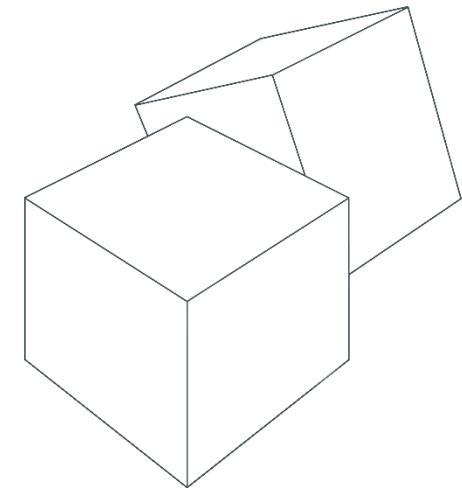


SFB 1288:

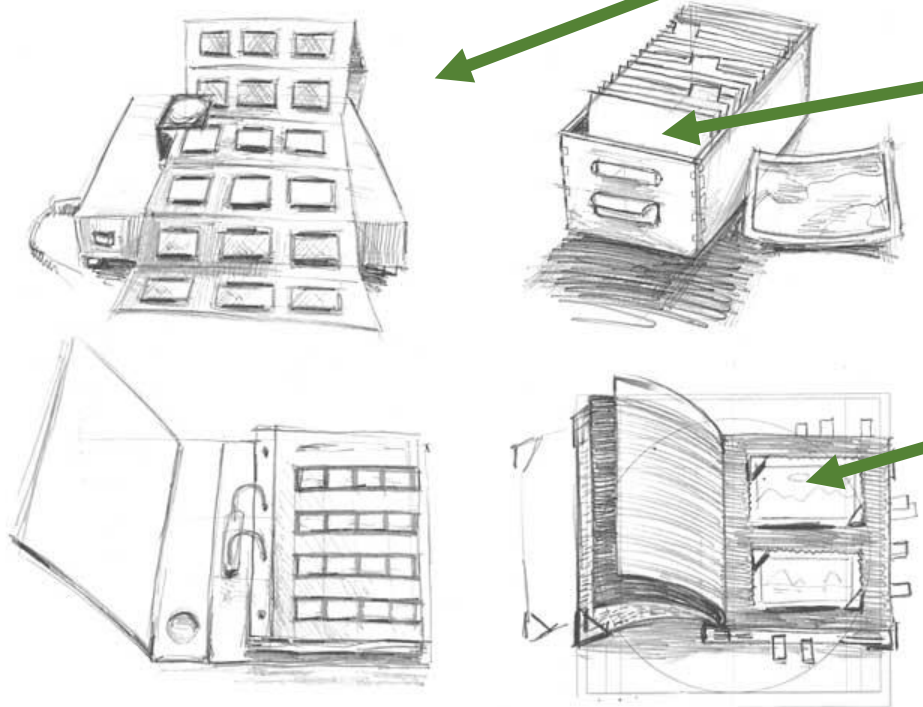
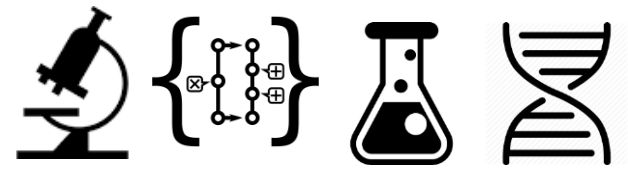
DATENINFRASTRUKTUR & DIGITAL  
HUMANITIES

02. FEBRUAR 2018



# WAS SIND „GEISTESWISSENSCHAFTLICHE“ FORSCHUNGSDATEN?

- Nicht so einfach zu definieren wie in den Naturwissenschaften
- Eher „Quelle“ oder „Forschungsliteratur“
- Charakter des **Vorläufigen** und **Unfertigen**
- Begriffs „Forschungsdaten“ hängt elementar mit der Digitalisierung zusammen, gebunden an:
  - (1) **Prozessierbarkeit**: die Möglichkeit digitalen Arbeitens
  - (2) **Methode, bzw. Transformation**

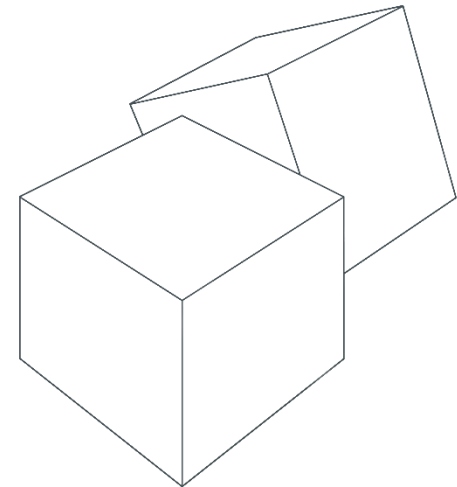


**Bildsuche:**  
Herstellung eines digitalen Faksimiles (*image*)

**linguistische Textanalyse:** Text in strukturierter Form (OCR, Transkription)

die **Analyse der Vernetzung** von Personen:  
Aufbereitung als LOD

Buzetti (2009): **Data is the representation of information in a form that can be processed by a machine.**





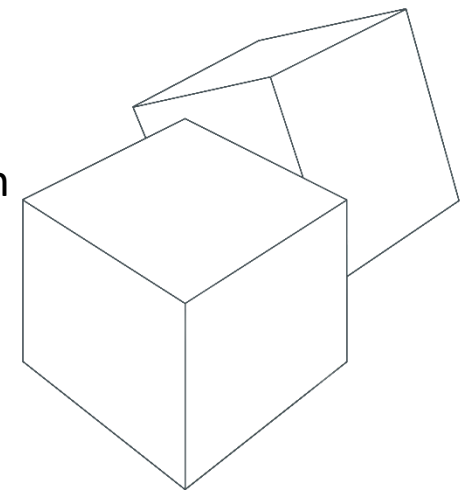
## DH-Aktivitäten (infrastrukturellen Belange)

- Vorgeschaltete Verfahren der Digitalisierung
- Speicherung von Forschungsdaten
- Entwicklung von Publikationsworkflows



## DH-Verfahren („DH als Forschungsmethodik“)

- Ziel: *die Untersuchung der Daten für den Erkenntnisgewinn*
- Finden von Antworten auf Forschungsfragen durch digitale Verfahren



DH is a research discipline that bridges the gap between engineering and humanities. DH is not simply about engineering software for humanities scholars. Instead, DH is where disciplines from engineering and humanities combine their expertise in dealing with information per se. (And obviously, one potential outcome of such a joint research may be tools that can be used in the humanities.)

*Tim Weyrich*

## DIGITAL HUMANITIES

Digital humanities is the synergy of the humanities and digital technologies, resulting in potentially powerful new ways of researching, analysing, synthesising and presenting humanities scholarship.

*Rachel Murphy*

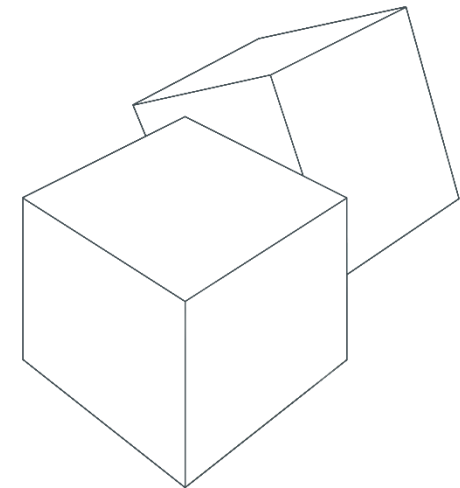
DH is a term that covers all scholarly methods, practices, endeavors, and challenges lying at the intersection of Humanities and Information science and technology.

*Aris Xanthos*

DH are part of the Humanities. Humanists ask research questions. From time to time it makes sense to answer them by using digital means and methods. It's not a decision between digital or analog. Digital Humanities has to match with scope and focus of the research project

*thomas*

Quelle: <http://whatisdigitalhumanities.com/>



DH is a research discipline that bridges the gap between engineering and humanities. DH is not simply about engineering software for humanities scholars. Instead, DH is where disciplines from engineering and humanities combine their expertise in dealing with information per se. (And obviously, one potential outcome of such a joint research may be tools that can be used in the humanities.)

*Tim Weyrich*

## DIGITAL HUMANITIES

Digital humanities is the synergy of the humanities and digital technologies, resulting in potentially powerful new ways of researching, analysing, synthesising and presenting humanities scholarship.

*Rachel Murphy*

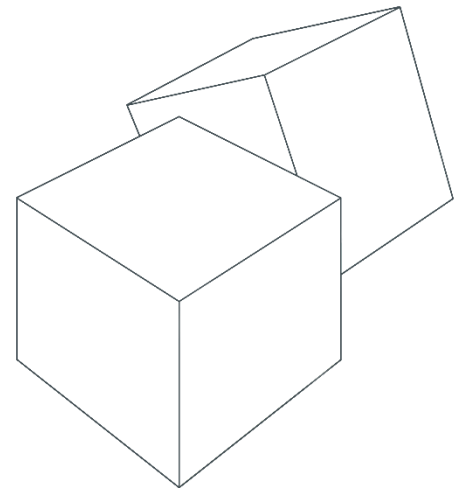
DH is a term that covers all scholarly methods, practices, endeavors, and challenges lying at the intersection of Humanities and Information science and technology.

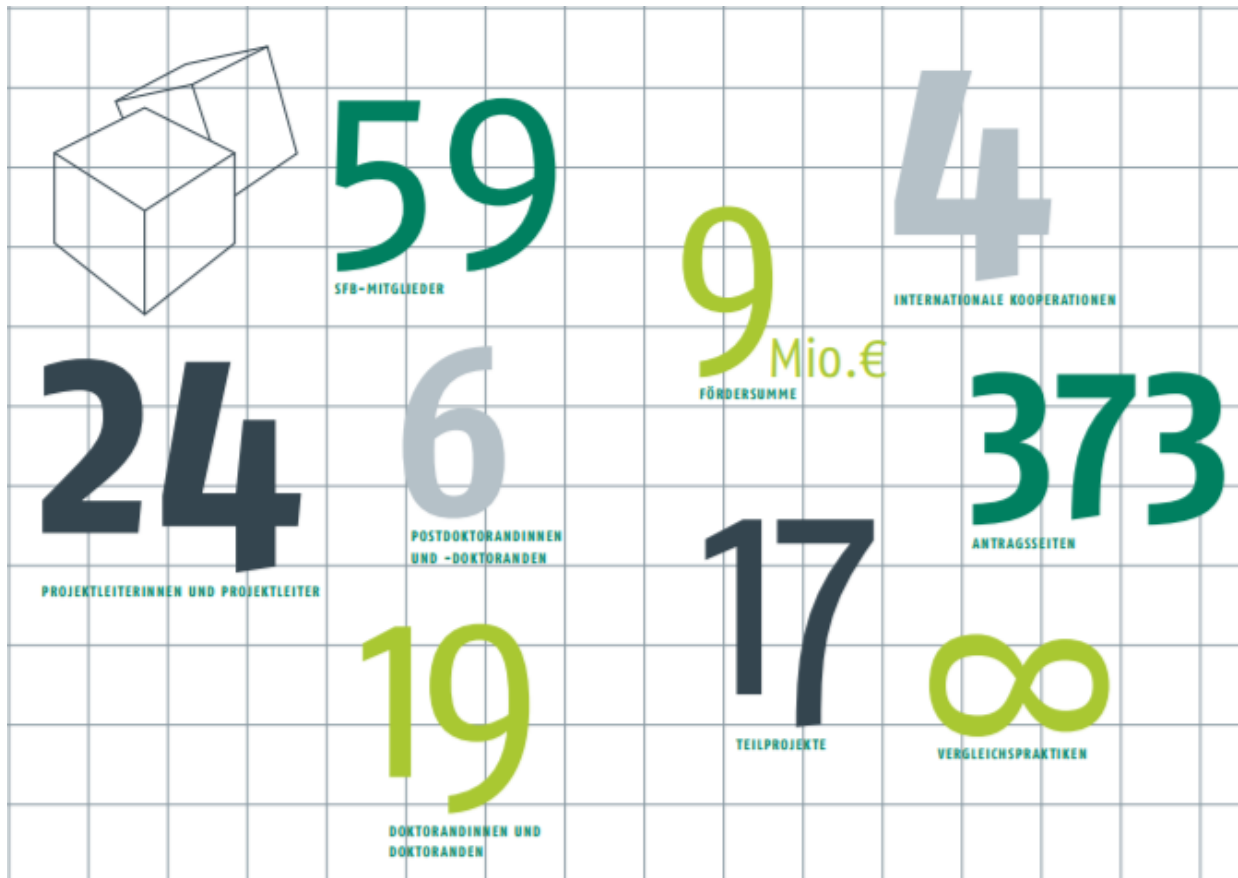
*Aris Xanthos*

DH are part of the Humanities. Humanists ask research questions. From time to time it makes sense to answer them by using digital means and methods. It's not a decision between digital or analog. Digital Humanities has to match with scope and focus of the research project

*thomas*

Quelle: <http://whatisdigitalhumanities.com/>





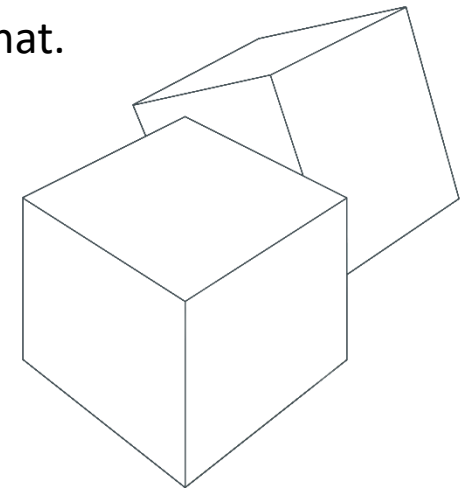
INF Projekt im SFB 1288  
(seit Januar 2017)  
„Praktiken des Vergleichens“

SFB 1288 nähert sich dem Vergleichen  
aus einer historischen Perspektive

**Was tun Akteure, wenn sie vergleichen?**

SFB will alles über die Anfänge herausfinden und wie sich das Vergleichen im Lauf der Geschichte geändert hat.

