instead of 1D model)

Vision III: time-dependency

- Provide access to the latest and best in primary data, data products and software but make available **full time history** of changes and easy reconstruction of any point in time
- Strive for automated analysis and provide **self-updating data products** which automatically assimilate additional data having become available since original paper
- For steps which cannot be automated: commit effort to keep data-product updated, or **isolate manual steps** and provide clear procedures for updating

Summary

Data should be available and discoverable.

Data products should be re-usable, reproducible, and capable of assimilating new data.

The full time history must be available for both as well as for the software used.

Sharing of data and software must be rewarded.

Peter Bartelheimer
Daten in Virtuellen Forschungsumgebungen

Symposium Forschungsdateninfrastrukturen
Helmholtz-Zentrum Potsdam, GFZ
22. Januar 2013

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 1

■ **Datenbereitstellung, Virtuelle Forschungsumgebung –
zwei Komponenten der IT-Forschungsinfrastruktur**

■ **Forschungsdateninfrastruktur**

- Ziel: Mikrodaten aus Bevölkerungsumfragen und Verwaltungsregistern berechtigten wissenschaftlichen Nutzer/innen verfügbar zu machen
- Leitbild: Zugriff, (Nach-) Nutzung, Kontrolle netzbasiert unmittelbar (»seamless«) zugänglich machen
- Problem Datengeheimnis: Daten sind unterschiedlich stark anonymisiert – Kontrolle von Zutritt, Zugang, Zugriff, Weitergabe und Output

■ **Virtuelle Forschungsumgebung (VFU)**

- Ziel: Wissenschaftliche Arbeitsprozesse vernetzen
- Leitbild: Kooperative Forschung unabhängig vom Ort zu gleicher Zeit ohne Ressourcen- und Zugangsprobleme
- Problem Workflow: Individualisierte Arbeitsprozesse analysieren und typisieren, um sie zu unterstützen

■ **Entwicklungsbedarf auf beiden Seiten der Datenschnittstelle**

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 2

■ **Das Projekt: Virtuelle Forschungsumgebung für die
sozioökonomische Berichterstattung (VFU soeb 3)**

■ **Projektpartner für die Entwicklung der VFU**

- SOFI mit IT-Partnern GESIS, GWDG, SUB
- Forschungsdatenzentren (FDZ-IAB, FDZ-RV, FDZ-SOEP)

■ **Funktionen der VFU in typischen Workflows der Datennutzung**

- IT-Portal mit Funktionen zur Unterstützung des gesamten Workflows bei der Nutzung sozial- u. wirtschaftswissenschaftlicher Mikrodaten durch
 - Kollaborationstools
 - Datenmanagement
 - Benutzer- und Rollenmanagement
 - Datenzugang
 - Metadaten- und Syntaxeditoren

■ **Derzeit: Einführungs- und Entwicklungsphase**

- Juni 2012 bis Februar 2014

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 3

■ Anwendungsfall: Dritter Bericht zur sozio- ökonomischen Entwicklung in Deutschland (soeb 3)

■ Verbundvorhaben 2013 - 2015 in Vorbereitung (Antragslage)

- Bis zu 24 beteiligte wissenschaftliche Einrichtungen mit bis zu 29 Arbeitspaketen (AP)
- Nutzung von 66 verschiedenen Forschungsdatensätzen, Datenbanken
- Neun der beim RatSWD akkreditierten FDZ als Datengeber
- An VFU-Entwicklung und Verbund beteiligt: FDZ-IAB, FDZ-RV, FDZ-SOEP

■ Schwerpunkte der Datennutzung

- Neue sozial- und wirtschaftswissenschaftliche Panel-Datensätze
- 13 Datensätze kollaborativ genutzt, davon sieben intensiv (> zwei AP)

■ Projektbegleitende Weiterentwicklung der VFU

- Operative Nutzung im Verbund
- Teilnahme von Verbundpartnern in Nutzungsstudien

■ Ziel: Bereitstellung einer »nachhaltigen« VFU-Plattform

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 4

■ Welche Daten (»Objektklassen«) und Metadaten?

■ Unveränderliche Daten

- Forschungsdaten (Ausgangsdaten)
 - Datensätze und Dokumentation (»Studien«)
 - Metadaten auf Studien-, Datensatz und Variablebene
- Statistikprogramme (z.B. R), statistische Packages (z.B. ado-Files)

■ Im Workflow erzeugte Daten

- Syntaxdateien, Syntax-Memos, Metadaten zu Syntax
- Output-Dateien, Metadaten zu Outputs, z.B. Log-Dateien
 - Tabellen, Grafiken, Textdateien (diverse Datenformate)

■ Speicherorte, Orte der Nutzung außerhalb der VFU

- Forschungsdatenzentren, Daten(service)einrichtungen
- Internet, Intranet
- Lokale Workstation, Intranet

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 5

■ An der Datenschnittstelle: Wie kommen Forschungsdaten in die VFU?

■ Bereitstellung für lokale Nutzung

- Public / Scientific Use Files -

■ Bereitstellung für Rechnen in Dateneinrichtungen

- Gastarbeitsplätze: interaktiver lesender Zugriff auf Originaldaten, Outputkontrolle
- Kontrolliertes Fernrechnen, zwei Varianten
 - Versenden von Syntax (»job submission«, »Remote Execution«), Übermittlung kontrollierter Outputs
 - »Remote Data Access«: interaktiver lesender Zugriff über gesicherte Verbindung (»safe centers«), Outputkontrolle

■ Jedes FDZ gestaltet Fernrechnen anders

■ Anforderungen an Authentifizierung, Rechte- u. Rollenverwaltung

- Zugangskontrolle, Zugriffskontrolle, Weitergabekontrolle

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 6

■ VFU muss alle Datenzugangswege unterstützen

- **Rechnen in der VFU als Variante externer Bereitstellung**
 - Public Use Files (PUF) / Scientific Use Files – nach Vereinbarung
 - Nutzung der Metadaten und Syntax in der VFU beim Rechnen in der VFU – mit R oder lizenzierten Statistikprogrammen
- **»Job Submission«**
 - Einheitliche Nutzungserfahrung durch Integration der verschiedenen Schnittstellen in VFU-Portlet?
 - Nutzung der Metadaten und Syntax in der VFU für Job-Erstellung
 - Historisierung der Jobs anhand von Metadaten
- **»Remote Data Access«**
 - Anmeldung über Aufruf der VFU an gesicherten Standorten oder zertifizierten Forschungseinrichtungen
 - Nutzung der Metadaten und Syntax in der VFU beim interaktiven Fernrechnen

