

Fünf Thesen zum Thema Nachhaltigkeit

Die Sicherstellung der Verfügbarkeit von (Text-)Daten als Aufgabe von Langzeitarchivierung
Erfahrungsbericht aus einem nationalen Forschungsdateninfrastrukturprojekt

Johannes Stigler
Zentrum für Informationsmodellierung
Karl-Franzens-Universität Graz

Data Center for the Humanities der Universität zu Köln (DCH)
Cologne Center for eHumanities (CCeH)
in Kooperation mit der Landesinitiative NFDI der Digitalen Hochschule NRW
und der Nordrhein-Westfälischen Akademie der Wissenschaften und der Künste

17. September 2018



These I

Der Status quo zum Thema Nachhaltigkeit resultiert aus den Prämissen der Forschungsförderung

Infrastrukturprojekte im Bereich der Digitalen Geisteswissenschaften werden (in Österreich) nicht gefördert. Das Rad wird ständig neu erfunden. Dort wo sie gefördert werden, sind sie zu weit weg von den Bedarfen spezifischer Domänen^a.

^aData Management Pläne sind in Hinblick auf Nachhaltigkeit m.E. nur eine völlig unbefriedigende Zwischenlösung



Versuch von Begriffsklärungen

„Nachhaltigkeit“

- Aus der Perspektive von Nutzer*innen
Die Publikation von Digitalen Editionen ist ebenso einfach wie heute die eines PDF/A-Dokumentes in einem Repositorium und ist nicht an einen spezifischen Projektkontext gebunden
- Aus der Perspektive von Systemadministrator*innen
Es gibt (a) eine Modellierungssprache für die Funktionalitäten von Digitalen Editionen und (b) eine standardisierte API.

„Kuratierung“

- Algorithmusgesteuerte Bestandserhaltung in einem sich in einem ständigen Wandel befindlichen IT-Ökosystem
- Erweiterung bestehender Projekte durch neue Funktionalitäten und neue Formen der Datenrepräsentation (Tools zur Visualisierung oder Analyse von (Text-)Daten, u.a.)



These II

Digitale Langzeitarchive sind Publikationsinstanzen für Digitalen Editionen

Projektkuratierung erfolgt DANN nicht; muß also eingebettet sein in Prozesse der nachhaltigen Bestandserhaltung. Der Betrieb eines Langzeitarchivs erfordert eine technische und organisatorische Infrastruktur, sowie fachliches Knowhow rund um die Domäne aus der die zu verwaltenden Daten stammen^a.

^aDanke für die ausführliche Klärung diese Tatbestandes im Grundsatzpapier der Arbeitsgruppe Datenzentren im Verband DHd



These III

Kuratierbarkeit ist ein notwendiges Strukturmerkmal nachhaltiger Repräsentationsform Digitaler Editionen

Die Möglichkeit zur Kuratierung ist konzeptionell bei der Entstehung von Daten und deren Repräsentationsweisen im Design von Editionsframeworks zu berücksichtigen.

These IV

Objektorientierung ist ein notwendiges Strukturmerkmal nachhaltiger Repräsentationsform Digitaler Editionen

Ein Framework, das geeignet ist, (Text-)Daten aus der Domäne der Digitalen Geisteswissenschaften aufzunehmen, ist kein System zur Verwaltung von Dateien, sondern erfordert eine objektorientierte Organisation seiner Inhalte gemäß dem Prinzip der Trennung von Inhalt und (Re-)Präsentation.

These V

Modellierung ist ein notwendiges Strukturmerkmal nachhaltiger Repräsentationsform Digitaler Editionen

Zur Sicherstellung von Nachhaltigkeit bedarf es standardisierter daten- (und nicht technologie)zentrierter Lösungsansätze in den Digitalen Geisteswissenschaften. Auch Repräsentationsschichten Digitaler Editionen sollten letztendlich modelliert und nicht programmiert werden.

Gemeinsam sind wir stark

KONDE – Kompetenznetzwerk „Digitale Edition“

Projektziele

- die Erarbeitung eines inhaltlichen und strategischen Konzepts zur Bündelung der Kompetenzen und den Aufbau einer nationalen digitalen Infrastruktur für Editionsprojekte
- die Definition von Workflows und Standards für die systematische und institutionenübergreifende Digitalisierung und Bereitstellung von Quellenmaterial
- die Erarbeitung einer Best-Practice Lösung für eine eng in LZA-Repositorien eingebundene Publikationsplattform für DE
- die Entwicklung von Konzepten für eine Übernahme des im Projekt entwickelten Referenzmodells für DE durch einschlägige Gedächtnisinstitutionen
- die (Weiter-)Entwicklung von Werkzeugen für die Verarbeitung des digitalen Materials mittels (semi-)automatischer Verarbeitungsschritte
- den Aufbau und die Etablierung einer einschlägigen Fortbildungsinfrastruktur zur Vermittlung der im Projekt erarbeiteten Standards