Der interdisziplinäre, aus **Geschichts-** und **Literaturwissenschaft, Philosophie, Kunstgeschichte, Politik- und Rechtswissenschaft** bestehende Forschungsverbund, fragt wie sich die historisch variablen Praktiken des Vergleichens zu Routinen, Regeln, Habitus, Institutionen und Diskursen fügen und so Strukturen schaffen, aber auch Dynamiken mittlerer Reichweite oder übergreifenden Wandel anstoßen können.



## INF Projekt im SFB 1288:

### „Dateninfrastruktur und Digital Humanities“

#### 1.3 Teilprojekt **Informationsinfrastruktur**

In einem Sonderforschungsbereich dient ein Teilprojekt Informationsinfrastruktur vor allem dem systematischen Management der im Kontext des Sonderforschungsbereichs relevanten Daten. Unter Forschungsdaten werden alle Ergebnisse und Bezugsquellen des Forschungsprozesses (u.a. auch Software oder Forschungsobjekte bzw. Proben) verstanden, die im Projekt erhoben, ausgewertet und/oder entwickelt werden. Ein auf diese Daten bezogener Einsatz sowie die Erprobung oder Entwicklung neuer wissenschaftlicher Kommunikationsformen ist ebenfalls möglich. Es können die Entwicklung



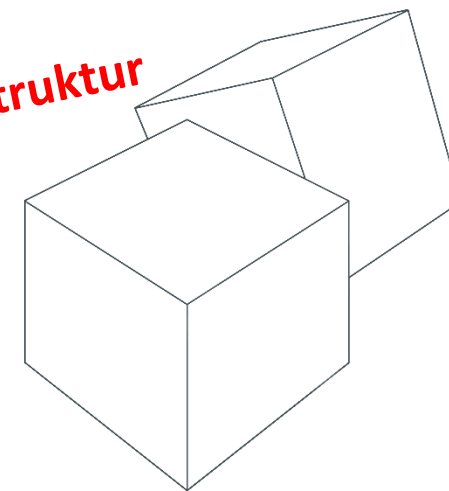
Wissenschaftliche Sy-  
Postanschrift: 53170 Bonn  
228 885-2777 · postmaster@dfg.de · www.dfg.de

Damit soll erreicht werden, dass sich im Sonderforschungsbereich wissenschaftliche Synergien durch gemeinsame Datenplattformen und/oder Kommunikationsforen sowie eine effiziente Datennutzung einstellen.

Grundsätzlich wird ein professionelles Management der Daten, die im Sonderforschungsbereich erhoben, ausgewertet und/oder entwickelt werden, erwartet. In der Regel soll daher mit den einschlägig ausgewiesenen Informationseinrichtungen am Standort zusammengearbeitet werden (z.B. Bibliotheken, Rechenzentren oder Biobanken der antragstellenden Hochschule). Der Nutzung bereits vorhandener Repositorien, Werkzeuge und Techniken ist gegenüber der Entwicklung neuer Instrumente in der Regel der Vorzug einzuräumen, siehe auch:

[www.re3data.org](http://www.re3data.org)

**Pilot für die Erweiterung der  
hochschulweiten  
Forschungsdateninfrastruktur**





## INF Projekt im SFB 1288:

Interdisziplinär besetztes Team, bestehend aus:



**Dr. Silke Schwandt**  
Digital Humanities  
(Fak. Für Geschichte)



**Dr. Johanna Vompras**  
Dateninfrastruktur / FDM  
(Universitätsbibliothek)



**Anna Maria Neubert**  
Digital Humanities



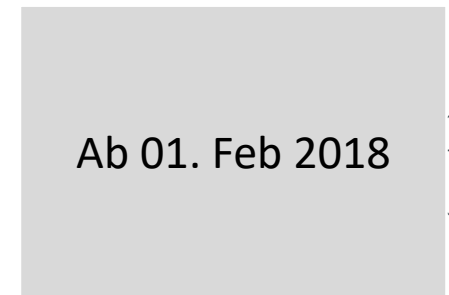
**Dr. Madis Rumming**  
Dateninfrastruktur  
/ FDM



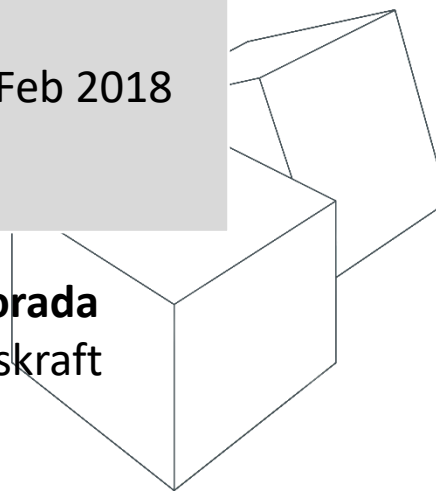
**Helene Schlicht**  
Digital Humanities



**Leonard Gödde**  
wiss. Hilfskraft

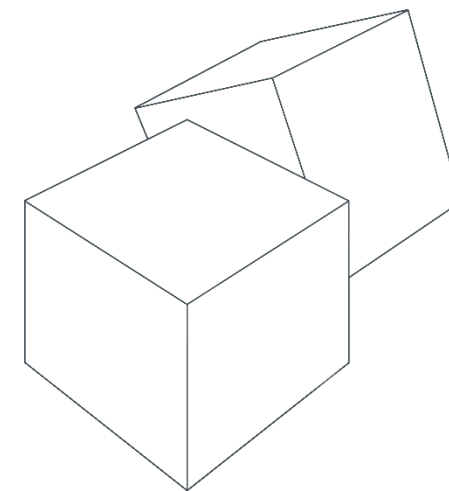
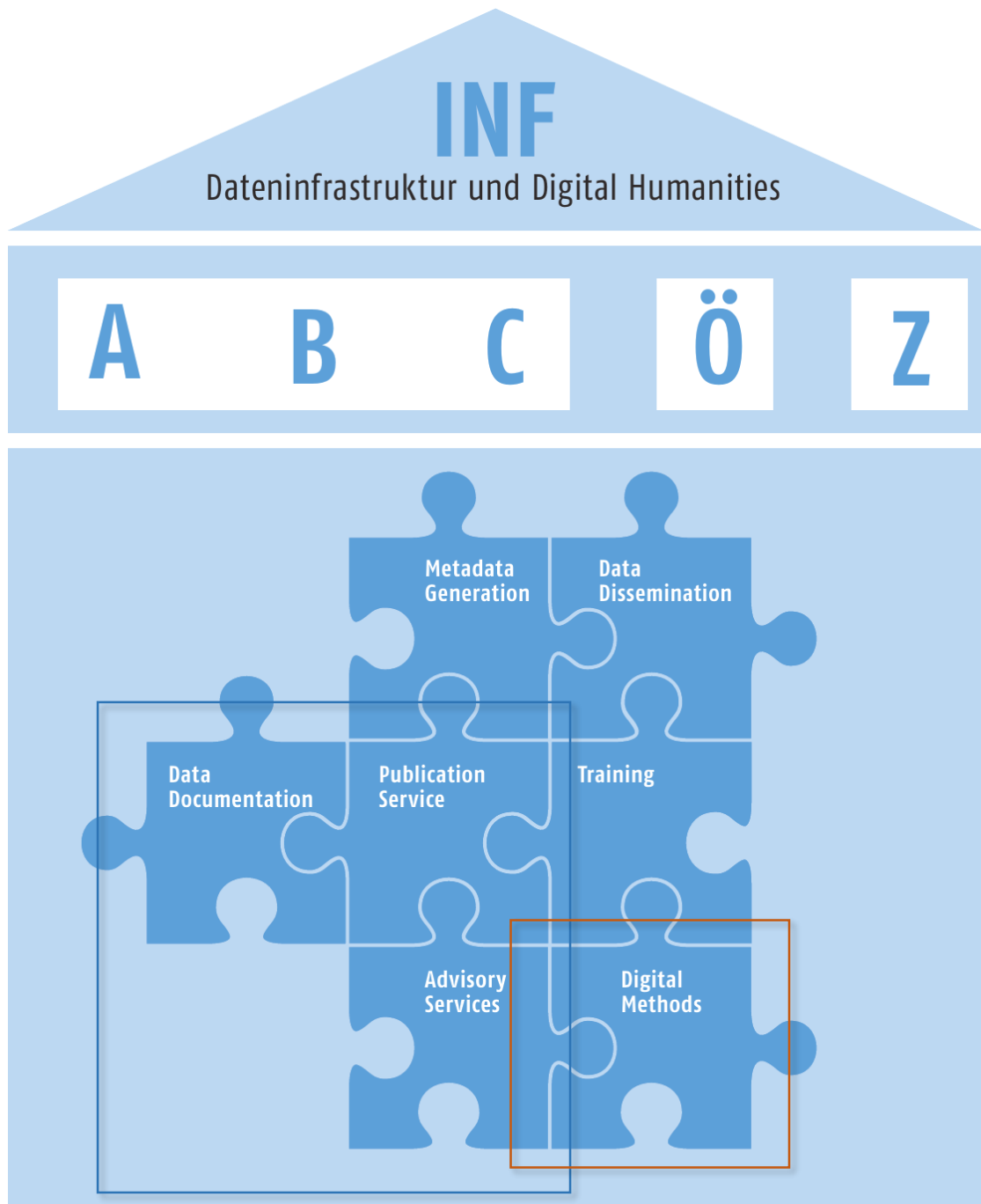


