

Sammelband

Seminar
Big Data Applications

Sommersemester 2013

Achim Streit, Rainer Stotzka

- 1) Steinbuch Centre for Computing
- 2) Institut für Prozessdatenverarbeitung und Elektronik



Editorial

Im Sommersemester 2013 wurde am Karlsruher Institut für Technologie das Seminar „Big Data Applications“ durchgeführt. Fachlich und formal begleitet wurde das Seminar von Prof. Dr. Achim Streit (Steinbuch Centre for Computing), Danah Tonne und Dr. Rainer Stotzka (beide Institut für Prozessdatenverarbeitung und Elektronik). Studierende der Studiengänge Informatik und Informationswirtschaft haben teilgenommen, indem sie sich in ein Big Data Thema eingearbeitet, einen Vortrag gehalten und eine Ausarbeitung angefertigt haben.

In diesem Sammelband finden Sie folgende Ausarbeitungen:

1. Meryem Cömert:
Digitalisierungs- und Bereitstellungsprojekte in Deutschland
2. Dominik Sauter:
Bitstream Preservation: Verfahren und Algorithmen
3. Frederic Markert:
Schafft Replikation Sicherheit? Kosten vs. Nutzen
4. Anjela Mayer:
Empfehlungen für die Langzeitarchivierung: Nestor, OAIS – TRAC, Data Seal of Approval
5. Mert Turan:
Metadaten und -standards: Welche Informationen sind wichtig?
6. Carsten Griesheimer:
Persistent Identifier Systems: Wie erzeuge ich PIDs?
7. Nico Kopp:
Ist das Open Archival Information System (OAIS) für Big Data geeignet?
8. Gül Kavak:
ESciDoc: Geeignet für Big Data?

Die oben genannten Autorinnen und Autoren sind alleine inhaltlich und formal für ihre Ausarbeitung verantwortlich und haben ihre schriftliche Einverständnis für die Veröffentlichung gegeben.

Beim Studium dieser spannenden Themen wünsche ich Ihnen viel Spaß und hoffe, Sie hiermit für „Big Data“ begeistern zu können.

Karlsruhe, 20. August 2013

Ihr Rainer Stotzka

Seminararbeit

Big Data Applications

Digitalisierungs- und Bereitstellungsprojekte in Deutschland

vorgelegt von
Meryem Cömert

Inhaltsverzeichnis

1 Motivation	3
2 Was ist Digitalisierung?	4
3 Ziele der Digitalisierung deutscher Bestände	4
4 Bilddigitalisierung	5
5 Aufnahmetechniken	6
6 Textdigitalisierung	7
7 OCR-Erkennung	7
8 Metadaten	8
9 Langzeitarchivierung	8
9 Kriterien zur Publikation: Deutsche Digitale Bibliotheken	9
10 Staatliche Kunsthalle Karlsruhe	10
11 Digitalisierte Altbestände der KIT-Bibliothek	
11	
12 Europeana	
12	
13 Die Deutsche Digitale Bibliothek	
12	
14 Die Deutsche Forschungsgemeinschaft (DFG)	13
15 Schlussbetrachtung	14
16 Quellenverzeichnis	15

Motivation

„Seek and you will find“ [Mathew 7:7]

versus

„If you can't find it, it doesn't exist“

Im Zeitalter der fortgeschrittenen Digitalisierung ist es immer einfacher geworden sich aus dem Überangebot an Informationen das Notwendige herauszusuchen. Dennoch ist es in dem für uns relevanten Kontext der „Big Data“ nicht auszuschließen, dass Daten gegebenenfalls noch nicht vorhanden oder nur schwer zu erzeugen sind. Nun ergeben sich folgende Fragestellungen:

- Ist wirklich alles digital Unauffindbare inexistent?
- Welche Schwierigkeiten ergeben sich beim Digitalisierungsprozess?
- Was möchte man mit der Digitalisierung erreichen?
- Und wie breit ist das Spektrum der digitalisierten Bestände in Deutschland eigentlich?

Auf diese und ähnliche Fragen möchte ich im Laufe dieser Seminararbeit eine Antwort geben.

Ziel dieser Arbeit soll außerdem sein, Techniken und Methodiken der Digitalisierung zu untersuchen und zu bewerten.

Was ist Digitalisierung?

Unter Digitalisierung versteht man allgemein die Überführung von physikalischen Einheiten in die digitale Form. So wird die Umwandlung von Informationen wie Ton, Bild oder Text in Zahlenwerte zum Zwecke ihrer elektronischen Bearbeitung, Speicherung oder Übertragung unter diesem Begriff zusammengefasst (Hans-Bredow Institut, 2006: 95). Physikalischer Einheiten sind hierbei Bestände wie beispielsweise Bücher, Bilder, Skulpturen und andere Materialien. Diese müssen sinnvoll und auf geeignete Weise organisiert werden.

Vor dem Digitalisierungsprozess muss die Verfügbarkeit der physikalischen Einheiten überprüft werden. Liegt das jeweilige Material bereits digitalisiert und in guter Qualität vor, soll darauf hingewiesen werden; denn dadurch können finanzielle, sowie zeitliche Kosten eingespart werden. Ebenfalls sollte eine geeignete Gruppe an Arbeitskräften für die Durchführung zur Verfügung stehen. Eine der Hauptkriterien, die eine perfekte Digitalisierung (ausführliche Erklärung auf Seite 9) erfüllen muss, ist die Schonung des eventuell fragilen Originals. Hierbei spielt die jeweilige Beschaffenheit des Materials eine große Rolle. Durch die Zuhilfenahme bestimmter Aufzeichnungstechniken, die im Folgenden erläutert werden, wird ein digitales Profil erstellt. Abschließend folgt die Langzeitarchivierung (siehe Seite 8), die die Bereitstellung der Medien für die Öffentlichkeit über Jahrhunderte hinweg gewährleisten soll. Um eine geordnete Übersicht zu garantieren werden nach der vollständigen Digitalisierung Metadaten zugeordnet (Seite 8f).

Ziele der Digitalisierung deutscher Bestände

Durch den immer schneller voranschreitenden Prozess der Globalisierung hat die Digitalisierung an Bedeutung gewonnen. Dieser Entwicklung verdanken wir wichtige, positive Wirkungen, wie etwa verbesserte internationale Zusammenarbeit, dichte technische Vernetzungen und schnellen Informationsaustausch. Dies stellt uns aber gleichzeitig vor die Herausforderung, mit wachsenden Datenmengen und den daraus resultierenden Folgen und Problemen umzugehen.

Zu diesem Zweck sollen wissenschaftliche Prozesse katalysiert werden und bei qualitativ hochwertiger Wiedergabe des Originals eine Plattform geschaffen werden, worauf die Forschung zeit- und ortsunabhängig Zugriff haben soll. Das „Wissen“ in Form von Daten wird dynamischer organisiert, Bibliotheken und Museen mit deutschem Kulturerbe werden

global zugänglich gemacht, sowie mit anderen verflochten (Deutsche Digitale Bibliothek 2012).

Bis diese Ziele in ausreichendem Maße umgesetzt werden können, liegt jedoch noch ein weiter Weg vor uns. Mit welchen Schwierigkeiten dieser Weg verbunden ist möchte ich im Folgenden veranschaulichen.

Bilddigitalisierung

Im privaten, wie auch im wissenschaftlichen Bereich hat die digitale Fotografie die analoge Fotografie längst überholt (Mumenthaler 2005: 44). Es ist jedoch von großer Wichtigkeit Altbestände der analogen Fotografie zu digitalisieren, um historische Werke länger verfügbar zu machen.

Um ein Bild in bester Qualität zu digitalisieren, sind die Farbtiefe, sowie die Auflösung ausschlaggebend (Mumenthaler 2005: 44). Jedes Bild ist eine Zusammenfügung von kleinsten Quadranten, den Pixeln. Ein Pixel misst den Wert in Bits, wodurch die Farbtiefe bestimmt wird. Ein Mehr an Farbtiefe bedeutet eine präzisere Darstellung des Bildes. Ist bei Graustufenbildern der Wert des Bits gleich eins, so sind zwei unterschiedliche Farben, also schwarz und weiß möglich. Bei acht Bit sind bereits 256 unterschiedliche Farbstufen von schwarz nach weiß erfassbar.

Die Bilddigitalisierung wird in verschiedene Modi unterteilt. Farbbilder werden für die Bildschirmanzeige im RGB-Modus (rot-grün-blau-Modus) angegeben. Verständlicherweise spielt hier die Farbtiefe eine größere Rolle. Ist der Wert des Bits eines Farbbildes im RGB-Modus beispielsweise gleich 16, so sind 48 Bit Bilder mit zwei hoch 48 unterschiedlichen Farbtönen möglich. Trivial betrachtet können bei einer höheren Anzahl an Bits schärfere Bilder erhalten werden. An dieser Stelle ist es sinnvoll, sich Gedanken über die Auflösung und den daraus resultierenden Speicherbedarf zu machen. Die Anzahl aller Pixel legt die Auflösung des Bildes fest. Damit das menschliche Auge ein Bild als „scharf“ empfindet ist eine Mindestauflösung von 300 dpi („dots per inch“) erforderlich (DFG-Praxisregeln zur Digitalisierung 2013: 9).



Quelle: <http://tinyurl.com/dpi300>

Schauen wir uns nun den benötigten Speicherbedarf für dieses Verfahren an. Betrachtet man, beispielsweise, ein Bild mit den Maßen zehn mal zehn Zentimeter, einer Auflösung von den oben genannten 300 dpi, sowie einer Farbtiefe von 24 Bit (das würden ca. 16.77 Millionen unterschiedlichen Farben im RGB-Modus entsprechen), so sind in diesem Bild 1181x1181, also 1.394.761 Pixel vorhanden. Der insgesamt erforderliche Speicherbedarf für dieses Bild betrüge 4,2 Megabyte. Wenn wir also annehmen, dass eine Million DIN A4 Seiten mit diesen Werten digitalisiert werden müssen, so stellen wir fest, dass sich eine Datenmenge von 27 Terabyte ergibt. Es lässt sich also festhalten, dass es an immensen Speicherplatz bedarf, den es gilt stabil zu organisieren.

Aufnahmetechniken

Elementare Bestandteile der Digitalisierung sind die verschiedenen Aufnahmetechniken. Die Methoden des Zeilenscan und des Flächensensor sind hierbei die bekanntesten Beispiele.

Bei der Methode des Zeilenscan wird die Bildinformation vom Scanner (Flachbett- oder Filmscanner) linienförmig abgelesen (Mumenthaler 2005: 44). Anders sieht es beim Flächensensor aus, bei diesem die gesamte Bildinformation gleichzeitig erfasst wird (Mumenthaler 2005: 44).

Die jeweilige Bildvorlage bestimmt das Verfahren. Um qualitativ hochwertig zu digitalisieren, liegt der Fokus liegt auf dem jeweiligen optischen System, sowie auf der Leistungsfähigkeit; zudem sind die Präzision des Mikrochips und die Zuverlässigkeit der Mechanik von großer Relevanz (Mumenthaler 2005: 44). Abschließend wird eine Bitmap (Bilddatei) erstellt, welche entsprechend archiviert werden muss (siehe Seite 8).

Nun gilt es jedoch festzuhalten, dass ein Digitalisierungsvorgang nicht die Erfassung sämtlicher Daten gewährleisten kann. Dies lässt sich vereinfacht darstellen, als dass die Kopie eines Bildes auch nur ein Abbild und nicht mehr das Original an sich ist. Dem kommt hinzu, dass die Bilder in komprimierter Form abgespeichert werden.

Hierzu gibt es verschiedene Formate: GIF (Graphics Interchange Format) und PNG (Portable Network Graphics) formatieren die Bilder vollständig, das bedeutet, dass die Bilder auf dem Bildschirm komplett wiederhergestellt werden (Wenk 2004: 3). Bei der Speicherung in JPEG (Joint Photographic Expert Group) Formaten hingegen sind Verluste unumgänglich; auf manche Bildinformationen können auch nicht mehr zugegriffen werden. Der Vorteil bei JPEG Speicherung ist allerdings der geringe Platzbedarf (Wenk 2004: 3).

Textdigitalisierung

Bei der Digitalisierung von Texten sind weniger hohe Qualitätsanforderungen notwendig, als bei einer Bilddigitalisierung, obwohl die Verfahren die Gleichen sind.

Um einen Text zu digitalisieren wird zunächst gemäß dem Bilddigitalisierungsverfahren eine Bitmap erstellt. Um die Daten, die nun in der entsprechenden Bilddatei vorhanden sind, in maschinenlesbare, sowie durchsuchbare Dateien umzuwandeln, wird die OCR-Erkennung (Optical Character Recognition) genutzt (Mumenthaler 2005: 46). Hierbei werden innerhalb des Bitmaps Einzelzeichen anhand bestimmter Schemata (Wiederholungen, Ähnlichkeiten etc.) erkannt. Nach diesen Wahrscheinlichkeitsannahmen identifiziert das Programm die Buchstaben und gibt diese wieder. Abschließend werden den Buchstaben Zahlenwerte gemäß üblicher Textkodierungen (ASCII, Unicode) zugeordnet. Probleme, die sich bei OCR ergeben, sind eventuelle Fehlerkennungen. Zur Verbesserung tragen spezielle Softwares wie FineReader und Omnipage bei (Mumenthaler 2005: 46). Die manuelle Korrektur ist dennoch nicht auszuschließen.

Um die etwaige Größe des Speicherbedarfs zu erfassen, ist es hinreichend zu wissen, dass für eine gewöhnliche DIN A4 Seite mit 63 Zeilen zu je 80 Zeichen nur ca. fünf Kilobyte erforderlich sind. Bei ebenfalls einer Million DIN A4 Seiten entspricht das einer Größe von gerade mal 0.0745058 Terabyte.

OCR-Erkennung

Das OCR-Verfahren zählt zu den ersten Pilotprojekten, mit der, mit einer ca. 99%igen Wahrscheinlichkeit die Buchstaben genau erkannt werden können (DFG-Praxisregeln 2013: 30). Das Verfahren läuft in Stufen ab, die aufeinander aufbauen.

Als Erstes wird das Bild bei der Binarisierung in ein bitonales Format übersetzt, wobei die Qualität des jeweiligen Bildes entscheidenden Einfluss auf diesen Prozess hat (DFG-Praxisregeln 2013: 36).

Infolgedessen wird die Segmentierung durchgeführt, die die Identifikation von Koordinaten versucht, um Textbereiche von Illustrationen zu trennen (DFG-Praxisregeln 2013: 36). Im dritten Schritt folgt der sogenannte „Pattern Matching Process“, welcher meistens durch ICR (Intelligent Character Recognition) unterstützt wird (DFG-Praxisregeln 2013: 36). Im Grunde genommen können diese drei Schritte getrennt voneinander vorgenommen

werden. Durch OCR soll es dem Endnutzer möglich sein, den Text zu durchsuchen und direkten Zugriff drauf zu haben. ALTO Standards der Library of Congress (ein XML-Schema um Daten zu beschreiben) werden empfohlen, um die OCR-Daten geeignet zu nutzen (DFG-Praxisregeln 2013: 36).

Metadaten

Die Beschreibung eines Objekts geschieht durch die Metadaten; diese enthalten also detailreiche Informationen über die jeweiligen Informationen. Von Interesse für die Organisation immenser Datenmengen, sowie für die Übersicht sind die bibliographischen Metadaten, wie Autor, Titel etc.. Die Metadatenvergabe für das Bestands-, sowie das digitale Objekt erfolgt unter einer strengen Qualitätskontrolle, die die Vollständigkeit und Korrektheit überprüft, um eine makellose Präsentation sicherzustellen. Ohne ordentliche Metadaten ist keine zweckmäßige Digitalisierung möglich.

Des Weiteren unterscheidet man zwischen administrativen und beschreibenden Metadaten (Mumenthaler 2005: 54). Die Beschriftung von Datenträgern, Signaturen und Standorten zählen zu den administrativen, diejenigen, die auf einer Liste, in einer Datenbank oder im Header eines TIFF (Tagged Image File Format)-Bildes vorzufinden sind, zu den beschreibenden (Mumenthaler 2005: 54).

Langzeitarchivierung

Digitale Daten zeichnen sich dadurch aus, verlustfrei kopiert werden zu können (Mumenthaler 2005: 49). Bei digitaler Kopie werden nur die Zahlen eins und null kopiert, während bei der analogen Kopie eventuelle Störungen, wie Alterungsprozesse des Originals, die Qualität der Kopie einschränken können. Trotz der einfachen Handhabung digitaler Daten müssen wir jedoch unterstreichen, dass das Hauptproblem der Digitalisierung zwar nicht bei sukzessivem Verfall, sondern totalem Informationsverlust liegt. Sind Daten einmal verloren, so ist es unmöglich, diese wieder zu erhalten.

Die begrenzte Lebensdauer der Datenträger beeinträchtigt das reibungslose Archivierungsverfahren ebenfalls. Dafür gibt es jedoch keine eindeutige Lösung. Die Anforderungen an den Datenträger sind Robustheit, Zuverlässigkeit, Standardisierung und eine möglichst hohe Speicherkapazität (Mumenthaler 2005: 52). Außerdem sollte eine Abhängigkeit von einem Lesegerät ausgeschlossen sein, sodass die Kompatibilität zu

zukünftigen Geräten möglich ist (Mumenthaler 2005: 52). Während die laufende Festplatte eine allgemeine Lebensdauer von maximal zehn Jahren hat (Feddern 2007), liegt die Lebensdauer eines USB-Sticks bei ca. 30 Jahren (Gilbert 2003) und bei optimaler Lagerung eines Mikrofilms geht man von einer Haltbarkeit von bis zu 500 Jahren aus (Fachverband für Multimediale Informationsverarbeitung e.V. 2011). Andere Speichermedien, wie CDs und DVDs fallen aus der Liste geeigneter Medien heraus, da diese über eine zu niedrige Speicherkapazität verfügen und für „Big Data“ nicht in Frage kommen.

Wichtig für die Langzeitdatenlagerung ist zudem die Datenpflege unter ständiger Qualitätsprüfung und Sicherung der Metadaten (Mumenthaler 2005: 53). Um einen Erfolg der Langzeitarchivierung zu beobachten ist die Schaffung technischer und wirtschaftlicher Rahmenbedingungen zu fokussieren (Handbuch der Digitalisierung 2010: 79). In Deutschland konzentrieren sich vor allem das „Network of Expertise in long-term Storage and availability of digital“ („nestor“), als auch die Deutsche Nationalbibliothek auf diese Problematik. Im ersteren Kompetenznetzwerk gibt es deutschlandweite Arbeitsgruppen, die mit nationalen als auch internationalen Erfahrungen an den Problemfeldern der Langzeitarchivierung arbeiten um diese zu verbessern (Wetzsenstein 2010: 19).

Kriterien zur Publikation: Deutsche Digitale Bibliotheken

Die wichtige Frage, die nun beantwortet werden muss, ist die nach den wichtigsten Kriterien und Ansprüchen, die eine einwandfreie, deutsche, digitale Bibliothek ausmacht.

Die Antwort hierauf gibt uns die DFG (Deutsche Forschungsgemeinschaft), die als einer der wichtigsten Förderer für die deutsche Grundlagenforschung agiert (Stifterverband für die Deutsche Wissenschaft e.V. 2012). Für diese Kriterien hat die DFG bestimmte Praxisregeln ausgelegt (Information für die Wissenschaft Nr. 07 | 6. Februar 2013).

Die Strategien der jeweiligen Institutionen bestimmen, wie die digitalen Inhalte organisiert werden. Die Erfolgsfaktoren hängen von den jeweiligen organisatorischen, sowie wirtschaftlichen Rahmenbedingungen ab (DFG-Praxisregeln – Digitalisierung 2013: 36).

Als Referenzmodell soll das Open Archival Information System (OAIS) angewandt werden (DFG-Praxisregeln – Digitalisierung 2013: 36). Zudem müssen gesetzliche Rahmenbedingungen beachtet werden.

Was für die perfekte Digitalisierung von hoher Bedeutung ist, ist die Wahl geeigneter Metadaten, um so viel Verlust wie möglich zu vermeiden. Außerdem ist wichtig, dass trotz

digitalisierter Bestände die physikalischen Originale nicht verloren gehen. Dies stellt einen zuzüglichen Aufwand dar, der nebst mehrerem Terabyte, die je nach Projekt nötig sind, finanziert werden muss.

Zur Datenhaltung sind Trägermaterialien wie Bandlaufwerke (streamer) und Festplatten von Wichtigkeit (DFG-Praxisregeln – Digitalisierung 2013: 36). Sind häufige Zugriffe auf den digitalen Master notwendig, so sind Bandarchivsysteme keine optimale Lösung, da eventuell ein Zwischenspeicher zugefügt werden muss (DFG-Praxisregeln – Digitalisierung 2013: 36). Festplattensysteme eignen sich daher eher für die Datenhaltung, da sie einen schnelleren und unkomplizierteren Zugriff auf die Daten erlauben.

Die Organisation übergroßer Datenmengen im Rechenzentrumsbetrieb erfolgt in Form der „Network Attached Storage“ (NAS)-Systeme, des „Storage Attached Network“ (SAN) oder der „Content Addressed Storage“ (CAS).

Zu einer weiteren Anforderung der DFG-geforderten Projekte zählt die Aufgabe, nachvollziehbare Aussagen zur institutionellen Langzeitsicherung und -archivierung, abzugeben (DFG-Praxisregeln – Digitalisierung 2013: 36).

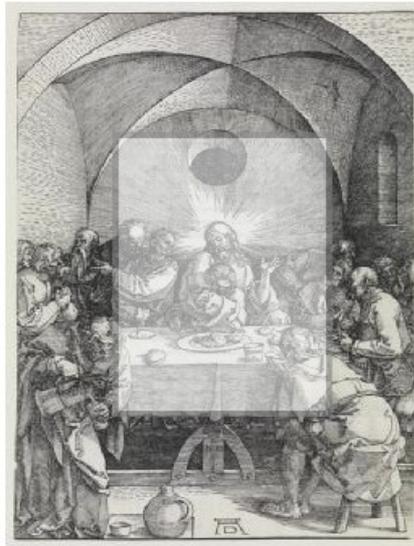
Staatliche Kunsthalle Karlsruhe

Aufgrund hoher Komplexität und etlicher digitalisierter Bestände in Deutschland habe ich mir als Beispiel für die Bilddigitalisierung die Staatliche Kunsthalle Karlsruhe ausgesucht. Seit 2011 wurden 1100 der Objekte, also mehr als ein Drittel des Gesamtbestandes der Kunsthalle, aus der Gemälde- und Skulpturensammlung zu der Online-Mediathek hinzugefügt, zu der man öffentlichen Zugang hat. Mithilfe der Metadaten wie Künstler- und Bildnamen ist es dem Benutzer möglich einen Online-Besuch durchzuführen. In Zukunft sollen mehr Informationen wie beispielsweise Kurztexpte zu den Werken ergänzt werden. Bemerkenswert ist, dass der Zugang zu den Bildern mit einer jeweiligen Druckfunktion sowie genauen Ansicht mit Lupeneffekt optimiert wurde.

Albrecht Dürer

1471 - 1528

Das letzte Abendmahl



Druckansicht



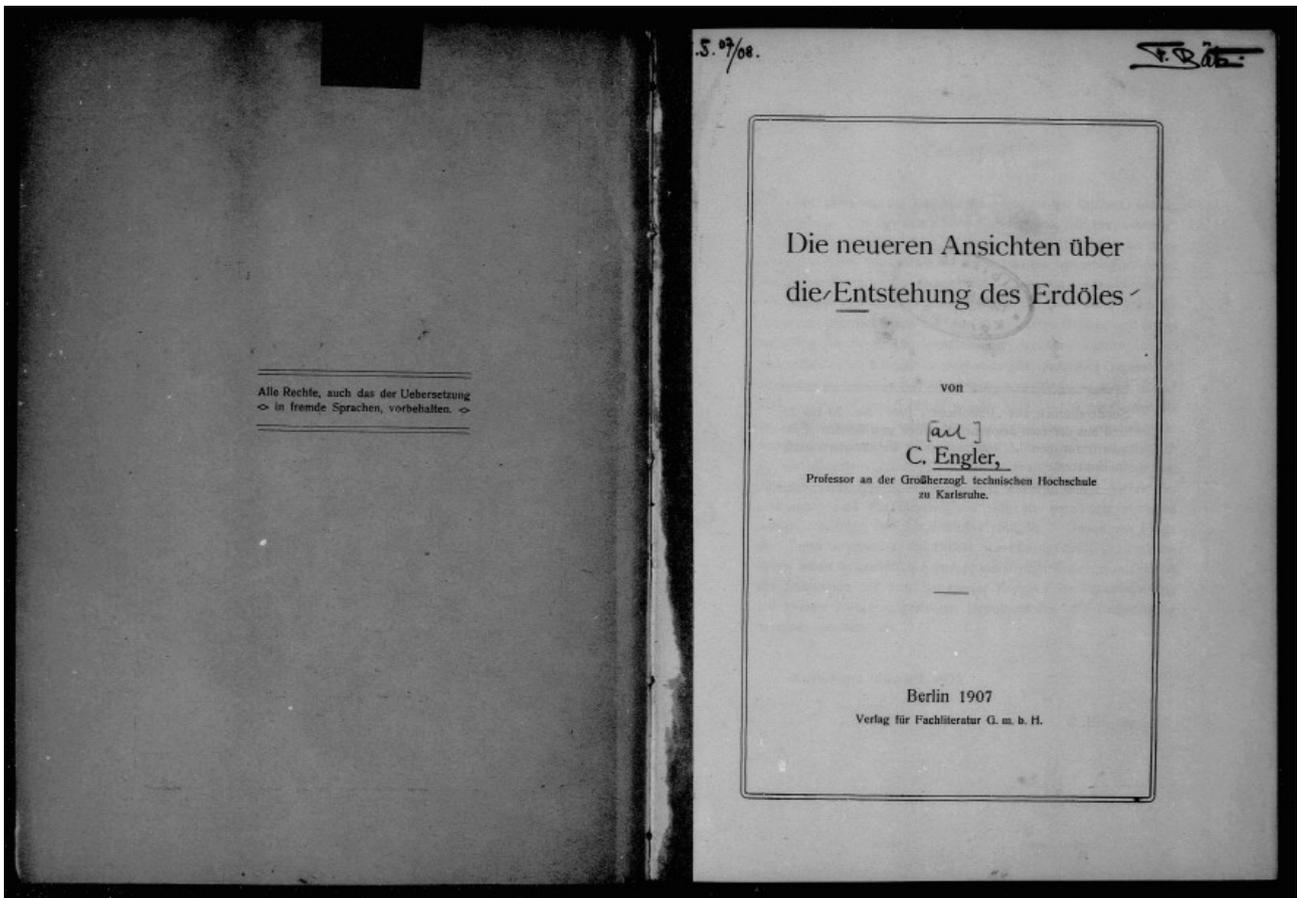
Bildbestellung

Quelle: <http://swbexpo.bsz-bw.de/skk/detail.jsp?id=696D3FE0452582C085C1B5B52BCE1CAF&img=1>

Digitalisierte Altbestände der KIT-Bibliothek

Zahlreiche berühmte deutsche Akademiker wie beispielsweise Josef Durm (Architekt und Hochschullehrer) oder Carl Engler (Chemiker und Professor) studierten erfolgreich an der Universität Karlsruhe (TH). Sowohl während als auch nach ihrer Bildungslaufbahn brachten diese etliche besondere Werke hervor. Bemerkenswert ist, dass auf ihren Erkenntnissen die Basis der heutigen industriellen Infrastruktur Baden-Württembergs beruht (KIT- Bibliothek 2013).

Da ihre Bücher zu den Altbeständen der KIT-Bibliothek zählen ist ein sukzessiver Verfall dieser Werke nicht auszuschließen, weshalb sie digitalisiert werden müssen. Auf der Website der KIT-Bibliothek kann auf die jeweiligen E-Books zugegriffen werden, welche sogar mit einer Downloadfunktion ausgestattet sind. Die Suchfunktionen bieten Präzisierungen in Form von Erscheinungsbereichen, Materialarten (Bücher, Zeitschriften etc.), Sprachen und Suchorten an. Als jeweilige Metadaten sind unter anderem Abrufzeichen, Urheber, ISBN und Jahr erkennbar. Sobald eine Download-Funktion in Anspruch genommen wird lässt sich erkennen, dass die Bücher anhand einer Bilddigitalisierung erfasst worden sind. Eine OCR-Bearbeitung mit entsprechender durchsuchbarer PDF-Datei ist bei den meisten Objekten nicht verfügbar.