KONDE – Kompetenznetzwerk „Digitale Edition“

Zehn Partner

- Universität Klagenfurt (AAU)
- Adalbert-Stifter Institut des Landes Oberösterreich (ASI)
- Universität Innsbruck (IBK)
- Universität Graz (KFU)
- Kunstuniversität Graz (KUG)
- Österreichische Akademie der Wissenschaften (ÖAW)
- Österreichische Nationalbibliothek (ÖNB)
- Universität Salzburg (SBG)
- Technische Universität, Graz (TUG)
- Universität Wien (VIE)



KONDE – Kompetenznetzwerk „Digitale Edition“

Zehn Arbeitsgruppen

- Quelledigitalisierung und -beschreibung
- Transkription und Textauszeichnung
- Korpuslinguistische Analyse und Textmining
- Webpräsentationsformen und -interfaces
- Alternative Formen der Textcodierung
- Netzwerkanalyse und Datamining
- Archivierung
- Textkritik und -kommentar
- Strategiewarbeitsgruppe
- Hybrid-Edition



KONDE – Kompetenznetzwerk „Digitale Edition“

Strategien zur Umsetzung der Projektziele

- Realisierung von digitalen Editionsprojekten
- Arbeitsgruppen zu Schwerpunktthemen
- Workshops
- Weißbuch
- Öffentlichkeitsarbeit



Objektorientierung: Ein Beispiel

- Jedes Objekt im digitalen Archiv hat einen komplexen Typ: Handschrift, museales Artefakt, Ontologie etc.)
- Heißt, jedem Objekt ist eine Klassendefinition zugeordnet ...
 - mit einer vordefinierten (aber grundsätzlich erweiterbaren) Menge an Datenströmen
 - einer vordefinierten Menge an Workflows zur Auslieferung der Daten
 - einer vordefinierten Menge an Workflows, die beim Ingest oder der Änderung der Daten ausgelöst werden
 - und einem Regelwerk für die Erstellung von neuen Objekten (des jeweiligen Typs)



Objektorientierung: Ein Beispiel

Eine Handschrift aus der digitalen Edition der Basler Jahrrechnungen^{ab}

^a<http://gams.uni-graz.at/osrbas.1535>

^b<http://gams.uni-graz.at/context:srbas>

The screenshot shows a web interface for the digital edition of the 'Jahrrechnungen der Stadt Basel 1535 bis 1610'. It features a navigation menu on the left with categories like 'Einnahmen' and 'Empfangen'. The main content area displays a list of entries, including 'Jahrrechnung a festo Johannis Baptistae anno xv'xxxv usque ad festum Johannis Baptistae anno xv'xxxvi'. A search bar and a 'Suche' button are visible at the top right.

Objektorientierung: Ein Beispiel

Das Transkript der Handschrift ist gemäß den Richtlinien der TEI kodiert

Beim Ingest des TEI-Dokumentes in das DLZA ...

- werden auch alle Faksimiles hochgeladen
- wird das TEI-Dokument mit Metadaten aus externen Quellen (geonames.org, GND, SKOS Ontologien, etc.) angereichert
- werden semantische Aussagen (RDF) aus dem TEI-Dokument extrahiert, im Objekt gespeichert und in eine semantische Datenbank übernommen
- werden deskriptive Metadaten (DCMI, EDM u.a.) aus dem TEI-Dokument extrahiert und im Objekt gespeichert
- wird – basierend auf den Strukturinformationen im TEI-Dokument – ein IIIF-Manifest erzeugt und im Objekt gespeichert

On-the-fly, wenn ein Objektinhalt angefragt wird ...

- werden HTML, PDF, \LaTeX PDF Repräsentationen des Objektes erzeugt und ausgeliefert
- werden Faksimiles im DFG- oder Mirador-Viewer angezeigt
- wird eine Datenmatrix aus dem Textkorpus nach einem Regelwerk erstellt und bereitgestellt



Objektorientierung: Ein Beispiel

Das TEI Objekt in der Anzeige des Mirador-Viewers^a

^a<http://gams.uni-graz.at/osrbas.1535/sdef:IIIF/getMirador>

The screenshot shows the Mirador viewer displaying a digital edition of a manuscript page. The viewer shows two views of the page: a thumbnail on the left and a larger view on the right. The manuscript text is clearly visible in the larger view.

Objektorientierung: Ein Beispiel

Die RDF-Daten des TEI Objektes als EXCEL Sheet^a

^a<http://gams.uni-graz.at/osrbas.1535/sdef:HSSF/get>

The screenshot shows an Excel spreadsheet with columns for 'Kategorie', 'Typ', 'Originaltext', 'Setting in Phrasing', 'Errechneter Betrag', and 'Komponente'. The data is organized in rows, showing various categories and their corresponding values.