**Stefan Porada**  
wiss. Hilfskraft

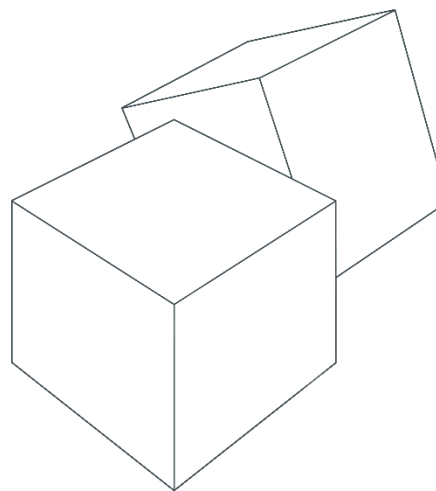
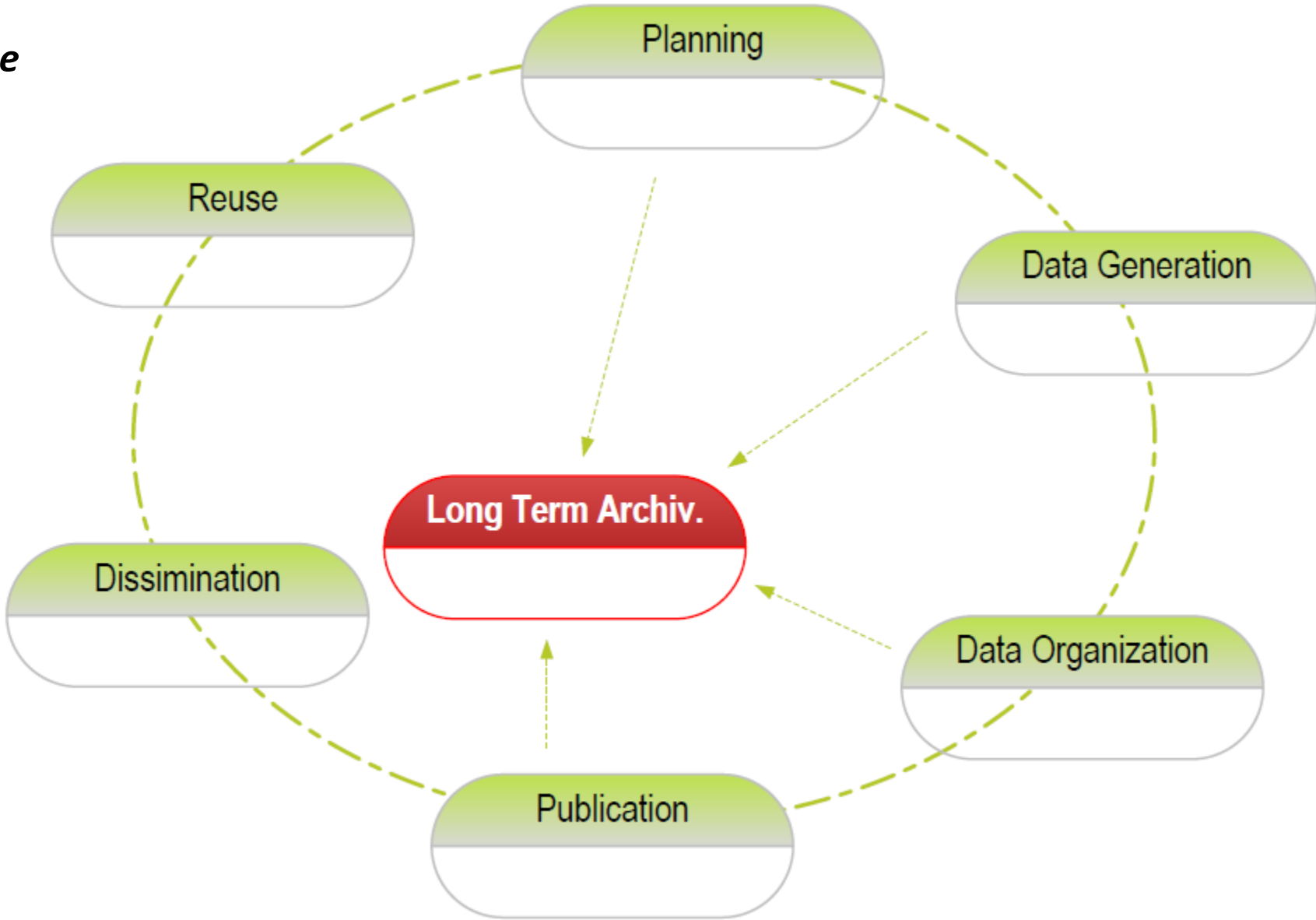


# INF Projekt im SFB 1288:

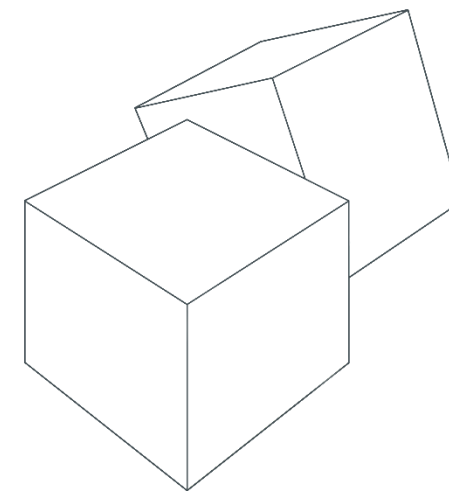
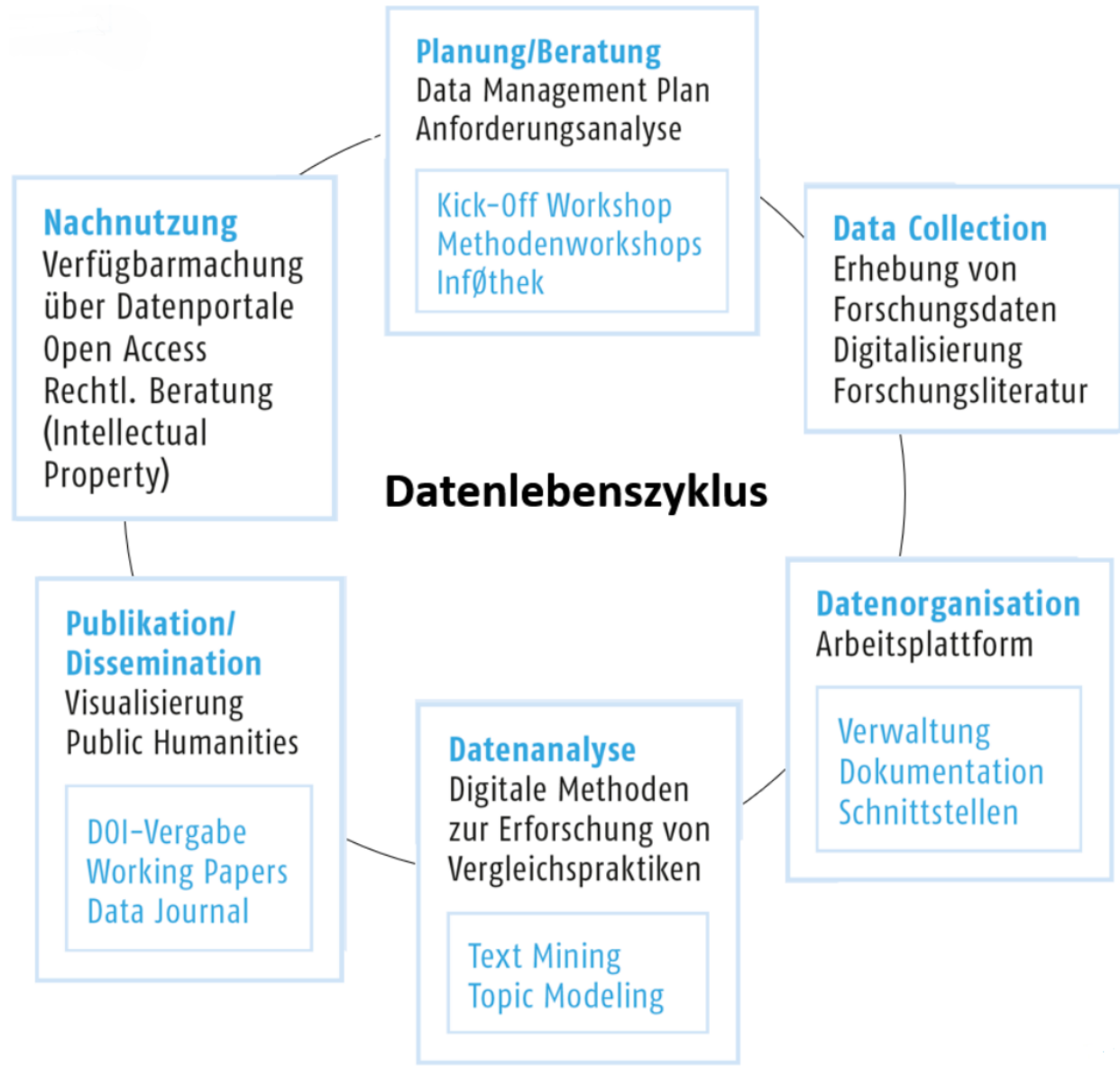
## Dateninfrastruktur und Digital Humanities



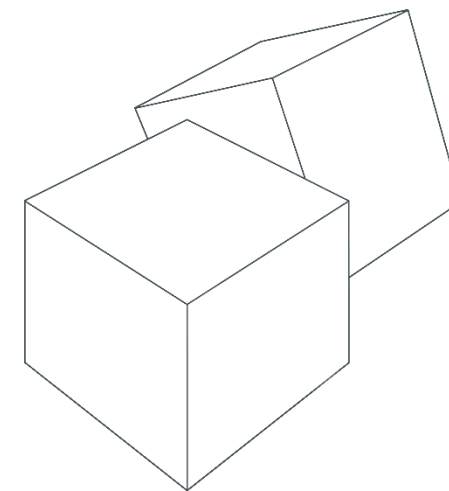
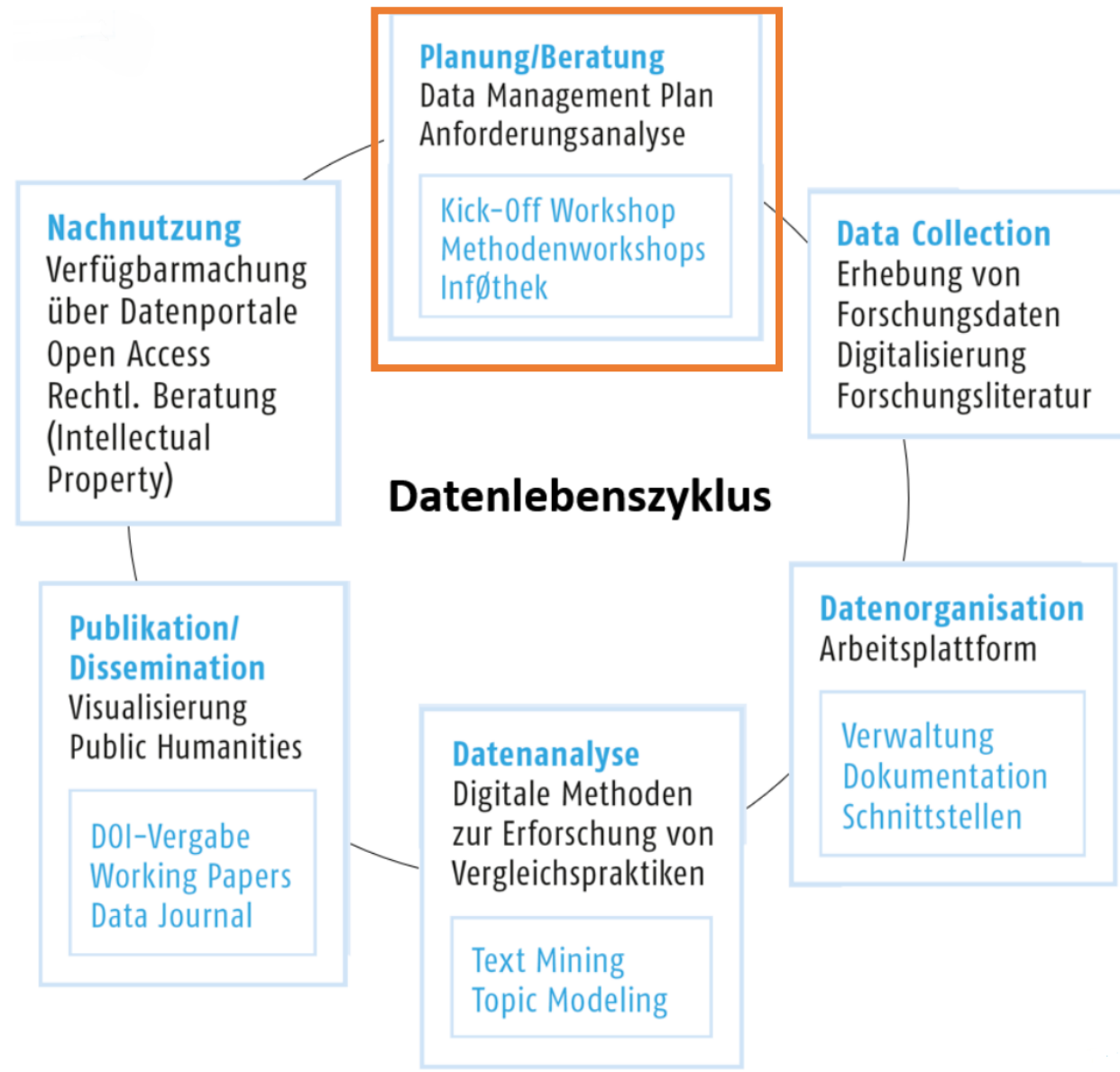
**Generischer  
Data Life Cycle**



Support der  
zahlreicher  
Stadien des  
*Data Life Cycle*  
in den Digital  
Humanities



Support der  
zahlreicher  
Stadien des  
*Data Life Cycle*  
in den Digital  
Humanities