Quelle: <http://tinyurl.com/carlengler>

Europeana

Im Europäischen Bereich wurde das Internetportal der digitalen Bibliothek Europeana erstellt. Hier sind Kulturgüter der Mitgliedsstaaten der EU zusammengestellt und zugänglich gemacht worden. Um die Forschung effizienter voranzutreiben fungiert Europeana multilingual. Die Digitalisate unterliegen keinem Copyright (Europe's Information Society 2007). Die Sammlung bei Europeana erfolgt lediglich durch die Metadaten der einzelnen europäischen kulturellen als auch wissenschaftlichen Institutionen. Durch Weiterleitung der relevanten Links auf Europeana gelangt man zu den jeweiligen Seiten. Die Metadatenvergabe muss dem Standard der Europeana Semantic Elements (ESE) entsprechen, um bessere Suchfunktionen zu ermöglichen (Europe's Information Society 2007).

Die Deutsche Digitale Bibliothek

Nachdem die Europäische Digitale Bibliothek Europeana entwickelt und Ende des Jahres 2008 freigeschaltet wurde, wurde die Initiative ergriffen ein deutsches Portal, nämlich die

Deutsche Digitale Bibliothek (DDB) zusammenzustellen. Auf diesem Portal sollen kulturelle und wissenschaftliche Werke zusammengetragen und verfügbar gemacht werden. Eine Vereinigung von über 30.000 Kultur- sowie Wissenschaftseinrichtungen Deutschlands steht in Aussicht. Mit 5,6 Millionen Objekten ging die DDB am 28. November 2012 online (Bundesregierung der Bundesrepublik Deutschland 2012). Die Finanzierung wurde bei der Initialisierung mit acht Millionen Euro vom Bund übernommen. Ab 2011 wurden vom Bund, von den Ländern und Kommunen 2,6 Millionen für fünf Jahre jährlich abgesichert (Deutscher Bibliotheksverband 2012).

Ein Ziel der DDB ist die Gewährleistung virtueller Museumsbesuche und der Austausch wissenschaftlicher Ergebnisse. Deutschland soll anhand der DDB Anschluss zur internationalen Wettbewerbsfähigkeit in Wissenschaft, Forschung und Bildung knüpfen (Deutsche Digitale Bibliothek: 2012). Signifikant dabei ist die sorgfältig filtrierte Information von fachkundigem Material, von dessen Authentizität der Nutzer sichergehen kann.

Hinsichtlich dieser Tatsache sollen praktische Metadaten bereitgestellt werden. Beispielsweise erhält man bei der Suche nach „Goethe's Faust“ nicht nur Literatur, sondern auch Illustrationen und Musikstücke. Dies ist eine große Herausforderung im Bereich der Big Data, da sämtliche Verfügbarkeiten gezielt organisiert sein müssen und auf dem Portal der DDB vernetzen werden zu können. Der freigeschaltete Zugang soll für jeden Nutzer zeitlich und lokal unbeschränkt sein. Im Hinblick auf die Zukunft soll dies kontinuierlich verbessert werden. Durch die DDB soll Deutschland seinen Platz in der weltweiten digitalen Zusammenarbeit finden.

Die Deutsche Forschungsgemeinschaft (DFG)

Zu den Aufgaben der DFG zählt die finanzielle Unterstützung von Forschungsaufgaben der Zusammenarbeit (Deutsche Forschungsgemeinschaft 2010). Sie beschreibt sich als Selbstverwaltungsorganisation, dessen Kernaufgabe darin besteht, die besten Forschungsvorhaben von Wissenschaftlern an Instituten und Hochschulen auszuwählen und finanziell zu unterstützen ((Deutsche Forschungsgemeinschaft 2010). Sie möchte leistungsfähige Informationssysteme für die Forschung unter überregionalen Gesichtspunkten aufbauen, deren Ergebnisse für die Wissenschaft frei und dauerhaft zugänglich sein sollen. Die Ausbesserung von Informations-Infrastrukturen in Deutschland steht dabei im Mittelpunkt (DFG-Praxisregeln – Digitalisierung 2013: 4).

Außerdem findet sich die DFG als zentrale nationale Förderorganisation in der Entwicklung des europäischen Forschungsraums wieder. Dabei sollen hohe Qualitätsstandards geschaffen werden,

die in dem wissenschaftlichen Wettbewerb angestrebt werden (Deutsche Forschungsgemeinschaft 2013). Die positive und produktive Vielfalt der verschiedenen nationalen Wissenschaftssysteme in Europa soll mit Forschungsstandort in Deutschland zur Gestaltung eines wissenschaftsorientierten Europäischen Forschungsraums genutzt werden. Ein weiteres Ziel sieht die DFG die grenzüberschreitende Kooperation um ein einheitliches europäisches Portal zur Verfügung zu stellen (Deutsche Forschungsgemeinschaft 2013).

Schlussbetrachtung

Bis heute hat sich die Digitalisierung zu einem hochkomplexen Phänomen entwickelt, dessen eindeutige Durchführung die internationale Gemeinschaft vor eine große Herausforderung stellt. Als Ausgangspunkt meiner Arbeit habe ich bewusst die Definition aus Matthäus 7:7 gewählt.

Betrachtet man die größten Internetportale digitalisierter Bestände der vergangenen 5 Jahre, so umfasst diese Analyse in meinen Augen viele zentrale Elemente des perfekten Digitalisierungsprozesses. Dabei ist in erster Linie ein Aspekt entscheidend – die weltweite Archivierung und Verbreitung von Daten und deren Organisation. Ist dem der Fall, so können nur noch wenige Gründe dagegen sprechen, dass eine Suche anhand ausgelegter Metadaten scheitert.

So müssen sich deutsche Institutionen immer wieder neu der Herausforderung stellen, auf dem internationalen Markt angemessen zu fungieren. Betrachtet man vergleichsweise den US-Konzern Google, der bekannter Weise alles Verfügbare digitalisiert, so stellt man fest, dass Deutschland mit seinen vergleichsweise kleinen Projekten weit hinten liegt (Lischka 2011). Nichtsdestotrotz lässt sich in meinen Augen ein klarer Fortschritt in der deutschen Digitalisierung erkennen.

Literaturverzeichnis

Bundesregierung der Bundesrepublik Deutschland 2012: *Stand des Projekts „Deutsche Digitale Bibliothek“*. In: „Antwort der Bundesregierung auf die Kleine Anfrage der Abgeordneten Dr. Petra Sitte, Jan, Korte, Dr. Rosemarie Hein, weiter Abgeordneter und der Fraktion DIE LINKE.“ Drucksache 17/9430

Online verfügbar unter:

dokumente.linksfraktion.net/drucksachen/26719_1709810.pdf [14.06.13]

Deutscher Bibliotheksverband 2012: *Deutsche Digitale Bibliothek*

Online Verfügbar unter:

<http://www.bibliotheksportal.de/themen/digitale-bibliothek/deutsche-digitale-bibliothek-ddb.html> [14.06.13]

Deutsche Digitale Bibliothek 2012: *Über uns*.

Online verfügbar unter:

<http://www.deutsche-digitale-bibliothek.de/content/about/> [14.06.13]

Deutsche Forschungsgemeinschaft 2010: *Aufgaben*.

Online verfügbar unter:

http://www.dfg.de/dfg_profil/aufgaben/index.html [14.06.13]

Deutsche Forschungsgemeinschaft 2013: *Digitalisierung und Erschließung der im deutschen Sprachraum erschienenen Drucke des 18. Jahrhunderts (VD 18) – Hauptphase*. In: Information für die Wissenschaft Nr. 27 | 3. Juni 2013.

Online verfügbar unter:

http://www.dfg.de/foerderung/info_wissenschaft/info_wissenschaft_13_27/index.html
[14.06.13]

Deutsche Forschungsgemeinschaft 2013: *DFG-Praxisregeln „Digitalisierung“* DFG-Vordruck 12.151 – 02/13.

Online verfügbar unter:

http://www.dfg.de/formulare/12_151/ [14.06.13]

Europe's Information Society 2007: *EURO-Photo: Disclosing the European Library on*

common visual historical heritage.

Online verfügbar unter:

http://ec.europa.eu/information_society/apps/projects/factsheet/index.cfm?project_ref=250468 [14.06.2013]

Fachverband für Multimediale Informationsverarbeitung e.V. 2011: *Archivierung auf Mikrofilm – eine bewährte Technologie mit Anschluss an die digitale Welt.*

Online verfügbar unter:

<http://www.fmi-ev.de/leistungen/mikrofilm.html> [14.06.13]

Feddern, Boi 2007: *Google-Studie zur Ausfallursache von Festplatten.*

Online verfügbar unter:

<http://www.heise.de/newsticker/meldung/Google-Studie-zur-Ausfallursache-von-Festplatten-147178.html> [14.06.2013]

Gilbert, Michael W. 2003: *Digital Media Life Expectancy and Care.*

Online verfügbar unter:

http://web.archive.org/web/20031222194846/http://www.oit.umass.edu/publications/at_oit/Archive/fall98/media.html [14.06.13]

Hans-Bredow Institut 2006: *Medien von A bis Z. Wiesbaden: VS-Verlag für Sozialwissenschaften.*

Zusammenfassung online verfügbar unter:

<http://www2.leuphana.de/medienkulturwiki/medienkulturwiki2/index.php/Digitalisierung>
[14.06.13]

KIT – Universität des Landes Baden-Württemberg und nationales Forschungszentrum in der Helmholtz-Gemeinschaft 2013: *Digitalisierte Altbestände*

Online verfügbar unter:

<http://www.bibliothek.kit.edu/cms/digitalisierte-altbestaende.php> [14.06.2013]

Lischka, Konrad 2011: *Digitale Bibliotheken: Der Staat spart, Google digitalisiert.* In: Spiegelonline: Netzwelt. 26.03.2011.

Online verfügbar unter:

<http://www.spiegel.de/netzwelt/netzpolitik/digitale-bibliotheken-der-staat-spart-google->

[digitalisiert-a-753229.html](#) [14.06.13]

Mumenthaler, Rudolf 2005: *Auf dem Weg zur digitalen Bibliothek: Strategien für die ETH-Bibliothek im 21. Jahrhundert* / Hrsg.: Gysling, Corinne und Neubauer, Wolfram. Zürich: ETH-Bibliothek (Schriftenreihe B der ETH-Bibliothek- Bibliothekswesen; 7)

Stifterverband für die Deutsche Wissenschaft e.V. 2012: *DFG-Förderung*

Online verfügbar unter:

<http://www.laendercheck-wissenschaft.de/drittmittel/dfg-foerderung/index.html> [14.06.13]

Wenk, B. 2004: *Digitalisierung von Bildern und Grafiken*. 24.11.04.

Online verfügbar unter:

<http://bscw.fh-htwchur.ch/pub/bscw.cgi/d53334/BilderUndGrafiken.pdf> [14.06.13]

Bitstream Preservation: Verfahren und Algorithmen

Proseminar-Ausarbeitung von

Dominik Sauter

Am

Steinbuch Centre for Computing (SCC)
Institute for Data Processing and Electronics (IPE)

16. Juni 2013

Inhaltsverzeichnis

1	Einleitung	1
1.1	Einordnung	2
2	Bitstream Preservation	3
2.1	Risikofaktoren	3
2.2	Die Kernstrategien	4
2.3	Modellbildung	5
3	Verfahren und Algorithmen	7
3.1	Fehlertoleranz	7
3.2	Prüfsummen	8
3.3	Datenintegritätsprüfung bei Speicher-APIs	9
3.4	Einsatz fehlerkorrigierender Codes	10
4	Fazit	13
	Literaturverzeichnis	15

1. Einleitung

In unserer heutigen Zeit liegt die meiste Information weltweit in Form von digitalen Daten vor.[Wik13a] Bilder, Musikstücke und Filme, aber auch Firmendaten und Daten, die im wissenschaftlichen Bereich anfallen, wie Messergebnisse aus Experimenten und wissenschaftliche Aufsätze, liegen oft ausschließlich digital vor.

Aber wie sicher sind eigentlich diese unzähligen Daten? Lässt sich der Lieblingssong auf CD auch noch in fünf, zehn oder zwanzig Jahren problemlos abspielen? Können wichtige Messdaten aus teuren wissenschaftlichen Experimenten, für deren sinnvolle Auswertung die Technologie vielleicht erst in 30 Jahren zur Verfügung steht, zu diesem Zeitpunkt überhaupt noch gelesen werden?[BSR⁺05]

Bei einem Vergleich der Lebensdauer digitaler und „analoger“ Speichermedien ergibt sich bereits eine gewisse Problematik:

„Analoge“ Speichermedien	Erwartete Lebensdauer
Bücher/Handschriften (aus säurefreiem Papier und mit säurefreier und nicht eisenhaltiger Tinte)	mehrere hundert Jahre (gesichert)
Keramiktafeln	5000 Jahre (gesichert), 10.000e Jahre (vermutet)
Steinzeugtafeln mit aufgebranntem keramischem Farbdruck	100.000e Jahre wenn erosionsgeschützt (vermutet)
Digitale Speichermedien	Erwartete Lebensdauer
CD (gepresst)	unter Idealbedingungen geschätzt 50-80 Jahre
Magnetbänder	mindestens 30 Jahre (gesichert)
Festplatten als Archivmedien (ohne Betrieb)	10–30 Jahre
Festplatten im laufenden Betrieb	2–10 Jahre (in Abhängigkeit der tägl. Betriebsdauer), durchschnittl. 5 Jahre

Abbildung 1.1: Vergleich der Lebensdauer „analoger“ und digitaler Speichermedien (Daten von [Wik13c])

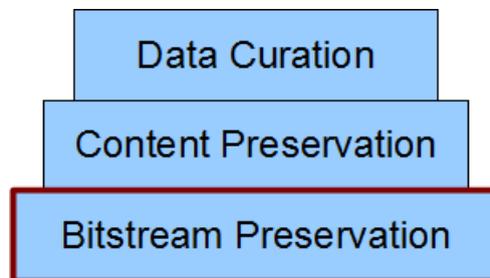
So kann z.B. eine Festplatte nicht etwa wie ein Buch einfach mit wichtigen Daten beschrieben und ins Regal gestellt werden, in der Erwartung, dass diese nach 50 oder 100 Jahren noch vorhanden sind. Dies wird, falls die Festplatte zu diesem Zeitpunkt überhaupt noch funktionstüchtig sein sollte, höchstwahrscheinlich nicht der Fall sein.

Dieser Aspekt der sehr begrenzten Lebensdauer digitaler Speichermedien, sowie weitere Aspekte verlangen schließlich nach einer fundierten Auseinandersetzung mit dem Thema Langzeiterhaltung von digitalen Daten.

In dieser Ausarbeitung soll nun zunächst eine kurze Einordnung von Bitstream Preservation in den Kontext der Langzeiterhaltung erfolgen. Daraufhin wird im zweiten Kapitel näher auf die Risikofaktoren für Bitstream Preservation und deren Lösungsstrategien eingegangen. Im dritten Kapitel werden schließlich heute gängige Verfahren und Algorithmen vorgestellt, die sich mit Fehlererkennung und -behebung beschäftigen. Den Abschluss bildet ein Fazit über das Thema.

1.1 Einordnung

Bei der Langzeiterhaltung von digitalen Daten wird für gewöhnlich zwischen drei verschiedenen Ebenen unterschieden:[grdc]



Diese Ausarbeitung behandelt die unterste Ebene, die sog. „Bitstream Preservation“, d.h. die physische Erhaltung der Bits auf den entsprechenden Medien. Darüber befindet sich die „Content Preservation“-Ebene, welche sich mit der Erhaltung der Lesbarkeit der Daten durch entsprechende Software und Dateiformate beschäftigt.[grda] Die oberste Ebene wird als „Data Curation“ bezeichnet und behandelt die Erhaltung der Bedeutung der Daten in ihrem Kontext.[grdb]

2. Bitstream Preservation

Im Folgenden sollen zuerst die Risikofaktoren und Gefahren für Bitstream Preservation genannt werden, um danach die heute bewährten Strategien vorzustellen, wie die Erhaltung von Bits sichergestellt werden kann. Zuletzt folgt ein Kapitel zum Thema Modellbildung von Speichersystemen.

2.1 Risikofaktoren

Um Lösungsansätze und Strategien für Bitstream Preservation zu entwickeln, muss zuerst ergründet werden, welche Gefahren und Risikofaktoren für die Erhaltung von Bits existieren. Erfolgt dies aus einer ganzheitlichen Sicht, so ergeben sich viele unterschiedliche Gefahren, welche im Folgenden aufgelistet werden:

- **technisches Versagen**

Wie eingangs erwähnt, besitzen digitale Speichermedien nur eine sehr begrenzte Lebensdauer. Materialverschleiß führt zur Unlesbarkeit einzelner Bits bis hin zum Ausfall des gesamten Mediums. Der Datenverlust muss dabei nicht kontinuierlich erfolgen, sondern kann auch durch ein plötzliches Ereignis eintreten, wie z.B. durch einen Head-Crash bei einer Festplatte. Des Weiteren können auch alle anderen Komponenten eines Systems von Ausfall betroffen sein. D.h. Hardware genauso wie Software und andere Dienste, auch solche von Dritt-Anbietern. Bei letzteren ergibt sich eine besondere Problematik, da diese nicht im eigenen Einflussbereich liegen und möglicherweise anderen Sicherheitskonzepten unterworfen sind.[BSR⁺05]

- **Großkatastrophen**

Auch höhere Gewalt in Form von Großkatastrophen kann zu Datenverlust führen und muss berücksichtigt werden. Dies können Naturkatastrophen, wie z.B. Erdbeben, Wirbelstürme, Überschwemmungen oder Brände sein, aber auch andere Arten von Katastrophen, wie Kriege oder Terroranschläge. Meist werden solche Katastrophen von anderen hier genannten Gefahren begleitet.[BSR⁺05]

- **menschliche Fehler**

Die zu erhaltenden Bits befinden sich stets in Menschenhand, und Menschen können Fehler machen. Diese Fehler wiederum können zu Datenverlust führen. Ein Beispiel dafür ist versehentliches oder absichtliches Löschen von Daten, die später wieder gebraucht werden. Es gibt aber auch noch zahlreiche andere menschliche Fehler, die z.B. den Umgang mit Hardware oder Infrastruktur betreffen. Insgesamt gesehen ist menschliches Versagen ein wesentlicher Faktor für Datenverlust.[BSR⁺05]

- **Angriffe**

Natürlich können digitale Daten auch Opfer von gezielten Angriffen werden. Darunter kann z.B. Zensur bzw. Modifikation der Daten verstanden werden, aber auch Diebstahl oder einfach nur mutwillige Zerstörung der Medien. Diese Angriffe können dabei sowohl von außen als auch von innen kommen, legal oder illegal nach der jeweils geltenden Rechtsprechung sein und sich über kurze oder lange Zeiträume erstrecken. Typischerweise sind sie jedoch von langsam zersetzender Natur. Die Motive reichen dabei von politischen Gründen über finanzielle bis hin zu privaten Gründen, etwa der Unzufriedenheit eines Mitarbeiters. Besonderes Augenmerk verdient die Überlegung, dass selbst ein (netzwerktechnisch) vollkommen von der Außenwelt isoliertes System durch das Wirken von „Insidern“ wieder angreifbar wird.[BSR⁺05]

- **Gefahren durch die Organisation, die für die Langzeiterhaltung zuständig ist**

Bei der Langzeiterhaltung von digitalen Daten und speziell Bitstream Preservation wird meistens von einer Organisation ausgegangen, die diese Aufgabe übernimmt, und nicht von einer Einzelperson. Somit können von dieser Organisation auch entsprechende Gefahren ausgehen. Im einfachsten Fall dadurch, dass sich die Organisation irgendwann auflöst, z.B. aufgrund von Bankrott. Es muss daher eine „Ausstiegsstrategie“ geben, die angibt, wie die Daten auf eine andere Organisation transferiert, oder ggf. auch vernichtet werden können.[BSR⁺05]

- **finanzielles Versagen**

Digitale Daten sind sehr anfällig für Unterbrechungen im Geldfluss. Können z.B. die laufenden Kosten für Strom, Kühlung, Miete, Löhne, Erneuerung der Infrastruktur usw. nicht mehr gedeckt werden, kann es schnell zu Datenverlust kommen. Problematisch ist, dass das zur Verfügung stehende Budget normalerweise gewissen Schwankungen unterworfen ist. Zudem ist es generell schwer vorherzusagen, wie sich die Kosten über längere Zeit entwickeln. Da auch die späteren Nutzer der Daten oft noch gar nicht feststehen oder existieren, ist es auch schwierig überhaupt eine Investition in Langzeiterhaltung zu motivieren. Das alles kann dazu führen, dass schon im Vorhinein weniger Daten zur Erhaltung ausgewählt werden (müssen) und nicht alles was es wert wäre, auch erhalten werden kann.[BSR⁺05]

Daneben gibt es noch zwei Arten von Risikofaktoren, die eher impliziterer Natur sind. Dies sind zum einen versteckte und zum anderen zusammenhängende (korrelierte) Fehler. Diese beiden Fehlerarten können jeweils eine Ausprägung durch eine der oben genannten Gefahren finden. So kann es z.B. versteckte technische Fehler geben oder menschliche Fehler, welche unentdeckt bleiben, sowie Angriffe, die nicht sofort erkannt werden usw. Gleiches gilt für zusammenhängende Fehler. So kann es z.B. in Folge einer Naturkatastrophe zu technischem Versagen, menschlichen Fehlern, finanziellem Versagen usw. kommen.[BSR⁺05]

Alles in allem ist jedoch die größte Gefahr für die digitale Langzeiterhaltung und somit auch für Bitstream Preservation ein zu knappes finanzielles Budget. Sobald die Kosten nicht mehr gedeckt werden können, besteht die akute Gefahr von Datenverlust.[BSR⁺05]

2.2 Die Kernstrategien

Nach [Ros10b] besteht der aktuelle Wissenstand, wie die Sicherheit von Bits erhöht werden kann, aus folgenden drei Strategien:

1. **Mehr Kopien**

Replikation ist die fundamentale Sicherheitsstrategie. Je mehr Kopien eines Datenobjekts angelegt werden, desto sicherer ist dieses gegen Ausfall von Kopien geschützt.

2. Größere Unabhängigkeit der Kopien

Kopien allein genügen allerdings nicht. Damit möglichst wenig Kopien dem selben Fehler oder Ereignis zum Opfer fallen können, müssen die Kopien möglichst unabhängig voneinander sein. Und zwar unabhängig bezüglich vieler verschiedener Gesichtspunkte, wie etwa geographisch (z.B. Kopien in mehreren Ländern), organisatorisch (z.B. Kopien nicht alle im Einflussbereich der gleichen Personen), technologisch (z.B. nicht nur Geräte des gleichen Herstellers) usw.[vF11]

3. Größere Häufigkeit in der die Kopien auf Fehler überprüft werden

Da digitale Speichermedien nur eine sehr begrenzte Lebensdauer haben, werden die Kopien mit der Zeit fehlerhaft oder können ganz ausfallen. Das bedeutet, dass sie auch regelmäßig auf Fehler überprüft und ggf. repariert oder ausgetauscht werden müssen. Je häufiger dies geschieht, desto schneller werden Fehler erkannt und desto sicherer sind die Kopien wiederum.

Bei allen drei Strategien wird darauf hingewiesen, dass eine Vergrößerung der Datenmenge eine Vergrößerung der Kosten für die jeweilige Strategie bedeutet und dadurch ein der Strategie gegenläufiger, negativer Effekt entsteht. So steigen z.B. bei der ersten Strategie durch eine größere Datenmenge die Kosten pro Kopie, was bei einem begrenzten Budget zu insgesamt weniger Kopien und damit geringerer Sicherheit führt.[Ros10b] Daraus folgt letztlich, dass diese Strategien bei großen Datenmengen auf ein Kostenproblem stoßen, welches sich nicht so einfach lösen lässt.[Ros10a]

2.3 Modellbildung

Modelle von Speicherdienstsystemen können dabei helfen, verschiedene Speichertechnologien und -techniken untereinander zu bewerten. Dabei muss jedoch darauf geachtet werden, dass die Modelle jeweils die gleiche Datengrundlage verwenden, um wirklich vergleichbare Ergebnisse zu erhalten. Typische Metriken sind die MTTF (Mean Time To Failure), d.h. die durchschnittliche Zeit bis zum Ausfall, die MTDDL (Mean Time To Data Loss), die durchschnittliche Zeit bis Datenverlust auftritt und die UBER (Unrecoverable Bit Error Rate), die Wahrscheinlichkeit, dass Bits nicht mehr korrekt gelesen werden können. Dabei macht die MTDDL jedoch keine Angaben über die Größenordnung des Datenverlusts, während die UBER keine zeitliche Dimension besitzt.[Ros10b]

Problematisch bei der Modellbildung ist, dass in der Regel bestimmte Annahmen gemacht werden, die sich in der Realität nur als bedingt oder überhaupt nicht zutreffend herausstellen. So werden im Modell des Pergamum-Projekts an der University of California, Santa Cruz, welches als State of the Art bezeichnet werden kann, etwa unkorrelierte Fehler, bugfreie Software und keine Gefahr durch Benutzerfehler oder Angriffe von außen angenommen. Dies sind aber alles Punkte, welche bei realen Speicherdiensten zu Datenverlusten führen. Anstatt eine Obergrenze für den möglichen Datenverlust zu geben, können diese Modelle daher nur die Untergrenze schätzen.[Ros10b]

Dies ist auch der Grund, warum viele Herstellerangaben von Speichermedien und -diensten sehr optimistisch ausfallen. Die Angaben werden in den meisten Fällen nämlich nicht durch entsprechende Experimente ermittelt, sondern lediglich theoretisch anhand von Modellen vorausgesagt.[Ros10b]

3. Verfahren und Algorithmen

In diesem Kapitel soll zunächst eine neue Sichtweise bezüglich des Auftretens von Fehlern motiviert werden, um danach heute gängige Verfahren und Algorithmen zur Fehlererkennung und -behebung bei Speichersystemen vorzustellen. Darunter Prüfsummen und deren Einsatz in Speicher-APIs, sowie der Einsatz fehlerkorrigierender Codes, speziell MDS Erasure Codes.

3.1 Fehlertoleranz

Nach [Ros10b] sind die Anforderungen, was die Datenmengen und die Zeitspannen angeht, für die sie erhalten werden sollen, heutzutage schlichtweg zu groß, als dass zuverlässige Prognosen getroffen werden könnten, dass bestimmte Speichersysteme dafür ausreichen. Modelle können nicht die ganze Bandbreite an Fehlerursachen und ihre Abhängigkeiten untereinander abdecken. Und selbst einfache Experimente würden zu lange dauern, oder wären zu teuer, oder sogar beides.

Das Anlegen von Kopien auf verschiedenen Speichersystemen und deren regelmäßiges Überprüfen auf Fehler verringert zwar die Wahrscheinlichkeit von Datenverlust, aber ganz ausgeschlossen werden kann er dadurch nicht. Zudem können zusätzliche Zugriffe auf die Speichermedien die Fehlerwahrscheinlichkeit erhöhen. Auch können Bugs in der Prüf-Software selbst zu Datenverlust führen. Des Weiteren ist auch keine technische Revolution in naher Zukunft zu erwarten, die das Problem lösen würde.[Ros10b]

Trotz der häufigen Erwartung, dass digitaler Speicher zuverlässig sei, sind in der Praxis heutzutage Datenverluste oder zumindest teilweise Beschädigungen der Daten an der Tagesordnung. Dieses Problem der Fehleranfälligkeit ist dabei selbst nur eine Facette im Kontext eines sich anbahnenden größeren Problems.[Ros10b]

Werden heute große High-Performance-Anwendungen im Petascale-Bereich mit langer Laufzeit betrachtet, so ist deren Fehleranfälligkeit so hoch, dass ein kompletter Neustart bei Eintritt eines Fehlers nicht mehr machbar ist. Es werden daher Neustart-Schemata mithilfe von Checkpunkten eingebaut, welche aber meist ebenso komplex wie teuer sind.[Ros10b] Jedoch wird auch dieser Ansatz bald nicht mehr funktionieren, da vorausgesagt wird, dass in der kommenden Exascale-Generation während eines Checkpoint-Neustarts bereits der nächste Fehler auftreten kann. Die Zeit für einen Checkpoint-Neustart wird also größer sein, als die Mean Time To Failure (MTTF) eines solchen Systems. Infolgedessen müssen Programmierer ihre Programme sehr viel sensibler in Hinsicht auf Fehler in ihrer Umgebung machen, da es nicht mehr möglich sein wird, alle Fehler im Vorhinein

abzufangen.[Ros10b]

Daher plädiert der Autor von [Ros10b] für ein fehlerbewussteres Design von Speichersystemen.