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 7

■ Hinter der Datenschnittstelle: Workflow und VFU

- **Aufarbeiten des Forschungsstands**
 - Arbeit mit Literatur: Recherche, Exzerpte, Bibliografie
- **Arbeit mit Forschungsdatensätzen**
 - Planung, Organisation, Dokumentation und Ausführung
 - Datensatzgenerierung
 - Recodierung und Analysen
 - Ergebnispräsentation
 - Datensicherung
- **Upload und Download**
 - von Dokumenten, Syntax, Outputs
- **Editieren**
 - von Syntax und Metadaten
- **Ergebnisdokumentation, Publikation (Open Access)**

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 8

■ Hinter der Datenschnittstelle: Wie kommen Syntaxdateien in die VFU?

- **Syntaxeditor als technische Entwicklungsaufgabe**
 - Funktionen des Erstellen und Nutzung von Syntax in der VFU
 - Upload / Download zwischen lokaler Workstation und VFU
 - Versionierung und Dateivergleich
 - Ggf. Volltextsuche
- **Soziale Aspekte – Arbeitsweise in der VFU**
 - Wer lädt welche Syntax hoch und wann?
 - Lauffähige (vollständige) Syntax oder Syntaxbausteine
 - Kontrolle über eigenes Arbeitsprodukt (geistiges Eigentum)
 - Rechte- und Rollenverwaltung: dateibezogene Rechte
 - Verhaltensanforderungen
 - z.B. Namenskonventionen, »kontrollierte Vokabulare«
 - Integrität zwischen lokaler Workstation und VFU?

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 9

■ **Hinter der Datenschnittstelle: Wie kommen Metadaten in die VFU?**

■ **Metadatenchema als konzeptionelle Entwicklungsaufgabe**

- Welcher Metadatenstandard (z.B. DDI 3.1, DDI 2.5)
- Für Metadaten zu Syntax kein etablierter Standard
- Dokumentationstiefe: Einschluss von Variablenebene?

■ **Metadateneditor als technische Entwicklungsaufgabe**

- Integration in VFU
 - Metadatenverwaltung in SQL-Datenbank
 - Verknüpfung mit »Objekten« in Dateiverwaltung
- Suchfunktion: strukturierte Suche und Volltextsuche
- Metadaten: Eingabe und Extraktion: z.B. aus Syntax, aus Datensätzen

■ **Soziale Aspekte – Arbeitsweise in der VFU**

- Verhaltensanforderungen und individuelle Arbeitsweise
 - Z.B. Erfassungsaufwand, Dokumentationstiefe, kontrollierte Vokabulare

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 10

■ **Zum Schluss**

■ **Zwei Perspektiven auf VFU**

- Kritisch für Akzeptanz datenhaltender Einrichtungen: Kontrollanforderungen
- Kritisch für Akzeptanz durch Nutzer/innen: intuitive individuelle Anwendung

■ **Nicht alle Daten stehen in der VFU**

- Leistungsversprechen der VFU: einheitliche Arbeitsumgebung
- Schrittweise mehr Daten integrieren

■ **Entwicklungsbedarf auf beiden Seiten der Datenschnittstelle**

- »Forschende FDZ« in VFU-Entwicklung und Nutzungsstudien integrieren

■ **»Koevolution« technischer und sozialer Entwicklung**

- Individualität der Forschungspraxis unterstützen
- Entwicklung als transdisziplinäre Herausforderung
- Selbstreflexion sozialwissenschaftlicher Datennutzung:
 - Analyse und Typisierung des Workflow ≠ Standardisierung

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 11

■ **Mehr ...**

- <http://www.soeb.de>

Dr. Peter Bartelheimer, 22. Januar 2013

SOFI | Soziologisches Forschungsinstitut Göttingen
an der Georg-August-Universität

Folie 12

Session 3 – Data Curation Continuum – Teil 2

„Persistente Domäne ,PID, DLZA, Zertifikate“ – Reiner Mauer (GESIS)



gesis
Leibniz-Institut für Sozialwissenschaften

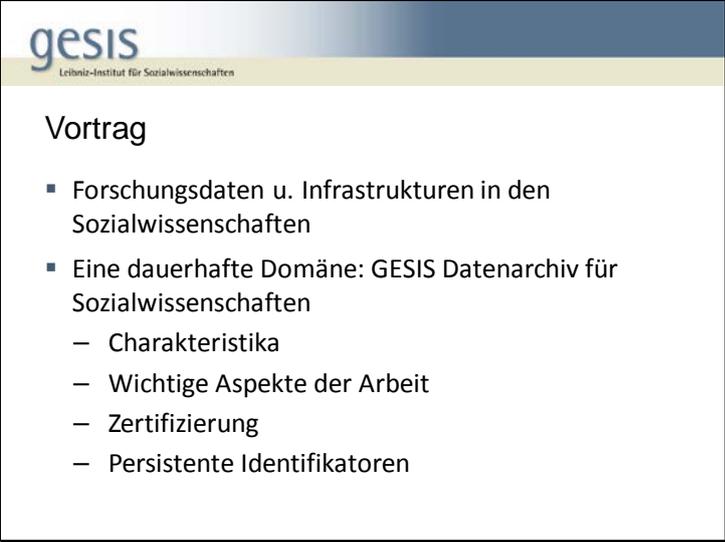
**Persistente Domäne:
PID, dLZA, Zertifikate**

Dauerhafte Zugänglichkeit von
Forschungsdaten

Reiner Mauer
GESIS – Leibniz-Institut für Sozialwissenschaften,
Datenarchiv für Sozialwissenschaften
Symposium Forschungsdaten-Infrastrukturen, Potsdam, 22.01.2013

Mitglied der
Leibniz-
Leibniz-Gemeinschaft

Folie 1



gesis
Leibniz-Institut für Sozialwissenschaften

Vortrag

- Forschungsdaten u. Infrastrukturen in den Sozialwissenschaften
- Eine dauerhafte Domäne: GESIS Datenarchiv für Sozialwissenschaften
 - Charakteristika
 - Wichtige Aspekte der Arbeit
 - Zertifizierung
 - Persistente Identifikatoren