## INF Projekt im SFB 1288: Formen der Zusammenarbeit

„*INFØthek*“: bilaterales Vorgehen

- Austausch über dateninfrastrukturelle Belange & Nutzung von Tools
- Analyse Ihrer **Forschungsprozesse** und **Arbeitsweisen** im Projekt
- Evaluation potenzieller **Schnittstellen zu INF** und **Anforderungen an digitale Methoden**

### Planung/Beratung

Data Management Plan  
Anforderungsanalyse

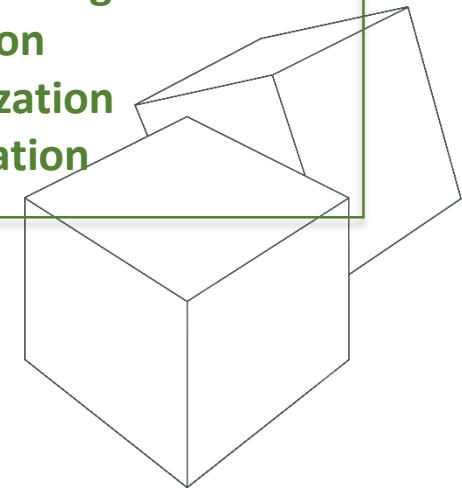
Kick-Off Workshop  
Methodenworkshops  
InfØthek

„Vergleichen“ im  
*Digital Age*

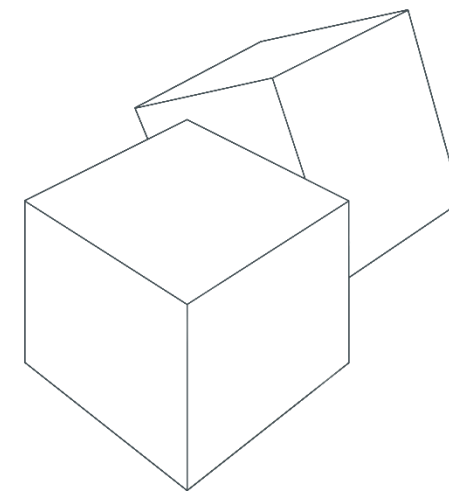
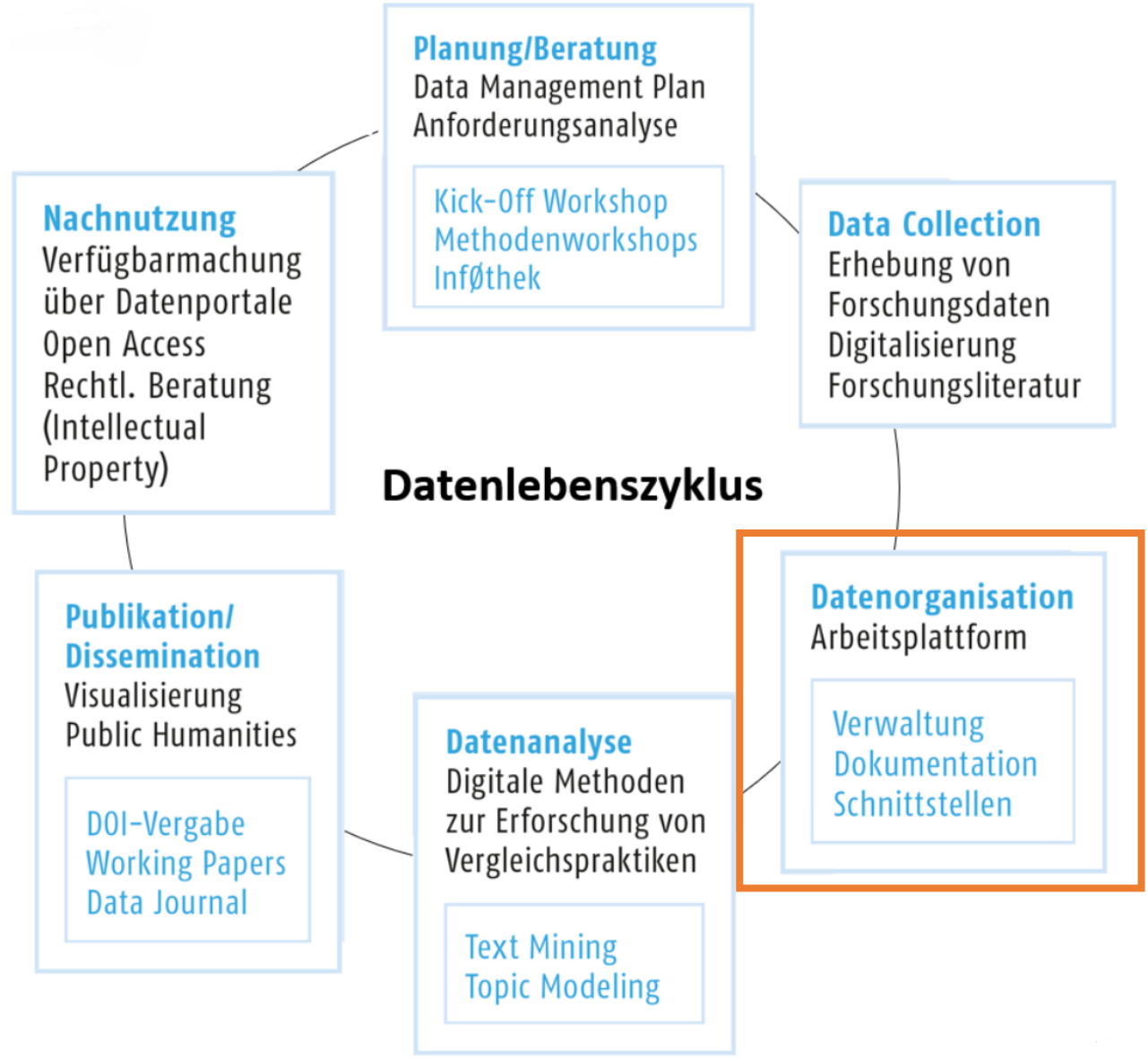
- Discovering
- Annotating
- Comparing
- Referring
- Sampling
- Illustrating
- Representing

### *Infrastruktursicht*

retrieval  
metadata enrichment  
similarity computation  
entity linkage  
selection  
visualization  
publication



Support der  
zahlreicher  
Stadien des  
*Data Life Cycle*  
in den Digital  
Humanities



# Kollaborationsplattform (Projektmanagement):

Hauptseite Projekte Kalender Urlaub Meine Seite SFB 1288 Webseite Administration

SFB 1288 / Suche  Zu

ein Konto Abmelden

Projekte

Geschlossene Projekte anzeigen

Projekte [+ Neues Projekt](#) | [Alle Tickets anzeigen](#) | [Aufgewendete Zeit a](#)

**A01 Streitkräftevergleich**

**B03: Weltvergleich und Weltwissen**

- B03: Weltvergleich und Weltwissen - Teilstudie 1 (Erhart)
- B03: Weltvergleich und Weltwissen - Teilstudie 2 (Kramer)

**C02 Nonkommensurabel?**

- Brainstorming
- C02 Teilprojekt 1 (Arlinghaus)
- C02 Teilprojekt 2 (Erhart)

**des Vergleichs**

**Geschichte**

**Gesamt-SFB**

- Antrag (Förderphase 1)
- Corporate Design
- Formulare
- GastwissenschaftlerInnen
- Grundordnung des SFB
- INFØthek
- Konferenzen
- Auftaktkonferenz Oktober 2017
- Mitgliederversammlung
- Newsletter - intern

**★ INF**

- ★ Digitalisierung
- ★ Digi\_HiWi
- ★ INF\_A01 Streitkräftevergleiche
- ★ INF\_A02 Vergleichen in der Konkurrenz
- ★ INF\_A03 (Welt-)Ordnungen und Zukunftsentwürfe
- ★ INF\_A04 Vergleichende Kulturdiskurse
- ★ INF\_B01 Interkulturelle Rechtsprechung
- ★ INF\_B02 Interamerikanische Vergleichspraktiken

**INF-Ö**

**Lehre Digital Humanities**

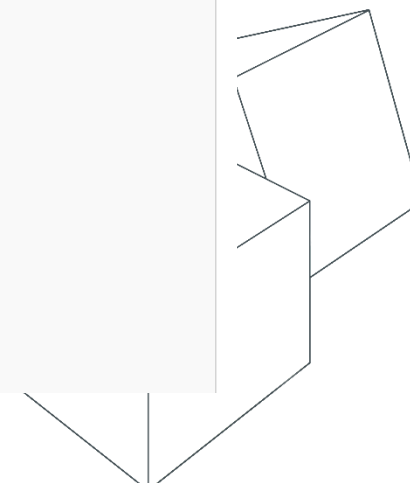
- BGHS Methods Class Sommersemester 2017

**Datenorganisation**  
Arbeitsplattform

Verwaltung  
Dokumentation  
Schnittstellen

Alle Projekte anzeigen

<https://redmine.sfb1288.uni-bielefeld.de/projects/new>





# Kollaborationsplattform (Projektmanagement):

Hauptseite Projekte Kalender Urlaub Meine Seite SFB 1288 Webseite Administration Angemeldet als admin Mein Kon...

Suche  Zu einem Projekt springen...

SFB 1288 /

Projekte

Geschlossene Projekte anzeigen

Anwenden

Projekte [+ Neues Projekt](#) | [Alle Tickets anzeigen](#) | [Aufgewendete Zeit aller Projekte anzeigen](#) | [Aktivitäten aller Proj...](#)

A01 Streitkr...

Hauptseite Projekte Kalender Urlaub Meine Seite SFB 1288 Webseite Administration Suche

SFB 1288 / INF

+ Übersicht **Aktivität** Tickets Gantt-Diagramm Kalender News DMS Wiki Foren Dateien Konfiguration

Aktivität

Tickets

Changesets

News

Dokumente

Dateien

Wiki-Bearbeitungen

Forenbeiträge

Benötigte Zeit

Dokumentenzugriffe

Dokumenteversion

Unterprojekte

Anwenden

Gesamt-SFB

Antrag (Förderp...

Corporate Desig...

Formulare

Gastwissens...

Grundordnung d...

INFØthek

Konferenzen

Auftaktkonfer...

Mitgliederversan...

Newsletter - inte...

<https://redmine.sfb1288.uni-bielefeld.de/projects/new>

Aktivität

von 16.09.2017 bis 15.10.2017

10.10.2017

🔗 14:13 Digi\_HiWi - Digitalisierung #89 (Neu): C01 Arbeitspaket#1

Die zu bearbeitenden Dateien des 1. Arbeitspakets (C01) finden Sie hier: <https://uni-bielefeld.sciebo.de/index.ph...>

Madis Rummig

🔗 12:39 Digi\_HiWi - Digitalisierung #88 (Neu): A04 Arbeitspaket#1

Die zu bearbeitenden Dateien des 1. Arbeitspakets (A04) finden Sie hier: <https://uni-bielefeld.sciebo.de/index.ph...>

Madis Rummig

🔗 12:38 Digi\_HiWi - Digitalisierung #87 (Neu): B03 Arbeitspaket#1

Die zu bearbeitenden Dateien des 1. Arbeitspakets (B03) finden Sie hier: <https://uni-bielefeld.sciebo.de/index.ph...>

Madis Rummig

🔗 12:29 Digi\_HiWi - Digitalisierung #85 (Neu): B01 Arbeitspaket#1

Die zu bearbeitenden Dateien des 1. Arbeitspakets (B01) finden Sie hier: <https://uni-bielefeld.sciebo.de/index.php/s/...>

Madis Rummig

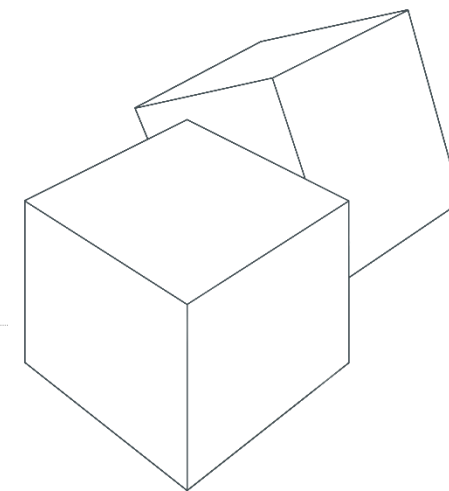
07.10.2017

🔗 10:18 Digitalisierung - Termine\_INF Digi #10 (Erledigt): Nachfrage Sabrina: WHKs und Oxygen