3.2 Prüfsummen

Prüfsummen dienen zur Sicherstellung der Datenintegrität bei Datenübertragung und -speicherung. Daher kommen sie auch in den meisten Fehlerkorrekturverfahren vor. Anstatt einfacher Prüfsummen werden dabei meist kryptographisch besonders starke Verfahren verwendet, um auch vor gezielter Datenmanipulation geschützt zu sein. Beispiele dafür sind sog. Einweg-Hash-Funktionen, bei denen Kollisionen sehr schwer ermittelbar sind.[Wik13e][Wik13b]

Eine solche kryptographische Hashfunktion ist der weit verbreitete Message-Digest Algorithm 5 (MD5). Dieser bildet jede Eingabe auf einen 128-Bit-Hashwert ab (siehe Abb. 3.1). Er wird verwendet, um Datenobjekte vor und nach einer Transaktion durch Vergleich der Hashwerte auf Integrität zu überprüfen, z.B. bei Datei-Downloads.[Wik13d] Unter Linux ist meistens das Programm „md5sum“ als Teil der coreutils (GNU core utilities) standardmäßig installiert und erlaubt das Erzeugen und Überprüfen von MD5-Hashes.[arc]

Eingabe	MD5-Hashwert (hexadezimal)
Big Data	3a9c1c8521b12754fc7582be72383e28
Big data	ca11a878b6e44cf23927de0ed69a5610

Abbildung 3.1: Beispiele von MD5-Hashwerten: bei Änderung nur eines Zeichens in der Eingabe ergibt sich ein vollkommen anderer Hashwert (berechnet mit [mir])

Heutzutage gilt MD5 jedoch nicht mehr als kryptographisch sicher, da ohne größeren Aufwand Kollisionen gefunden werden können. Dennoch ist MD5 noch relativ weit verbreitet und im Einsatz. Ein Hauptgrund dafür dürfte seine Performance sein, so ist MD5 ca. drei mal schneller zu berechnen als der vergleichbare Hash-Algorithmus SHA-1 (Secure Hash Algorithm). Dies kommt aber auch daher, dass kryptographisch sichere Algorithmen sich durch eine lange Berechnungsdauer auszeichnen. Sichere Algorithmen müssen oft langsam sein, da der hohe Berechnungsaufwand meistens die Grundlage für die Sicherheit bildet. Des Weiteren sind Algorithmen in Software-Systemen im Allgemeinen nicht besonders einfach auszutauschen. Austauschkosten und Abwärtskompatibilität sind dabei oft wichtige Faktoren. Letztendlich muss auch gesagt werden, dass MD5 bis heute nur bezüglich der Kollisions-Angriffe (Abb. 3.2) gebrochen ist, und nicht bezüglich eines auch um Größenordnungen schwereren Preimage-Angriffs (Abb. 3.3). Das bedeutet, dass z.B. als Hashwerte gespeicherte Passwörter noch relativ sicher sind.[Wik13d][cod][yah][per]

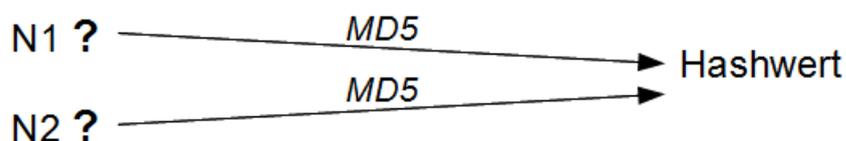


Abbildung 3.2: Kollisions-Angriff: Es werden zwei beliebige Nachrichten N1 und N2 gesucht, welche auf den gleichen Hashwert abgebildet werden.

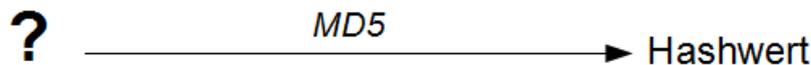


Abbildung 3.3: Preimage-Angriff: Es wird zu einem vorgegebenen Hashwert eine Eingabe gesucht, die auf diesen abgebildet wird.

3.3 Datenintegritätsprüfung bei Speicher-APIs

Bei Speicher-APIs gibt es bereits eine Möglichkeit, die Integrität von Daten nachzuprüfen, indem eine Anwendung dem Schreib- oder Ladebefehl einen Hashwert der Daten mitgibt. Damit kann das Speichersystem dann überprüfen, ob die Daten zwischen Anwendung und Speichergerät unbeschädigt übertragen wurden oder nicht. Datenpfade und Speicherpuffer sind auf diesem Weg nämlich erhebliche Fehlerquellen. So ist es bei der Amazon S3 (Simple Storage Service) REST API möglich, durch einen Content-MD5-Header (MD5 für „Message-Digest algorithm 5“) den Hashwert der Daten an den Speicherdienst mitzuschicken. Bei einer Antwort wird ein ETag-Header (entity tag) mit dem Hashwert der Daten, welchen der Speicherdienst berechnet hat, zurückgeschickt (Abb. 3.4). Auf diese Weise kann eine Anwendung den Hashwert, den sie zum Speicherdienst geschickt und selbst bei sich gespeichert hat, mit dem zurückgegebenen Hashwert abgleichen.[Ros10b]

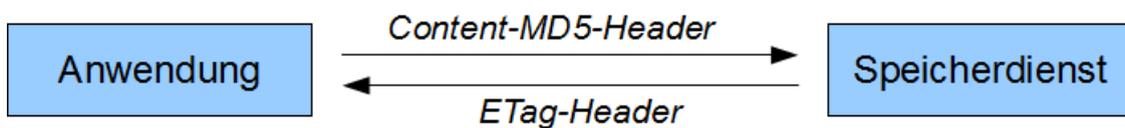


Abbildung 3.4: Austausch der Hashwerte durch die jeweiligen Header

Das Problem das sich hierbei ergibt ist aber, dass der Speicherdienst seinen Datenhashwert nicht notwendigerweise berechnen muss. Er erhält diesen ja von der Anwendung selbst und bräuchte ihn daher lediglich zu speichern und wieder zurückzugeben. Ein böser oder fehlerhafter Speicherdienst, der die Daten die er bekommt nicht einmal speichert, würde daher überhaupt nicht auffallen. Eine Methode um ein solches Verhalten aufzudecken wäre, die Daten vom Speicherdienst wieder anzufordern, und den Hashwert davon zu berechnen. Nun könnte dieser neu berechnete Hashwert mit dem vorher gespeicherten oder dem des Speicherdienstes (ETag-Header) verglichen werden. Stimmen die Hashwerte überein, sind die Daten unverändert, stimmen sie nicht überein, wurden die Daten korrumpiert. So scheint es auf den ersten Blick zumindest. Aber es darf nicht vergessen werden: Auch der Datenhashwert selber ist von Korruption nicht ausgeschlossen. Insgesamt ergeben sich somit vier Fälle (siehe Abb. 3.5). Es muss daher erkennbar sein, ob der Datenhashwert verändert wurde oder nicht.[Ros10b]

Datenhashwert	Übereinstimmung	Keine Übereinstimmung
Unverändert	Daten O.K.	Daten korrumpiert
Verändert	Absichtliche Änderung	Daten und/oder Hashwert korrumpiert

Abbildung 3.5: Die vier Fälle beim Vergleich der Hashwerte (Quelle: [Ros10b])

Des Weiteren muss ein Speicherdienst auch beweisen können, dass er seine Daten korrekt gespeichert hat. Ein solcher Beweis wäre z.B. nach einem Hashwert eines zufälligen

Bitstrings (Nonce) und des gespeicherten Datenobjekts zu verlangen (Abb. 3.6). Da der Bitstring zufällig gewählt ist, kann er vom Speicherdienst nicht vorhergesehen werden und zwingt ihn so den Hashwert tatsächlich mithilfe der Daten zu berechnen. Das bedeutet aber wiederum, dass er die Daten korrekt gespeichert haben muss.[Ros10b]

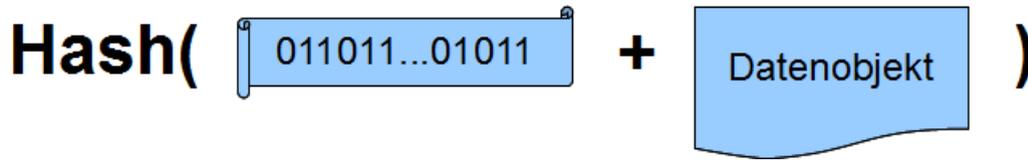


Abbildung 3.6: Berechnung des Hashwerts von Nonce und Datenobjekt

Zusammenfassend kann gesagt werden, dass die bereits vorhandenen Möglichkeiten für Integritätsprüfungen bei Speicher-APIs Schritte in die richtige Richtung sind. Von optimalen Lösungen sind diese jedoch noch weit entfernt.[Ros10b]

3.4 Einsatz fehlerkorrigierender Codes

Replikation als Sicherheitsstrategie ist relativ einfach umzusetzen, erzeugt jedoch einen hohen Speicher-Overhead. Dieser schlägt sich nicht nur in primären Kosten durch die Speichermedien, wie z.B. Festplatten, nieder, sondern auch in Form von Mehrkosten durch erhöhten Energieverbrauch, sowie zusätzlichen Wartungs- und Managementaufwand. Alternative Techniken bieten dabei ein besseres Redundanzverhältnis und sind anpassbar bezüglich der Zuverlässigkeit. Ein Beispiel dafür sind fehlerkorrigierende Codes, speziell sog. MDS (Maximum Distance Separable) Erasure Codes. Diese bieten bei geringem Speicher-Overhead eine hohe Fehlertoleranz. Dabei werden aus den Originaldaten redundante Blöcke berechnet, welche im Falle von Datenverlust die Originaldaten wiederherstellen können.[Pet11]

Eine Idee ist es dann, Replikation und MDS Erasure Codes als Sicherheitstechniken zu kombinieren. Das Ziel dabei ist, Speicher-Overhead zu verringern, ohne die Fehlerwahrscheinlichkeit zu stark zu erhöhen und damit die Zuverlässigkeit zu stark zu beeinträchtigen. Durch die Kombination beider Techniken wird eine dynamische Anpassung von Speicher-Overhead und Zuverlässigkeit ermöglicht. Beim Speichervorgang einer Datei (siehe Abb. 3.7) werden dabei gleich große Datenblöcke erstellt und in sog. „Stripes“ zusammengefasst. Diese werden anschließend, je nach Vorgabe, entweder nur repliziert, oder es werden auch noch zusätzliche redundante Blöcke berechnet. Bei Datenverlust können dann die Kopien oder die berechneten redundanten Blöcke eines Stripes dazu verwendet werden, um die Originaldaten wiederherzustellen.[Pet11]

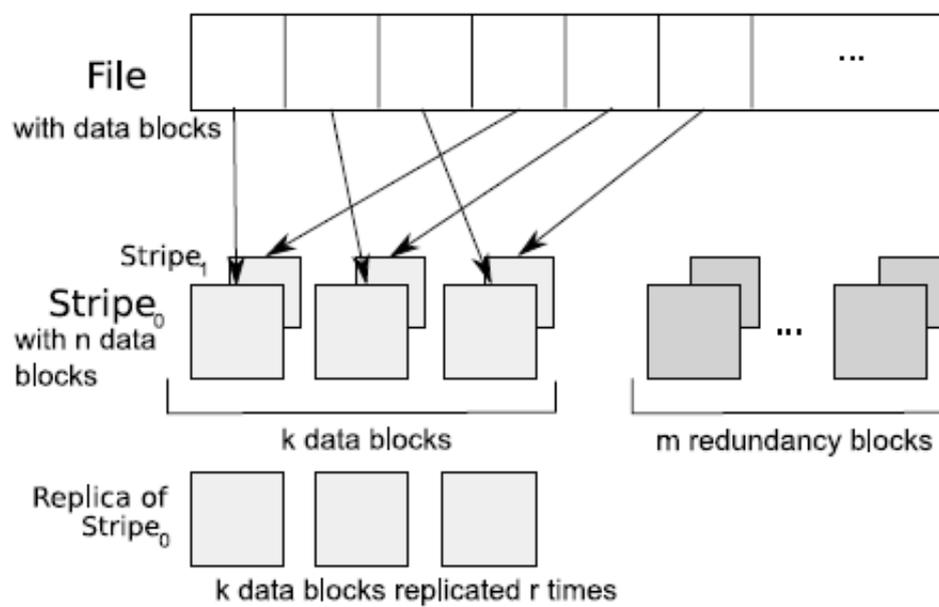


Abbildung 3.7: Speichervorgang einer Datei mit $(k + m)$ MDS Erasure Codes und r -facher Replikation (Quelle: [Pet11])

4. Fazit

Zusammenfassend kann gesagt werden, dass Bitstream Preservation heutzutage immer noch kein vollständig gelöstes Problem ist. Dies hängt aber auch mit den unterschiedlichen Ansichten bezüglich eines „vollständig gelösten“ Status' von Bitstream Preservation zusammen. Daher darf bezweifelt werden, ob die Anerkennung eines solchen Status' je erreicht werden kann.[Ros10a] Obwohl es bewährte Strategien gibt, um die Sicherheit von Bits zu erhöhen, kann die Wahrscheinlichkeit für Datenverlust und -beschädigung lediglich verringert werden. Ein vollständiges Ausschließen von Datenverlust wird nie möglich sein. In diesem Punkt gibt es auch keinen Unterschied zu „analogen“ Speichermedien. Auch diese sind von Verfall und letztendlich Datenverlust betroffen (vgl. Abb. 4.1). Insofern dürfen die Ansprüche an digitalen Speicher auch nicht zu hoch angesetzt werden.

Des Weiteren sind Modelle von Speichersystemen zwar geeignet, um verschiedene Technologien und Techniken untereinander zu vergleichen, aber häufig nur bedingt brauchbar, um realitätsnahe Vorhersagen zu treffen. Dies ist natürlich vor allem der hohen Komplexität von Speichersystemen mit unzähligen Abhängigkeiten geschuldet, welche sich sehr schwer modellieren lassen. Nichtsdestotrotz ist der Punkt Modellbildung ein Ansatzpunkt, um das Problem Bitstream Preservation weiter in den Griff zu bekommen.

Zudem wäre die Installation eines großen anonymen Erfahrungsarchivs für Speichersysteme, wie in [BSR⁺05] vorgeschlagen, ein sinnvoller Schritt. Damit könnten die Erfahrungen mit den von Natur aus nur selten auftretenden Zwischenfällen und Problemen bei Speichersystemen gesammelt werden. Durch die anonyme Form wäre es insbesondere auch möglich, dass prekäre Fälle von Datenverlust gefahrlos dokumentiert werden. So könnte eine gemeinsame Erfahrungsbasis geschaffen werden, welche dabei helfen würde, die Zuverlässigkeit von Speichersystemen in Zukunft weiter zu erhöhen.

Abschließend kann gesagt werden, dass Bitstream Preservation in der Praxis heutzutage in erster Linie ein Kostenproblem ist. Dieses besteht aus einer grundsätzlichen Abwägung zwischen Kosten und dadurch gewonnener Sicherheit. Bei vorgegebenem Budget bedeutet das konkret eine Abwägung verschiedener Sicherheitstechniken und deren Kombination, um die maximale Gesamtsicherheit zu erreichen.

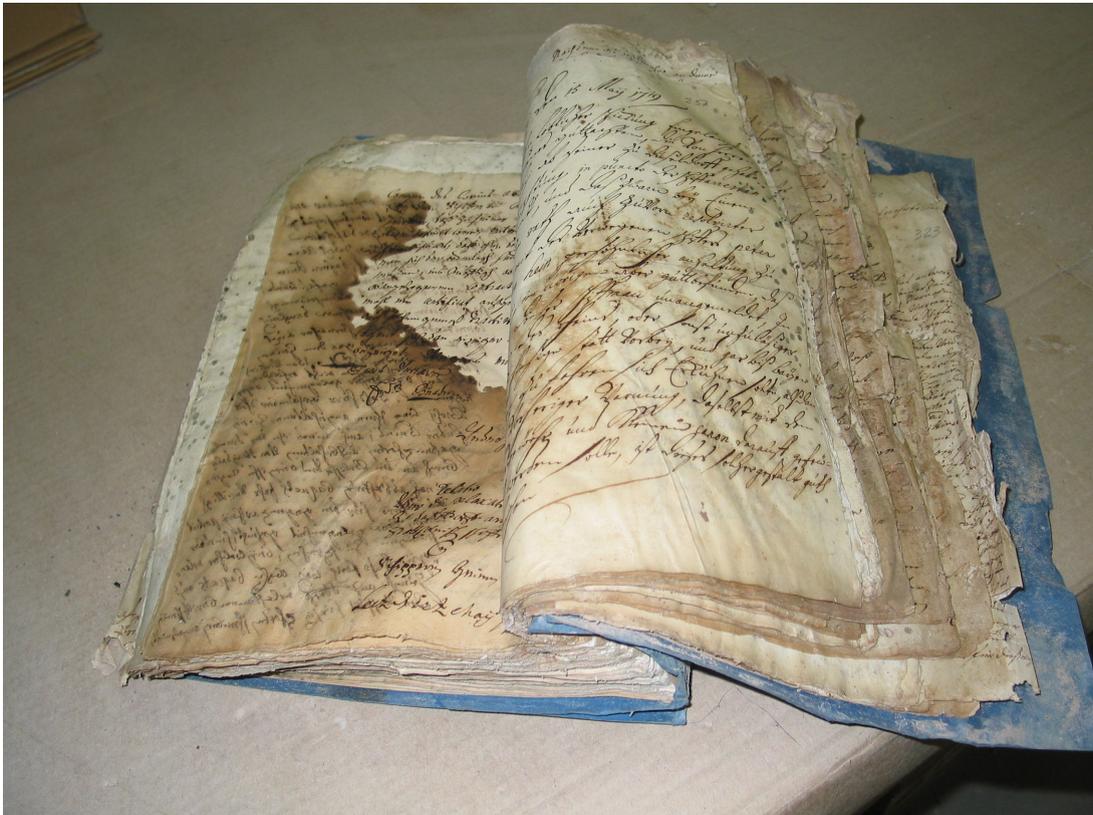


Abbildung 4.1: Datenverlust bei "analogen" Speichermedien: ein durch Nässe stark beschädigtes Buch nach dem Einsturz des Kölner Stadtarchivs im Jahr 2009 (Quelle: [min])

Literaturverzeichnis

- [arc] *Arch Linux Wiki: md5sum*. <https://wiki.archlinux.de/index.php?title=Md5sum&oldid=14492>, Abruf: 26.05.2013.
- [BSR⁺05] Mary Baker, Mehul Shah, David S. H. Rosenthal, Mema Roussopoulos, Petros Maniatis, T. J. Giuli und Prashanth Bungale: *A Fresh Look at the Reliability of Long-term Digital Storage*. arXiv:cs/0508130, August 2005. <http://arxiv.org/abs/cs/0508130>, besucht: 2013-04-14.
- [cod] *Coding Horror: Speed Hashing*. <http://www.codinghorror.com/blog/2012/04/speed-hashing.html>, Abruf: 29.05.2013.
- [grda] *GRDI2020 wiki: Content Preservation*. http://grdi2020.eu/mediawiki/index.php/Content_Preservation, Abruf: 17.05.2013.
- [grdb] *GRDI2020 wiki: Data Curation*. http://grdi2020.eu/mediawiki/index.php/Data_Curation, Abruf: 17.05.2013.
- [grdc] *GRDI2020 wiki: Data Curation and Preservation*. http://grdi2020.eu/mediawiki/index.php/Data_Curation_and_Preservation, Abruf: 17.05.2013.
- [min] *Bild von www.minden.de*. <http://www.minden.de/medien/bilder/archiveinsturz1%5B1%5D.jpg>, Abruf: 29.05.2013.
- [mir] *Miracle Salad: md5 Hash Generator*. <http://www.miraclesalad.com/webtools/md5.php>, Abruf: 29.05.2013.
- [per] *PerlMonks: What is MD5 Hashing and Why is it Important?* http://www.perlmonks.org/?node_id=145704, Abruf: 29.05.2013.
- [Pet11] Kathrin Peter: *Reliability study of coding schemes for wide-area distributed storage systems*, 2011. 19th International Euromicro Conference on Parallel, Distributed and Network-Based Processing.
- [Ros10a] David S. H. Rosenthal: *Bit Preservation: A Solved Problem?* International Journal of Digital Curation, 5(1):134–148, Juni 2010, ISSN 1746-8256. <http://www.ijdc.net/index.php/ijdc/article/view/151>, besucht: 2013-04-14.
- [Ros10b] David S. H. Rosenthal: *Keeping bits safe: how hard can it be?* Commun. ACM, 53(11):47–55, November 2010, ISSN 0001-0782. <http://doi.acm.org/10.1145/1839676.1839692>, besucht: 2013-04-14.
- [vF11] Arbeitspaket 3: Langzeitarchivierung von Forschungsdaten: *Bitstream Preservation: Bewertungskriterien für Speicherdienste*, 2011. <http://www.wissgrid.de/workgroups/ap3/2011-03-08--bitstream-preservation.pdf>, Version 0.5.

- [Wik13a] Wikipedia: *Digitale Revolution* — *Wikipedia, Die freie Enzyklopädie*, 2013. http://de.wikipedia.org/w/index.php?title=Digitale_Revolution&oldid=118874657, [Online; Stand 10. Juni 2013]; Abruf: 26.05.2013.
- [Wik13b] Wikipedia: *Hashfunktion* — *Wikipedia, Die freie Enzyklopädie*, 2013. <http://de.wikipedia.org/w/index.php?title=Hashfunktion&oldid=118963538>, [Online; Stand 11. Juni 2013]; Abruf: 26.05.2013.
- [Wik13c] Wikipedia: *Langzeitarchivierung* — *Wikipedia, Die freie Enzyklopädie*, 2013. <http://de.wikipedia.org/w/index.php?title=Langzeitarchivierung&oldid=118424934>, [Online; Stand 10. Juni 2013]; Abruf: 26.05.2013.
- [Wik13d] Wikipedia: *Message-Digest Algorithm 5* — *Wikipedia, Die freie Enzyklopädie*, 2013. http://de.wikipedia.org/w/index.php?title=Message-Digest_Algorithm_5&oldid=119373501, [Online; Stand 11. Juni 2013]; Abruf: 26.05.2013.
- [Wik13e] Wikipedia: *Prüfsumme* — *Wikipedia, Die freie Enzyklopädie*, 2013. <http://de.wikipedia.org/w/index.php?title=Pr%C3%BCfsumme&oldid=116737568>, [Online; Stand 11. Juni 2013]; Abruf: 26.05.2013.
- [yah] *The Gnutella Developer Forum (GDF): "What are the advantages of md5 and sha1?"*. http://groups.yahoo.com/group/the_gdf/message/20755, Abruf: 29.05.2013.

Erklärung

Ich versichere, dass ich die Arbeit ohne fremde Hilfe und ohne Benutzung anderer als der angegebenen Quellen angefertigt habe, und dass die Arbeit in gleicher oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen hat und von dieser als Teil einer Prüfungsleistung angenommen wurde. Alle Ausführungen, die wörtlich oder sinngemäß übernommen wurden, sind als solche gekennzeichnet.

Karlsruhe, den 16. Juni 2013

(Dominik Sauter)

Karlsruher Institut für Technologie (KIT)
Steinbuch Centre for Computing (SCC)

Dr. Rainer Stotzka
Prof. Dr. Achim Streit
Danah Tonne

Schafft Replikation Sicherheit?
Kosten vs. Nutzen

Proseminar „BigData“
Sommersemester 2013

Jan-Frederic Markert
Matrikelnummer 1580236

Inhaltsverzeichnis

I. Einleitung.....	3
1. Motivation.....	3
2. Begriffserläuterung: Schafft Replikation Sicherheit?.....	3
3. Abgrenzung.....	3
II. Datensicherung.....	4
1. Erwartungshaltung.....	4
2. Technologien.....	4
III. Replikation.....	6
1. Aktueller Stand.....	6
2. Kosten.....	7
3. Nutzen	8
4. Auseinandersetzung Kosten vs. Nutzen.....	10
IV. Aus der Praxis	11
V. Quellen.....	12

STEINBUCH CENTRE FOR COMPUTING – SCC
INSTITUTE FOR DATA PROCESSING AND ELECTRONICS – IPE



I. Einleitung

1. *Motivation*

Die Erhaltung von Informationen war schon seit jeher essentiell für die Menschheit. Die Medien haben sich jedoch mit der Zeit gewandelt: Waren es zu Beginn in Stein gemeißelte Hieroglyphen, später auf Papier geschriebene Buchstaben und anschließend der Foto-Film oder die Schallplatte, so ist es heutzutage das Magnetband beziehungsweise die elektronische Speicherung in Form von Bits.

Während diese Entwicklung bewirkt, dass die Erzeugung, Manipulation und Weitergabe von Daten immer weiter vereinfacht werden, nimmt die Robustheit der Medien im Gegenzug stetig ab.

Dies führt zu einer Spannung, denn die digitale Informationsverarbeitung nimmt inzwischen eine zentrale Rolle in unserer Gesellschaft ein. Die Digitalisierung ist auf dem Vormarsch und durchdringt alle Bereiche gegenwärtigen Lebens - immer mehr Daten werden erhoben und gespeichert.

Mit dem Begriff BigData wird die Problematik adressiert, dass die anfallenden Datenmengen zum Teil ungekannte Größen annehmen, die zur Verarbeitung nötigen Technologien und Werkzeuge jedoch mit diesem Wachstum nur in beschränktem Maße schritthalten.

Im Folgenden betrachte ich, inwiefern Replikation einen Nutzen bringt bei der langfristigen Erhaltung solcher großen Datenmengen und in welchem Verhältnis dazu die Kosten stehen.

Datenverlust kann verschiedene Konsequenzen mit sich bringen. Es entstehen meist in erster Linie ökonomische Verluste, allerdings können auch unwiederbringliche Informationen verloren gehen oder es notwendig werden, eine Daten-Erhebung oder ein Experiment wiederholen zu müssen.

2. *Begriffserläuterung: Schafft Replikation Sicherheit?*

Unter Replikation wird verstanden, Daten mehrfach zu speichern - meist werden hierzu verschiedene Medien eingesetzt.

Datensicherheit bedeutet, dass die zu erhaltenden Daten zu einem späteren Zeitpunkt wieder so ausgelesen werden können, wie sie einst auf das Medium geschrieben wurden.

3. *Abgrenzung*

Nicht eingegangen wird in dieser Arbeit auf die logische, bzw. funktionale Bewahrung digitaler Inhalte, sondern rein auf die Ebene der Erhaltung der blanken Bits.

Unter logischer Bewahrung von Informationen wird verstanden, dass nicht nur die Medium und die darauf enthaltenen Daten an sich erhalten werden, sondern auch gewährleistet wird, dass die Informationen selbst anschließend immer noch gelesen und verstanden werden können.

Digitale Datenformate werden in vergleichsweise kurzen Abständen abgelöst, Software wird durch neue ersetzt und auch Betriebssysteme kommen und gehen. In diesem unvorhersehbaren Wandel eine beständige Grundlage zur Informationsdarstellung zu erhalten oder zumindest die ursprünglich notwendigen Werkzeuge weiterhin verfügbar zu machen, ist die Aufgabe der logischen Bewahrung.

II. Datensicherung

1. Erwartungshaltung

Computer bestehen aus festverdrahteten Komponenten und arbeiten nach deterministischen Algorithmen. Die Ausführung von Computerprogrammen soll bei fixen Umständen und auf Basis identischer Daten ein vorhersehbares Resultat erbringen.

Wenn man nun diese allgemeine Betrachtung direkt auf die Anwendung der Replikation überträgt, so sollte das Kopieren von Bits fehlerfrei vonstatten gehen können. Die Tatsache jedoch, dass Bits korrekt kopiert werden können, führt oft zu der Annahme, dass dem auch immer so sein muss.

Als Folge dessen wird digitaler Speicher als zuverlässig angesehen, obwohl die Zuverlässigkeit von vielen Faktoren abhängt und beschränkt ist, und nicht etwa allgemein und uneingeschränkt gültig ist.

2. Technologien

Digitale Daten können auf allen Medien gespeichert werden, durch die binäre Zeichen darstellbar sind. Jedoch bieten sich je nach Anwendungsgebiet unterschiedliche Medien an.

Im Bereich von daten-intensiven Anwendungen existiert nur eine geringe Streuung auf verschiedene Medien: hier gibt es eine starke Fokussierung auf Festplatten und Magnetband.