Folie 2

Forschungsdaten in den Sozialwissenschaften

- Sozialwissenschaften: Sammelbegriff für viele Fächer
→ vielfältige Datentypen (Statistik, Audio, Video, Text, ...)
- Sekundär-/ Nachnutzung v. Forschungsdaten relativ verbreitet (Soziologie, Politologie)
- Lange Tradition von Datenarchiven
erste Gründung 1940er-Jahre in USA; 1960 Zentralarchiv für Empirische Sozialforschung (jetzt GESIS Datenarchiv für Sozialwissenschaften)
- gegenwärtig allein in Europa 21 Datenarchive im
CESSDA-Verbund organisiert (www.cessda.org)

Folie 3

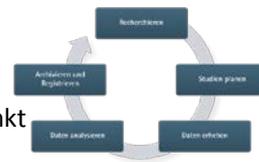
Forschungsdaten in den Sozialwissenschaften (2)

- Gründung des **Rats für Sozial- u. Wirtschaftsdaten (RatSWD)** in 2004 u. paralleler Aufbau von **Forschungsdatenzentren** verbreitern Datenzugänglichkeit deutlich
- SoWi-Dateninfrastrukturlandschaft entwickelt sich insgesamt sehr dynamisch, aber
 - positive Entwicklung führt auch zu gewisser Unübersichtlichkeit im Datenzugang
 - viele Daten existieren (?) weiterhin im ‚Verborgenen‘
 - Wachsende Anzahl an Akteuren in der ‚Access‘-Domäne, aber persistente Domäne (noch) unterentwickelt bzgl. Breite u. Tiefe der Abdeckung

Folie 4

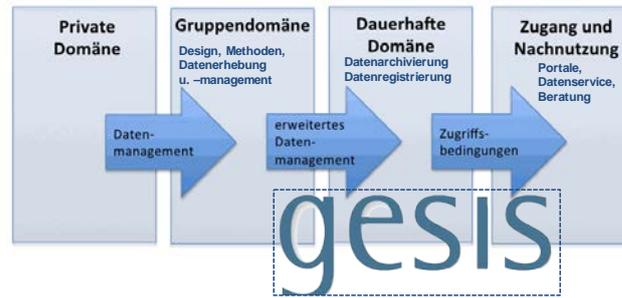
Was ist GESIS?

- Infrastruktureinrichtung für die Sozialwissenschaften, Schwerpunkt empirische Sozialforschung
 - **Forschungsdaten** (Erhebung, Archivierung, Registrierung, Analyse)
 - **Integrierte Fachinformationen** (Daten, Literatur, Projekte)
 - **Methodenberatung u. -entwicklung**
 - **Wissensvermittlung** (Summer School, Workshops, Seminare etc.)
- Basis der Serviceangebote ist eigene Forschung
- Langzeitarchivierung von Forschungsdaten explizit als Zielsetzung in der Satzung verankert (§ 2, 2c)



Folie 5

GESIS – Aktivitäten im Curation Continuum



[Grafik: WissGrid, Leitfaden zum Forschungsdaten-Management, eigene Textergänzungen (blau)]

Folie 6

GESIS Datenarchiv: persistente Domäne

- 50+ Jahre Archivierung von Forschungsdaten
 - Kontinuität trotz (oder wegen?) institutionellem Wandel gewahrt
 - technologischer Wandel extrem; veränderte Anforderungen u. Nutzungsszenarien, organisatorischer u. rechtlicher Rahmen
- Fokus auf Kuratieren u. Bereitstellen für Nachnutzung
 - Kuratieren dient nicht nur dem Erhalt, auch d. Verbesserung der Nutzbarkeit u. dem Schaffen von Mehrwert
 - Langzeitarchivierung lange Zeit eher implizite Aufgabenstellung
- Fachwissenschaftlicher/ disziplinärer Ansatz
 - eigene Forschung u. Einbettung in Forschungszusammenhänge sind Teil des Selbstverständnisses (Ausdifferenzierung der Funktionen hat allerdings stark zugenommen)

Folie 7

Datenarchiv: Persistenz, LZA (2)

- Archivierungseinheiten (Studien) bzw. AIPs sind
 - ... **heterogen**: bestehen aus mehreren bis vielen (im Extremfall mehreren hundert) Objekten, die in ihrer Zusammensetzung variieren (ein od. mehrere Datensätze, Messinstrumente, Metadaten, begleitende Materialien ...)
 - ... **nicht statisch**: Daten u. Metadaten werden im Archiv (kontinuierlich) verändert bzw. erzeugt (korrigiert, erweitert um neue Datenpunkte, integriert, aufgewertet)
- **LZA an wichtigen Stellen von manuellen/ intellektuellen Prozessen bestimmt** (bspw. Ingest inkl. Eingangskontrolle und Datenbereinigung, Datenaufbereitung u. -dokumentation)

Folie 8

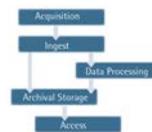
Datenarchiv: Persistenz, LZA (3)

- eher **kleine Daten**
 - vergleichsweise geringe Speicheranforderungen
- kleine Daten, **aber komplexe Inhalte**
 - Aufbereitung, Dokumentation u. Metadaten spielen wichtige Rolle beim Kuratieren; sonst keine (sachgerechte) Nachnutzung möglich
 - SoWi verfügt mit **DDI** über einen komplexen, internationalen Metadatenstandard  (www.ddialliance.org/)

Folie 9

Datenarchivierung@GESIS: Funktional

- **Akquisition** (aktiv, passiv)
- **Aufnahme ins Archiv** (Ingest)
- **Datenaufbereitung u. –dokumentation** (Standard für alle, Added-value für ausgewählte Studien)
- **Langzeitarchivierung**
- **Datenservice (Access):** Beratung, Datenzugang (Download, Online-Analyse, manuelle Bereitstellung)
- **Datenregistrierung (da|ra):** Vergabe von DOIs im DataCite Verbund

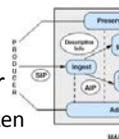


10

Folie 10

Ingest / Aufnahme ins Archiv (1)