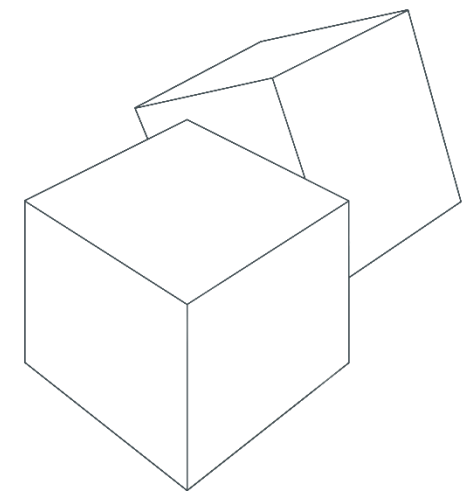
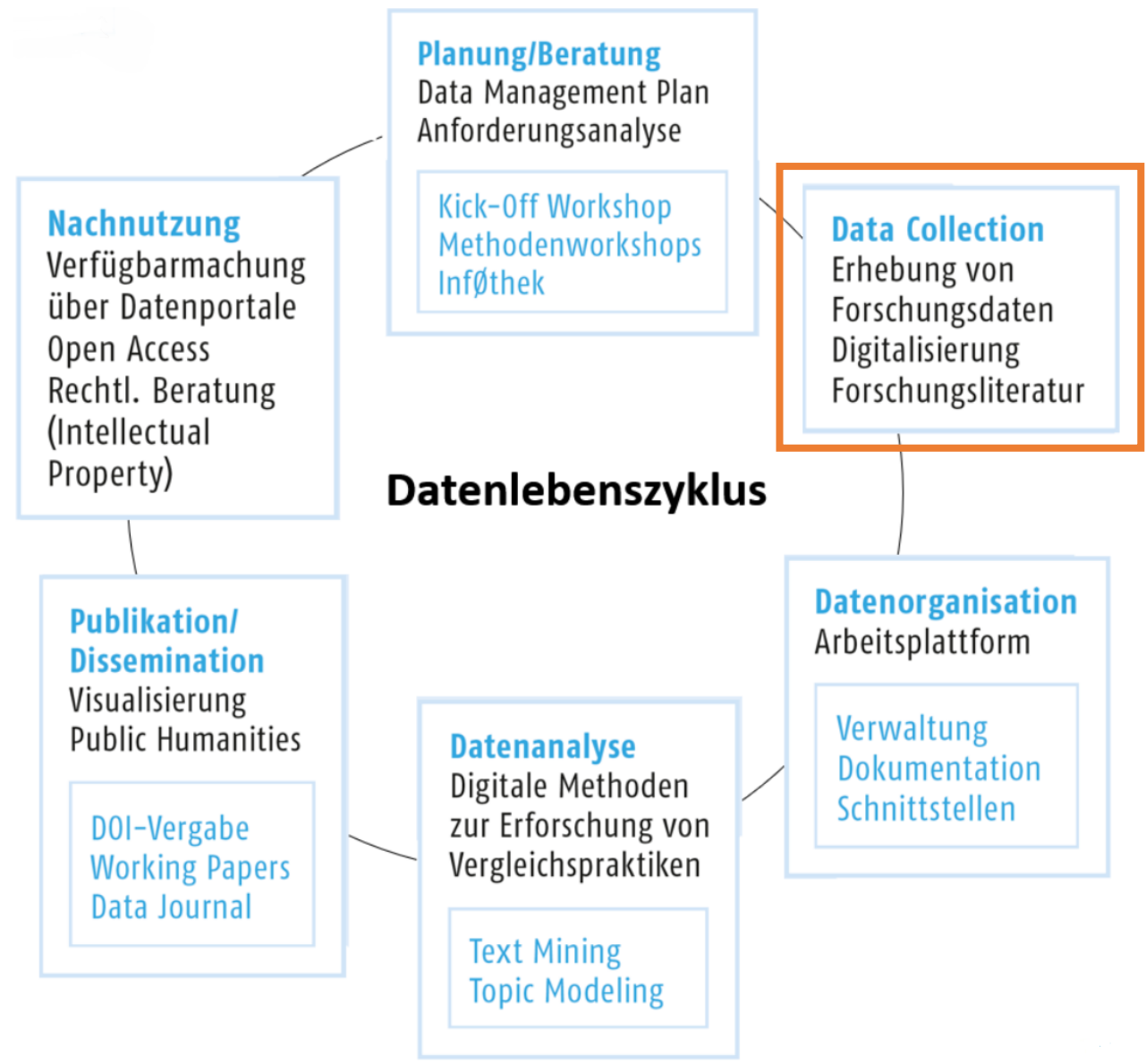
Madis Rummig

Datenorganisation  
Arbeitsplattform

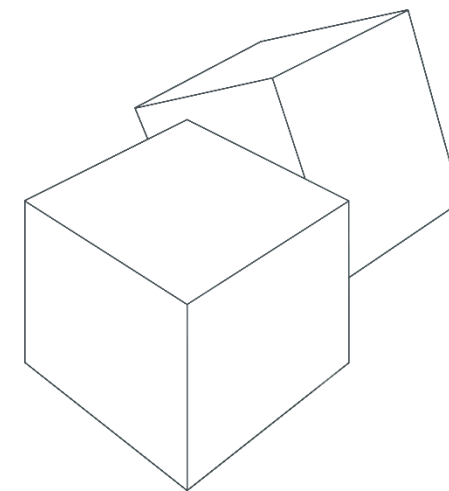
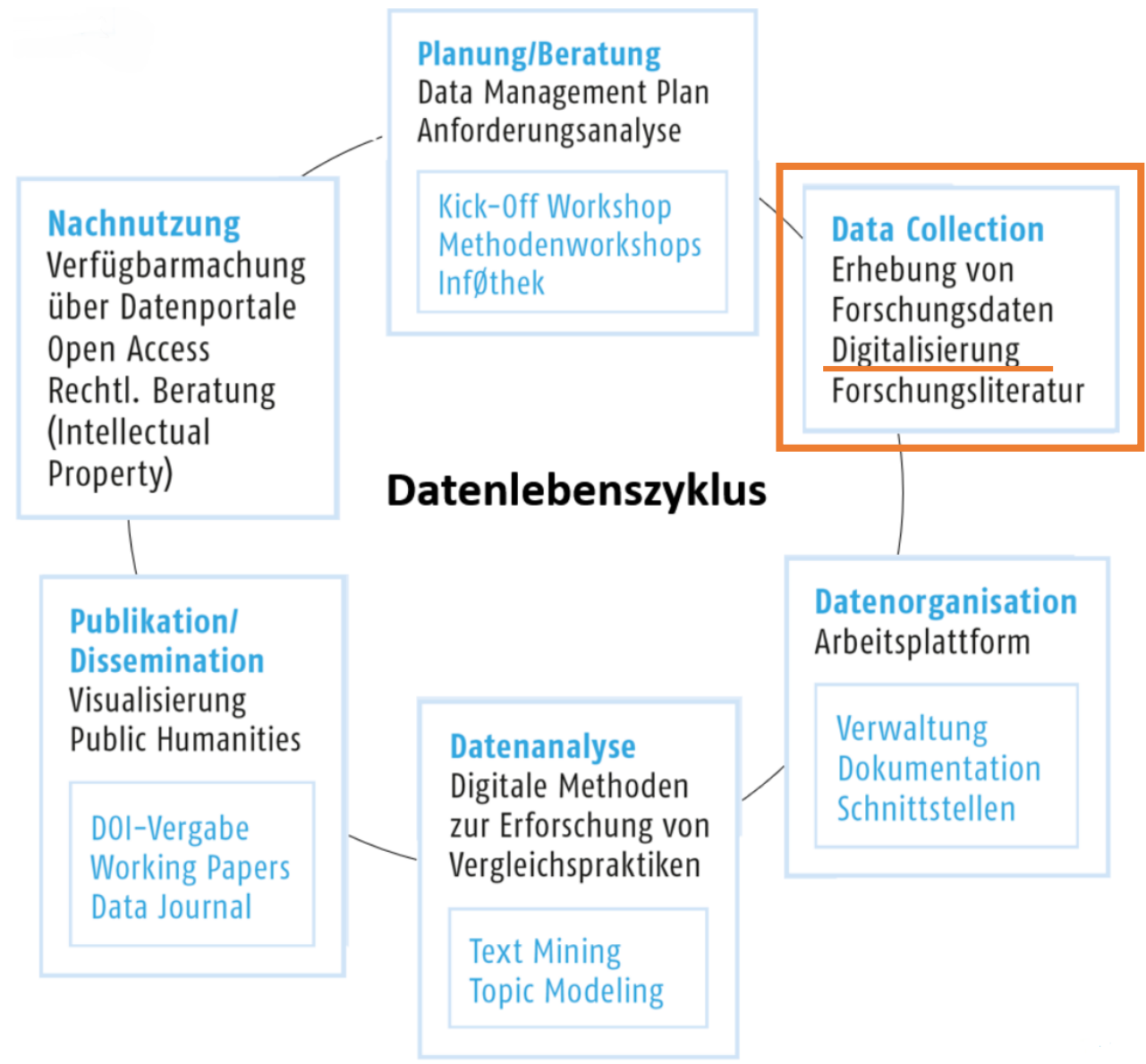
Verwaltung  
Dokumentation  
Schnittstellen



Support der  
zahlreicher  
Stadien des  
*Data Life Cycle*  
in den Digital  
Humanities



Support der  
zahlreicher  
Stadien des  
*Data Life Cycle*  
in den Digital  
Humanities



Data Collection  
Erhebung von  
Forschungsdaten  
Digitalisierung  
Forschungsliteratur

# Digitalisierung

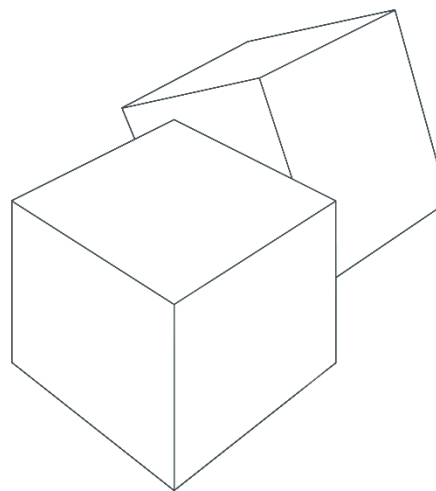
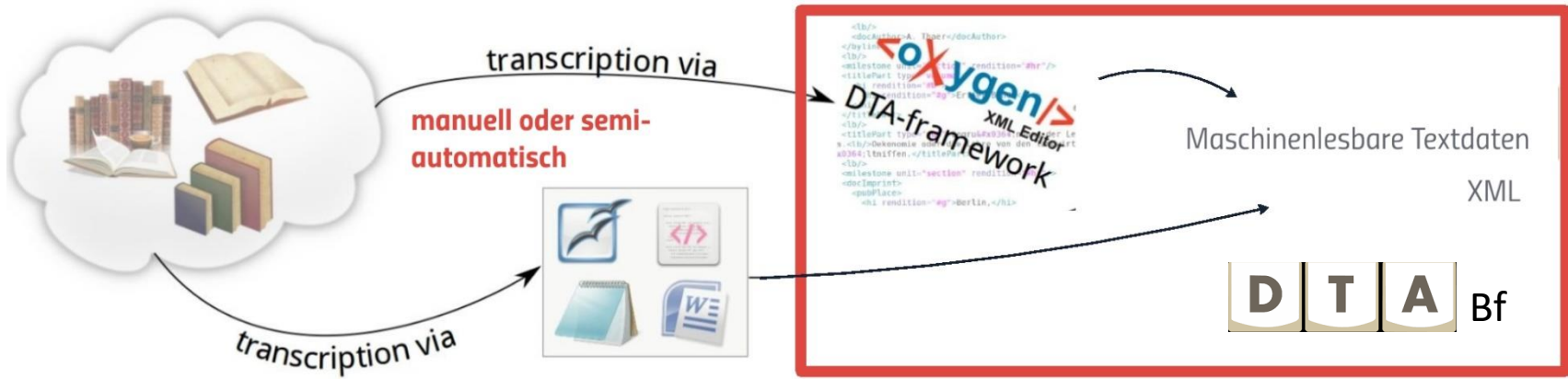
## Etablierung von Workflows