Festplatten, die auch im Heimanwenderbereich eine sehr große Bedeutung besitzen, zeichnen sich durch geringe Herstellungskosten, teilweise auf Grund der hohen Verbreitung, und die komfortable Bedienung mit wahlfreiem Zugriff und hohen Zugriffsraten aus.

Magnetbänder bieten eine hohe Lebensdauer, große Kapazitäten und eine günstige Anschaffung. Die Tatsache, dass die Bänder unhandlicher sind und nur ein sequentieller Zugriff möglich ist, ist im Zusammenhang der Langzeitarchivierung weniger von Bedeutung.

RAID

RAID (Redundant Array of Independent Disks, dt. Redundante Anordnung unabhängiger Festplatten) dient dazu, die Zuverlässigkeit von Festplatten-Speichersystemen durch Hinzunahme von Replikationen zu erhöhen. Dazu werden die einzelnen Festplatten so zu einem logischen Laufwerk zusammengeschlossen, sodass eine Redundanz unter den gespeicherten Daten entsteht – die Wahrscheinlichkeit steigt, dass ein Datenverlust, bedingt durch einen Plattenausfall, kompensiert werden kann.

Darüber hinaus können mit RAID auch die Zugriffsraten gesteigert werden, da beim Lesevorgang meist mehrere Festplatten beteiligt sind und somit nicht ein Flaschenhals vorhanden ist.

Es existieren verschiedene RAID-Level, welche sich in der Ausprägung der Redundanz und anderer Eigenschaften zum Teil erheblich unterscheiden. Hier aufgelistet sind die am meisten verwendeten:

- RAID 0: Beschleunigung, keine Redundanz (für Datensicherung ungeeignet)
 - Daten werden ähnlich dem Reißverschlussprinzip auf den Festplatten verteilt
- RAID 1: Beschleunigung, volle Redundanz
 - Alle Festplatte enthalten gleiche Daten (schlechte Ausnutzung des gegebenen Speichers)
 - Alle bis auf eine Festplatte können ohne Verlust ausfallen
- RAID 5: RAID 1 mit Paritätsinformation
 - Paritätsinformationen verteilt auf die einzelnen Festplatten
 - Ausfall einer Festplatte kann verkraftet werden
- RAID 6: RAID 5 mit zusätzlicher Paritätsinformation
 - Wiederherstellung komplexer als bei RAID 5
 - Ausfall von 2 Festplatten kann verkraftet werden

SMART

Bei SMART (Self-Monitoring Analysis and Reporting Technology; System zur Selbstüberwachung, Analyse und Statusmeldung) handelt es sich um einen Standard zum Protokollieren von Statusparametern von Festplatten. Er dient vor allem zum frühzeitigen Aufdecken von drohenden Defekten oder Ausfällen.

Unter anderem werden folgende Daten gesammelt:

- Seek Error Rate Fehler beim Lesen
- Raw Read Error Rate Fehler beim Lesen
- Scan Error Rate Fehler beim Überprüfen der Plattenoberfläche
- Drive Temperature Laufwerks-Temperatur
- Reallocated Sector Count Anzahl der verbrauchten Reserve-Sektoren

SMART und RAID

Bei in Software umgesetzten RAID-Systemen können die SMART-Werte ohne Einschränkungen genutzt werden. Allerdings kommt es bei Hardware-RAID auf den RAID-Controller an, ob SMART überhaupt angeboten wird, und wenn, ob gewisse Hersteller unterstützt werden oder nicht.

III. Replikation

1. Aktueller Stand

Um die Zuverlässigkeit eines Datensicherungssystems zu erhöhen, bietet es sich zum einen an, die Ausfallwahrscheinlichkeit der zugrunde liegenden Speichermedien zu minimieren, und zum anderen die die Replikation zu erhöhen.

Bei Anwendung von Replikation wird grundsätzlich nach drei Richtlinien gehandelt:

1. Je mehr Kopien, desto sicherer

Mit dem Wachsen der Datenmengen steigen auch die Kosten pro Kopie und somit verringert sich die Anzahl erschwinglicher Replikationen.

2. Je unabhängiger die Kopien, desto sicherer

Aufgrund der steigenden Größe der Daten verringert sich die Anzahl verfügbarer Speicheroptionen. Dadurch basieren immer mehr Kopien auf der gleichen Technologie und das durchschnittliche Maß an Unabhängigkeit nimmt ab. Unabhängigkeit lässt sich bei den Organisationen, den Speicher-Herstellern, bei Betriebssystemen und bei der geographischen Lage der Speicherzentren erreichen.

3. Je öfter die Kopien überprüft werden, desto sicherer

Durch die größer werdenden Datenmengen werden auch die einzelnen Überprüfungen der Datensätze aufwendiger und teurer, und somit nimmt auch deren Häufigkeit ab. Die Überprüfungen werden meist anhand von vorher errechneten Prüfsummen durchgeführt.

Im Bereich von BigData ist Replikation ein fester Bestandteil der Sicherheitsstrategie, jedoch ist die Frage nach dem konkreten Einsatzumfang nicht allgemein zu beantworten.

2. Kosten

Maßgebliche Kostentreiber der Replikation sind die Speichermedien.

Im zeitlichen Verlauf lässt sich die Kapazität von angebotenen Speicher bei konstant gehaltenen Kosten mit einem ungefähr exponentiellen Wachstum beschreiben (siehe Abbildung 1. logarithmisch fallende Kosten pro Gigabyte).

Diese Tatsache führt zu der Annahme, dass Replikation immer preiswerter werden müsste und sich die Kostenfrage in naher oder ferner Zukunft erübrigen müsste. Demgegenüber stehen jedoch ebenfalls exponentiell steigende Datenaufkommen.

Weiterhin bleibt unberücksichtigt, dass ein großer Anteil der anfallenden Kosten nicht durch die reine Produktion beziehungsweise Anschaffung anfallen, sondern von weiteren Faktoren, wie etwa Betrieb, Wartung und Verwaltung, abhängen.

Die für die Anschaffung und den laufenden Betrieb anfallenden Kosten sind jedoch alles in allem relativ gut im Vorhinein kalkulierbar.

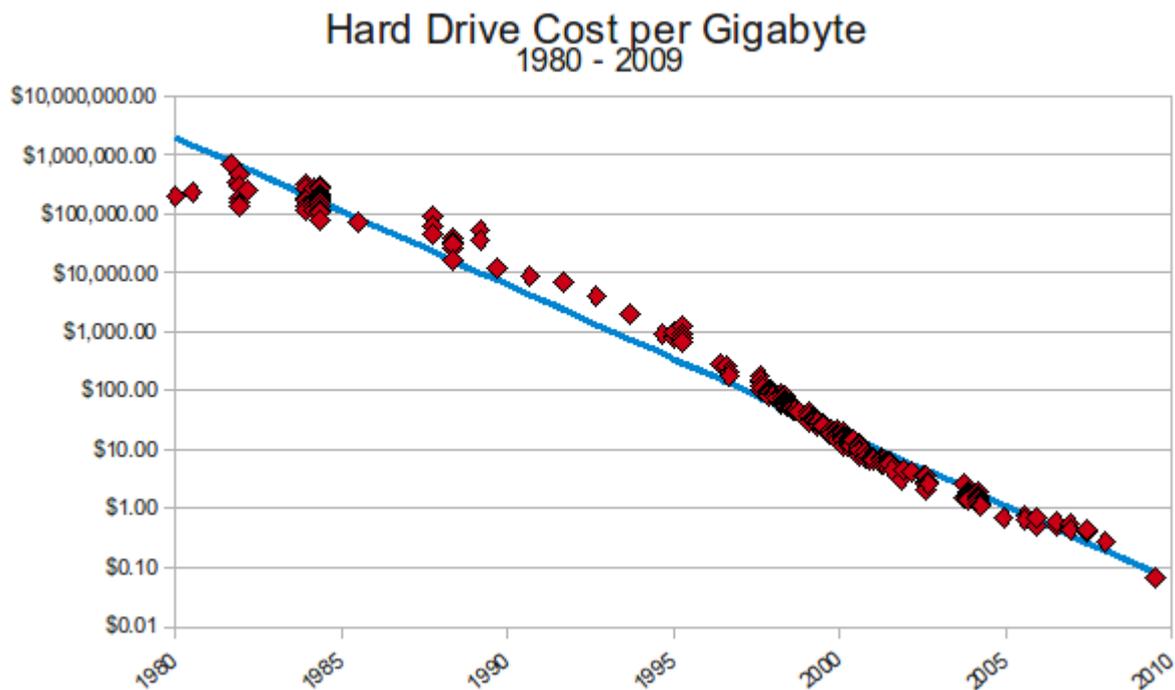


Abbildung 1: Kapazität und Kosten in Relation (<http://www.mkomo.com/cost-per-gigabyte>)

3. Nutzen

Die Quantifizierung des Nutzens, der durch Replikation gestiftet wird, ist nicht ganz so einfach wie bei den Kosten. Hierzu ist zunächst eine formellere Definition dieses Nutzens erforderlich.

Bit Preservation

Bit preservation beschreibt das Problem, einzelne Bits und in erweiterten Sinne digitale Daten für einen gewissen Zeitraum beziehungsweise auf unbestimmte Dauer in ihrem ursprünglichen Zustand zu erhalten.

Zur Messung der Zuverlässigkeit benötigt man eine Metrik, und an diese sind gewisse grundsätzliche Anforderungen zu stellen:

- Berechenbar
 - Vernünftige Methode vorhanden, um das Maß zu berechnen
- Aussagekräftig
 - Direkter Zusammenhang mit der Zuverlässigkeit des eingesetzten Speichersystems
- Verständlich
 - Das Maß und seine Einheiten müssen verständlich sein für alle möglichen Adressaten
- Vergleichbar
 - Gewährleistung von Vergleichbarkeit über unterschiedliche Skalen, Architekturen und zugrunde liegende Technologien

MTTDL

MTTDL (Mean Time To Data Loss) war für einen langen Zeitraum das Standard-Maß für die Verlässlichkeit in Speichersystemen. Damit wird die erwartete Zeit beschrieben, die verstreichen müsste, damit ausreichend viele Versagen vorliegen, um mindestens einen nicht rekonstruierbaren Datenblock vorzuweisen.

Die MTTDL gibt jedoch keine Auskunft darüber, wie viele Daten bei einem durchschnittlichen Auftreten von Datenverlust tatsächlich verloren gehen, sondern lediglich wann ein solches Ereignis eintritt. Daher wird die MTTDL inzwischen nicht mehr als optimales Maß für Leistungsfähigkeit von Bit-Erhaltung.

MTTDL basiert auf einer einfachen Formel, die mittels eines Markov-Modells ausgerechnet werden kann und von der Fehler- und der Reparationsrate und der Anzahl verwendeter Festplatten abhängt.

Gemessen an den vier eingangs genannten Kriterien für ein Zuverlässigkeits-Maß wird nur die Berechenbarkeit zufriedenstellend erfüllt.

Bit Half-Life

Mit der Bit-Halbwertszeit wird die Zeit beschrieben, die vergeht, bis eine Wahrscheinlichkeit von 50% erreicht ist für das Ereignis, dass der Zustand eines Bits kippt.

Dieses Maß ist explizit losgelöst von der Ausfallsicherheit einzelner Speichersysteme und beschränkt sich auf die Daten selbst. Damit ist die bit half-life informativer als die MTTDL, schließlich wird nicht nur die Zeit bis zu Verlust, sondern auch der Umfang adressiert.

Jedoch bleiben auch hier noch Einschränkungen bestehen, so etwa ist die bit half-life wie die MTTDL eine statistische Abschätzung und birgt somit eine Unbestimmbarkeit. Außerdem kann nicht davon ausgegangen werden, dass Bits unabhängig voneinander kippen – somit treten Fehler wahrscheinlicher auf als unter Unabhängigkeit.

Auch hier zeigt sich, dass nur das Kriterium der Berechenbarkeit ausreichend erfüllt wird.

NOMDL

NOMDL (Normalized Magnitude Of Data Loss, dt. normalisierte Größenordnung von Datenverlust) versucht, alle vier Kriterien zu erfüllen.

Bei diesem Ansatz wird die Größenordnung eines Datenverlusts innerhalb einer Einsatzzeit in den Mittelpunkt gestellt. Dies hat zweierlei Vorteile: Es lässt sich für beliebige Einsatzdauern die erwartete Menge an verlorenen Daten bestimmen, und die Einheiten – Bytes und Jahre – sind leicht verständlich. Um die Vergleichbarkeit über eingesetzte Systeme hinweg zu gewährleisten, wird eine Normalisierung vorgenommen bezüglich der verwendbaren Speicherkapazität in einem RAID4-Array.

Darüber hinaus ist die NOMDL berechenbar mittels einer Monte Carlo Simulation.

NOMDL erfüllt somit alle vier an ein Zuverlässigkeitsmaß gestellten Anforderungen.

Quantifizierung

Mit der zunehmenden Zuverlässigkeit gegenwärtiger Speichersysteme nimmt die erforderliche Anzahl zu betrachtender Bits und die Betrachtungsdauer zu, um die Bit-Halbwertszeit ausreichend genau abschätzen zu können.

Abgesehen von den hohen wirtschaftlichen Hürden, die sich hier auftun, wäre auch die Praktikabilität in Frage zu stellen. Und darüber hinaus würde ein solches Experiment zu lange dauern, als dass dessen Ergebnisse noch wirklich sinnvoll zu verwerten wären.

4. Auseinandersetzung Kosten vs. Nutzen

Das Argument, dass eine Lösung von bit preservation nicht experimentell nachweislich ist, verdeutlicht, dass bit preservation per se nicht lösbar sein kann. Eine solche Zuverlässigkeit ließe sich nicht bestätigen.

Zu dieser natürlichen Nicht-Erfüllbarkeit einer vollständigen Zuverlässigkeit kommen noch weitere Faktoren, etwa „silent data corruption“. Damit bezeichnet man das Auftreten von Änderungen an gespeicherten Daten, ohne dass eine Erklärung gefunden oder ein Fehler festgestellt werden könnte. Das bedeutet, dass selbst unter der unmöglichen Annahme optimaler Betriebsbedingungen und jeglicher Fehlerfreiheit Datenverlusten unvermeidlich sind.

Als Annäherung an diesen Umstand sei die Tatsache angeführt, dass selbst auf atomarer Ebene Verfall ein natürliches Ereignis ist und ein gewisser nicht-deterministischer Faktor nicht auszuschließen ist.

Somit bleibt nur eine immer bessere Fehlerbeschränkung, um die Zuverlässigkeit zu erhöhen.

Da diese aber nicht messbar ist, stellt sich die Frage wie eine Gegenüberstellung von Kosten und Nutzen überhaupt bewerkstelligt werden kann.

Erreichbarkeit

Großen Einfluss bei der Konzeption von Speichersysteme in der BigData ist Erreichbarkeit: Müssen die Daten ständig erreichbar, oder werden die Daten erst einmal nur für spätere Verwendung gespeichert? Müssen die Daten unmittelbar verfügbar sein oder ist eine gewisse Latenz vertretbar?

Wird ständig auf die Daten zugegriffen und ist dieser Zugriff nicht stark eingegrenzt auf bestimmte Datensätze, so ist vom Einsatz von Magnetbändern abzusehen. Liegt der Schwerpunkt bei der Datensicherung jedoch darin, die Daten für (eventuelle) zukünftige Verwendung zu erhalten, oder ist eine gewisse Verzögerung unproblematisch, so verhält es sich schon wieder anders.

Kosten für Datenverlust

Ein weiterer Faktor, der bei der Konzeption einen großen Einfluss besitzt, ist die ökonomische Relevanz von Datenverlusten. Verursachen bereits minimale Datenverluste oder können gewisse Verluste ohne wesentliche Konsequenzen kompensiert werden?

Stellt man sich ein wissenschaftliches Forschungsprojekt in der Physik vor, in der ein gewisser Zusammenhang zwischen einzelnen Einflussgrößen betrachtet wird, so wird ein einzelner Datensatz kaum einen Unterschied machen im Resultat – die Masse macht's. In diesem Fall wäre ein gewisser Datenverlust durchaus verschmerzbar.

Demgegenüber sind auch Szenarien denkbar, bei denen einzelne Datensätze selbst erhalten werden sollen und ein Verlust nicht ohne weiteres kompensiert werden kann. Ein Beispiel wäre die Archivierung von Kunstobjekten.

Zur Verfügung stehendes Kapital

Verbesserte Zuverlässigkeit kostet Geld, auch wenn der Zusammenhang nicht so leicht quantifizierbar ist. Im konkreten Anwendungsfall bedeutet das, dass Projekte mit engerem Kapitalrahmen eher weniger Replikationen einsetzen können.

IV. Aus der Praxis

Google-Studie zur Festplattenzuverlässigkeit

Die meisten Studien über die Zuverlässigkeit von Festplatten stammen von den Herstellern selbst und sind somit generell mit Vorsicht zu genießen – zumal die Ergebnisse oft durch Hochrechnungen und Schnellalterungstests gewonnen werden, sodass reale Szenarien nur zu einem gewissen Grad abgebildet werden.

Google hat als End-Nutzer mit einer große Basis an echten Betriebsdaten eine vielversprechende Ausgangssituation für eine Studie zu Plattenzuverlässigkeit. So konnten im Zuge einer Studie Informationen von mehr als 100.000 Festplatten ausgewertet werden, wobei alle maßgeblichen Hersteller abgedeckt wurden.

Unter Berücksichtigung von Daten wie der Umgebungstemperatur, Aktivitätslevel und SMART-Parameter sind folgende Kernpunkte hervorgetreten:

- Es existiert nur eine geringe Korrelation zwischen Ausfallrate und erhöhter Temperatur oder Aktivität
- Gewisse SMART-Parameter wie Abtastfehler und Umverteilungshäufigkeit haben großen Einfluss auf die Ausfallwahrscheinlichkeit.
- Auf Basis von SMART alleine können Festplattenausfälle nicht ausreichend vorhergesagt werden, da vor einigen Ausfällen gar keine Indikation bestand.

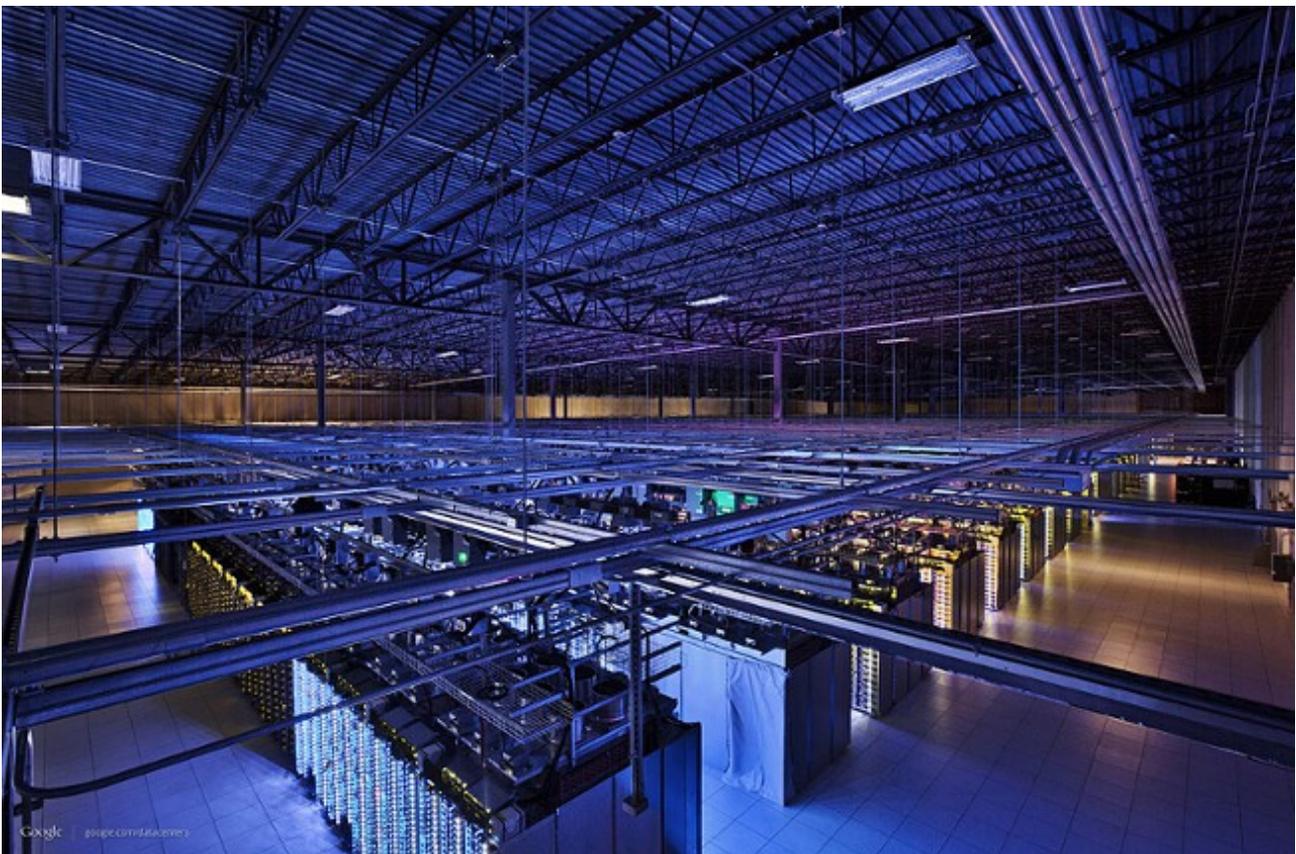


Abbildung 2: Eines der vielen enormen Rechenzentren von Google ([Quelle](#))

V. Quellen

Rosenthal, David S. H.: „Keeping Bits Safe: How Hard Can It Be?“

In: Queue Magazine, Volume 8, Issue 10 (Oktober 2010).

<http://queue.acm.org/detail.cfm?id=1866298>

Rosenthal, David S. H.: „Bit Preservation: A Solved Problem?“

In: The International Journal of Digital Curation, Volume 5, Issue 1 (2010).

<http://www.ijdc.net/index.php/ijdc/article/view/151>

Greenan, Kevin M.; Plank, James S.; Wylie, Jaj J.:

„Mean time to meaningless: MTTL, Markov models, and storage reliability“

In: Hot Storage '10, 2nd Workshop on Hot Topics in Storage and File Systems (2010)

http://static.usenix.org/event/hotstorage10/tech/full_papers/Greenan.pdf

Anderson, Martha: „B is for Bit Preservation“

In: The Library of Congress, Blog on Digital Preservation

<http://blogs.loc.gov/digitalpreservation/2011/09/b-is-for-bit-preservation/>

Lee, K.H.; Slattery, O.; Tang, X.; Lu, R.; McCrary, V.:

„The State of the Art and Practice in Digital Preservation“

In: Journal of Research of the National Institute of Standards and Technology

<http://archive.org/details/jresv107n1p93>

Google Inc.; Pinheiro, E.; Weber, W.-D.; Barroso, L. A.:

„Failure Trends in a Large Disk Drive Population“

In: Proceedings of the 5th USENIX Conference on File and Storage Technologies (FAST'07)

http://labs.google.com/papers/disk_failures.pdf

Proseminar „Big Data“

Langzeitarchivierung von Big Forschungsdaten

Vorgelegt von: Anjela Mayer

am 24. Juni 2013, am KIT

Betreuer: Dr. Rainer Stotzka

Inhaltsverzeichnis

Inhaltsverzeichnis	2
Abkürzungsverzeichnis	3
1 Motivation	4
2 Einführung	5
3 Aktuelle Herausforderungen	7
4 Systematische Lösungsansätze	9
4.1 Deutsche Forschungsgemeinschaft	9
4.2 Data Seal of Approval	11
4.3 Open Archival Information System (OAIS).....	13
4.4 Nestor	16
5 Fazit	17
Literaturverzeichnis	18

Abkürzungsverzeichnis

ALICE	A L arge I on C ollider E xperiment
AOD	A nalysis O bject D ata
ATLAS	A T oroidal L H C A pparatu S
CERN	C onseil E uropéen pour la R echerche N ucléaire
CMS	C ompact- M uon- S olenoid- E xperiment
DFG	D eutsche F orschungs G emeinschaft
DSA	D ata S eal of A pproval
LHCb	L arge H adron C ollider b eauty
LZA	L ang Z eit A rchivierung

Motivation

Die Langzeitarchivierung digitaler Daten hat in den letzten Jahren immer mehr an Bedeutung gewonnen. Dabei bringt sie viele Vorteile mit sich, z.B. die Langzeitverfügbarkeit von Daten. Andererseits gibt es auch viele Probleme die Daten über einen langen Zeitraum verfügbar zu halten und ihre Richtigkeit zu garantieren.

Vor allem in der Forschung ist die Archivierung und Langzeitverfügbarkeit der Daten von enormer Wichtigkeit.

Zu einer guten wissenschaftlichen Praxis gehört die Verifikation der Forschungsarbeiten und -ergebnisse. Die Daten, die zu den bestimmten Ergebnissen geführt haben und auf die sich die Forschungsarbeiten stützen, sollten archiviert und über einen langen Zeitraum erreichbar sein. Dies ermöglicht, dass wissenschaftliche Arbeiten jederzeit auf ihre Korrektheit überprüft werden können, was wiederum die Qualität der Wissenschaft sichert.

Des Weiteren ermöglicht die Langzeitarchivierung eine produktive Nachnutzung der Daten und Kontinuität der Forschung auch über mehrere Generationen hinweg. Langzeitstudien, wie z.B. in der Klimaforschung oder Astrophysik, können betrieben werden wobei die vor langer Zeit archivierten Daten ausschlaggebend für neue Erkenntnisse sind.

In der theoretischen Physik werden häufig neue, bessere Methoden entwickelt um Forschungsdaten auszuwerten. Auch hier ist das Interesse besonders hoch auf Daten aus vergangenen Experimenten zuzugreifen und diese anhand der neuen Methoden neu zu interpretieren.

Häufig entstehen in der Forschung Daten unter besonders seltenen Umständen oder gehen aus Experimenten hervor, die nicht wiederholt werden können. Diese nicht reproduzierbaren Daten müssen ebenfalls archiviert und verfügbar gemacht werden.

Neben seltenen Daten gibt es auch Forschungsexperimente mit einem sehr hohen finanziellen Aufwand. Diese Experimente können aufgrund der hohen Kosten nicht beliebig oft wiederholt werden und schon gar nicht viele Institutionen können sich die Technologie und den Aufbau leisten. Indem diese teuren Experimentdaten archiviert und verfügbar gemacht werden, müssen sie nicht immer wieder reproduziert werden. Außerdem lassen sich die Kosten für den Aufbau und die Instandhaltung der Experimente sparen. Andererseits können Institutionen und Forscher auf die Daten zugreifen, die sonst dazu finanziell nicht in der Lage wären.

1 Einführung

Diese Arbeit beschäftigt sich mit der Langzeitarchivierung von „Big“ Forschungsdaten.

Mit „Big Data“ werden besonders große Datenmengen bezeichnet, die mit herkömmlichen Datenbanken und Verwaltungs-Tools nicht bearbeitet werden können.

Der Begriff Forschungsdaten ist immer in Bezug zur jeweiligen Fachdisziplin zu setzen.

Das können zum Beispiel in der Astrophysik Messdaten aus Teleskopen sein, oder in der Klimaforschung Messdaten, die von Satelliten erfasst werden.