- Vorbereitung d. Archivierung mit Datengeber
 - Ziel ist es möglichst früh im Lebenszyklus anzusetzen
 - Klärung bzgl. Umfang, Formaten, Aufbereitungs- u. Dokumentationszielen, Zugangsklasse
 - Abschluss Archivierungsvertrag
- Eingangskontrolle SIP und Basis-Aufbereitung
 - Technische Kontrollen (Formate, Lesbarkeit, Virenfreiheit ...)
 - Vollständigkeit u. Nutzbarkeit
 - Konsistenz der Daten
 - **ggf. korrigierende Aufbereitungen**



11

Folie 11

Access / Datenbereitstellung

- 40.000+ Datenweitergaben in 2012
- Weit überwiegend akademische Nutzung: Lehre und Forschung
- Je nach Angebot 30%-70% internationale Nutzer
- Datenzugang über Online-Portale und (individuelle) Bereitstellung
- (Meta)Datenportale: Retrieval, Download, Online-Analyse
- Zugang zu Datenbeständen internationaler Partnerarchive (ICPSR, CESSDA, IFDO)



15

Folie 15

Langzeitarchivierung: Substanzerhaltung

- physischer Schutz vor Ort
- räumlich getrennte und redundante Datenhaltung
- Diversität eingesetzter Speichertechnik
- regelmäßige Medienmigration bzw. Refreshment



16

Folie 16

Langzeitarchivierung: Erhalt von Nutzbarkeit und Interpretierbarkeit

- Bedrohung der digitalen Bestände durch technischen Fortschritt ist eine Konstante in der 50jährigen Archivarbeit
 - Maßnahmen, die auf Erhalt des bitstreams abzielen nicht ausreichend
- Erhalt wird hauptsächlich durch Migrationsstrategien erreicht (zur Überbrückung auch Emulation bzw. Virtualisierung)



17

Folie 17

Langzeitarchivierung: Erhalt von Nutzbarkeit und Interpretierbarkeit (2)

Maßnahmen:

- Dokumentation und Metadaten, Metadaten, Metadaten,
- Definierte und standardisierte Formate (Ingest): erleichtert Monitoring, senkt Migrationsbedarf, verringert Ressourcenbedarf bei Migration
- Verfolgen d. technischen Entwicklung (insbesondere Software und damit verbundene Dateiformate)
- Formatmigrationen
nur wenn Gefahr der Beeinträchtigung der Nutzbarkeit oder wenn mit der Migration so große Vorteile für die Nutzung oder die Archivarbeit einhergehen, dass der Aufwand zu rechtfertigen ist

18

Folie 18

Zertifizierung: „Vertrauenswürdige Archive“

Vertrauen in der dauerhaften Domäne von großer Bedeutung:

- Datenproduzenten u. -nutzer, Kooperationspartner, Forschungsförderer müssen darauf vertrauen und sich davon überzeugen können, dass ein Langzeitarchiv in der Lage ist, dauerhaft
 - Authentizität,
 - Integrität,
 - Verfügbarkeit und Zugänglichkeit,
 - sowie Interpretierbarkeit der Archivobjekte zu erhalten

→ Kurz u. mittelfristig evtl. kein Problem, aber dauerhaft ...?

Folie 19

Zertifizierung (2)

- Vertrauen muss erworben werden, aber **Vertrauenswürdigkeit** kann durch eine Zertifizierung nachgewiesen werden (wird mittelfristig ein ‚Muss‘ in der persistenten Domäne)
- Zertifizierung **auch und gerade für etablierte Institutionen relevant**, wie bspw. das GESIS Datenarchiv
 - Grundlegende Abläufe und Verfahren wurden lange vor OAIS, DIN 31644 od. ISO 16363 etabliert
 - Zertifizierung hilft bei der Identifikation v. Schwachstellen u. kann zur Verbesserung der Abläufe u. Services beitragen

Folie 20

Zertifizierung (3)

- **Optionen:**
 - Data Seal of Approval (DSA): Selbst-Audit mit Review
 - DIN 31644: Kriterien für vertrauenswürdige digitale Langzeitarchive
 - ISO 16363: Audit and certification of trustworthy digital repositories
- 3-stufiges Verfahren **“European framework for audit and certification of digital repositories”**:
 - **Basis:** DSA
 - **Erweitert:** DSA plus extern begutachtetes self-audit (ISO od. DIN)
 - **Formal:** DSA plus externes Audit und Zertifizierung (ISO od. DIN)

Folie 21

Zertifizierung (4): GESIS Datenarchiv

- Zunächst Basiszertifizierung (gerade angelaufen); anschließend ‚internes‘ Selbst-Audit DIN
- Weitere Schritte derzeit noch nicht entschieden, sehr wahrscheinlich ‚erweiterte‘ Zertifizierung:
 - entsprechende Angebote für die ‚erweiterte‘ u. ‚formale‘ Zertifizierung sind in Entwicklung (z.B. **nestor** für DIN)
 - kritische Entscheidungsgröße ist der dafür notwendige Ressourceneinsatz

Folie 22

Persistente Identifikatoren (PID): Was ist das?

- Ein Identifier (Name), der **eindeutig** auf **Dauer** mit einem **digitalen Objekt** verknüpft ist
- Zielsetzung: digitale Objekte langfristig identifizierbar, referenzierbar und verfügbar zu halten
- Verschiedene technische Systeme, verschiedene Policies, verschiedene Vor- u. Nachteile
 - Zu möglichen Risiken u. Nebenwirkungen fragen Sie Ihren Arzt oder Apotheker

Folie 23

Persistente Identifikatoren (PID): Wozu?

(1) Zitation: Daten werden selten od. ungenügend zitiert

- Probleme:
 - Forschungsergebnisse sind nicht replizierbar (weil Daten nicht auffindbar; Unsicherheit bzgl. Versionen u. Quellen)
 - Zitationen weisen **Reputation** zu; es **fehlen Anreize, Daten zu publizieren und damit zu teilen**
- Zitation mit PIDs eindeutig; erlauben schnellen Zugriff; können für bibliometrische/ szientometrische Analysen sowie Impact-Messung verwendet werden

(2) Verlinkung (Daten mit Daten, Daten mit Publikationen, Daten mit ...)