- OCR
- Transkription / Richtlinien
- Dokumentationsprache

## Etablierung von Tools

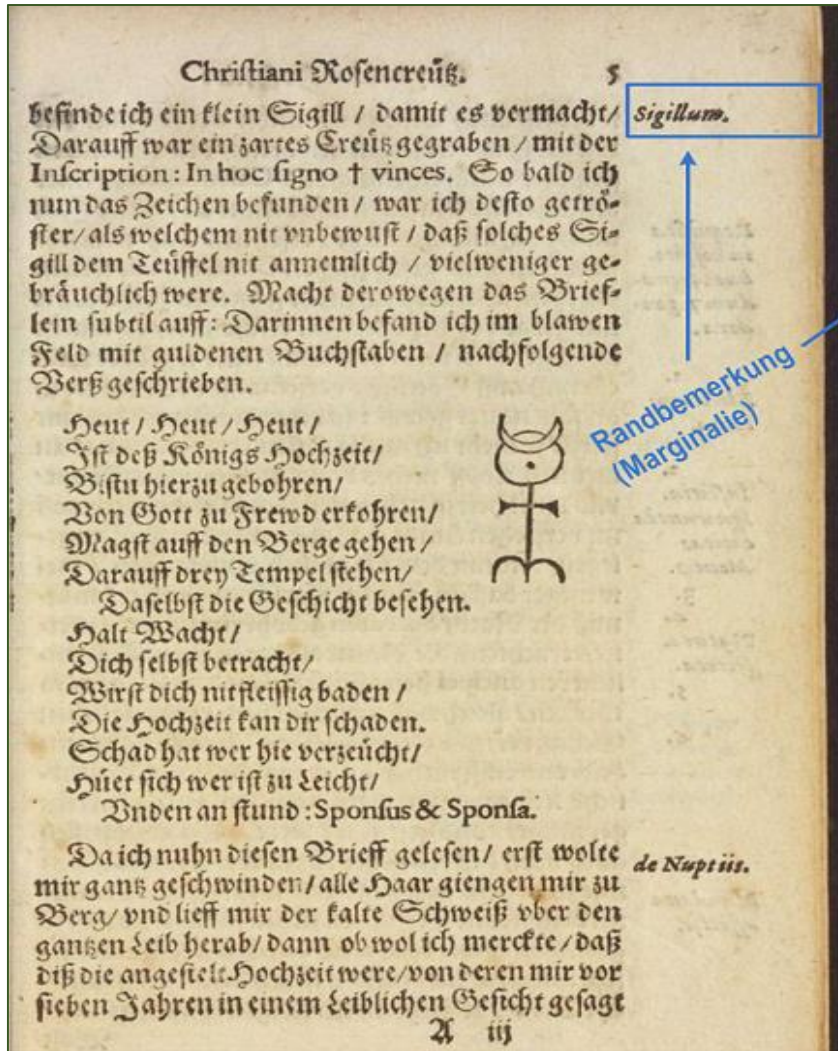
- Software (z.B. PoCoTo)
- Oxygen Framework
- DTA-Bf (TEI-P5-basiert)

TEI Standard: Text Encoding Initiative



# Digitalisierung

Textauszeichnung mit DTA-Bf (umfasst 133 TEI-P5 Elemente, XML-basiert)



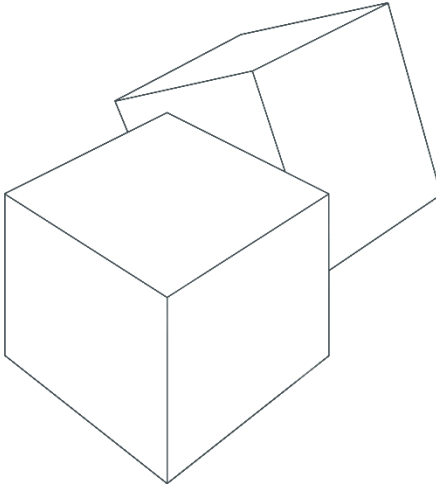
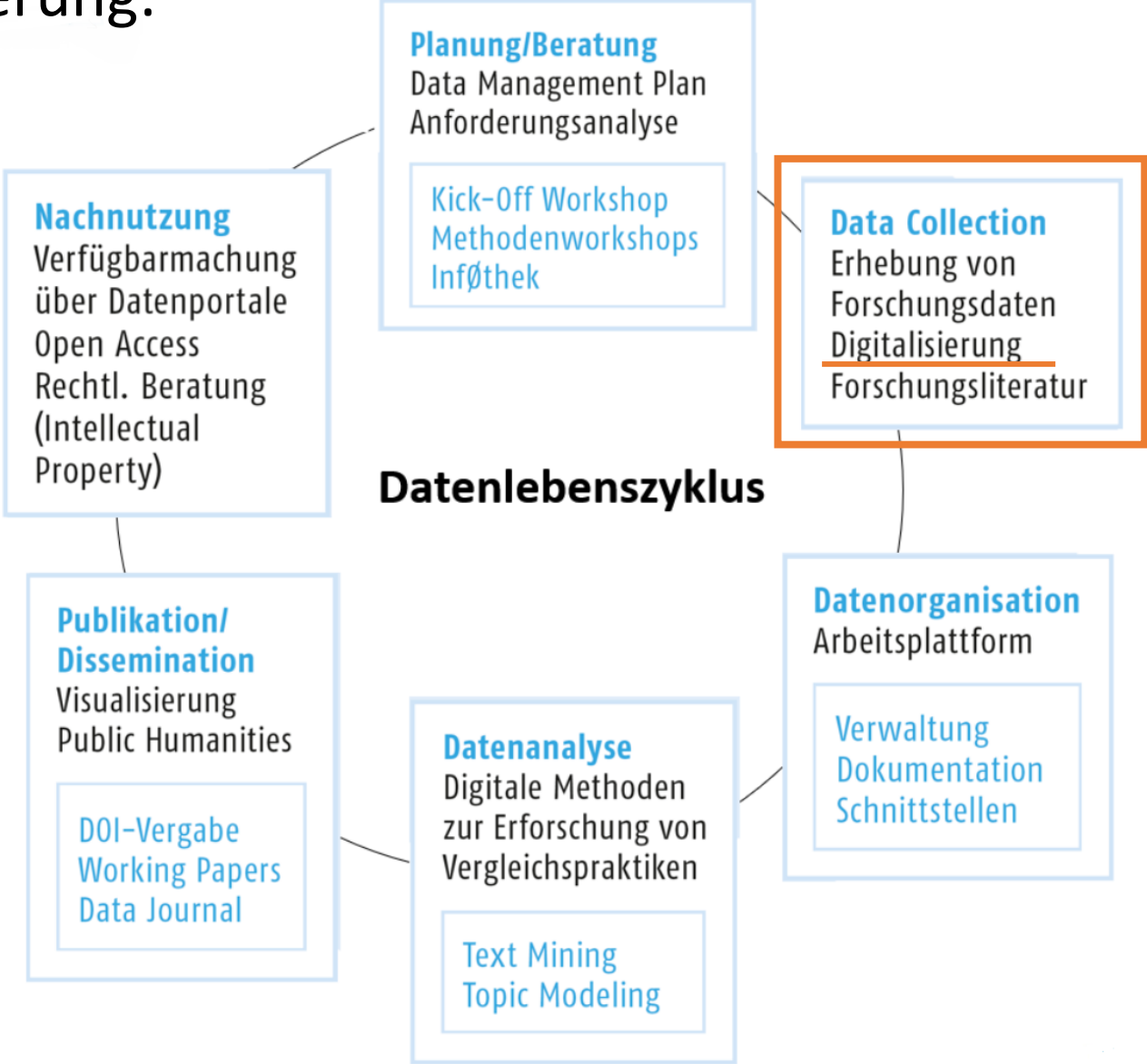
Sigillum.

Randbemerkung  
(Marginalie)

```
<TEI>
  <text>
    <body>
      <div n="1">
        <div n="2">
          <p><pb facs="#f0009" n="5"/><fw place="top" type="header"><hi
rendition="#aq">Chri&#x017F;tiani</hi> Ro&#x017F;encreu&#x0364;tz.</fw><lb/>
befinde ich ein klein Sigill/ damit es vermacht/<note place="right"><hi
rendition="#i"><hi rendition="#aq">Sigillum.</hi></hi></note><lb/>
Darauff war ein zartes Creu&#x0364;tz gegraben/ mit der<lb/><hi
rendition="#aq">In&#x017F;cription: In hoc &#x017F;igno &#x2020; vinces</hi>.
So bald ich<lb/>
nun das Zeichen befunden/ war ich de&#x017F;to getro&#x0364;-<lb/>
&#x017F;ter/ als welchem nit vnbewu&#x017F;t/ daß &#x017F;olches Si-<lb/>
gill dem Teu&#x0364;ffel nit annemlich/ vielweniger ge-<lb/>
bra&#x0364;uchlich were. Macht derowegen das Brief-<lb/>
lein &#x017F;ubtil auff: Darinnen befand ich im blawen<lb/>
Feld mit guldenen Buch&#x017F;taben/ nachfolgende<lb/>
Verß ge&#x017F;chrieben.</p><lb/>
          <lg type="poem">
            <lg n="1">
              <l>Heut/ Heut/ Heut/ <figure/></l><lb/>
              <l>I&#x017F;t deß Ko&#x0364;nigs Hochzeit/</l><lb/>
              <l>Bi&#x017F;tu hierzu geböhren/</l><lb/>
              <l>Von Gott zu Frewd erköhren/</l><lb/>
              <l>Mag&#x017F;t auff den Berge gehen/</l><lb/>
              <l>Darauff drey Tempel &#x017F;tehen/</l><lb/>
              <l>Da&#x017F;elb&#x017F;t die Ge&#x017F;chicht be&#x017F;ehen.</l>
            </lg><lb/>
            <lg n="2">
              <l>Halt Wacht/</l><lb/>
              <l>Dich &#x017F;elb&#x017F;t betrachte/</l><lb/>
              <l>Wir&#x017F;t dich <choice>
<sic>nitfleißig&#x017F;ig&#x017F;ig</sic><corr>nit fleißig&#x017F;ig</corr>
</choice> baden/</l><lb/>
              <l>Die Hochzeit kan dir &#x017F;chaden.</l><lb/>
            </lg>
          </lg>
        </div>
      </div>
    </body>
  </text>
</TEI>
```

# Pilotphase Digitalisierung:

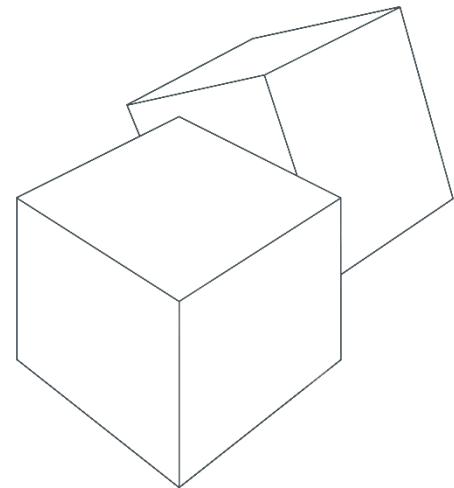
## Vorbereitung für maschinengestützte Korpusanalyse



# Pilotphase Digitalisierung: Vorbereitung für maschinengestützten Korpusanalyse

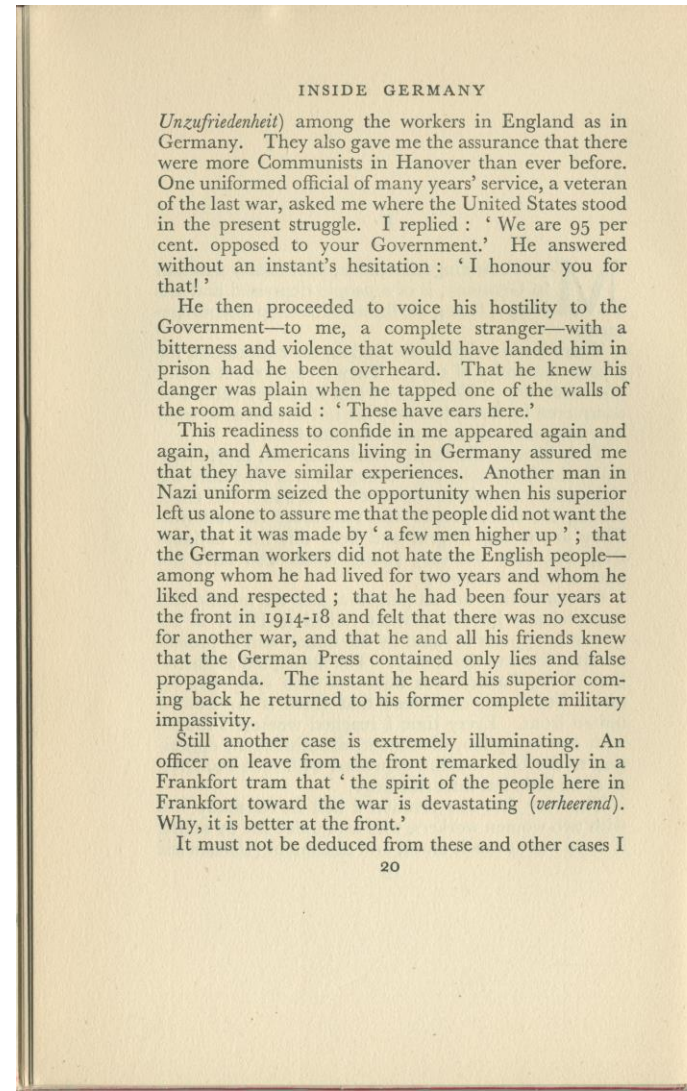
Korpus: Sammlung von (repräsentativen) Texten oder sprachlichen Äußerungen zum Zwecke wissenschaftlicher Untersuchungen.