In der Teilchenphysik entstehen „Big“ Forschungsdaten bei Experimenten mit dem Teilchenbeschleuniger. Ein gutes Beispiel hierfür ist der Large Hadron Collider (LHC) am CERN, der Europäischen Organisation für Kernforschung (frz.: *Conseil Européen pour la Recherche Nucléaire*):

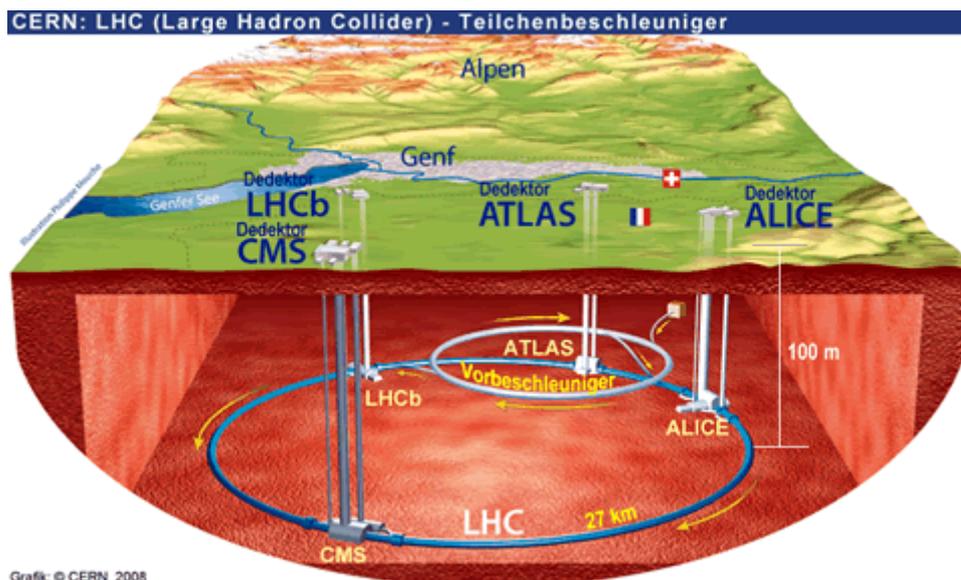


Abbildung 1: Large Hadron Collider (LHC) am CERN¹

Der LHC ist der zurzeit größte und energetischste Teilchenbeschleuniger. Da solche Projekte sehr hohe Investitionssummen erfordern und auch hohe Betriebskosten haben, sind sie nur in Kollaborationen durchführbar. Am LHC werden von weltweiten Kollaborationen die vier Experimente ATLAS, CMS, LHCb und ALICE betrieben. Doch wie entstehen nun Forschungsdaten?

¹ Abbildung 1 Quelle: http://raonline.ch/pages/edu/nw2/phys_research020102b2.html

Jedes Experiment verfügt über einen eigenen Detektor. An den Detektoren werden während des Betriebes sogenannte Rohdaten gemessen. Dabei entstehen riesige Mengen an Daten. Diese werden mindestens für die Dauer des Experimentes vor Ort, am CERN, gespeichert. Im Jahr 2010 wurden am CERN 13 PetaBytes an Daten gesammelt und weltweit zur Analyse verteilt [LZA, Seite 259].

Doch bevor aus den Daten für die Analyse und Langzeitarchivierung geeignete Forschungsdaten entstehen, müssen diese zuvor mehrere Schritte durchlaufen. Diese Rohdaten in Kombination mit den Daten, die den Messdetektor beschreiben ergeben „RECO“-Daten und werden am CERN im Tier-0 Rechenzentrum zwischengespeichert und danach zur Datenanalyse an Tier-1 Rechenzentren weltweit verteilt [LZA, Seite 263].

Da die „RECO“-Daten für wissenschaftliche Analysen wegen ihrer Größe immer noch ungeeignet sind, werden diese zuvor in Analysis Object Data (AOD) Formate konvertiert [LZA, Seite 263]. Dabei werden für das jeweilige Experiment uninteressante Daten weggelassen, sodass im Endeffekt unterschiedliche AOD Datensätze entstehen. Diese AOD Datensätze können dann schließlich analysiert werden.



Abbildung 2: Verarbeitungsschritte von Forschungsdaten (Quelle: Eigene Abbildung)

Ein derartiger Teilchenbeschleuniger erfordert Investitionen von mehreren Milliarden Euro. Zusätzlich sind die Betriebs- und Instandhaltungskosten ebenfalls sehr hoch. Aus diesen Gründen handelt es sich meistens um einmalige Experimente. Dabei kommt es zu sehr seltenen physikalischen Phänomenen. Diese Faktoren machen die dabei entstehenden Daten für die Wissenschaft sehr wertvoll und es ist von höchster Priorität sie für eine lange Zeit verfügbar zu machen.

Die Sicherung und Langzeitverfügbarkeit solcher Daten sind die Aufgaben der digitalen Langzeitarchivierung.

Für die digitale Langzeitarchivierung existiert noch keine genaue Definition. Hier verstehen wir darunter die Bewahrung digitaler Daten, für einen bestimmten Zeitraum und über technologische und soziokulturelle Wandlungsprozesse hinaus. Solche Wandlungsprozesse können beispielsweise Medien, Formate oder Benutzergruppen verändern.

Die Langzeitarchivierung soll also sichern, dass auf die Daten auch in Zukunft zugegriffen werden kann und diese für die Nachnutzung bereithalten.

2 Aktuelle Herausforderungen

Bei der Langzeitarchivierung von „Big“ Forschungsdaten gilt es unterschiedliche Herausforderungen zu meistern.

Aufgrund der enormen Datenmengen, die bei wissenschaftlichen Experimenten, wie beim LHC Beschleuniger, anfallen, ist es unmöglich alle Daten zu speichern. Eine Herausforderung ist die richtige Auswahl an Daten zu treffen. Dabei muss beachtet werden, welche Daten für zukünftige Forscher relevant sein könnten.

Um die Daten zu einem späteren Zeitpunkt auswerten und richtig interpretieren zu können, müssen außerdem noch zusätzliche Daten gespeichert werden, die die technischen und organisatorischen Rahmenbedingungen beschreiben. Zum Beispiel ist es bei den LHC Experimenten besonders wichtig die Daten über den jeweiligen Detektor zu speichern. Außerdem werden in Metadaten alle nötigen Informationen über die Daten selber gespeichert, die für die spätere Analyse notwendig sind.

Zu weiteren Herausforderungen zählen die Wandlungsprozesse. Vor allem die Technologie verändert sich ständig. Damit die Daten auch in Zukunft verfügbar sind, muss auch die technologische Umgebung gesichert werden, wie z.B. Analyse- und Rekonstruktionssoftware für AOD Formate.

Seitens des Langzeitarchivs, sollten gewisse Standards festgelegt werden, wie beispielsweise das Format, in dem die Daten archiviert werden oder die Art und Weise wie die Daten ins Archiv eingespeist bzw. ausgelesen werden. Außerdem muss an die Sicherheit und Korrektheit der im Archiv gespeicherten Daten gedacht werden.

Die aus finanzieller Sicht bedeutendste Frage ist: Wie lange sollen die Daten archiviert werden? Oder anders formuliert: Für wie lange reichen die finanziellen Mittel aus? Je mehr Daten gesichert werden sollen und je länger der Zeitraum in dem sie verfügbar sein sollen, desto mehr Kosten fallen an. Reichen die finanziellen Mittel für das Archiv bzw. die Rechenzentren nicht mehr aus, können die Daten nicht mehr zur Verfügung stehen.

Bei der Langzeitverfügbarkeit der Daten spielen Forschungskollaborationen eine große Rolle. Für die Dauer des Experimentes und die anschließende Analysephase ist die jeweilige Kollaboration im Besitz der erhobenen Forschungsdaten und ist für ihre Langzeitarchivierung und Veröffentlichung zuständig. Nach dem Ende der Kollaboration und in der eigentlichen Phase der Langzeitarchivierung verbleiben die Daten selbstverständlich weiterhin an den Kollaborationsinstituten. Es ist jedoch nicht bestimmt wer nach

der Beendigung der Kollaboration für diese Daten zuständig ist und sich weiterhin um ihre Veröffentlichung und Sicherung kümmert.

Eine andere Herausforderung bei der Langzeitarchivierung ist das Vertrauen in digitale Langzeitarchive, da es keine genauen Richtlinien und Standards für die Gestaltung eines Langzeitarchivs gibt. Das erschwert die Bewertung der Seriosität und Qualität eines Langzeitarchivs.

Für wertvolle Forschungsdaten muss das Archiv auf jeden Fall zuverlässig sein. Es sollte die Daten für den bestimmten Zeitraum archiviert haben und zur Verfügung stellen.

Außerdem sollte es über Mechanismen und Methoden verfügen, die für die Korrektheit der Daten sorgen.

Erste Lösungsansätze gewisse Standards für die Langzeitarchivierung und den richtigen Umgang mit wissenschaftlichen Daten zu etablieren gibt es erst seit dem Jahr 1998 mit der Veröffentlichung des Dokumentes „Sicherung guter wissenschaftlicher Praxis“.²

² Vorschläge zur Sicherung guter wissenschaftlicher Praxis, Empfehlungen der Kommission „Selbstkontrolle in der Wissenschaft“ (1998),
http://www.dfg.de/download/pdf/dfg_im_profil/reden_stellungnahmen/download/empfehlung_wiss_praxis_0198.pdf

3 Systematische Lösungsansätze

Da sich die Wissenschaft noch nicht allzu lange mit dem Thema Langzeitarchivierung auseinandersetzt, gibt es dabei noch keine etablierten Methoden oder Standards, wie die unterschiedlichen Herausforderungen zu bewältigen sind.

Im Folgenden werden verschiedene Lösungsansätze vorgestellt, die Anforderungen und Empfehlungen an die digitalen Langzeitarchive beschreiben.

3.1 Deutsche Forschungsgemeinschaft

Die Deutsche Forschungsgemeinschaft (DFG) ist eine Einrichtung, die sich mit der Förderung der Wissenschaft und Forschung beschäftigt.³

Im Jahr 1998 veröffentlichte die DFG die Denkschrift „Sicherung guter wissenschaftlicher Praxis“. Dieses Dokument beinhaltet 16 Empfehlungen, die eine gute und korrekte wissenschaftliche Arbeit definieren. Für die Langzeitarchivierung war vor allem Empfehlung 7 ausschlaggebend.

Empfehlung 7: „Primärdaten als Grundlagen für Veröffentlichungen sollen auf haltbaren und gesicherten Trägern in der Institution, wo sie entstanden sind, für **zehn Jahre** aufbewahrt werden.“ [WissP, Seite 12]

Dabei sollen auch für die Reproduzierbarkeit der Daten relevante Aufzeichnungen mitgesichert werden. Angezweifelte Forschungsergebnisse sollen somit auch an einem anderen Ort nachvollziehbar sein.

Vor dieser Veröffentlichung spielte der Umgang mit wissenschaftlichen Daten keine große Rolle. Erst mit dem Bekanntwerden von Fällen wissenschaftlichen Fehlverhaltens, gewann die Zugänglichkeit an wissenschaftliche Daten an Bedeutung.

Diese Empfehlung soll vor allem gegen das Verschwinden von Originaldaten, auf die sich die wissenschaftliche Arbeit stützt, absichern.

Das Abhandenkommen von Originaldaten aus einem Labor wird in diesem Dokument als grob fahrlässiges Verhalten bezeichnet. Außerdem wird verdeutlicht, dass solche Verstöße gegen die Grundregeln der wissenschaftlichen Sorgfalt unterbunden werden müssen [WissP, Seite 13].

³ DFG – Deutsche Forschungsgemeinschaft: <http://www.dfg.de/>

Für die Langzeitarchivierung war diese Veröffentlichung der erste Schritt in Richtung Standardisierung. Allerdings ist sie für „Big“ Forschungsdaten jedoch nicht realisierbar. Es wäre natürlich optimal, wenn alle Daten gesichert und über eine lange Zeit verfügbar gemacht werden könnten, doch wegen der enorm großen Datenmengen ist es für „Big“ Forschungsdaten finanziell unmöglich.

Deshalb werden bei den LHC Experimenten am CERN, die Originaldaten mindestens für die Dauer des Experiments gespeichert.

Hypothetisch betrachtet wären für eine Speicherung von jährlich ca. 13 PetaBytes an Experimentdaten und über 10 Jahre hinweg, Kapazitäten von 130 PetaBytes nötig.

Außerdem müssen die Daten für die eigentliche Analyse immer noch in AOD Datensätze konvertiert werden um eine geeignete Datengröße zu erreichen.

Die Speicherung der Rohdaten wäre praktisch und finanziell keine sinnvolle Lösung.

3.2 Data Seal of Approval

Das Data Seal of Approval (DSA)⁴ ist ein Qualitätssiegel für die Langzeitarchivierung. Es richtet sich unter anderem an Forschungsinstitutionen und garantiert, dass die Daten auch in Zukunft in hoher Qualität und Zuverlässigkeit verfügbar sind.



Abbildung 3: Data Seal of Approval⁵

Der Siegel lässt sich über die DSA Webseite beantragen⁶, woraufhin das Archiv auf die Kriterien des Leitfadens geprüft wird. Der DSA Leitfaden⁷ beinhaltet 16 Kriterien, die an das Archiv und an den Archivnutzer gerichtet sind. Im Folgenden sind die Kriterien zusammengefasst.

Archivnutzer

Die Benutzerkriterien richten sich einerseits an den Produzenten der Daten, der sie im Repository ablegen will, und andererseits an den Konsumenten, der auf die Daten zugreifen und diese benutzen möchte.

Der Datenproduzent ist dazu verpflichtet seine Daten in einem zulässigen Format im Repository abzulegen [DSA, Punkt 2.]. Außerdem müssen diese Daten ausreichend dokumentiert werden, z.B. mit Angaben zum Urheber, Metadaten und Rahmenbedingungen unter welchen die Daten entstanden sind [DSA, Punkt 1.].

Der Datenkonsument muss die Nutzungsbedingungen erfüllen und außerdem die Lizenzvereinbarungen des Repositoriums akzeptieren [DSA, Punkte 14.-16.].

⁴ DSA – Data Seal of Approval: <http://datasealofapproval.org/>

⁵ Abbildung 3 Quelle: <http://datasealofapproval.org/>

⁶ Apply for DSA: <https://assessment.datasealofapproval.org/apply/>

⁷ DSA Guidelines: <http://datasealofapproval.org/?q=node/35>

Repositorium

Bei dem Repositorium muss ein Plan für die Langzeitarchivierung der digitalen Daten existieren, worin die notwendigen organisatorischen, technischen, finanziellen und rechtlichen Methoden und Mittel festgelegt sind.

Außerdem garantiert das Repositorium die Integrität von den digitalen Objekten und den zugehörigen Metadaten [DSA, Punkte 4.-13.].

Wie schon bei den Empfehlungen der „guten wissenschaftlichen Praxis“ sind die Kriterien auch hier sinnvoll für die Langzeitarchivierung, jedoch an manchen Stellen nicht für „Big“ Forschungsdaten geeignet.

Auf erste Schwierigkeiten trifft man bereits bei dem für das Repositorium zulässigen Format. In den Rechenzentren am CERN werden die Experimentdaten in AOD Datensätze konvertiert. Dabei werden die, für das jeweilige Experiment, unwichtigen Daten weggelassen. Da die unterschiedlichen Experimente andere Präferenzen an den Daten haben, entstehen aus denselben Rohdaten, aufgrund der variierenden Konvertierung, AOD Datensätze in unterschiedlichen Formaten. Diese werden vor Ort in Tier-0 Rechenzentren am CERN gespeichert. Somit sind dort Datensätze in unterschiedlichen Formaten archiviert.

Zur Analyse werden die Daten weltweit auf Tier-1 Rechenzentren von unterschiedlichen Institutionen verteilt [LZA, Seite 263]. Da es jedoch keine allgemeinen organisatorischen Standards für Langzeitarchive und Langzeitarchivierung gibt, können die Pläne zur Langzeitarchivierung stark voneinander abweichen.

Eine weitere Problematik in der Teilchenphysik betrifft die Benutzergruppen. Die Datenproduzenten sind trivialerweise die Mitglieder der Kollaboration, die an dem jeweiligen Experiment beteiligt sind. Die Datennutzer sind meistens jedoch auch nicht mehr als nur die Mitglieder der Kollaborationen. Es ist üblich, dass die Daten während der Experiment- und Analysephasen nicht veröffentlicht werden. In der Langzeitarchivierungsphase, wenn die Experimente und meisten Analysen schon abgeschlossen wurden und wenn sich manche Kollaborationen bereits beginnen aufzulösen, werden die Daten jedoch ebenfalls nicht veröffentlicht und zugänglich für weitere Benutzergruppen gemacht. Ein Grund dafür könnte neben dem fehlenden Interesse auch die fehlende Zuständigkeit für die Veröffentlichung der Daten sein.

3.3 Open Archival Information System (OAIS)

Das Open Archival Information System (OAIS) ⁸ ist ein Referenzmodell für ein Archivinformationssystem. Es wurde 2002 vom Consultative Committee of Space Data Systems (CCSDS) ⁹ zuerst als Empfehlung und ab 2003 auch als ISO-Standard 14721:2003 veröffentlicht. Die aktuellste Version des OAIS Referenzmodells ist ISO 14721:2012.

Das OAIS Modell ist eines der einflussreichsten Dokumente für die Langzeitarchivierung. Es behandelt nicht nur den technischen Aufbau des Archivs, sondern auch dessen Organisation. OAIS beschreibt, wie die vom Produzenten erzeugten digitalen Daten in das Archivsystem gelangen sollen und wie der Konsument auf die im Archiv gespeicherten Daten zugreifen kann. Es beinhaltet außerdem einen Überblick über die Bearbeitungsschritte, die das Archiv für die langfristige Verfügbarkeit der Daten vornehmen soll. OAIS beschreibt jedoch nicht die konkrete Implementierung des Repositoriums.

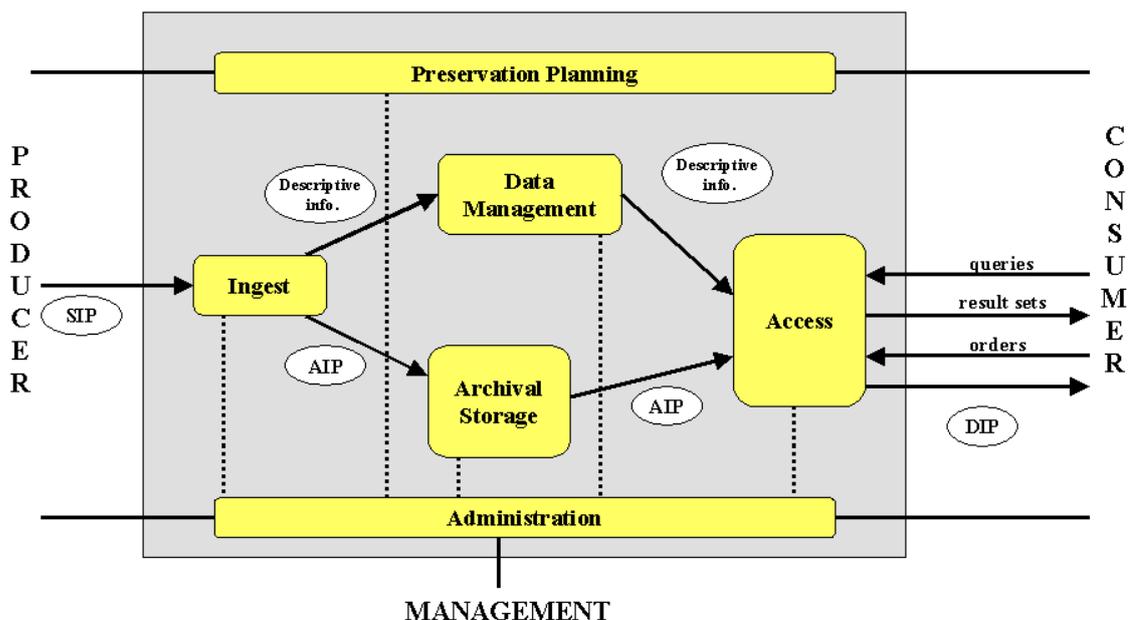


Abbildung 4: Open Archival Information System ¹⁰

Der Hintergrund für die Entwicklung des Modells war die Feststellung, dass Daten nach längerer Zeit in Archiven nicht mehr lesbar sein könnten.

Das OAIS beinhaltet Empfehlungen, die es ermöglichen digitale Daten langfristig zu archivieren. Langfristig bedeutet über technologische und soziokulturelle Wandlungs-

⁸ OAIS – Reference Model for an Open Archival Information System: <http://public.ccsds.org/publications/archive/650x0m2.pdf>

⁹ CCSDS – Consultative Committee of Space Data Systems: www.ccsds.org

¹⁰ Abbildung 4 Quelle: <http://www.ukoln.ac.uk/metadata/publications/ijlim-2003/>

prozesse hinaus. Es muss also die Auswirkungen des technologischen Wandels, wie z.B. die Unterstützung neuer Dateiformate und Datenträger berücksichtigen und darüber hinaus anpassungsfähig gegenüber wechselnden Benutzergruppen sein.

OAIS TRAC

TRAC (eng.: *Trustworthy Repositories Audit & Certification*)¹¹ ist ein ISO Standard für die Zertifizierung von OAIS Repositorien. Es wurde von OCLC (eng.: *Online Computer Library Center*), RLG (eng.: *Research Libraries Group*) und von der NARA (eng.: *National Archives and Records Administration*) entwickelt.

TRAC beinhaltet eine Checkliste für ein zuverlässiges OAIS Langzeitarchiv. Diese Checkliste gliedert sich in drei Untergruppen:

- A) **Organisational Infrastructure** – Organisatorischer Rahmen des digitalen Langzeitarchivs, der finanzielle und personelle Ressourcen, rechtliche Bedingungen und die Zieldefinition des digitalen Langzeitarchivs beinhaltet.
- B) **Digital Object Management** – Der Umgang mit Objekten beinhaltet organisatorische sowie auch technische Aspekte. In diesem Abschnitt werden die Funktionen, Prozesse und Methoden beschrieben, über die das Langzeitarchiv verfügen muss, für die korrekte Aufnahme, Management und Zugänglichkeit digitaler Objekte.
- C) **Technologies, Technical Infrastructure & Security** – Hier werden technischen Kriterien der Infrastruktur des Langzeitarchivs und die notwendigen Sicherheitsvorkehrungen beschrieben.

[TRAC]

TRAC bildet die Basis für den ISO 16363¹² Standard, der 2012 im „Trusted Digital Repository“ Dokument veröffentlicht wurde.

Das OAIS Modell und die TRAC Kriterien sind ziemlich allgemein definiert. Dadurch finden sich hier viele nützliche Kriterien für die digitale Langzeitarchivierung von „Big“ Forschungsdaten. Da diese Dokumente jedoch nicht speziell für Big Data verfasst wurden, können diese nicht ohne weitere Änderungen übernommen werden. Vor allem in dem letzten Abschnitt der TRAC Checkliste, wo es um die Planung der Sicherung der digitalen Daten geht müssten Änderungen vorgenommen werden. Zum Beispiel sind

¹¹ TRAC – Trustworthy Repositories Audit & Certification: Criteria and Checklist, http://www.crl.edu/sites/default/files/attachments/pages/trac_0.pdf

¹² ISO 16363: 2012: http://www.iso.org/iso/catalogue_detail.htm?csnumber=56510

dort Backups der digitalen Daten an unterschiedlichen Orten vorgesehen [TRAC, Seite 44]. Je mehr Replikationen der Daten existieren, desto sicherer sind sie. Bei sehr großen Datenmengen, wie den Forschungsdaten aus dem LHC Teilchenbeschleuniger wird Replikation aufgrund der Datenmengen jedoch erschwert und bedarf spezieller Lösungsansätze.

3.4 Nestor

Nestor¹³ ist ein deutsches Kompetenznetzwerk, welches sich seit 2003 mit dem Thema Langzeitarchivierung beschäftigt. Dieses Netzwerk verbindet spartenübergreifend betroffene Institutionen, Experten und aktive Projektnehmer, die mit der Langzeitarchivierung arbeiten. Nestor fördert vor allem den Austausch von Informationen und die Entwicklung der für die Langzeitarchivierung notwendigen Standards. Es ist die größte deutschsprachige Informationsplattform für Fragen rund um digitale Langzeitarchivierung.

Nestor stellt viele wichtige Publikationen zum Thema Langzeitarchivierung zur Verfügung, wie z.B. das nestor-Handbuch¹⁴, eine Enzyklopädie zum Einstieg in die Langzeitarchivierung.

Außerdem bietet Nestor auch Übersetzungen internationaler Publikationen zu dem Thema an. Beim „Kriterienkatalog vertrauenswürdige digitale Langzeitarchive“ handelt es sich um die deutsche Version der OAIS TRAC Checkliste.

Aus diesem Kriterienkatalog ist auch der deutsche DIN 31644:2012¹⁵ Standard entstanden.

Äquivalent zum TRAC beinhaltet auch der Kriterienkatalog eine Checkliste mit 34 Kriterien, die die Zuverlässigkeit des digitalen Langzeitarchivs sicherstellen sollen.



Abbildung 5: Nestor-Siegel¹⁶

Über die Nestor Webseite kann, ähnlich wie beim DSA Siegel, der Nestor-Siegel für das digitale Langzeitarchiv beantragt werden. Grundlage dafür ist die DIN 31644 Norm.

¹³ Nestor: <http://www.langzeitarchivierung.de>

¹⁴ DIN 31644:2012: <http://www.beuth.de/de/norm/din-31644/147058907>

¹⁵ Nestor-Handbuch: Eine kleine Enzyklopädie der digitalen Langzeitarchivierung, Version 2.3 – 2010
<http://www.beuth.de/de/norm/din-31644/147058907>

¹⁶ Abbildung 5 Quelle: http://www.langzeitarchivierung.de/Subsites/nestor/DE/nestor-Siegel/siegel_node.html

4 Fazit

Diese vorgestellten Lösungsansätze lassen sich größtenteils, aufgrund der allgemeinen Formulierung, auch auf die Langzeitarchivierung von „Big“ Forschungsdaten übertragen. An ein paar Stellen wird jedoch deutlich, dass diese nicht als Lösungen für Big Data verfasst wurden. Es wäre finanziell unmöglich alle Forschungsdaten für einen Zeitraum von 10 Jahren zu sichern und verfügbar zu machen. Vor allem wenn es in die konkrete, praktische Umsetzung geht, benötigt die Langzeitarchivierung von Big Data andere Methoden, wie zum Beispiel bei der Sicherung der digitalen Forschungsdaten.

Außerdem handelt es sich bei den vorgestellten Lösungsansätzen lediglich um Empfehlungen. Wie die digitalen Langzeitarchive letztendlich an den Forschungsinstitutionen realisiert werden, bleibt jedem selbst überlassen. Da es lange Zeit keine allgemein geltenden Richtlinien für den Umgang mit Forschungsdaten gab, wird in jeder wissenschaftlichen Fachdisziplin mit der Langzeitarchivierung digitaler Daten anders umgegangen. Es gibt keine fachübergreifenden Standards.

Andererseits sind solche Standards nicht einfach realisierbar, da zum Beispiel jede Forschungsgruppe die Daten in unterschiedlichen Dateiformaten abspeichert.

Um das Wissen zu fördern und eine sinnvolle Zusammenarbeit zu gestalten sollte jedoch weiterhin nach gemeinsamen Standards gesucht werden, die eine fachinterne aber auch fachübergreifende Kooperation ermöglichen.

Ein weiterer wichtiger Punkt ist die Nachnutzung der Forschungsdaten. Es muss unbedingt die Zuständigkeit für die Veröffentlichung der Forschungsdaten und für den „Open Access“ geklärt werden. Wenn auf diese wertvollen und teuren Daten nicht mehr zugegriffen wird, weil sie nicht frei zugänglich sind, ist es ein großer Verlust des wissenschaftlichen Potentials. Es wird geschätzt dass ca. 10% der wissenschaftlichen Publikationen während der Langzeitarchivierungsphase entstehen [LZA, Seite 259].

Literaturverzeichnis

LZA (2012): Nestor Langzeitarchivierung von Forschungsdaten: Eine Bestandsaufnahme, Version 1.0 – 2012, Hrsg. von: Heike Neuroth, Stefan Strathmann, Achim Oßwald, Regine Scheffel, Jens Klump, Jens Ludwig (Quelle: http://nestor.sub.uni-goettingen.de/bestandsaufnahme/nestor_lza_forschungsdaten_bestandsaufnahme.pdf)

WissP (1998): Vorschläge zur Sicherung guter wissenschaftlicher Praxis, Empfehlungen der Kommission „Selbstkontrolle in der Wissenschaft“ (1998), (Quelle: http://www.dfg.de/download/pdf/dfg_im_profil/reden_stellungnahmen/download/empfehlung_wiss_praxis_0198.pdf)

DSA: Data Seal of Approval Guidelines (Quelle: <http://datasealofapproval.org/?q=node/35>)

TRAC (2007): Trustworthy Repositories Audit & Certification: Criteria and Checklist, OCLC and CRL (Quelle: http://www.crl.edu/sites/default/files/attachments/pages/trac_0.pdf)



Metadaten – Welche Informationen sind wichtig?

Proseminar – Big Data Applications

Mert Turan

Matrikelnummer: 1499532

Studiengang: Informationswirtschaft B. Sc.