Folie 24

PID: Fragen zur Implementierung

- Welches System?
Technologie, Stabilität, Kosten, Zusatzdienste, Metadaten, Registrierungsprozess
- Was soll mit welcher Granularität registriert werden?
- Was geschieht, wenn das Objekt verändert wird? Wann wird ein neuer PID vergeben?
- Wie soll mit Objekten umgegangen werden, die auch an anderen Orten/ bei anderen Institutionen vorgehalten werden?
- Wohin löst der PID auf? Direkt auf das Objekt, landing page ...?

Folie 25

Herausforderungen

- Dateninfrastruktur zunehmend verteilter
→ Stärkere Bedeutung strukturierende u. kooperative Dienste
- Bedeutung komplexer Forschungsdesigns und neuer Datenformen wächst
→ LZA muss mit immer komplexeren Forschungsdaten umgehen;
- Verlinkung/Integration von Daten aus unterschiedlichen Datenquellen u. Disziplinen nimmt zu
→ höhere Anforderungen an Daten u. Metadaten; rechtliche Fragen, Lizenzfragen; Qualität: Aktualisierung, Versionierung, autoritative Versionen

Folie 26

„Eternity is really long, especially near the end“

[Woody Allen]

Vielen Dank!

Kontakt: Reiner Mauer, reiner.mauer@gesis.org

Folie 27

GESIS Datenarchiv: Datenquellen

Daten wurden / werden

- ... **selbst bzw. unter Beteiligung von GESIS erhoben:**
z. B. ALLBUS, EVS, ISSP, GLES
- ... **akquiriert** (externe Daten):
Großteil des Archiv-Bestands, wie z.B. Politbarometer, Eurobarometer, Mediaanalyse, Reiseanalyse, viele Einzelstudien)
- ... **entwickelt / produziert / transformiert**
Produktion zeit- u./od. ländervergleichender Datensätze:
Eurobarometer, EVS, ISSP, Politbarometer, Daten der historischen Statistik (HISTAT)

Folie 28

„Zugangsdomäne ,Portale, Best Practices“ - Dr. Hans Pfeiffenberger (AWI)

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

Portale, Best Practices – Forschungsdaten sichtbar und nachnutzbar machen

Hans Pfeiffenberger
Alfred-Wegener-Institute for Polar and Marine Research, Helmholtz Association - Germany

Symposium FDI, 2013-01-22, Potsdam

Earth System Science Data
Volume 3, Number 1, 2013

AWI

Folie 1

PHILOSOPHICAL TRANSACTIONS: GIVING SOME ACCOUNT OF THE PRESENT UNDERTAKINGS, STUDIES, AND LABOURS OF THE INGENIOUS IN MANY CONSIDERABLE PARTS OF THE WORLD.

Vol. I.
For Anno 1665, and 1666.

In the SAVOY, Printed by T. N. for John Martyn at the Bell, a little without Temple-Bar, and James Allestry in Duck-Lane, Printers to the Royal Society.

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 2

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

2012: Nature CC & ESSD; data reduction at global scale

Earth Syst. Sci. Data Discuss., 5, 1107–1157, 2012
www.earth-syst-sci-data-discuss.net/5/1107/2012/
doi:10.5194/essd-5-1107-2012
© Author(s) 2012. CC Attribution 3.0 License.

Earth System Science Data

ESSDD
5, 1107–1157, 2012

The global carbon budget 1959–2011
C. Le Quéré et al.

This discussion paper (as has been under review for the journal Earth System Science Data (ESSD)). Please refer to the corresponding final paper in ESSD if available.

The global carbon budget 1959–2011

C. Le Quéré¹, R. J. Andres², T. Boden³, T. Conway³, R. A. Houghton⁴, J. I. House⁵, G. Marland⁶, G. P. Peters⁶, G. van der Werf⁶, A. Ahlström⁷, R. M. Andrew⁸, L. Bopp⁹, J. G. Canadell¹⁰, P. Ciais¹⁰, S. C. Doney¹², C. Enright¹, P. Friedlingstein¹³, C. Huntingford¹⁴, A. K. Jain¹⁵, C. Jourdain¹⁶, E. Kato¹⁸, R. F. Keeling¹⁷, K. Klein Goldewijk¹⁸, S. Levis¹⁹, P. Levy¹⁹, M. Lomas⁹, B. Poulter²⁰, M. R. Reuspeck¹, J. Schwyling²⁰, S. Sitch²¹, B. D. Stocker²², N. Viovy¹⁰, S. Zaehle²³, and N. Zeng²⁴

¹Tyndall Centre for Climate Change Research, University of East Anglia, Norwich Research Park, Norwich, NR4 7TJ, UK
²Carbon Dioxide Information Analysis Center (CDIAC), Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA
³National Oceanic & Atmospheric Administration, Earth System Research Laboratory (NOAA/ESRL), Boulder, Colorado 80535, USA
⁴Woods Hole Research Centre (WHRC), Falmouth, Massachusetts 02540, USA
⁵Cabol Institute, Dept of Geography, University of Bristol, UK

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 3

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS



“Without the infrastructure that helps scientists manage their data in a convenient and efficient way, no culture of data sharing will evolve.”

Stefan Winkler-Nees
(Deutsche Forschungs-Gemeinschaft, DFG)

7 UNIVERSITÄT POTSDAM APA ALLIANCE FOR POLAR AND MARINE RESEARCH HELMHOLTZ ASSOCIATION Opportunities for Data Exchange

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 4

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS



“[Researchers would prefer] just one point of access to all data, which would be simple to use and ‘fool proof.’”