Kategorie	Typ	Originaltext	Setting in Phrasing	Errechneter Betrag	Komponente
1 Kategorie	sk:entry	Prima angaria pro locis lb		209500	Kategorie
2 Kategorie	sk:entry	Secunda angaria pro locis lb		209500	Kategorie
3 Kategorie	sk:entry	Tertia angaria pro locis lb		180750	Kategorie
4 Kategorie	sk:entry	Quarta angaria pro locis lb		227750	Kategorie
5 Kategorie	sk:entry	Quinta angaria pro locis lb		176250	Kategorie
6 Kategorie	sk:entry	Sexta angaria pro locis lb		213440	Kategorie
7 Kategorie	sk:entry	Septima angaria pro locis lb		279640	Kategorie
8 Kategorie	sk:entry	Octava angaria pro locis lb		227800	Kategorie
9 Kategorie	sk:entry	Nonima angaria pro locis lb		255560	Kategorie
10 Kategorie	sk:entry	Decima angaria pro locis lb		255560	Kategorie
11 Vom mullerkornungelt	sk:total	Suma mullerkornungelt		1019400	1011840 Kategorie
12 Kategorie	sk:entry	Prima angaria pro locis lb		209500	Kategorie
13 Kategorie	sk:entry	Secunda angaria pro locis lb		209500	Kategorie
14 Kategorie	sk:entry	Tertia angaria pro locis lb		180750	Kategorie
15 Kategorie	sk:entry	Quarta angaria pro locis lb		227750	Kategorie
16 Kategorie	sk:entry	Quinta angaria pro locis lb		176250	Kategorie
17 Kategorie	sk:entry	Sexta angaria pro locis lb		213440	Kategorie
18 Kategorie	sk:entry	Septima angaria pro locis lb		279640	Kategorie
19 Kategorie	sk:entry	Octava angaria pro locis lb		227800	Kategorie
20 Kategorie	sk:entry	Nonima angaria pro locis lb		255560	Kategorie
21 Vom mullerkornungelt	sk:total	Suma mullerkornungelt		1019400	1011840 Kategorie
22 Kategorie	sk:entry	Prima angaria pro locis lb		209500	Kategorie
23 Kategorie	sk:entry	Secunda angaria pro locis lb		209500	Kategorie
24 Kategorie	sk:entry	Tertia angaria pro locis lb		180750	Kategorie
25 Kategorie	sk:entry	Quarta angaria pro locis lb		227750	Kategorie
26 Kategorie	sk:entry	Quinta angaria pro locis lb		176250	Kategorie
27 Kategorie	sk:entry	Sexta angaria pro locis lb		213440	Kategorie
28 Kategorie	sk:entry	Septima angaria pro locis lb		279640	Kategorie
29 Kategorie	sk:entry	Octava angaria pro locis lb		227800	Kategorie
30 Kategorie	sk:entry	Nonima angaria pro locis lb		255560	Kategorie
31 Vom mullerkornungelt	sk:total	Suma mullerkornungelt		1019400	1011840 Kategorie
32 Kategorie	sk:entry	Prima angaria pro locis lb		209500	Kategorie
33 Kategorie	sk:entry	Secunda angaria pro locis lb		209500	Kategorie
34 Kategorie	sk:entry	Tertia angaria pro locis lb		180750	Kategorie
35 Kategorie	sk:entry	Quarta angaria pro locis lb		227750	Kategorie
36 Kategorie	sk:entry	Quinta angaria pro locis lb		176250	Kategorie
37 Kategorie	sk:entry	Sexta angaria pro locis lb		213440	Kategorie
38 Kategorie	sk:entry	Septima angaria pro locis lb		279640	Kategorie
39 Kategorie	sk:entry	Octava angaria pro locis lb		227800	Kategorie
40 Kategorie	sk:entry	Nonima angaria pro locis lb		255560	Kategorie
41 Vom mullerkornungelt	sk:total	Suma mullerkornungelt		1019400	1011840 Kategorie
42 Kategorie	sk:entry	Prima angaria pro locis lb		209500	Kategorie
43 Kategorie	sk:entry	Secunda angaria pro locis lb		209500	Kategorie
44 Kategorie	sk:entry	Tertia angaria pro locis lb		180750	Kategorie
45 Kategorie	sk:entry	Quarta angaria pro locis lb		227750	Kategorie
46 Kategorie	sk:entry	Quinta angaria pro locis lb		176250	Kategorie
47 Kategorie	sk:entry	Sexta angaria pro locis lb		213440	Kategorie
48 Kategorie	sk:entry	Septima angaria pro locis lb		279640	Kategorie
49 Kategorie	sk:entry	Octava angaria pro locis lb		227800	Kategorie
50 Kategorie	sk:entry	Nonima angaria pro locis lb		255560	Kategorie
51 Vom mullerkornungelt	sk:total	Suma mullerkornungelt		1019400	1011840 Kategorie