Betreuer:

Dr. Rainer Stotzka

Motivation

Heut zu Tage gibt es eine sehr hohe Anzahl an Daten, die man in Archiven findet. Wie ist es nun möglich, in einem Archiv mit Unmengen an Dokumenten, ein bestimmtes davon zu finden? Um die Suche zu erleichtern und überhaupt zugänglich zu machen, müssen Schlagwörter und Beschreibungen vorliegen. Der Beginn des „Digitalen Zeitalters“ begann am Anfang des Jahres 2000¹. Seitdem ist die Digitalisierung, Texterfassung (und –erkennung) und die Suche in traditionellen textbasierten Medien bereits auf hohem Niveau möglich. Um Dokumente leichter zu finden, müssen diese beschrieben und in Textform vorliegen. Was kann nun in einem Dokument beschrieben werden? Hierzu gehören beispielsweise der Inhalt, technische Parameter, Produktionsprozesse etc.

Die vorliegende Arbeit widmet sich der Frage, wie Daten beschrieben werden, was die Gründe hierfür sind und wo diese eingesetzt werden.

¹ <https://de.wikipedia.org/wiki/Digitalisierung>

Inhaltsverzeichnis

1 Einführung	1
2 Arten von Metadaten.....	3
2.1 Welche Arten von Metadaten gibt es?.....	3
2.2 Wie kann man Metadaten kategorisieren?.....	4
3 Zu welchem Zweck werden Metadaten eingesetzt?	5
3.1 Erfassung von Metadaten.....	5
4 Metadaten in Big Data	5
5 Metadaten in wissenschaftliche Experimente	7
6 Wie werden Metadaten in der Wirtschaft eingesetzt?.....	8
7 Beispiel von Metadaten in Fotos	8
8 Fazit	10

1 Einführung

Die meist verbreitete Definition für Metadaten ist: „Daten über Daten“², wird jedoch dabei häufig aufgrund ihrer geringen Differenzierung des Themengebietes, als unzureichend angesehen.

Jedes Objekt besitzt „Hintergrundinformationen“. Metadaten geben Informationen über bestimmte Objekte wie z.B. über Gemälden, Bücher, Tondokumente aber auch Datenbanken aus, um z.B. die Recherche zu erleichtern. Tim Berners-Lee, ein britischer Physiker und Informatiker, sagt über Metadaten im Web: “Metadata is machine understandable information about web resources or other things“.³

Wenn man sich ein Buch kauft, so ist es auch für den Leser wichtig, von wem das Buch verfasst wurde (*Autor*), um welche *Genre* es sich handelt, wie viele Seiten das Buch besitzt (*Seitenanzahl*) u.v.m. Diese Informationen bezeichnet man als Metadaten. Nicht nur für den Nutzer sind diese Daten wichtig, sie werden auch für andere Zwecke benötigt. So helfen Metadaten beispielsweise bei der Buchrecherche in einer Bibliothek oder bei der Suche nach Gemälden in Inventar(system)en.

Metadaten umfassen aber nicht nur Informationen für die bessere Recherche, sondern auch zu anderen Aspekten, die für die Nutzung, Verwaltung und Pflege der Objekte notwendig sind. Genauso wie wertvolle Gegenstände für Menschen wichtig sind und deshalb aufgehoben und gepflegt werden, müssen Metadaten auch in bestimmten Archiven gespeichert und verwaltet werden (Siehe unterer Abschnitt). Metadaten sind nicht nur einfache Beschreibungen von Daten, sie ermöglichen auch den physischen Zugang zu den Objekten und sind unveränderbar. Um den Wert des Objektes und die Suche nach einzelnen Objekten beispielsweise in einer Bibliothek zu erleichtern, sollten Metadaten ausführlich und umfassend sein. Dadurch können auch Daten desselben Typs leichter differenziert und gefunden werden. Durch eine umfassende und differenzierte Beschreibung von Metadaten ist

² http://de.wikipedia.org/wiki/Wikipedia:WikiProjekt_Metadaten

³ <http://www.w3.org/DesignIssues/Metadata.html>

Metadaten

die Recherche für den Suchenden flexibler und damit auch der wissenschaftliche Wert des Objektes.

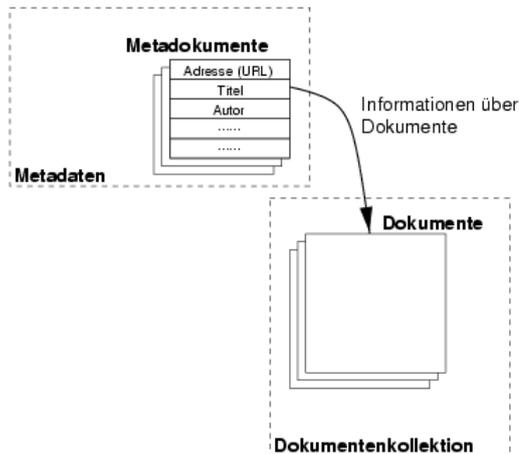


Abbildung 1⁴

Mit der Abbildung 2 soll noch einmal deutlich gemacht werden, wozu Metadaten dienen und wie Daten ohne Hilfe anderer Daten aussehen würden. Jedoch kann man nicht alle Daten, die Informationen über andere Daten geben, als Metadaten bezeichnen: Geben Daten nützliche Informationen, - dazu gehört, dass die Metadaten z.B. die Auswertung bzw. die Verwaltung anderer Daten erleichtern – so bezeichnet man diese als Metadaten. Wie schon oben erwähnt, wurden die meisten Metadaten für Bibliotheken entworfen. In so einer Bibliothek wären beispielsweise das Veröffentlichungsdatum, der Autor und der Verlag eines Buches die Metadaten. Was nicht zu den Metadaten gehören würde, wären beispielsweise die Farbe der Schrift oder die Qualität des Deckblatts, obwohl es sich bei beiden Informationen um Daten über die verwalteten Daten, nämlich die der Bücher, handelt (Ein genaueres Beispiel im Abschnitt „Beispiele zu Metadaten in Fotos“ mehr). Zusammenfassend lässt sich somit sagen, dass es sich bei Metadaten um Informationen handelt, die nützlich, informativ und zugänglich für den Nutzer sind. Metadaten erleichtern somit die Verwaltung und Auswertung von Daten.

⁴http://www.swisseduc.ch/informatik/internet/internet_recherche/informationsbeschaffung_im_internet/Metadaten.html

Grundsätzlich lässt sich sagen, dass eine Verwaltung von Daten, die keine Metadaten besitzen, ohne Bedeutung ist.

2 Arten von Metadaten

2.1 Welche Arten von Metadaten gibt es?

Metadaten können hinsichtlich vieler Perspektiven⁵ unterschieden werden:

Datenerfassungsmethoden

Metadaten werden automatisch durch das System bzw. durch den Computer erfasst (generiert), oder durch den Menschen manuell.

Datenquellen

Interne Daten, die bei der Produktion entstehen. Dazu gehören z.B. Produzenten, Entstehungsdaten etc.

Externe Daten, die bei der Produktion neuer Produkte entstehen.

Status

Dynamische Metadaten (können sich während der Nutzung des Objektes ändern).

Zu diesen gehören:

- Kurzfristige Metadaten
- Langfristige Metadaten
- Statische Metadaten (Autor, Titel, Laufzeit...)
- Unstrukturierte Metadaten

⁵http://www.archaeobooksalkriese.de/nestor/index.php?option=com_content&view=article&id=62&Itemid=70

2.2 Wie kann man Metadaten kategorisieren?

Hinzu kommt, dass Metadaten in mehrere Kategorien⁶ unterteilt werden können:

Inhaltliche Metadaten

Der Content des digitalen Objektes wird durch inhaltliche Metadaten beschrieben. Inhaltliche Metadaten wären bspw. *Autor, Titel, Schlagworte*, um Objekte leichter zu finden

Strukturelle Metadaten

Der Zusammenhang von Daten mit anderen Daten wird beschrieben. „Bei einer Audio – Datei kann somit der Zusammenhang mit einer Powerpoint – Präsentation festgehalten werden“

Administrative Metadaten

Hierzu gehören beispielsweise der Speicherort, Dateiname und rechtliche Aspekte (Urhebernachweis)

Technische Metadaten

Hierzu gehören Dateiformate, Angabe der Version (von dem Objekt), Datenmenge, Softwarenamen etc....

Metadaten zur Langzeitbewahrung

Es existieren Informationen, die besonders für digitale Langzeitbewahrung benötigt werden. Zu diesen Informationen könnte beispielsweise das Datum des Kopierens einer Datei auf andere Medien wie z.B. auf Archive. Auf solche Informationen ist die Nutzbarkeit angewiesen.

⁶ http://files.dnb.de/nestor/sheets/06_metadaten.pdf

3 Zu welchem Zweck werden Metadaten eingesetzt?

Metadaten sind für alle wissensbasierte Systeme, für das Internet und auch für Data Warehouse Systeme wichtig und werden somit größtenteils auch in diesem Gebiet eingesetzt. Als Ziel ist dabei die *Verbesserung der Qualität* der Daten. Außerdem sollen Metadaten es dem Nutzer bei der Selektion von Daten erleichtern, notwendige Daten einfacher zu finden und beim Verständnis von relevanten Daten unterstützen. Heutzutage gibt es eine Vielzahl von Unternehmen, die sich mit dem Einsatz von Metadaten im Internet beschäftigen. Die „Dublin Core Metadata Initiative“ der OCLC⁷ wäre hierfür ein gutes Beispiel.

3.1 Erfassung von Metadaten⁸

Es gibt unterschiedliche Varianten, Metadaten zu erfassen und bereitstellen.

Technische Daten (siehe „Abschnitt 1.3.1“) werden automatisch festgelegt. Andere müssen recherchiert oder durch „geschultes Personal“ erstellt werden. Zu diesen gehören einfache, inhaltliche Beschreibungen, Einordnung und rechtliche Aspekte. Außerdem ist festzustellen, welche Informationen für die Langzeitbewahrung benötigt werden, um zuordnen zu können, welche Metadaten welchen Originalen angehören. Somit ist es wichtig bei Digitalisierungen, möglichst viele Metadaten zu identifizieren und diese richtig einzusetzen.

4 Metadaten in Big Data

Große Unternehmen, Museen, Bibliotheken etc. besitzen „Big Data“, also „besonders große Datenmengen“, die nur durch bestimmte Arten von Software wie z.B. „Verarbeitung vieler Spalten innerhalb eines Datensatzes“ verarbeitet werden kön-

⁷ Mehr zur Dublin Core Metadata Initiative unter „<http://www.oclc.org/de/de/default.htm>“

⁸ http://files.dnb.de/nesstor/sheets/06_metadaten.pdf

nen.⁹ Bei Unmengen an Daten kommt es oft vor, dass Daten verloren gehen bzw. nicht mehr gefunden werden.¹⁰ Auch kleine Unternehmen wie z.B. kleine Banken, besitzen auf ihren Servern Informationen, die eine Speicherkapazität von mehreren Terrabytes benötigen, die man auch schon als „Big – Data“ ansehen kann. Ohne bestimmte Prozesse lassen sich diese Daten nicht analysieren. So sehen Unternehmen seit wenigen Jahren die Gefahr, dass eine Menge von Daten so rasch anwachsen wird, sodass diese bald gar keinen Überblick mehr haben werden. Der Big – Data Bestand steigt nämlich rasant an. Die Autoren der Studie „The Future of Big Data“¹¹ behaupten, „dass es den meisten Unternehmen und Behörden an dem nötigen Methodenwissen und den Tools fehlt, das Potenzial, das dieser Big – Data – Bestand birgt, auszuschöpfen.“ Damit Unternehmen unter keinem Datenverlust leiden, sind Metadaten eines der wichtigsten Kernpunkte, die es gibt. „Sie sind das A und O für das Verwalten von Daten.“ Wie können nun Metadaten in Big Data aussehen? Abbildung 3 soll verdeutlichen, dass es recht kompliziert werden kann, welche Daten durch andere Daten beschrieben werden sollen. Das Model soll zeigen, dass bei jeder Studie, die in Instituten durchgeführt werden, alle Daten gespeichert werden und in einer Datenkollektion gehalten werden müssen. Eine Studie wird dem Model nach in 3 Arten unterteilt: Experimente, Messungen und Simulationen, in denen Tausende, gar Millionen von weiteren Daten vorhanden sind, wodurch auch die Studien erleichtert und somit in Archiven gehalten werden können.

⁹ http://de.wikipedia.org/wiki/Big_Data

¹⁰ <http://www.computerwoche.de/a/big-data-fluch-segen-oder-einfach-viel-arbeit,2527644>

¹¹ <http://pewinternet.org/Reports/2012/Future-of-Big-Data.aspx>

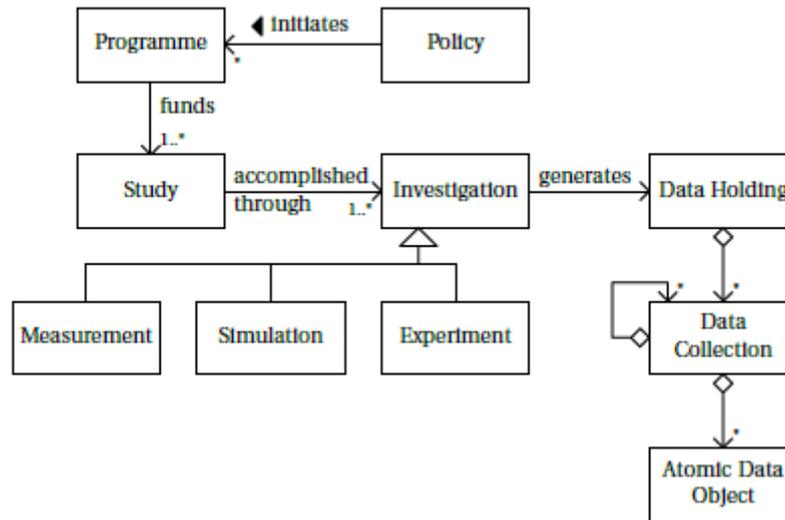


Abbildung 3: High Level entity model from the CCLRC Scientific Metadata Model ¹²

5 Metadaten in wissenschaftliche Experimente

Metadaten sind wichtige Bestandteile eines jeden Objektes und somit auch wichtig für alle Nutzer, Suchmaschinen oder Bibliotheken. Ebenso wichtig sind sie jedoch auch für wissenschaftliche Experimente:

Führt man ein Experiment durch, so braucht man für das Objekt, mit dem man das Experiment durchführen möchte, jede brauchbare Information. Es ist wichtig, zu wissen, wann das Objekt erstellt/kreiert wurde, woraus es besteht usw. Alle Eigenschaften sind für solche Experimente wichtig. Mit Metadaten kann man außerdem mehrere, sich sehr ähnelnde Objekte unterscheiden, da – wie schon erwähnt – Metadaten Informationen über Herstellungsdatum, Größe, Versionen etc. geben.

¹²“Scientific Data Application Profile Scoping Study Report” von Alexander Ball, UKOLN, University of Bath Vers. 1.1

6 Wie werden Metadaten in der Wirtschaft eingesetzt?

Woher wissen viele Unternehmen, auf was bestimmte Kunden gezielt Wert setzen? Großunternehmen besitzen Informations- und Kommunikationssysteme, die dafür zuständig sind, durch Metadaten explizit an Kunden näher zu kommen. Durch das Internet wird einem Unternehmen auch die Möglichkeit gegeben, Kundenkontakte zu pflegen, das wir als CRM (engl.: Customer Relationship Management) bezeichnen und Werbung gezielt an Kunden zu verteilen (Marketing:Targeting). Als Beispiel kann hier „Amazon“ erwähnt werden. Der Kontakt zu Kunden wird sehr eng gehalten: Wenn der Suchende ein Produkt gefunden hat, werden auf der jeweiligen Produktseite weitere ähnliche Artikel zu diesem Produkt angezeigt, welches man als Cross Selling bezeichnet.¹³ Die Grundlage für Cross-Selling sind Metadaten. Je nach freier Auswahl des Kunden können ihm auch per Mail für ihn interessante Artikel vorgeschlagen werden. Bewertungen für Produkte besitzen auch einen hohen Wert. Durch Produktempfehlungen (negative oder positive Kritik) kann der Umsatz des jeweiligen Produktes je nach Empfehlungswahrscheinlichkeit entweder sinken oder stark zunehmen, da durch „Fremder“ Beurteilung, verschiedene Meinungen zum Produkt vorhanden sind. Ohne Metadaten der jeweiligen Kunden wären solche Strategien gar nicht möglich. Somit sind diese auch in der Wirtschaft sehr wichtig, da durch sie verschiedene Kaufverhalten von unterschiedlichen Kunden abgebildet werden können.¹⁴

7 Beispiel von Metadaten in Fotos

Wie wir schon wissen, sind Metadaten Informationen bzw. Beschreibungen zum Inhalt eines Objekts. Wenn wir konkret Metadaten von Bildern ansprechen, muss man sich auch mit Begriffen wie „IPTC“, „Exif“ oder „XMP“ vertraut machen. Was bedeuten diese?

¹³ Mehr zu „Cross Selling“: <http://en.wikipedia.org/wiki/Cross-selling>

¹⁴ <http://www.informatik.uni-oldenburg.de/~iug11/ge/website/handel.html>

Metadaten

Exif (Engl.: Exchangeable Image File Format) sind Daten und bedeutet, dass Fotoaufnahmegeräte wie Kameras, Scanner oder auch Smartphones Informationen direkt auf dem Gerät in der Bilddatei schon automatisch abspeichern. Zu diesen Informationen gehören unter anderem:

- Blende
- Bildgröße (Pixelangabe)
- Auflösung
- Kamerahersteller
- Kameramodell

Auch können damit Informationen wie „Copyright“ oder „Bildbeschreibung“ später ergänzt werden.

Heutzutage werden Exif - Daten von allen bekannten Bildbearbeitungsprogrammen unterstützt und können auch direkt im Programm verändert werden (jedoch nur wenige wie „Copyright“ oder „Bildname“).

IPTC (Engl.: International Press Telecommunications Council) sind ähnliche Daten wie Exifdaten, die von der Firma „Adobe“ entwickelt wurde. Sie unterscheiden sich nur dadurch, dass IPTC Daten auch den Inhalt des Bildes beschreiben. Wenn man also zum Beispiel ein Foto im Urlaub geschossen hat, kann man dem Bild zeitgleich eine Beschreibung geben. Zu den IPTC Informationen gehören unter anderem:

- Überschrift
- Schlüsselwörter
- Kategorien
- Bildbeschreibung

Metadaten

Weiter unterscheidet sich IPTC dadurch, dass die Daten direkt kategorisiert werden können und mit „speziellen – vor allem für Bildagenturen wichtig – Informationen zu Urheber und Copyright versehen.“¹⁵

Sowohl bei Exif-, als auch bei IPTC – Daten können die Metadaten nur in JPG, TIF und PSD – Dateien abgespeichert werden.

XMP (Engl.: Extensible Metadata Platform = Erweiterbare Plattform für Metainformationen) ist auch, wie die vorigen auch, eine Entwicklung von Adobe. Es ist eine Art Zusammenfassung von IPTC und Exif: Es werden die Informationen von Exif und IPTC abgespeichert und zeitgleich ist eine eigene Beschreibungen zum Bild zu machen, gestattet. Metadaten können bei XMP jedoch auch in andere Formate wie JPG, TIF oder PSD abgespeichert werden. Ein wichtiger Vorteil von XMP Daten ist die „Erweiterbarkeit der Daten. So können zusätzlich zu den vorgegebenen Informationen neue hinzugefügt werden, wenn sie für einen bestimmten Arbeitsablauf notwendig sind.“¹⁶

8 Fazit

Aus der Arbeit heraus ist deutlich erkennbar, dass Metadaten überall in der Technologie vorhanden sind. Wie sollten Konsumenten Suchmaschinen ohne Metadaten nutzen? Wie sollten Unternehmen ohne Metadaten Marketing machen? Diese und viele andere Fragen würden auftauchen, wenn keine Metadaten vorhanden wären. Metadaten setzen sich überall durch. Das alles zeigt, dass Metadaten, ohne dass wir sie wirklich sehen bzw. erkennen, uns alles im Web, in Softwares, in Unternehmen, Museen, Bibliotheken und in vielen anderen Bereichen existieren und uns vieles um Einiges erleichtern.

¹⁵ http://www.optimal-foto.de/documents/Fototipp_0802.pdf

¹⁶ <http://blog.panodapter.com/was-ist-exif-iptc-xmp-metadate.html>

Seminar Big Data Applications

Persistent Identifier Systems: Wie verwende ich PIDs?

Carsten Griesheimer

Sommersemester 2013

Dr. Rainer Stotzka, Danah Tonne

Fakultät für Informatik

Karlsruher Institut für Technologie

Inhaltsverzeichnis

1	Motivation	3
2	Grundlage: Persistent Identifier (PID)	3
3	Handle System	3
3.1	Auflösung	4
3.2	Metadaten	5
3.3	Suche	5
3.4	Erstellung von PIDs	5
4	Digital Object Identifier (DOI)	6
4.1	Format	6
5	Anwendungen	6
6	Fazit	7
6.1	Vergleich	7
6.2	Anwendung für Big Data	8

1 Motivation

Das Internet hat ein Problem. Eine Grundfunktion ist das Verlinken von Webseiten. Da sich URLs von Webseiten u.a. wegen Softwareupdates, Restrukturierung oder gar Löschung ständig ändern, laufen immer wieder Verweise ins Leere. Das ist für einen privaten Internetnutzer ärgerlich, aber für den professionellen und wissenschaftlichen Gebrauch ein ernst zu nehmendes Problem. Bereits 1998 versuchte das World Wide Web Consortium (W3C) mit der Kampagne „Cool URIs don't change“ Webseiten-Betreiber für dieses Problem zu sensibilisieren [14]. Besonders im Zuge der Open-Access-Bewegung und dem steigenden Bedarf Video- und Tonaufnahmen und Datensätze aus Experimenten online bereitzustellen und zu referenzieren gewinnt dieses Problem an Brisanz [15]. Um dieses Problem zu lösen gibt es einige Lösungsansätze. Einer davon ist das Handle System.

2 Grundlage: Persisten Identifiere (PID)

Eine PID ist eine eindeutige ID mit der digitale Objekte identifiziert werden. Sie hat bestenfalls eine unbegrenzte Lebensdauer und ändert sich nicht.

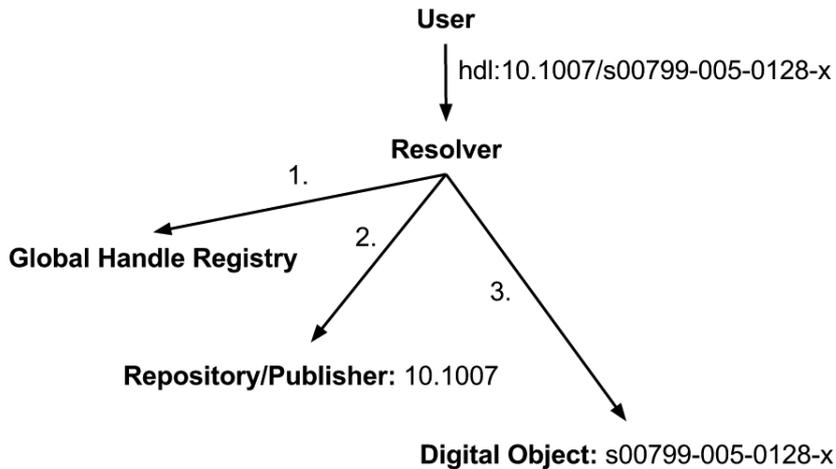
Ein digitales Objekt kann jeder Art von Datei sein, die man auf einem Computer speichern und abrufen kann. Häufig sind es Bilder, Filme oder Texte.

3 Handle System

Das Handle System stellt eine Spezifikation dar, welche die Erstellung, Verwaltung und Auflösung von Persisten Identifiers (PID) für digitale Objekte beschreibt. Eine PID bzw. ein Handle ist eine eindeutige ID, zu der Metadaten gespeichert werden. Die Metadaten enthalten Informationen über das referenzierte digitale Objekt und wie auf dieses zugegriffen werden kann. Ein Anwendungsfall ist das Auflösen einer PID zu einer URL. Durch die Trennung von Bezeichner und Speicherort und der langfristigen und eindeutigen Speicherung der PIDs wird ein Auflösen garantiert. Ändert sich zum Beispiel der Speicherort einer Datei, müssen lediglich die Metadaten der PID angepasst werden. Die PID bleibt unverändert gültig. Das System findet Anwendung in Bereichen in denen Daten veröffentlicht werden und eine langfristige Referenzierbarkeit garantiert werden muss. Das betrifft u.a. Verlage, Universitäten und Bibliotheken.

3.1 Auflösung

Um eine PID aufzulösen, sendet man eine Anfrage mit der aufzulösenden PID an einen Resolver. Dieser fragt bei der globalen Registrierungsstelle nach, welcher lokale Anbieter zuständig ist (1). Eine Anfrage an diesen liefert den Ort des Archivs (2), welche letztendlich die Metadaten des gesuchten Objekts enthält (3). Im Normalfall enthalten die Metadaten die aktuelle URL, mit welcher man letztendlich an die gesuchten Daten gelangt.



Title	A framework for distributed digital object services
Publisher	Springer-Verlag
URL	http://link.springer.com/article/10.1007/s00799-005-0128-x

Das System wurde mit dem Gedanken entworfen „Single Point of Failures“ zu vermeiden. Verschieden Server übernehmen einzelne Aufgaben der Infrastruktur und können beliebig repliziert werden. Zusätzlich können Caching-Instanzen hinzugefügt werden und Anfragenlast Abzufangen.

3.2 Metadaten

Für jede PID werden Typ/Wert-Paare gespeichert. Beispielsweise wäre der Typ einer URL „0.TYPE/URL“ und der Wert entsprechend dem Verweis. Welche Typen zur Verfügung stehen wird eindeutig vom Handle System definiert. Wird ein Typ gebraucht, der nicht definiert ist, bietet das System einen Typ-Registrierungs-Service an, mit dem neue Typen definiert werden können [8]. Es ist jedoch möglich nicht registrierte Typen zu verwenden, sollte aber vermieden werden. Für viele gebräuchliche Informationen, wie Titel, Verlag, Veröffentlichungsdatum, etc. gibt es bereits Typen. Neben diesen gibt es Typen, die vom System ausschließlich für administrative Zwecke verwendet werden. Jedes Paar enthält zudem Informationen über das Cach-Verhalten, Zugriffsberechtigungen und dem letzten Änderungszeitpunkt.

Mit der freien Definierbarkeit der Metadaten lassen sich eine Vielzahl an Funktionen realisieren. Beispielsweise ist es möglich alternative Speicherorte anzugeben oder auch verschiedene URLs, falls eine Datei in verschiedenen Formaten vorliegt.

Index	Type	Timestamp	Data
1	URL	Thu Oct 05 2006 09:26:15 EDT	http://www.sciencemag.org/cgi/doi/10.1126/science.169.3946.635
2	700050	Wed Oct 25 2006 17:13:02 EDT	200610250852
100	HS ADMIN	Thu Oct 05 2006 09:26:15 EDT	handle=0.na/10.1126; index=200; [delete hdl,read val,modify val,del val,add val,modify admin,del admin,add admin,list]

3.3 Suche

Das Handle System an sich bietet keine Such-Funktionalität. Jedoch bieten einige Anbieter, wie CorssRef ¹ oder DataCite ², Suchmaschine an, die auf die Metadaten des Systems zugreifen.

3.4 Erstellung von PIDs

Systeme die das Handle System implementieren bieten in der Regel eine plattformunabhängige Web-API an, über die angemeldete Mitglieder PIDs verwalten und erstellen können. Bei einigen Dienstleistern muss man einen Beitrag zahlen und nachweislich in Wissenschaft und Forschung tätig sein, um Mitglied zu werden. Verlage, die wissenschaftliche Artikel veröffentlichen, erzeugen meist automatisch für jeden Artikel eine PID, über die der Artikel referenziert werden kann. Es besteht neben der Nutzung eines Dienstleister, aber auch die Möglichkeit ein Handle System auf eigener Infrastruktur zu betreiben, indem man sich an frei verfügbaren Open-Source-Lösungen bedient. Dies bietet sich für Testumgebungen an.

¹<http://www.crossref.org/guestquery>

²<http://search.datacite.org/ui>

4 Digital Object Identifier (DOI)

Das Digital Object Identifier (DOI) System ist eine weit verbreitete Implementierung des Handle Systems. Es verwaltet über 84 Millionen Handels und über 215.000 DOI Präfixe. Das System bearbeitet jährlich 1 Milliarde DOI Abfragen.