But she suspects it is wishful thinking to ask for Google-like simplicity when one looks for **“chlorophyll data in the Atlantic at 200 meters depth”**

Karin Lochte
(Alfred Wegener Institute for Polar and Marine Research)

7 UNIVERSITÄT POTSDAM APA ALLIANCE FOR POLAR AND MARINE RESEARCH HELMHOLTZ ASSOCIATION Opportunities for Data Exchange

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 5

Earth System Science Data
The Data Publishing Journal

Home Online Library ESSD Online Library ESSD

ESSDD - Special Issue

MAREDAT **Towards a world atlas of marine plankton functional types**
Editor(s):

Database of diazotrophs in global ocean: abundances, biomass and nitrogen fixation rates
13 Feb 2012
V. W. Lutz, S. C. Dooney, L. A. Anderson, M. Benavides, A. Boile, S. Brunet, K. H. Brostrom, D. Blodje, D. G. Capone, E. J. Carpenter, Y. L. Chen, M. J. Church, J. E. Dore, L. J. Falcón, A. Fernández, R. A. Foster, K. Funuyó, F. Gómez, K. Gundersen, A. M. Hynes, D. M. Karl, S. Kitajima, R. J. Langlois, J. LaRoche, R. M. Letelier, E. Marañón, D. J. McGillicuddy Jr., P. H. Moseander, C. W. Moore, B. Nowirko-Capraloto, M. R. Ruellet-Vivien, J. A. Neundorfer, K. M. Grout, A. J. Poulton, P. Raimbault, A. P. Rees, L. Riemann, T. Shiozaki, A. Subramaniam, T. Tyrrell, K. A. Turk-Kubus, M. Varela, T. A. Villareal, E. A. Webb, A. E. White, J. Wu, and J. P. Zehr
Earth Syst. Sci. Data Discuss., 5, 47-106, 2012
⇒ Abstract ⇒ Discussion Paper (PDF, 3215 KB) ⇒ Supplement (79 KB)
⇒ Interactive Discussion (Closed, 4 Comments) ⇒ Final Revised Paper (ESSD)

A global diatom database – abundance, biovolume and biomass in the world ocean
18 Apr 2012
K. Leblanc, J. Arstegui, L. Armand, P. Assmy, B. Bekker, A. Boile, E. Brunet, V. Cornet, J. Gibson, M.-C. Gosselin, E. Koppenscha, H. Marañón, J. Pedoja, S. Plankovski, A. J. Poulton, B. Quéguener, R. Smeets, R. Stoeck, J. Stoeck, M. A. van Leeuwe, M. Varela, C. Wildcombe, and M. Yallop
Earth Syst. Sci. Data Discuss., 5, 147-189, 2012
⇒ Abstract ⇒ Discussion Paper (PDF, 1788 KB) ⇒ Interactive Discussion (Closed, 5 Comments) ⇒ Final Revised Paper (ESSD)

Distribution of known macrozooplankton abundance and biomass in the global ocean
20 Apr 2012
R. Mortarty, E. T. Dauterhals, C. Le Quééré, and M.-P. Gosselin
Earth Syst. Sci. Data Discuss., 5, 187-226, 2012
⇒ Abstract ⇒ Discussion Paper (PDF, 1345 KB) ⇒ Interactive Discussion (Open, 3 Comments)
Manuscript under review for ESSD

Picophytoplankton biomass distribution in the global ocean
27 Apr 2012
E. T. Dauterhals, W. K. W. Li, D. Valiela, M. W. Lomas, M. Landry, F. Partensky, D. M. Karl, O. Ulloa, L. Campbell

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 6

Perspectives – “Records Management”

- DoW APARSEN Task 2430
Provenance Interoperability and Reasoning
 - A single acquisition process may create thousands of images and some terabyte of data. The complex processes yielding massive intermediate data and multiple versions of final products, reprocessing with improved methods or corrected input give rise to a need for complex generic reasoning over provenance data ... such as inheritance ... merging metadata of intermediate steps, relevance, assessment, obsolescence control, "garbage collection" and ...
- Credibility, Simplification, Selection/ Relevance, ...

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 7

APARSEN WP26 Annotation, an Example

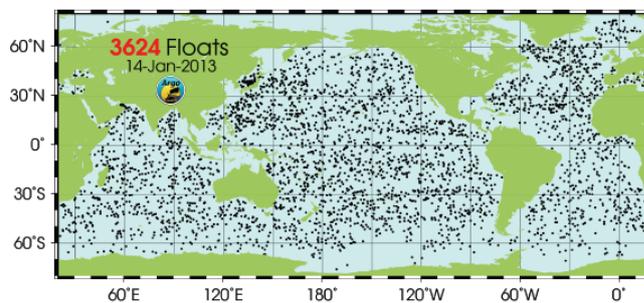
NSF CLIVAR S04P cruise report

Station /Cast	Sample No.	Quality Property	Code	Comment
21/1	134	ctds	3	Shiproll plus poor mixing in gradient cause much CTD signal oscillation during bottle stop. Code CTD salinity questionable, CTD is okay, just does not compare well with bottle salinity.
21/1	134	salt	2	Bottle salinity is low compared with CTD, agrees with adjoining stations. Variation in CTD at bottle trip. Salinity as well as oxygen and nutrients are acceptable.

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

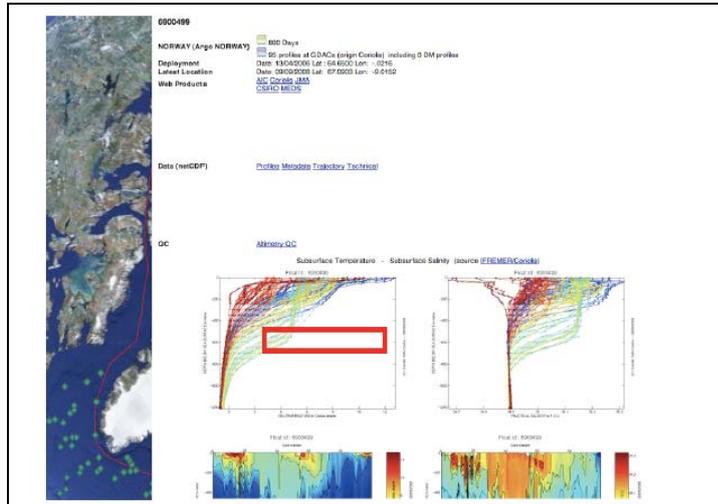
Folie 8

The biggest experiment in the world (not at CERN!)