Objektorientierung: Ein Beispiel

Ein auf Basis des gesamten TEI-Korpus generierter Summary Report

ZIM Nachhaltigkeits Digitaler Editionen 17. September 2018 13 / 18

Objektorientierung: Ein Beispiel

... und all das aus dem TEI Quellcode^a

^ahttp://gams.uni-graz.at/orsbas.1535/TEI_SOURCE

```

1
2 <TEI xmlns="http://www.tei-c.org/ns/1.0" xmlns:r="http://gams.uni-graz.at/rem/ns/1.0">
3 <teiHeader xml:id="d2e4">
4 <fileDesc xml:id="d2e6">
5 <titleStm xml:id="d2e8">
6 <title xml:id="d2e10">Jahrechnung Stadt Basel 1535/1536</title>
7 <respStm xml:id="d2e13">
8 <resp xml:id="d2e15">Transkription von</resp>
9 <persName xml:id="editor">Sonia Calvic</persName>
10 </respStm>
11 <respStm xml:id="d2e22">
12 <resp xml:id="d2e24">Korrigiert von</resp>
13 <persName xml:id="d2e27">Lukas Meili</persName>
14 </respStm>
15 <principal xml:id="d2e31">Burghartz, Susanna (Universität Basel)</principal>
16 </titleStm>
17 <publicationStm>
18 <distributor xml:id="d2e37">Zentrum für Informationsmodellierung - Austrian Centre for Digital
19 Humanities, Universität Graz</distributor>
20 <idno type="PID">o:orsbas.1535</idno>
21 </publicationStm>
22 <sourceDesc xml:id="d2e41">
23 <msDesc xml:id="d2e43">
24 <msIdentifier n="http://query.staatsarchiv.bs.ch/query/detail.aspx?ID=1137716" xml:id="d2e45">
25 <repository xml:id="d2e47">StAB8</repository>
26 <collection xml:id="d2e50">Finanz H</collection>
27 <idno xml:id="d2e53">92.1</idno>
28 </msIdentifier>
29 <msContents xml:id="d2e57">
30 <p xml:id="d2e59">
31 <origDate extent="52 Wochen" from="1535-06-26" to="1536-06-24" xml:id="d2e60"></origDate>
32 </p>
33 </msContents>
34 </msDesc>

```

Kuratierbarkeit: Ein Beispiel

GAMS – Geisteswissenschaftliches Asset Management System

- Eine Infrastruktur für unterschiedlichste Projekte (Editionen, Sprachkorpora, Epigraphische Sammlungen, Bildsammlungen u.a) mit je sehr spezifischen Funktionalitäten
- Mengengerüst des derzeitigen Datenpools: 90.000 Objekte mit rund 800.000 Datenströmen und 100.000.000 RDF-Triples
- Das älteste Objekt im Repository ist 14 Jahre alt

Kuratierbarkeit: Ein Beispiel

GAMS goes FEDORA 5.x

- Komplettaustausch der Framework-Komponenten ohne Änderungen an den Einzelprojekten
- Vollständig Docker-basierte Systemarchitektur aller Services
- Implementierung der Content Model Logistik auf Basis von API-X OSGI Modulen
- Die Verlagerung vieler, derzeit clientseitig implementierter Workflows in ein Server-basiertes Doorkeeper-Webservice (fcgate) ermöglicht die Thin Client Entwicklung in beliebigen Programmiersprachen

Kuratierbarkeit: Ein Beispiel

Summe der Erfahrungen aus der Testmigration

- Ausgangsbefund bei Testmigration mit Standardimage von FEDORA 4.7.5: Exponentiell ansteigende Ingest- und Requestzeiten
- Befund am Ende der Testphase: Stabile durchschnittliche Ingestzeit von rund 1.2 sec. pro Objekt und Requestzeiten im zweistelligen Millisekundenbereich.
- Schrauben an denen gedreht wurde: Tomcat-Konfiguration, PostgreSQL anstelle von Google Key-Value-Store, flachere Hierarchiestrukturen im Datenmodell bei gleichzeitiger Einführung von Zwischenschichten



Finis

Danke für Ihre Aufmerksamkeit!

<http://gams.uni-graz.at>
<http://fedora-commons.org>
<http://apache.org>
<http://tei-c.org>
<https://iiif.io>
<http://lido-schema.org>
<http://geonames.org>
 u.a.



Kuratierbarkeit: Ein Beispiel

Systemarchitektur unter FEDORA 5.x