4.1 Format

Die PID hat das Format <ID des Anbieters>/<ID des Dokuments>

Beispiel PID: 10.2478/s11533-013-0250-8

Beispiel Quellenangabe: Green, T. (2009). „We Need Publishing Standards for Data-sets and Data Tables“. Research Information. doi:10.1787/603233448430.

DOI.org stellt einen Resolver mit dem man einfach DOIs auflösen kann. Entweder ein Anwender gibt die DOI auf <http://dx.doi.org> in ein Eingabefeld ein oder er ruft die URL mit folgendem Schema auf: <http://dx.doi.org/<doi>>

Beispiel URL: <http://dx.doi.org/10.1787/603233448430>

Der DOI hat dabei den Vorteil menschenlesbar zu sein, d.h., aus einer URL kann leicht der DOI abgelesen werden oder er kann aus einer Quelle leicht abgeschrieben werden.

5 Anwendungen

Drei verbreitete Systeme, die das Handle System einsetzen, um digitale Objekte langfristig referenzierbar zu machen, sind CrossRef, DataCite und EPIC. CrossRef und DataCite setzen dabei DOIs ein.

CrossRef verwaltet vor allem geschriebenen Dokumenten. Bereits heute beteiligen sich bekannte Verlage wie der Springer Verlag [9] oder Bibliotheken wie die Universitätsbibliothek von Harvard [10].

DataCite hingegen verwaltet Forschungsdaten und arbeitet dabei mit Datencentern und Forschungseinrichtungen weltweit zusammen. Der Ansprechpartner für Deutsche Organisationen ist die Technische Informationsbibliothek (TIB), welche ein Mitgründer dieser Plattform ist [11].

European Persistent Identifier Consortium (EPIC) ist neben CrossRef und DataCite eine Alternative für europäische Forschungseinrichtungen. Der Dienst bietet die Möglichkeit PIDs zu erstellen, unabhängig davon ob diese später für Publikationen verwendet werden [12]. Dieses System empfiehlt sich, wenn beispielsweise bei Experimenten

große Datenmengen anfallen und zunächst unbekannt ist, welche Daten langfristig gespeichert werden sollen, aber es dennoch nötig ist, diese zu referenzieren. Später können diese PIDs in andere Systeme wie DataCite überführt werden.

6 Fazit

6.1 Vergleich

Andere Systeme sind nur schlecht mit dem Handle System vergleichbar, da sie meist unterschiedliche Ziele verfolgen und einen anderen Satz an Funktionen aufweisen.

Persistent uniform resource locator (PURL) ist eine URI, die auf einen Resolver verweist. Dieser leitet den Benutzer weiter zum referenzierten digitalen Objekt. Damit ähnelt das System auf den ersten Blick dem Handle System. Das System ist vor allem dadurch beschränkt, dass keine Metadaten gespeichert werden und es stark abhängig von den zugrunde liegenden Techniken und Protokollen ist [13]. Ein System sollte in Hinblick auf Langzeitspeicherung möglichst flexibel und unabhängig sein. Darüber hinaus bekommt ein Benutzer keinerlei Informationen, falls das Objekt gelöscht wurde. Im Handle System würden die Metadaten noch rudimentäre Aussagen über das Objekt machen.

OpenURL ist ebenfalls eine URI und referenziert ein Object. Aber im Gegensatz zu PURLs oder dem Handle System liegt der Fokus darauf zu einem Objekt verschiedene Speicherorte anzubieten. Man könnte mit diesem System auf ein Buch verweisen und die URI würde auf den Bestand einer beliebigen Bücherrei verweisen, die dieses Buch führt.

Uniform Resource Names (URN) sind PIDs, aber können nicht direkt aufgelöst werden. Sie identifizieren ein Objekt eindeutig, sagen aber nichts über den Speicherort aus. Es werden keine Metadaten mit einer URN verknüpft. URNs werden aber von anderen Systemen verwendet. So sind zum Beispiel Handles des Handle Systems eine spezielle Form einer URN.

Archival Resource Key (ARK) fast die drei zuvor genannten Systeme im Allgemeinen im Funktionsumfang zusammen, aber hat einen entscheidenden Nachteil. Es wird keine Persistenz garantiert.

ISBN ist für digitale Objekte im Internet, aufgrund der fehlenden Auflösung zum Speicherort, nicht geeignet.

6.2 Anwendung für Big Data

Das Handle System scheint aufgrund seiner Flexibilität und Robustheit gut geeignet um Daten langfristig online referenzierbar zu machen. Es gibt dem Anbieter von Forschungsdaten die Möglichkeit unabhängig von Änderungen an der Infrastruktur Daten anzubieten. Durch APIs lässt sich die Erstellung von PIDs automatisieren. Es stehen verschiedene Open Source Produkte zur Verfügung mit dem man ein eigenes Handle System aufsetzen kann. Falls aber der Wunsch besteht sich an einem großen, verbreiteten System zu beteiligen muss man sich entscheiden. Für Projekte bei denen viele Daten entstehen aber unklar ist welche langfristig gespeichert werden bietet sich EPIC an. Falls dies bereits geklärt ist und Daten der Öffentlichkeit zugänglich gemacht werden sollen, ist DataCite eine weitere Möglichkeit. Das dieses System DOIs verwendet eignet es sich besonders gut für die Verwendung durch Menschen. Probleme können durch langsame Reaktionszeiten von APIs entstehen, die bei großem Datenaufkommen beträchtliche Verzögerungen verursachen können. Die API des Handle Systems mit der Version zwei versucht diesem Problem entgegenzuwirken, indem es erlaubt wird in einem API aufruf mehrere Operationen auszuführen. Welches System man wählt, jedes ist besser als URL für Referenzen zu benutzen.

Literatur

- [1] R. Kahn und R. Wilensky, „A framework for distributed digital object services“, *Int J Digit Libr*, Bd. 6, Nr. 2, S. 115–123, Apr. 2006. doi:10.1007/s00799-005-0128-x
- [2] Handle System, Wikipedia, The Free Encyclopedia, http://en.wikipedia.org/w/index.php?title=Handle_System&oldid=543779077 (11.06.2013)
- [3] Digital object identifier, Wikipedia, The Free Encyclopedia, http://en.wikipedia.org/w/index.php?title=Digital_object_identifier&oldid=558582109 (11.06.2013)
- [4] Digital Object Identifier System Factsheets, <http://www.doi.org/factsheets> (10:21, 11.06.2013)
- [5] European Persistent Identifier Consortium, <http://www.pidconsortium.eu> (10:21, 11.06.2013)
- [6] High-Availability, Complementary Infrastructures for Persistent & Unique Identifiers for Data Objects & Published Collections based on Handle System, Ulrich Schwardmann, März 2012, http://www.doi.org/topics/EPIC_DataCite_March2012.pdf
- [7] Encyclopedia of Library and Information Sciences, Third Edition DOI: 10.1081/E-ELIS3- 120044418
- [8] Sam Sun, Sean Reilly, Larry Lannom: Handle System Namespace and Service Definition. RFC 3651, November 2003, <http://www.ietf.org/rfc/rfc3651.txt>, Abschnitt 3.1 Handle Value Set
- [9] Publishers and Societies, Crossref.org, <http://www.crossref.org/01company/06publishers.html> (01.07.2013)
- [10] Libraries, Crossref.org, <http://www.crossref.org/01company/07libraries.html> (01.07.2013)
- [11] German National Library of Science and Technology, DataCite, <http://www.datacite.org/TIB> (01.07.2013)
- [12] Goal, EPIC, <http://www.pidconsortium.eu/index.php?page=goal> (01.07.2013)
- [13] Factsheet - DOI® System and Persistent URLs (PURLs), DOI.org, http://www.doi.org/factsheets/DOI_PURL.html (01.07.2013)
- [14] Cool URIs don't change, Tim Berners-Lee, 1998, <http://www.w3.org/Provider/Style/URI> (01.07.2013)
- [15] Open Access, Wikipedia, Die freie Enzyklopädie. http://de.wikipedia.org/w/index.php?title=Open_Access&oldid=118132869 (01.07.2013)

Seminar Big Data Applications

Ist das Open Archival Information System (OAIS) für Big Data geeignet?

Eingereicht von

Nico Kopp

Matrikelnummer:

1629376

Sommersemester 2013

Inhaltsverzeichnis

1. Motivation	3
2. Das Open Archival Information System (OAIS).....	4
2.1 Die Entwicklung von OAIS.....	4
2.2 Die Umgebung des OAIS.....	5
2.3 Die Funktionsweise von OAIS	6
2.4 Der Einsatz von OAIS-Archiven	9
3. Ist OAIS für Big Data geeignet?.....	11
4. Fazit	12
5. Literatur	13

1. Motivation

Jährlich werden bis zu 30 Petabyte Speicher für neue Forschungsdaten benötigt, die vom Großen Hadronen-Speicherring (engl.: Large Hadron Collider, LHC) in CERN generiert werden [ORA]. Von diesen Daten müssen „[...] mehr als 70 Petabyte¹ an archivierten Altdaten aus der Forschung für eine schnelle Abfrage und Analyse für die kommenden Jahrzehnte verfügbar [...]“ [ORA] gehalten werden.

Auch in der Raumfahrt sind seit den 60er Jahren jede Menge an digitalen Daten angefallen.

Nun ergeben sich einige größere Probleme:

Wo sollen diese Daten gespeichert werden? Die Daten auf einem Speichermedium zu speichern und dieses Medium in einem geeigneten Raum zu verstauen und bei Bedarf wieder herauszuholen, funktioniert nicht, da, wie oben beschrieben, eine schnelle Abfrage auf Altdaten möglich sein soll. Also muss es irgendeine Art eines digitalen Archivs geben, in das die Daten hineingelegt und bei Bedarf schnell wieder herausgenommen werden können. Aber wie muss eine so große Menge an Daten archiviert werden, damit sie für Jahrzehnte verwendbar ist? Wie wird sichergestellt, dass die Daten in 10 bis 20 Jahren noch lesbar sind und das Datenformat noch bekannt ist?

Um diese Probleme zu lösen werde ich mich mit dem Open Archival Information System (OAIS) genauer beschäftigen und überprüfen, ob dieses Archivierungsmodell für eine Langzeitarchivierung von Forschungsdaten geeignet ist.

¹ 1 Petabyte entspricht 1.000 Terrabyte.

2. Das Open Archival Information System (OAIS)

2.1 Die Entwicklung von OAIS

In diesem Kapitel wird die Geschichte von OAIS erläutert. Außerdem wird gezeigt, dass es auch noch in der heutigen Zeit an Aufmerksamkeit verdient.

„Das als ISO 14721 verabschiedete Referenzmodell „Open Archival Information System – OAIS“ beschreibt ein digitales Langzeitarchiv als eine Organisation, in dem Menschen und Systeme mit der Aufgabenstellung zusammenwirken, digitale Informationen dauerhaft über einen langen Zeitraum zu erhalten und einer definierten Nutzerschaft verfügbar zu machen.“ [NES][Seite 4]. OAIS wurde 1997 von dem „Consultative Committee for Space Data Systems“ (CCSDS), einer Kooperation verschiedener Luft- und Raumfahrtorganisationen unter der Führung der NASA, entwickelt. An der Entwicklung waren außerdem die amerikanische nationale Archivverwaltung (engl.: National Archives and Records Administration, NARA) und die Research Libraris Group (RLG) beteiligt. Wie in der Motivation beschrieben, ist in dem Bereich der Raumfahrt eine große Menge an Daten angefallen, die in geeigneter Form archiviert werden muss.

Im Jahre 1999 wurde dann die erste Fassung des OAIS publiziert, 2001 unter dem ISO/DIS 14721 als internationaler Standard angenommen und schließlich im Jahr 2002 dem Normenwerk der Internationalen Standardorganisation hinzugefügt [NES][Seite 5]. Es ist erstaunlich, dass ein Modell so schnell zu einem internationaler Standard wird und noch erstaunlicher, dass es in so kurzer Zeit zu einer Norm für Langzeitarchive wird. Dies spricht dafür, dass die Entwicklung von OAIS gut durchdacht war.

2009 wurde eine überarbeitete Version des OAIS von der CCSDS veröffentlicht. Dies zeigt, dass dieses Modell nicht veraltet ist, sondern immer noch den Bedürfnissen eines aktuellen Archivs entspricht. Selbst zehn Jahre nach der ersten Publikation wird immer noch an diesem Referenzmodell gearbeitet.

2.2 Die Umgebung des OAIS

Im vorausgegangenen Kapitel wurde die Entstehung von OAIS dargelegt, nun wird in diesem Abschnitt erklärt, welche Teilnehmer in OAIS benötigt werden.

Im vorherigen Kapitel wurde im Zusammenhang mit OAIS von einem Referenzmodell gesprochen.

Ein Referenzmodell ist ein allgemeines Modell für eine allgemeine Problemstellung (hier in unserem Beispiel: ein Modell für die Langzeitarchivierung), mit dessen Hilfe eine konkrete Konstruktion eines Lösungsansatzes zu diesem Problem entwickelt werden kann. Das heißt OAIS ist keine konkrete Softwareimplementierung eines Langzeitarchivierungssystems für die nur noch die Hardware bereitgestellt werden muss, sondern OAIS beschreibt, welche Voraussetzungen Software und Hardware für ein erfolgreiches Langzeitarchiv erfüllen müssen.

Als Erstes stellt sich die Frage nach der Umgebung des Archivs: Wer also kommuniziert mit dem Archiv und auf welche Weise tut er dies, damit ein reibungsloser Zugriff, eine problemlose Wartung von Archivdaten sowie eine reibungslose Einspeisung von neuen Daten möglich ist. Dafür wurde im OAIS-Modell zwischen drei Benutzergruppen unterschieden (siehe Abbildung 1).

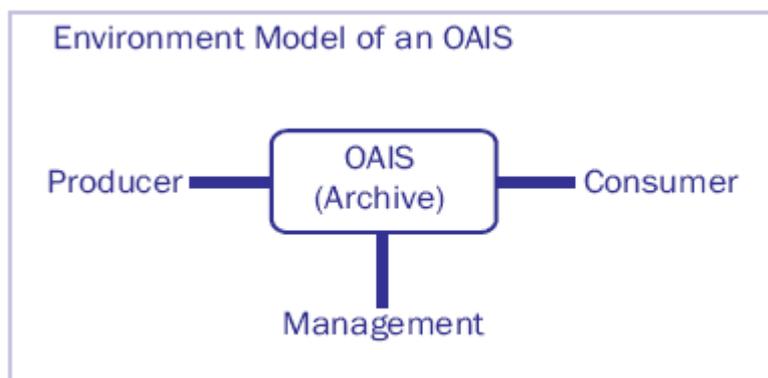


Figure 2: Environment Model of an OAIS Abbildung 1 [OAIS1]

Der erste Teilnehmer ist der Produzent (engl.: Producer), also derjenige, der Daten in das Archiv für die Langzeitarchivierung einreicht. Dies kann sowohl ein Benutzer direkt, als auch ein Klient sein, der die Kommunikation mit dem Archiv für den Benutzer übernimmt.

Der zweite Teilnehmer ist der Konsument (engl.: Consumer), also derjenige, der Lesezugriff auf die im Archiv langzeitarchivierten Daten anfordert. Dies kann, wie beim Producer, sowohl ein Benutzer direkt, als auch ein Klient sein. Ein Teilnehmer kann Producer und Consumer gleichzeitig sein. Ein Benutzer, der Daten einreicht, kann also/folglich Daten auch wieder einsehen.

Der letzte Teilnehmer ist das Management. Dieser ist für die Regelungen innerhalb des OAIS-Archivs zuständig. Einerseits stellt das Management die Trennung der Vorgänge, die automatisierbar sind, von denen, die von Menschen durchgeführt werden müssen, sicher. Andererseits ist es auch für die Qualität des Archivinhalts, also für die Lesbarkeit, Nutzbarkeit und Verständlichkeit der Daten, zuständig. Kurzum: Für all das, was Producer und Consumer für ein einfaches und sicheres Benutzen des Archivs benötigen.

2.3 Die Funktionsweise von OAIS

Im vorherigen Abschnitt wurde dargelegt, welche Teilnehmer für ein OAIS-Archiv benötigt werden. In diesem Abschnitt wird auf die genaue Funktionsweise eines OAIS-Archivs eingegangen.

Die Grundidee hinter OAIS ist, dass die Daten, die vom Producer in das Archiv eingereicht wurden, in einem gegen Fehler unanfälligen Format umgewandelt und archiviert werden. Sobald ein Consumer die Datei einsehen möchte, wird diese wieder in ein für den Menschen lesbares Format konvertiert. Das Format der Datei, die der Consumer vom Archiv erhält, muss allerdings nicht dasselbe Format haben, wie die ursprüngliche Datei des Producers. Der Grund hierfür wird später genauer erläutert.

Es wurde gerade schon angedeutet, dass in OAIS zwischen unterschiedlichen Informationspaketen unterschieden wird, genauer gesagt gibt es drei unterschiedliche Informationspakete:

Das erste Informationspaket ist das Submission Information Package (SIP), also das Paket aus Informationen, welches dem Archiv zur Langzeitaufbewahrung übergeben wird.

Das zweite Informationspaket ist das Archival Information Package (AIP), also das Informationspaket, welches das Archiv für die Speicherung der Daten verwendet. Das

Format eines AIP ist ein einheitlicher Archivstandard, der gegen Fehler (z.B. das Kippen eines Bits) unanfällig sein soll. Ein AIP ist also das aus einem oder mehreren SIPs generierte Informationspaket, welches in das Format des Archivstandards überführt wurde. Wie ein solcher Archivstandard auszusehen hat, ist allerdings nicht Thema dieser Arbeit.

Schließlich das letzte Informationspaket ist das Dissemination Information Package (DIP). Ein DIP ist eine aus einem AIP generierte Datei für den Consumer. Wie aus einem AIP ein DIP generiert wird, wird später genauer erläutert.

Nun da die wichtigsten Begriffe im OAIS-Referenzmodell erklärt wurden, kann die genaue Funktionsweise erläutert werden.

Die Funktionalität eines OAIS-Archivs kann in fünf Einheiten (engl.: Entities) unterteilt werden (vgl. blaue Kästchen in Abbildung 2):

1. Die Datenübernahme (engl.: Ingest)
2. Der Archivspeicher (engl.: Archival Storage)
3. Das Datenmanagement (engl.: Data Management)
4. Der Zugriff (engl.: Access)
5. Die Administration

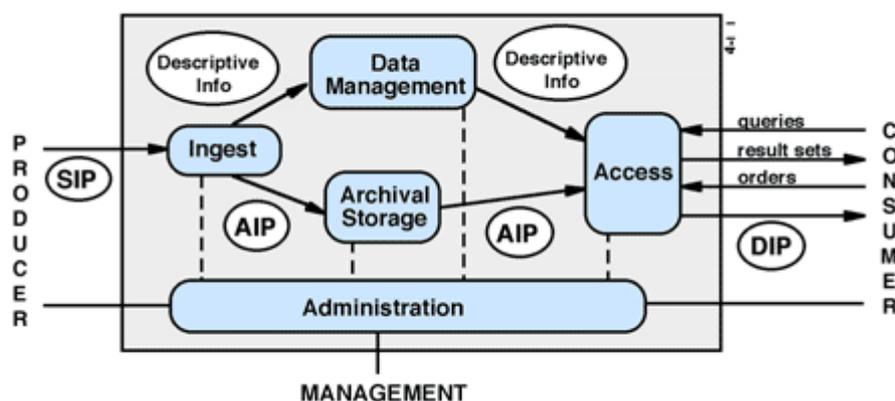


Abbildung 2 [OAIS2]

Jeder dieser fünf Entities hat für sich seinen eigenen Aufgabenbereich. So ist der Ingest dafür zuständig die SIPs des Producers entgegenzunehmen, diese auf Lesbarkeit, Fehler und Kontext zu prüfen, aus diesen dann die AIPs und beschreibende Informationen (engl.:

Descriptive Infos) zu generieren und schließlich die Descriptive Infos an das Data Management und die AIPs an den Archival Storage zu senden.

Eine Descriptive Info beinhaltet unter anderem Informationen darüber, welches Verfahren von der Ingest angewendet wurde, um das SIP in ein AIP umzuwandeln und welche Methode notwendig ist, um aus einem AIP ein DIP zu generieren. Außerdem enthält eine Descriptive Info Metadaten, also Daten, die Daten beschreiben. Diese werden benötigt, damit ein Consumer hinterher, wenn er eine Datei sucht, diese dann auch durch passende Eingaben findet.

Der Archival Storage ist, wie der Name schon sagt, dafür zuständig die AIPs so zu speichern, dass ein hohes Maß an Datensicherheit und Wiederauffindbarkeit gewährleistet wird. Welche Verfahren (z.B. Redundante Datenspeicherung) dafür am besten geeignet sind, würde den Umfang dieser Arbeit sprengen. Diese Fragestellung gehört in den Forschungsbereich der Bitstream Preservation, welcher in einer anderen Seminararbeit behandelt wurde.

In der dritten Entity, dem Data Management geht es um das Verwalten und kontinuierliche Erweitern/ Verbessern der Descriptive Infos, sowie um das Verwalten der Datenbanken, in welche die Descriptive Infos eingelagert werden. Damit soll garantiert werden, dass die Daten lesbar, verständlich und nutzbar bleiben. Sollte ein Datenformat eines DIPs irgendwann einmal veraltet sein, wird in dieser Entity sichergestellt, dass die Methode, die ein AIP in das veraltete Format umwandelt, durch eine neue Methode ersetzt wird. Die neue Methode generiert aus einem AIP ein DIP mit einem neueren, lesbareren Format.

Der Aufgabenbereich des Accesses besteht darin, Anfragen der Consumer entgegenzunehmen und diese dann auszuwerten. Das heißt, dass der Access Anfragen nach Daten entgegen nimmt und die Descriptive Infos und die AIPs der angeforderten Daten sucht. Danach wird mittels der Informationen der Descriptive Infos DIPs aus den AIPs generiert und schließlich dem Consumer die DIPs übergeben. Allerdings hat nicht jeder Consumer Zugriff auf jede Datei. In den Descriptive Infos befinden sich auch Information darüber, welche Benutzer(gruppen) vollen bzw. eingeschränkten Zugriff auf die Zieldaten besitzen. Dadurch erhält ein Consumer, der nur eingeschränkte Zugriffsrechte besitzt, nur so

viel von der Zieldatei, wie es ihm gestattet ist. Ein Consumer, der keine Zugriffsrechte hat, erhält dementsprechend auch kein DIP.

Die letzte Entity ist die Administration. Sie ist für all das im Archiv zuständig, was unter anderem nicht mit dem tagtäglichen Datenaustausch im Archiv zu tun hat. Die Administration stellt zum Beispiel Lösungsansätze bereit, wenn eine Datenübergabe an den Ingest misslungen ist und regelt auch die Nutzungsrechte innerhalb des Archivs. Außerdem ist die Administration quasi die „Wartung“ des Archivs, da sie dem Data Management und dem Archival Storage mitteilt, in welchen Abständen ein Datenbankupdate beziehungsweise ein Refreshing² der Datenträger durchgeführt werden soll.

Zusammenfassend lässt sich die Funktionsweise von einem OAIS-Archiv so beschreiben, dass eine Datei für die Langzeitarchivierung in ein sicheres Format konvergiert und gespeichert wird. Die Informationen zur Wiederherstellung in ein lesbares Format werden ebenfalls gespeichert, wobei die gespeicherten Daten in regelmäßigen Abständen gewartet und aktualisiert werden. Sollte ein Consumer schließlich eine Datei im Archiv benötigen, wird mittels der gespeicherten Informationen und der Archivdatei das gewünschte Informationspaket für den Consumer erstellt.

2.4 Der Einsatz von OAIS-Archiven

Im obigen Kapitel wurde genauer erläutert, wie ein Archiv nach dem OAIS-Modell zu funktionieren hat. In diesem Kapitel werden einige Einsatzgebiete von OAIS-Langzeitarchiven genannt.

Das erste digitale Archiv, das nach dem OAIS-Modell aufgebaut wurde, ist die Niederländische Nationalbibliothek in Den Haag. Im Jahr 2002 wurde mit Hilfe von IBM der erste Prototyp eines digitalen Langzeitarchivs nach dem OAIS-Standard entwickelt, der Publikationen zugänglich halten soll [NES][Seite 14]. Auch das australische Nationalarchiv orientiert sich im Rahmen des Pandora-Projektes an OAIS [HKI]. Dieses Archiv ist genauso wie das in Den Haag für Publikationen zuständig.

² Refreshing ist „die Überprüfung der verwendeten Datenträger auf ihre Lesbarkeit und die Verständlichkeit“ [NES][Seite 12].

Es ist zu erkennen, dass im Bereich der Archivierung von Publikationen OAIS häufig verwendet wird. Diese Tatsache ist auch einfach nachzuvollziehen, da Publikationen meist ein einheitliches Format (PDF) besitzen. Dadurch werden nur wenige unterschiedliche Funktionen benötigt, um die SIPs in AIPs und die AIPs in DIPs umzuwandeln. Außerdem sind Publikationen nicht sonderlich groß (einige MB). Das Konvertieren einer so kleinen Datei ist dementsprechend nicht so sonderlich zeitintensiv.

Ein anderes Einsatzgebiet für OAIS-Archive sind Regierungsakten. Das Britische und das Amerikanische Nationalarchiv benutzen beide Archive, die auf OAIS basieren [NES][Seite 14]. Die Gründe, warum sich OAIS dafür anbietet, sind die gleichen wie bei den Publikationen. OAIS garantiert eine relativ hohe Datensicherheit und da die zu archivierenden Daten nicht so groß und nicht so verschieden voneinander sind, ist der benötigte Rechen- bzw. Verwaltungsaufwand für das Konvertieren der Dateien nicht sonderlich groß.

Es ist also zu erkennen, dass OAIS gerade in den Bereichen beliebt ist, wo es gilt, viele kleine und ähnliche Dateien zu archivieren.

3. Ist OAIS für Big Data geeignet?

Im letzten Kapitel wurde gezeigt, dass OAIS häufig für kleine Dateien verwendet wird. Aber wie sieht es nun mit Big Data, genauer gesagt mit großen Forschungsdaten aus?

Erst einmal Big Data sind sehr große Daten, die mehrere Terrabyte oder sogar Petabyte groß sein können. Forschungsdaten haben kein einheitliches Format, es gibt unzählig viele verschiedene Formate von Forschungsdaten. So erzeugt der LHC in CERN Daten in einem anderen Format, als es beispielsweise das Pierre-Auger-Observatorium tut, das die kosmische Strahlung misst.

Ein OAIS-Langzeitarchiv für große Forschungsdaten müsste also für jedes Format Methoden besitzen, um die SIPs in AIPs und die AIPs schließlich wieder in DIPs umzuwandeln. Der Aufwand all diese Methoden zu erstellen und sie dann auch zu verwalten ist im Vergleich zum Nutzen, d.h. Dem Verhindern von Datenverlust, viel zu hoch. Als weiterer Punkt gegen OAIS und Big Data spricht die Tatsache, dass die Dateien in einem anderen Format archiviert werden als die Ursprungsdaten. Wenn also terrabytegroße Daten archiviert werden sollen, müssen diese erst in den Archivstandard umgewandelt werden. Je größer Daten sind, umso größer ist auch der Rechenaufwand für die Konvertierung und dementsprechend auch die Zeit, die benötigt wird. Das gleiche gilt dann ebenfalls für die Generierung der DIPs für die Consumer. Der Consumer müsste Stunden bzw. Tage warten bis er seine Daten bekommt, da diese für ihn erst generiert werden müssten. Ein Archiv bearbeitet ja nicht nur eine Anfrage, sondern gleich viele gleichzeitig. Dadurch würde die Anzahl der Anfragen, die vom Archiv abgearbeitet werden müssten, immer größer werden, was zur Folge hätte, dass der Vorgang, bis eine Anfrage beendet worden wäre, eine lange Zeit in Anspruch nehmen würde.

Zusammenfassend ist also zu sagen, dass meiner Meinung nach die Datensicherheit eines OAIS-Archivs nicht den Verwaltungs- und Rechenaufwand für Big Data kompensiert und OAIS deshalb für Big Data ungeeignet ist.

4. Fazit

Wie in der Motivation schon erläutert, muss eine Lösung für ein Archivierungssystem von Forschungsdaten entwickelt werden. Anfangs sah das OAIS-Referenzmodell als Lösungsansatz dafür gar nicht so schlecht aus. Allerdings stellte sich in der Folge heraus, dass die Stärken dieses Modells, nämlich die hohe Datensicherheit und das Sicherstellen der Lesbarkeit der Daten, Probleme bei großen und unterschiedlichen Daten mit sich ziehen.