H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 9



Folie 10

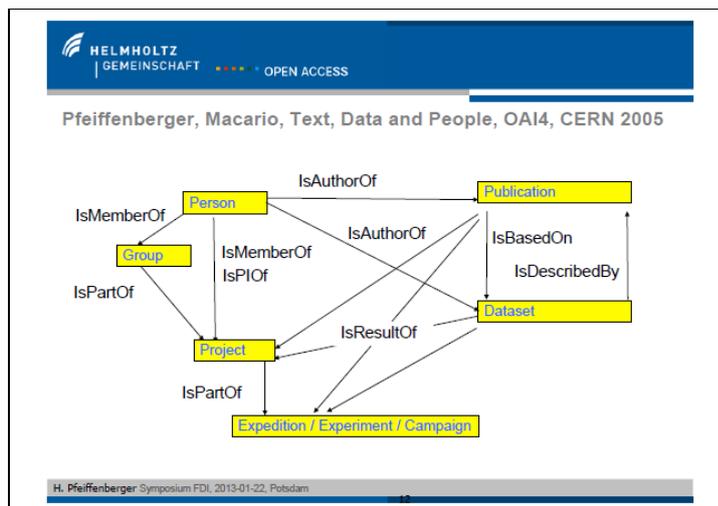
HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

An important, typical Experiment

- EISENEX / EIFEX : Two expeditions of "Polarstern" :
 With a few tons of iron fertilizer, south of Capetown
- EIFEX (2004):
 - 54 scientists and students from
 - 14 institutes and 3 companies from
 - 7 EU countries and South Africa
 - Oceanographers
 - Biologists
 - Chemists.....
- "Biogeochemistry"
- + Satellite observations !

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 11



Folie 12

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

eXpedition (in production, since 2005)

Related information: ["Reports on Polar and Marine Research"](#) (1982 to date)
[Primary data](#) (all polarstern datasets in PANGAEA)
[Handbook and scientific device documentation](#) (in deutsch)
[DShip](#) (Polarstern Data Acquisition System)
[VirtualPSI: Virtual Polarstern Tour](#)

Expedition	Date Port	Region Research	Publications & Primary Data	Details
ANT-XXI/3	21.01.2004 - 25.03.2004	Atlantic/Indian Ocean, Polar frontal zone Biology EIFEX	ePIC: Publications ePIC: Reports on Polar and Marine Research ePIC: Weekly reports PANGAEA: Stations PANGAEA: Datasets <i>(Note: Publications and datasets for recent cruises may not yet be available)</i> Meteorology	
ANT-XXI/4	27.03.2004 - 06.05.2004	Lazarev Sea Biology, Kribi, GLOBEC	ePIC: Publications ePIC: Reports on Polar and Marine Research	

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 13

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

eXpedition (Alpha)

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 14

HELMHOLTZ GEMEINSCHAFT OPEN ACCESS

PANGAEA - Elsevier

Purchase PDF (743 K) Export citation

Abstract Article Figures/tables References

Marine Microgeobotany
Volume 09, Issues 3-4, 20 February 2008, Pages 192-207

doi:10.1016/j.microm.2007.09.005 | How to Cite or Link Using DOI
Copyright © 2007 Elsevier B.V. All rights reserved.
Permissions & Reprints

Organic matter rain rates, oxygen availability, and vital effects from benthic foraminiferal $\delta^{13}\text{C}$ in the historic Skagerrak, North Sea

Sylvia Brückner and Andreas Mackensen

Walter Wegener Institute for Polar and Marine Research, Columbusstr. D-27568 Bremerhaven, Germany
Received 27 March 2007; revised 21 September 2007; accepted 24 September 2007; Available online 4 October 2007.

Abstract
The sediment cores 225514 and 225510 were recovered from 420 and 285 m water depth, respectively. They were investigated for their benthic foraminiferal $\delta^{13}\text{C}$ during the last 500 years.

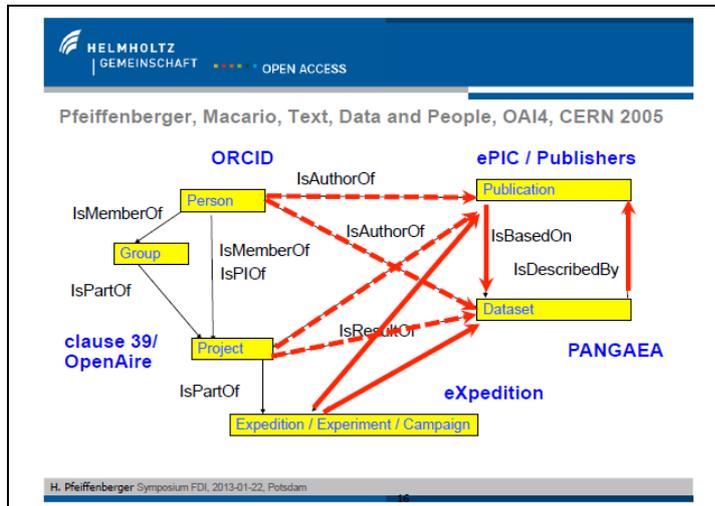
PANGAEA[®] - Supplementary Data
Stable carbon isotopic composition of benthic foraminifera from sediments of the Skagerrak...

Related Articles

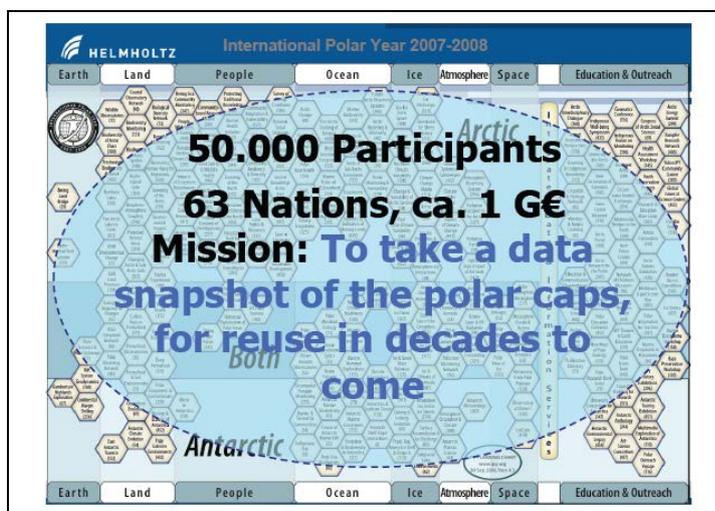
- The tropical carbon and productivity changes off north...
- Marine Microgeobotany
- Transport variability in living deep-sea benthic foraminifera...
- Earth Science Reviews
- Early Mesozoic benthic foraminiferal assemblages from...
- Marine Microgeobotany