- Digitale Methoden zur Verarbeitung textueller Daten setzen Maschinenlesbarkeit voraus
- Ziel: Erzeugung maschinenlesbarer Textformate durch Texterkennung in Digitalisaten



# Vorbereitende Data Collection zur Datenanalyse und (Daten)Publikation

- Digitalisat bisher: Digitale Photographien oder Scans
  - Einzelne Bilddateien
  - Eingebettet in eine PDF
- Digitalisat in den Digital Humanities: Annotiertes Textformat, vorzugsweise XML-basiert
- Digitalisat im SFB1288: Markup gemäß DTAbf oder TEI P5



```

<pb n="20" facs="0023.tif"/>
<fw type="header" place="top">INSIDE GERMANY</fw>
<lb/>
Unzufriedenheit) among the workers in England as in<lb/>
Germany. They also gave me the assurance that there<lb/>
were more Communists in Hanover than ever before.<lb/>
One uniformed official of many years' service, a veteran<lb/>
of the last war, asked me where the United States stood<lb/>
in the present struggle. I replied : 'We are 95 per<lb/>
cent. opposed to your Government.' He answered<lb/>
without an instant's hesitation : 'I honour you for<lb/>
that!'

```

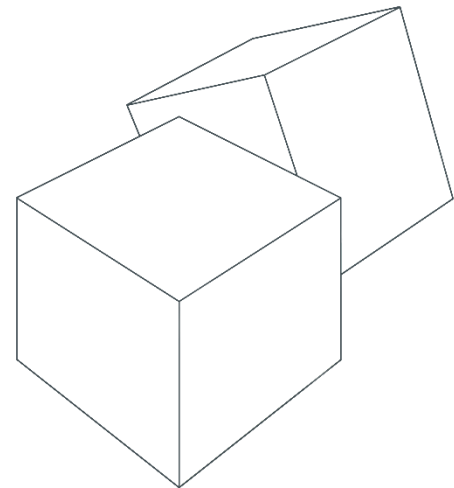


# DTabf: Deutsches Textarchiv Basisformat

- Entwickelt vom Deutschen Textarchiv
  - Deutschsprachige Texte von 1600-1900
  - 3319 annotierte Volldigitalisate
- DTabf: Erweiterung von TEI P5
  - Eliminierung von TEI P5 Ambiguität
  - Erfassung von Inhalt und formaler Struktur

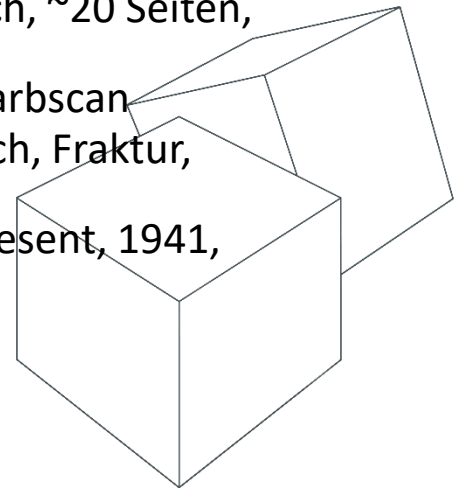


```
nicht sagen/ das Christus allein im Himmel vnd nicht bey vns<lb/>
<cit xml:id="bibl9" next="#quote12">
.....<bibl>
.....<note place="left">Pfal. 139.</note>
.....</bibl>
.....</cit>
auff Erden sey/ sintemahl er sitzet zu der Rechten der krafft Got-
tes/ welche sich nicht theilen lest/ sondern allenthalben ist/ vnd alles
erfüllet. Deñ also stehet geschrieben im 139 Pfalm.
<cit xml:id="quote12" prev="#bibl9">
.....<quote>Wo<lb/>
..... soll ich hingehn für deinem Geist/ vnd wo sol ich hinfliehen für
..... deinem Angesicht? Fuhre ich gehn Himmel/ so bistu da/ bettet<lb/>
..... ich mir in die Helle/ siehe so bistu auch da/ Nehme ich flügel der
..... Morgenröte/ vnd bliebe am eussersten Meer/ so würde mich doch<lb/>
..... deine Handt daselbst führen/ vnd deine Rechte mich halten:</quote>
</cit>
```

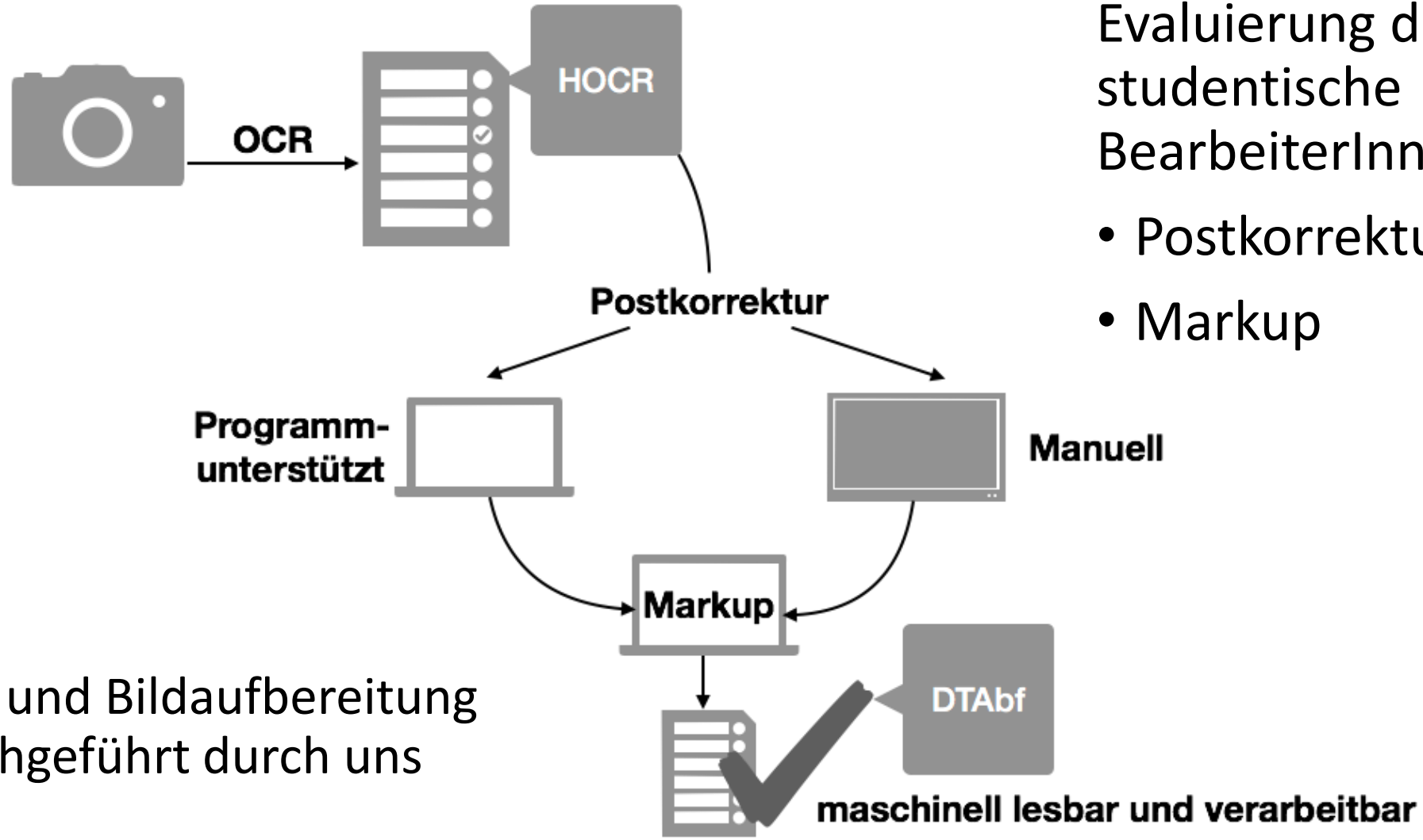


# Teilnehmende Projekte

- A04 — Kulturvergleiche durch Zeitschriften
  - Die Gegenwart 1915-1917, deutsch, Fraktur, 300+ Seiten, Microfiche (2003)
- B01 — Welches Recht gilt für wen?
  - Mémoire a consulter, et consultation, 1768, französisch, 75 Seiten, Farbfotografien
  - Voyage dans les mers de l'Inde, 1779, französisch, 143 Seiten, Farbscan
  - Mœuers et coutumes des indiens, 1987, französisch, 129 Seiten, SW-Scan
- B03 — Travel is the school of comparison
  - Entdeckungs-Reise in die Süd-See und nach der Herings-Straße zur Erforschung einer nordöstlichen Durchfahrt, 1821, deutsch, Fraktur, ~290 Seiten, SW-Scan
- B05 — Der englische Roman als Labor
  - Castle Rackrent, an hiberian tale, 1880, englisch, 224 Seiten, SW-Scan
- B06 — Rechtliche Vergleichsverbote
  - Ausgewählte UN-Protokolle, 1972-2012, englisch/ französisch, ~300 Seiten, SW-Scans (Schreibmaschine, PDF als Officeexport)
- C01 — Das vergleichende Sehen
  - An essay on the theory of painting, 1725, englisch, ~290 Seiten, Farbscan
- C03 — Eine Begriffsgeschichte des Vergleichens
  - Moskau 1937 Ein Reisebericht für meine Freunde, 1937, 152 Seiten, Farbscan
  - One year of hitlerism, 1934, englisch, ~20 Seiten, Farbscan
  - Inside Germany, 1939, 90 Seiten, Farbscan
  - Das System Mussolini, 1924, deutsch, Fraktur, ~140 Seiten, SW-Scan
  - Black Record: Germans past and present, 1941, englisch, ~70 Seiten, Farbscan



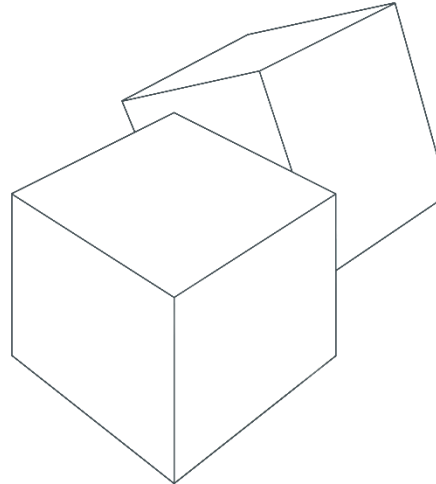
# Digitalisierungspipeline



Evaluierung durch studentische BearbeiterInnen:

- Postkorrektur
- Markup

OCR und Bildaufbereitung durchgeführt durch uns



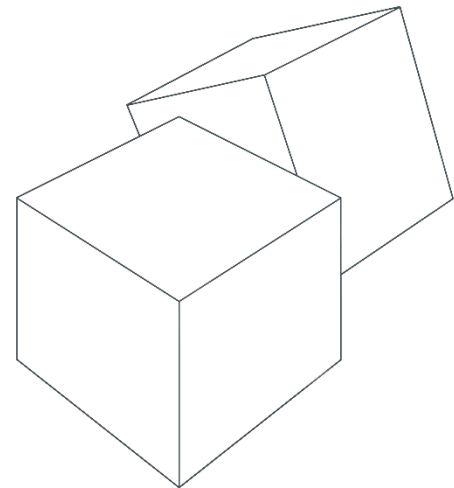
# Ziele der Pilotphase Digitalisierung: Ursprungsmedium



Welche Eingabeformate eignen sich für eine tiefergehende Digitalisierung?

Welche Mindestanforderungen müssen erfüllt sein?