Das heißt also, dass jedes Archiv genau für seinen Verwendungszweck geeignet ist. Ein Archiv, das für Publikationen geeignet ist, ist nicht unbedingt dafür geeignet Forschungsdaten zu archivieren.

Es kommt beim Entwickeln eines Archivs also immer darauf an, welche Art von Daten archiviert werden sollen. Somit gibt es kein „ultimatives“ Archivmodell.

5. Literatur

- [NES] Nestor. [Online] 2009. [Abruf am 06.11.2013] http://nestor.sub.uni-goettingen.de/handbuch/artikel/nestor_handbuch_artikel_366.pdf
- [ORA] Oracle. [Online] [Abruf am 11.06.13] <http://www.oracle.com/us/corporate/customers/customersearch/cern-storagetek-snapshot-1705027-de-ch.html>
- [OAIS1] Paradigm. [Online] [Abruf am 01.07.13] <http://www.paradigm.ac.uk/images/environmental-model.gif>
- [OAIS2] NASA. [Online] [Abruf am 01.07.13] http://nssdc.gsfc.nasa.gov/nssdc_news/dec00/oais_fig3.gif
- [HKI] Universität Köln. [Online] [Abruf am 02.07.13] http://www.hki.uni-koeln.de/sites/all/files/courses/4478/Referat_OAIS_Pr%C3%A4sentation.pdf

Karlsruher Institut für Technologie
Steinbuch Centre for Computing (SCC)
Institute for Data Processing and Electronics (IPE)
Prof. Dr. Achim Streit
Dr. Rainer Stotzka



eSciDoc – Geeignet für Big Data?

Seminararbeit

Abgabetermin:

15. Juli 2013

Vorgelegt von: Gül Kavak
Matrikelnummer: 1606460
Studiengang: Informationswirtschaft
Email: ufdkr@student.kit.edu
Betreuer: Prof. Dr. Achim Streit
Dr. Rainer Stotzka

Inhaltsverzeichnis

Einleitung.....	3
Was ist e-Science?	3
FIZ Karlsruhe und Max-Planck-Gesellschaft	3
FIZ Karlsruhe und Max-Planck-Gesellschaft	4
Was ist eSciDoc?.....	5
Das eSciDoc Schichtenmodell.....	6
eSciDoc-Solutions, Applications und Service.....	8
Ziele von eSciDoc.....	9
Vorteile und Nachteile von eSciDoc	10
Fazit	11
Fazit	12
Quellenverzeichnis	13

1. Einleitung

Die Sicherstellung von effizienter und nachhaltiger Organisation des Daten- und Wissensflusses stellt eine wichtige Anforderung wissenschaftlichen Arbeitens dar, in der heutigen in fast jeder Hinsicht globalisierten Welt. Um diese Anforderungen zu erfüllen wird heute an einer Informations-, Kommunikations-, und Publikationsplattform geforscht und gearbeitet. In der sog. e-Science („Enhanced Science“) wird eine innovative Infrastruktur hergestellt, die die Massendaten technisch gekonnt dokumentiert, per Internet veröffentlicht, ständig aktualisiert und Dienstleistungen zur Verfügung stellt. Diese neue Plattform wird eSciDoc genannt. Eine neue Form von Wissenschaftskommunikation und Zusammenarbeit von Forschern wird, durch die Internetbasis und der gigantischen Speichermöglichkeit weltweit ermöglicht.

2. Was ist e-Science?

Die e-Science zielt auf eine in weltweiter Zusammenarbeit entwickelnde Forschung mit Hilfe des Internets.

Die e-Science Infrastruktur beinhaltet alle nötigen Forschungsressourcen für einen Forschungsbereich und ermöglicht somit nicht nur die Kenntnisnahme auch die Verarbeitung und Entwicklung dieser Analysedaten und Ergebnisse, durch die Publikationen. So kommt es zu neuen Ideen, die die Forschung fortführt. Die e-Science vertritt damit nicht nur den technologischen Gesichtspunkt, sondern auch den sozialen und wissenschaftspolitischen Gesichtspunkt.

3. FIZ Karlsruhe und Max-Planck- Gesellschaft

Die eSciDoc Plattform wird mit einer Partnerschaft von der Max-Planck-Gesellschaft und des Fachinformationszentrum (FIZ) Karlsruhe umgesetzt und vom Bundesministerium für Bildung und Forschung gefördert bzw. sieben Jahre lang, zwischen 2004 und 2009, finanziert.



<http://www.fiz->

[karlsruhe.de/home.html?&L=skdiyeusug](http://www.fiz-karlsruhe.de/home.html?&L=skdiyeusug)

Das Fachinformationszentrum (FIZ)

Karlsruhe ist eine gemeinnützige GmbH. Diese stellt die größte außeruniversitäre Informati-

onsinfrastruktur-Einrichtung in Deutschland dar. Es bietet eine professionelle Ausübung der Wissenschaft bzw. der Wirtschaft mit Forschungs- und Patentinformation und eine Entwicklung von innovativen Dienstleistungen. Zu den Kunden des FIZ Karlsruhe zählen die internationalen Marktführer aus der Pharma- und Chemieindustrie, große Patentämter und Forschungseinrichtungen. Das FIZ Karlsruhe übernimmt außerdem die Entwicklung von innovativen e-Science-Solution und Services, welche die wichtigsten Anwendungen bei der eSciDoc ermöglicht.



MAX-PLANCK-GESELLSCHAFT

http://imprs.evolbio.mpg.de/images/Logo_s.gif

Die Max-Planck-Gesellschaft (MPG) ist Deutschlands erfolgreichste Forschungsorganisation mit den weltweit besten Forschungsinstitutionen. Sie vertritt 17 Nobelpreisträger in den Reihen ihrer Wissenschaftler. Die MPG hat jedes Jahr mehr als 15.000 Publikationen in international renommierten Fachzeitschriften. Das Jahresbudget beträgt ca. 1,4 Milliarden Euro pro Jahr. Derzeit betreiben 82 Max-Planck-Institute Grundlagenforschung in den Natur-, Bio-, und Sozialwissenschaften im Dienste der Allgemeinheit.

Welche Bedeutung hat eSciDoc für die Max-Planck-Gesellschaft?

eSciDoc hat eine strategische Bedeutung für die MPG. Mit dem Projekt ist eine langfristige interne Finanzierung gesichert. Durch die effektive Zusammenarbeit mit FIZ Karlsruhe, wird das ursprüngliche Ziel des Projektes, nämlich Kommunikation und Zusammenarbeit, verwirklicht. Die Publizierung bietet der MPG größere Entwickler Basis und damit mehr und schnellere Lösungen. Die MPG gewinnt mit diesem Projekt auch einen größeren Kundenstamm und gewinnt somit der für sich an Stabilität.

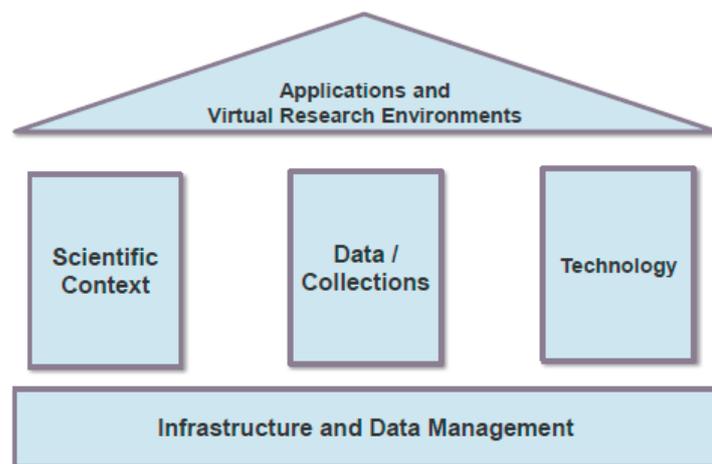
Die Max-Planck-Gesellschaft (MPG) und das FIZ Karlsruhe wollen die Sicherstellung vom ständigen Zugang zu den Forschungsergebnissen und Forschungsmaterialien des Max-Planck-Gesellschaft und die nahtlose Integration in eSciDoc, sowie die Integration in einer aufstrebenden, globalen wissenschaftlichen Wissens-Raum. Sie wollen die Bereitstellung von effektiven Möglichkeiten des Zugangs zu Informationen für die Wissenschaftler der Max-

Planck-Gesellschaft und ihrer Arbeitsgruppen. Außerdem wollen sie die Unterstützung der wissenschaftlichen Zusammenarbeit in Zukunft e-Science-Szenarien fördern.

Das Projekt soll nicht nur an bestimmten wissenschaftlichen Disziplinen genutzt werden, sondern es soll einen viel breiteren Anwendungsbereich haben. Das Ziel ist, eine generische Infrastruktur und eine Lösung für alle wissenschaftlichen Fachbereich der Max-Planck-Gesellschaft (Naturwissenschaften, Lebenswissenschaften, Sozial- und Geisteswissenschaften) zu liefern.

4. Was ist eSciDoc?

ESciDoc ist ein System, welches in Forschungseinrichtungen, Hochschulen, Instituten und Unternehmen genutzt wird. Es bietet eine generische Infrastruktur zur Speicherung und Verwaltung von Forschungsprimärdaten, Analysedaten und Publikationen. Die noch in Entwicklung steckende Plattform ermöglicht die Generierung, die Verarbeitung, die Verbreitung und die



<https://www.escidoc.org/pdf/escidoc-days-2011/day1-dreyer-overview.pdf>

Archivierung von wissenschaftlichen Kenntnissen. Jedoch nicht nur Endergebnisse des Forschungsprozesses, sondern alle Zwischenschritte: Primärdaten, experimentellen Daten etc. können dieser Verwaltung unterliegen. In der Zukunft soll dieses System für die Kommunikation der Forscher aus der ganzen Welt vorliegen. Alle im Verlauf dieses Prozesses anfallenden Daten (Dokumente, Primärdaten wie z. B. Spektren, Textkorpora, Mikroskop-Aufnahmen) werden mit Hilfe der eSciDoc-Infrastruktur nachhaltig gespeichert. Durch die weitgefächerte Zugriffsmöglichkeit der Forscher, wird die Nutzung von Forschungsergebnissen und Informationen optimiert.

Das eSciDoc-Schichtenmodell

Das eSciDoc-System besteht aus drei Teilbereichen:

1. *eSciDoc Core*: Die unterste Ebene bietet große Anzahl an allgemeine Funktionalitäten vor allem zur Speicherung und zur Versionierung der Daten zur Verfügung.



<http://www.uni-tuebingen.de/einrichtungen/zentrum-fuer-datenverarbeitung/projekte/bw-escit/projektuebersicht/escidoc.html>

2. *eSciDoc Services*: In der mittleren Schicht sind die Basis Dienstleistungen wie das Hinzufügen und Aktualisieren von Forschungsressourcen, das Verbinden individueller Forschungsressourcen mit dauerhaftem Identifizieren, das Suchen von Metadaten oder auch allgemeine Such- und Indexierungsmöglichkeiten werden angeboten.
3. *eSciDoc Solutions*: Die oberste Schicht bietet auf den Kundenbedarf zugeschnittene, angefertigte Anwendungen (Lösungen), die für individuelle Forschungsgruppen bestimmt sind.

Grundsätzlich fällt die eSciDoc-Software-Stack in drei Kategorien:

Nämlich in die Kategorien Anwendungen (Applikations), Dienstleistungen (Services) und die Infrastruktur. Anwendungen sind die Endanwender-Applikationen für Forscher und Bibliothekare (z.B. PubMan, Virr).

Die Dienstleistungen stellen Aggregationsfunktionen (oder Zusammenfassungsfunktionen) von den Dienststellen der eSciDoc-Infrastruktur zur Verfügung und bieten eine höhere Ebene und mehr anwendungsorientierte Funktionalität. eSciDoc bietet Dienstleistungen zur Spei-

cherung von Objekten, Suche und Indizierung, Statistiken, Berichte, anhaltende Identifizierung, Workflows, Aussortierung von wichtigem und unwichtigem und Transformation. Die Infrastruktur liefert die grundlegenden und allgemeinen Funktionen und ist eine Voraussetzung für alle Anwendungen und die meisten eSciDoc –Dienstleistungen. Die zur Verfügung gestellten Anwendungen können auch ohne die Zwischenschicht, Dienstleistungen, also mit einer Umgehung, aufgerufen werden.

Das eSciDoc System ist Service orientiert aufgebaut. Die Konstruktion ist eine sich an die Leistungsfähigkeit des Hardware anpassendes, wiederverwendbares und erweiterbares Service System. Diese Service orientierte Konstruktion fördert die mehrmalige Nutzung der vorhandenen Dienstleistungen. Ein eSciDoc Service kann von verschiedenen Institutionen und Einrichtungen sowohl lokal als auch ferngesteuert verwendet werden, so kommt es zu einer größeren e-Science-Infrastruktur. Diese Dienstleistungen können sowohl gekoppelt als auch getrennt verwendet werden. Die eSciDoc-Plattform sollte nicht als ein einzelnes Paket zum Download verstanden werden, sondern als eine Reihe von Diensten und Anwendungen (Applications), aus denen Sie die, die Sie interessiert, auswählen und implementieren. Damit die eSciDoc-Infrastruktur bestmöglich zusammenarbeiten kann, muss die Möglichkeit bestehen andere bestehende Systeme und Datenquellen zu integrieren. Dazu wären mehrere verschiedene Metadaten-Profile erforderlich. Neben den grundlegenden und weit verbreiteten Metadaten-Formaten, wie Dublin Core oder MODS, sind viel mehr Community-spezifische Metadaten-Formate nötig, die für bestimmte Anwendungen genutzt werden können. Diese müssen aber von der Infrastruktur unterstützt werden. Die entwickelten Dienste müssen eine hohe Flexibilität in dieser Hinsicht bieten, da es nicht Möglich ist zu wissen, ob diese Formate bei allen potenziellen eSciDoc Lösungen unterstützt werden.

Mit diesem agnostischen Ansatz, wurde und wird dieses Projekt entwickelt und umgesetzt. Das Projekt hängt ab von der Weiterentwicklung des Fedora und vom Repository-System, welches im Einsatz für grundlegendes Objekt-Management ist. Die Kerntechnologie basiert auf Java und Extensible Markup Language, abgekürzt XML. Anstatt den Aufbau der Infrastruktur "von Grund auf neu" anzugehen, entschied sich das eSciDoc Team zu vorhandenen Open-Source-Komponenten, da die Plattform von „Allen“ genutzt werden soll, soll die Allgemeinheit so weit wie möglich schon integriert sein. eSciDoc Dienstleistungen im Allgemeinen bieten sowohl SOAP- (ursprünglich für *Simple Object Access Protocol*) als auch

Representational State Transfer-, Akronym REST, Stil an. Dies ermöglicht die weitere Entwicklung von Lösungen, ohne dabei die Auswahl oder den Tausch von Programmiersprachen. So werden deren Umsetzung und damit die Beteiligung der verschiedenen Entwickler-Gruppen beschleunigt und erhöht. Deshalb ist eines der Grundprinzipien des eSciDoc-Projekts: Das Wiederverwenden und zwar bestehende Konzepte, Dienstleistungen und Implementierungen. Beispiele für Software, die in die eSciDoc-Infrastruktur integriert wurden:

- Der Repository-Architektur Fedora Commons zur Speicherung von Objekten
- SRW / U für die Suche
- Fedora Generische Suche Service für die Indizierung
- JHOVE für technische Metadaten-Extraktion
- Sun XACML Motor für politische Auswertung
- Shibboleth der Internet2-Projekt für verteilte Authentifizierung

eSciDoc wird nicht nur von der Technik und der Informatik angetrieben. Das Team vereint funktionale Anforderungen. Ingenieure, beschreiben die Szenarien und die erforderlichen Anwendungsfälle, sowie technologische Experten und Entwickler wirken stark mit ein. Alle Gruppen innerhalb des Teams sind fest in einem gut festgesetzten Kommunikations- und Abstimmungsprozess gekoppelt und müssen gut zusammenarbeiten, um ein Kompromiss für die verschiedenen Ansichten einzugehen und allgemein akzeptierte Ansätze, also Bedürfnisse der Akteure innerhalb der wissenschaftlichen Gemeinschaft, in Erfüllung zu bringen. Die Infrastruktur der E-Science-Umgebung benötigt Basistechnologien der Berechnung, Kommunikation, Software-Programme und ein angemessenes Datensicherungs-System, um große Daten zu verwalten. Für das Überstehen dieser Herausforderungen zur Realisierung einer generischen Infrastruktur, die zur gleichen Zeit in der Lage sein soll, fachspezifische Anwendungen auszuführen, wählte das Projekt eine Architektur, welches zu Vergleichen ist mit einem Baukastensystem, die die Kombination von Modulen erlaubt. Für die Entwicklung von Lösungen und fachspezifischen Anwendungen und Arbeitswelten, werden Dienstleistungen zusammengefügt, um die spezifischen Bedürfnisse zu erfüllen. Die ersten Anwendungen, die derzeit umgesetzt werden, decken die Bereiche der Veröffentlichung und des wissenschaftlichen Datenmanagements ab. Für die Arbeit mit diesen Datensammlungen muss eine flexible Versionierung möglich sein und sie müssen Workflow-Modelle haben, bevor sie weiter umgesetzt werden können, damit auch Dienste für anspruchsvollere Daten-

verarbeitungen vom Server unterstützt werden können. Da diese Daten stark spezifisch angelegt sind, werden sie auch in allgemeiner Weise dargelegt, um so standardisierte Schnittstellen zu bieten. Damit werden extern integrierbare Dienstleistungen zur Verfügung gestellt. In dieser Methode wird die Middleware genutzt. Definitionen von Middleware, die die Sache klarer darstellen sollen:

Middleware (englisch für Diensteschicht oder Zwischenanwendung) oder Vermittlungssoftware bezeichnet in der Informatik anwendungsneutrale Programme, die so zwischen Anwendungen vermitteln, dass die Komplexität dieser Applikationen und ihre Infrastruktur verborgen werden.¹ Oder (Klassische) Middleware ist eine Menge von wenig spezialisierten ("general-purpose"-) Diensten, die zwischen der Systemplattform (Hardware+Betriebssystem) und den Anwendungen angesiedelt sind und deren Verteilung unterstützen.²

So werden Dienste zur Verfügung gestellt, die die Verteilung von Anwendungen unterstützen und die plattformübergreifend sind. Außerdem werden Standard-Schnittstellen für Anwendungen und Benutzern mit einheitlichen Benutzeroberflächen zur Verfügung gestellt. Damit lassen sich vorhandene Systeme mit geeigneter Middleware ohne große Umorganisation an andere Systeme ankoppeln.

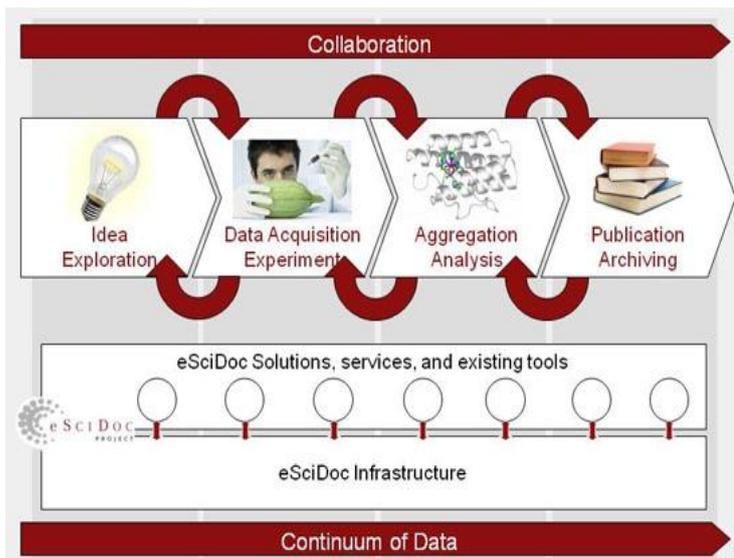
5. Ziele von eSciDoc

Die Ziele von eSciDoc sind der Aufbau einer e-Science-Infrastruktur, die Entwicklung einer integrierten Publikations- und Kommunikationsplattform, die Forschungsergebnisse und Materialien im e-Research Netzwerk integriert. So haben die Nutzer von eSciDoc einen umfassenden Zugriff auf Daten und Informationen. Damit wird die Zusammenarbeit und interdisziplinäre Forschung in zukünftigen e-Science-Szenarien stark unterstützt. Die Entwickler dieses Projekt möchten eine bessere Sichtbarkeit von Forschungseinrichtungen und Organisationen erreichen.

¹ Nach: W. Ruh u. a.: Enterprise Application Integration. Wiley, 2001

² Nach: <http://avs-www.informatik.uni-hamburg.de/teaching/ss-96/vortrag1/def.htm>

Die Zusammenarbeit funktioniert folgendermaßen: Wie bei herkömmlichen wissenschaftlichen Arbeiten, werden Hypothesen aufgestellt und Versuche gemacht, die diese Hypothesen stützen sollen. Die anfallenden Daten bei der Auswertung dieser Versuche werden protokolliert. All die Informationen, Protokolle,



Beobachtungen, Fragen, Hypothesen, Simulationen etc., die bei diesem Prozess entstehen, werden veröffentlicht. Also auf die eSciDoc- Plattform hochgeladen und dabei nachhaltig gespeichert und archiviert. Wiederum werden die archivierten Daten und Protokolle von Interessenten runtergeladen, analysiert und neue Hypothesen und Fragen werden aufgestellt. Dazu werden Experimente durchgeführt, eventuell die vorherigen Fragen beantwortet und neue Fragen gestellt. Diese werden ebenfalls hochgeladen und publiziert, sodass weitere Interessenten darauf agieren können und weitere Verknüpfungsstellen entstehen, die das Thema fortführen soll. Die hochgeladenen Daten, werden über die Dienstleistungen und Anwendungen weitergeführt, sodass sich die Interessenten auch nur den Teil beobachten oder analysieren, der sich fachspezifisch auf sie bezieht. Eine riesige Datenmenge sammelt sich dabei an, worauf die Forscher/Innen Zugriff haben können. Auf diese Weise können global neue Erkenntnisse gewonnen werden und zwar mit Hilfe dieser riesigen Infrastruktur.

6. Vorteile und Nachteile von eSciDoc

Dieses Projekt bringt wie alle Projekte sowohl Vorteile als auch Nachteile mit sich, jedoch welche Gesichtspunkte Ihnen schwerer fallen, muss jeder für sich selbst entscheiden.

Zunächst bietet das eSciDoc-Projekt eine riesige Speichermöglichkeit. Außerdem wird die Organisation der Daten z.B. für die Forschungscoordination erleichtert. Mit der Veröffentlichung der Forschungsdaten bleiben diese erhalten vor allem mit der richtigen Archivierung gehen diese Daten nicht im World Wide Web verloren, sondern sie können schnell und effi-

zient wiederverwendet werden. Somit werden stabile Identifizierungen vergeben. Diese können jederzeit wieder aktiviert werden für die Wiederverwendung von anderen. Eine Zusammenarbeit über institutionelle Grenzen wird auf dieser Weise ermöglicht. Durch die Aggregationsfunktion der Dienstleistungen, wird die Zusammenfassung ermöglicht und die Ordnung der Projekte in diesem Verzeichnis. Da die ESciDoc-Software kein festgelegtes Gesamtpaket ist, können die Interessenten auch nur ihren fachspezifischen Teil bzw. Bereich nutzen und auf ihrer Hardware installieren. Eine einfache Integration in bestehende Werkzeuge, Instrumente und Systeme stellt ebenfalls ein Vorteil dar. Als weiterer Vorteil kann die Modularität genannt werden. Die Anwendungen können wie in einem Baukastensystem kombiniert werden, somit wird Ungewolltes vermieden. Die Service orientierte Architektur fördert Erweiterbarkeit, Kapselung von Konzepten und die grundlegende Durchsuchung und Kontrolle der Komponenten. Sie unterstützt den vollständigen Objekt-Lebenszyklus und die Objekte behalten ihre eigenen Kontext Informationen, welches gut geeignet für die Langzeitarchivierung ist. Die Services verbinden und verbreiten Daten und Objekte. Die Infrastruktur sorgt für die Nachhaltigkeit von Daten und den (Kern-) Dienstleistungen.

Als Nachteil dieses Projektes kann genannt werden, dass die Installation sehr schwierig ist und viel Zeit in Anspruch nimmt. Außerdem wird viel Speicherplatz benötigt, welches auch von der möglichen Speicherkapazität des Projektes bedingt ist. Als weiteren Punkt kann die anspruchsvolle Versionierung in Betracht genommen werden. So kann doch nicht „jeder“ die eSciDoc-Software nutzen, sondern sie verlangt fachkundliche Kompetenz.

7. Fazit

Der Generaldirektor des Research Council (UK) hat diese Tatsache wie folgt zusammengefasst:

"e-Science ist über die globale Zusammenarbeit in wichtigen Bereichen der Wissenschaft und der nächste Generation der Infrastruktur, die es ermöglichen wird."

E-Science bedeutet Zusammenarbeit und Kommunikation während des gesamten Prozesses der Wissenserweiterung und Verbreitung.

E-Science verlangt nach innovativen Methoden, Services und IT-Infrastrukturen, die effektiv Forscher bei ihren täglichen Arbeitsprozessen unterstützt. Unmengen von experimentellen Primärdaten sind heute in wissenschaftlichen Repositorien, also in den alten Aktenschränken. Die aktuelle Forschung ist mehr und mehr gezwungen den direkten Zugriff auf Primärdaten zu ermöglichen. Allerdings sollte e-Science in einem weiteren Sinne verstanden werden. Direkter Zugang zum primären Daten, Publikationen und andere wissenschaftliche Materialien nicht unbedingt den Einsatz von großen Rechen- und Speicher-Kapazitäten.

Die Wissensgrundlage enthält alle Daten, Informationen, Wissen und Kompetenz, dass eine Organisation nutzt. Zu den wichtigen Schritten hin zu einer funktionalen wissenschaftlichen Wissen- Raum gehören die Standardisierung von Datenformaten, die Entwicklung von allgemein anerkanntes System von Informationen mit logischen Relationen und das Management von Uneinheitlichkeit in Fällen, in denen die Standardisierung nicht integrierbar ist.

Aktuelle Informationen Systemen oft konzentrieren sich auf strukturierte Daten. Allerdings ist ein Großteil des vorhandenen wissenschaftlichen Materials weniger strukturiert. Heutige Systeme sind oft Arten von Materialien beschränkt, die in traditionellen Bibliotheken verblieben sind. Jedoch sind auch andere Daten für die Fortführung von Wissenschaft erforderlich, nämlich: Primäre Daten, Ergebnisse aus Simulationen, informelle Ergebnisse, Erkenntnisse, Anmerkungen und so weiter

Auf die Frage, ob eSciDoc geeignet für große Datenmengen (BigData) ist, würde ich sagen: Ja, mit eSciDoc lassen sich große Daten verwalten, organisieren und nachhaltig speichern. Die Objekte behalten auch die eigenen Kontextinformationen, die Langzeitarchivierung wird bereitgestellt, die Datenqualität wird erhöht. Durch die Zuverlässigkeit, die Genauigkeit der Daten, die Datenqualität und Datenwiederherstellung gibt es eine wachsende Anzahl an Interessenten und Nutzern. Wenn die Vorteile bzw. die Nachteile beobachtet werden, kann ebenfalls erkannt werden, dass die Vorteile ziemlich stark überragen.

Quellenverzeichnis

E-Science 5. April 2013

- <http://de.wikipedia.org/wiki/E-Science> (15.Juli 2013)

e-SciDoc 9. März 2012

- <https://www.escidoc.org/> (15.Juli 2013)

FIZ Karlsruhe 10. Juli 2013

- <http://www.fiz-karlsruhe.de/index.php?id=15> (15.Juli 2013)

Max-Planck-Gesellschaft 2013

- <http://www.mpg.de/> (15.Juli 2013)

Middleware

- <http://avs-www.informatik.uni-hamburg.de/teaching/ss-96/vortrag1/def.htm> (15.Juli 2013)

Modularität 26. Mai 2013

- <http://de.wikipedia.org/wiki/Modularit%C3%A4t> (15.Juli 2013)

eSciDoc in der MPG

- http://www.agmb.de/papooopro/dokumente/upload/b1b64_Tschida_.pdf (15.Juli 2013)

eSciDoc 29.07.2010

- <http://www.uni-tuebingen.de/einrichtungen/zentrum-fuer-datenverarbeitung/projekte/bw-escit/projektuebersicht/escidoc.html> (15.Juli 2013)