H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

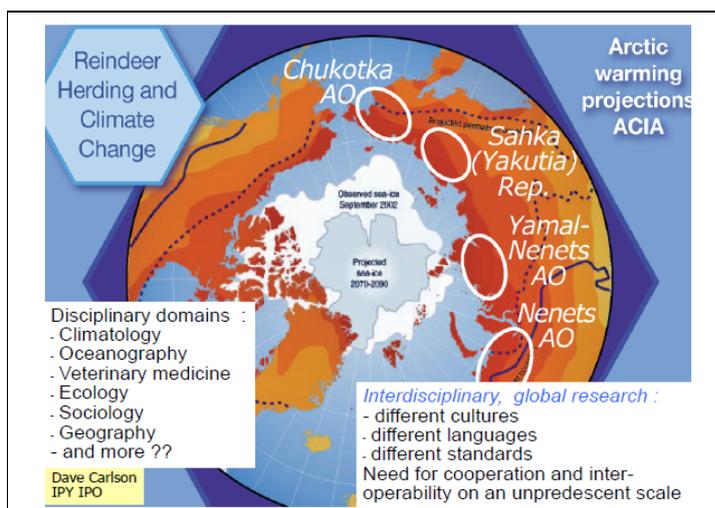
Folie 15



Folie 16



Folie 17



Folie 18

MaNIDA – Enabling Data-Intensive Marine Science



Global Change

- Assessing, understanding, and predicting environmental changes
- Human environmental impact

Hazards

- Risk analysis and support for disaster management
- Understanding environmental factors affecting human health

Resources

- Sustainable ecosystem management
- Energy from the ocean

Folie 19

Status of Marine Science Data in Germany



Extremely wide range of data sources

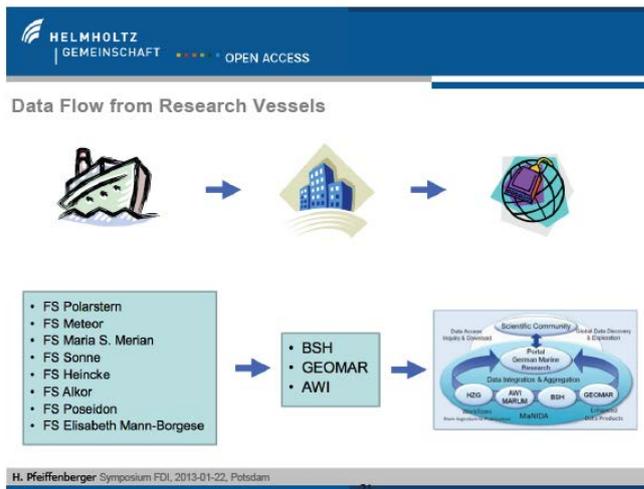
- Ship-based instruments
- Instruments in water column and at seafloor
- Air- and space-borne instruments
- Sensor networks (increasingly in the deep ocean)

- “snapshots“ (individual projects) and long-term monitoring

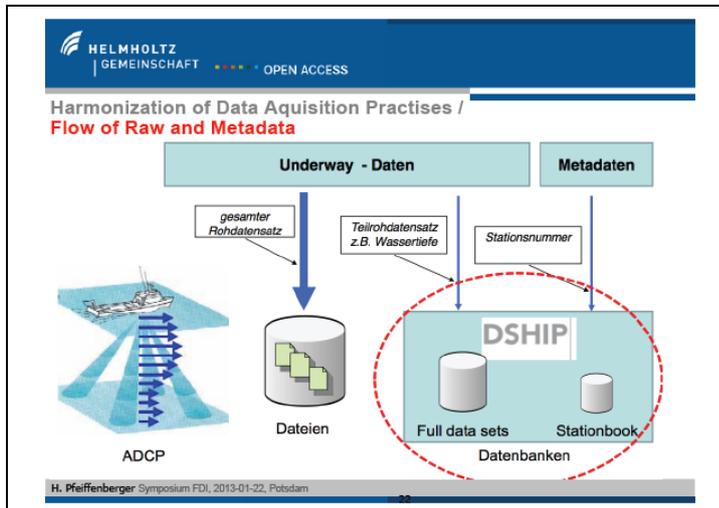
- stored in different data centers (e.g. AWI/PANGAEA, BSH/DOD, HZG/COSYNA)

→ Requirement : Data-intensive research through **coherent data portal with common access strategy**

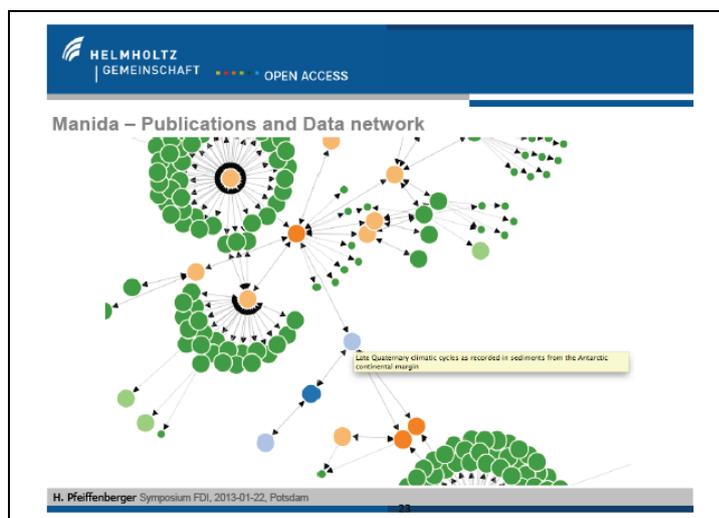
Folie 20



Folie 21



Folie 22



Folie 23

- HELMHOLTZ
GEMEINSCHAFT OPEN ACCESS
- ### Zusammenfassung
- **Wir helfen der Wissenschaft zurück auf dem Weg zur Reproduzierbarkeit**
 - Vielfältige Datenquellen
 - Verstreute Datenquellen
 - alle notwendigen Kontext/Abstammungsinformationen
 - Publikationen als "Hubs" / beste verfügbare Metadaten/ Kontextbehälter
 - **Modularisierung des Problems**
 - Vereinfachung/Separierung der Probleme
 - Rohdaten => qualitätsgesicherte Primärdaten => abgeleitete/ Datenprodukte
- H. Pfeiffenberger Symposium FDI, 2013-01-22, Potsdam

Folie 24

Thank you!

www.awi.de

www.manida.org

oa.helmholtz.de