- Auflösung, Detailfülle
- Farbe, Schwarz-Weiss
- Maximale Verzerrung im Bild, Verschmutzungsgrad

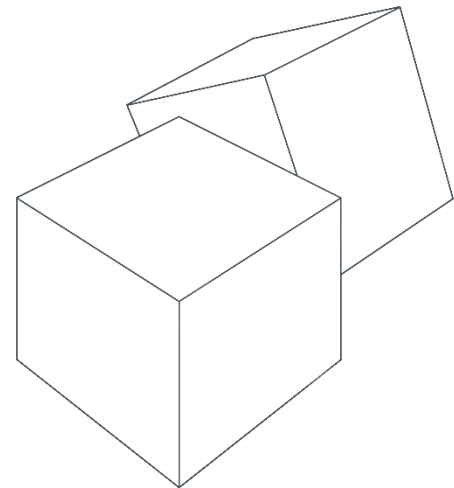


# Ziele der Pilotphase Digitalisierung: OCR

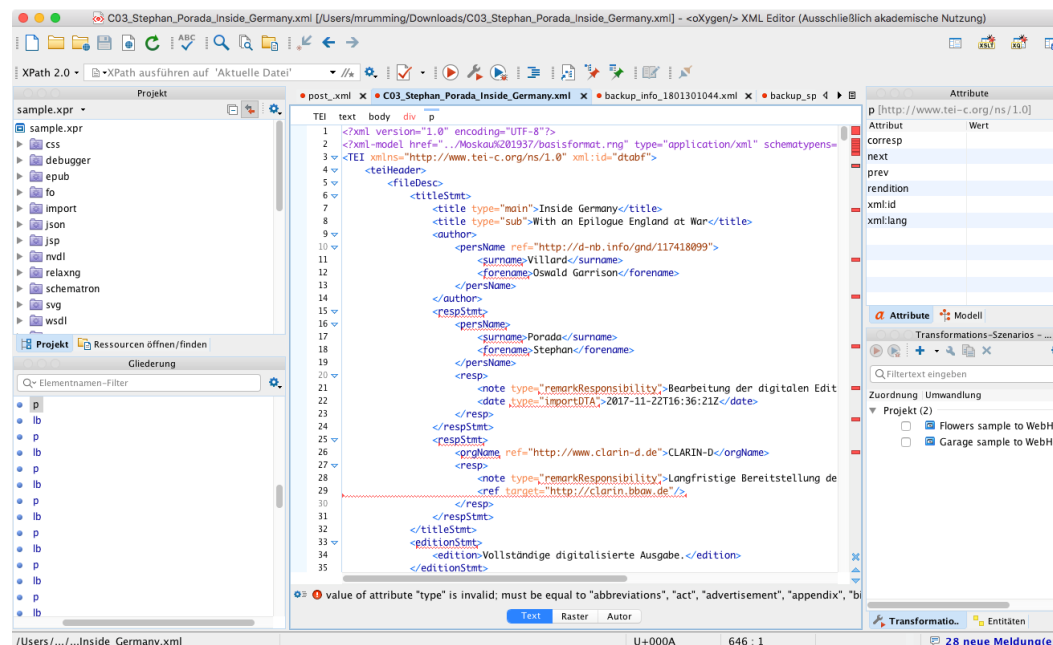
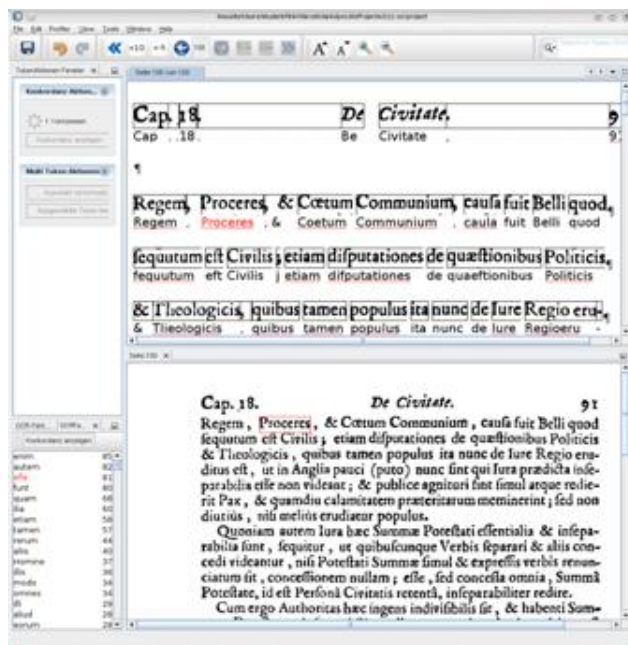
117	<b>Lachs</b>	xi7_Kchs
<p>Männlein aber sich hauptsächlich im Haupt-Fluß, oder in der Elbe zu halten pflegten. Es gedencket auch eben dieser Auctor aus einem alten Manuscripto, das An. 1432. ein so grosses Heer von Lachsen angekommen, daß sie bey nahe die Elbe nicht beherrbergen, und ein Fisch dem andern nicht ausweichen können, daher die Leute Hauffen Weise mit Netzen herzugelauffen, und die Fische erschlagen. Den Vortheil des Lachs-Fangs genüßet auch Schlesien von der Oder, und es sind von langen Jahren her ansehnliche Fangereyen längst der Oder, z. E. bey</p>	<p>Männlein<del>c</del>cher sich hauptsächlich im Haupt-Fluss, öderm der Gbe zu halten pflegten. Es gedencktt auch eben dieser Auctor aus einem alten Mannferipto, das An. 1431. ein 1o grosses Heer von Lachsen angekommen, daß sie bey nahe die Elbe nicht beherrbergen, und ein Fisih dem andern nicht auSweichm können, daher die Leute Haussen Weise mit Aexem bcr;ugelauffen, und die Zische erschlagen. Den Vortheil des LachS-Fangs gmüßet auch Schlesim von der Obtti und es sind von langen Jahren her ansehnliche Fangereyen längst der Oder, 5. &amp; bey</p>	

Welche open source OCR-Engine liefert die besten Ergebnisse?  
Welches freie Erkennungsmodell ist am besten?

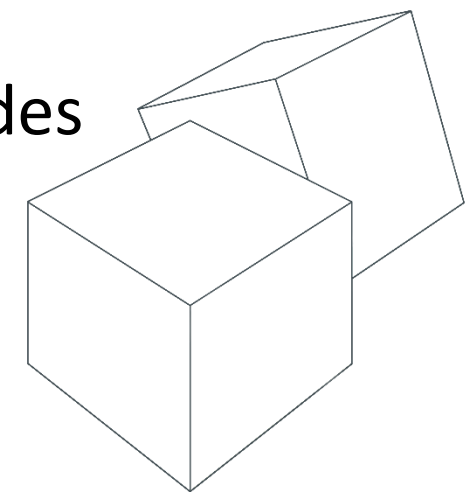
- ocropy: en, french\_balzac, de, de\_frak(2x), antiqua
- tesseract: deu, deu\_fra, fra, frm, eng, enm



# Ziele der Pilotphase Digitalisierung: Postkorrektur und Markup

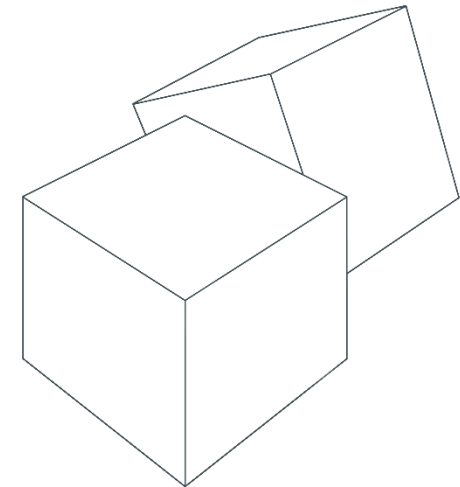
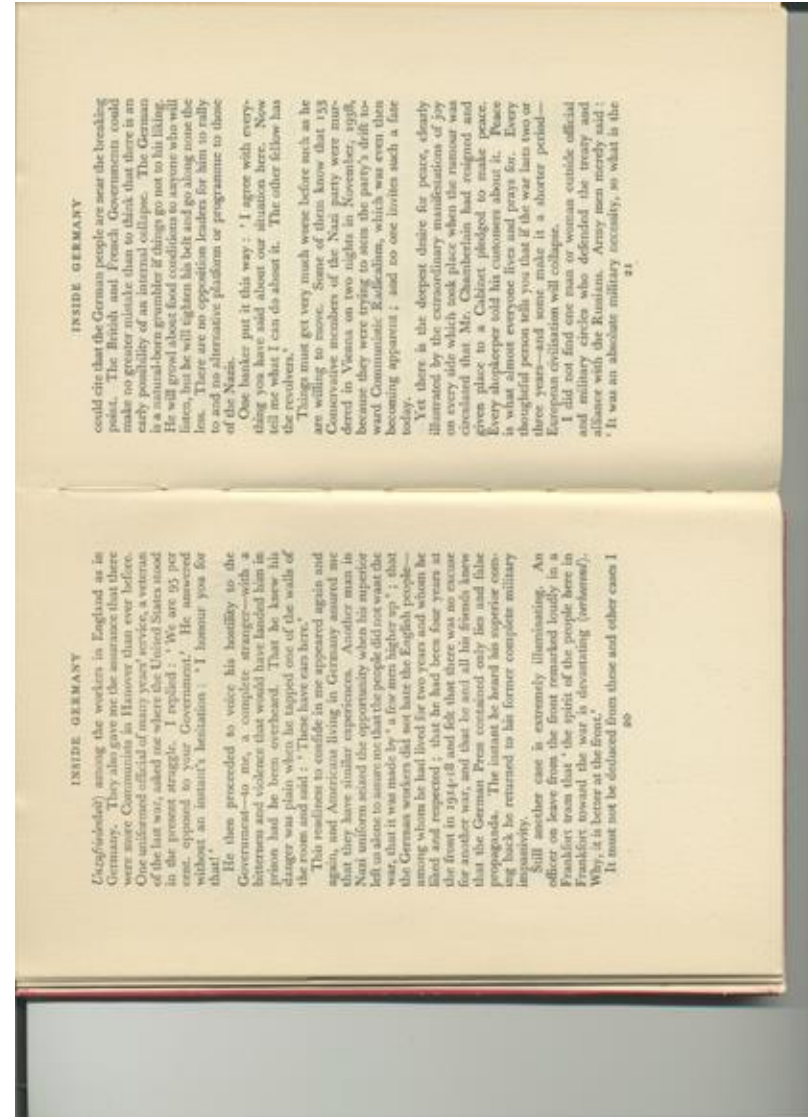


- Programmunterstützte Postkorrektur und danach anschließendes Markup
- Postkorrektur und Markup direkt in einem Schritt
- „Abtippen“ und Markup schneller als vorheriges OCR?



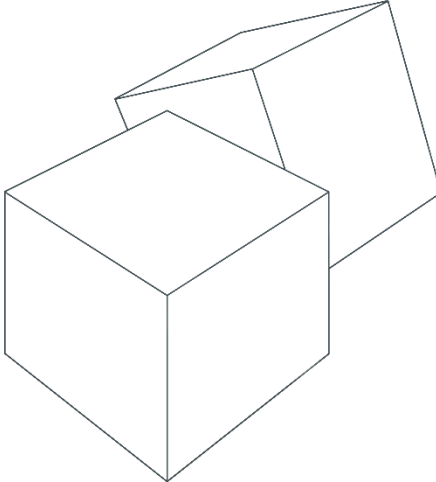
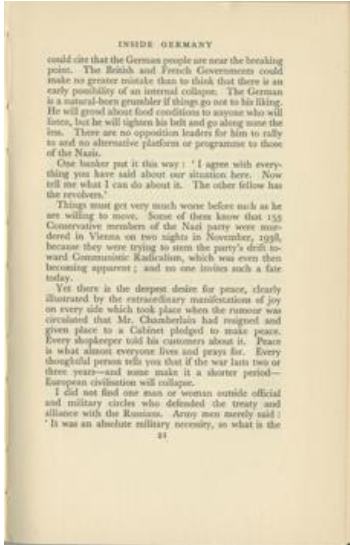
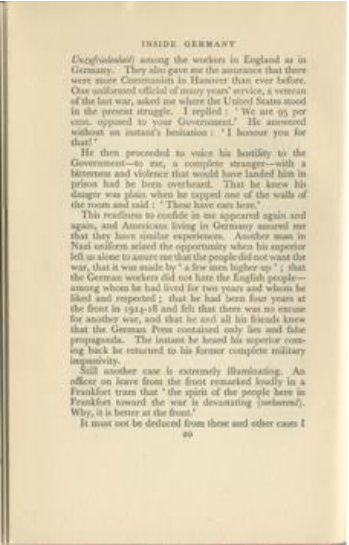
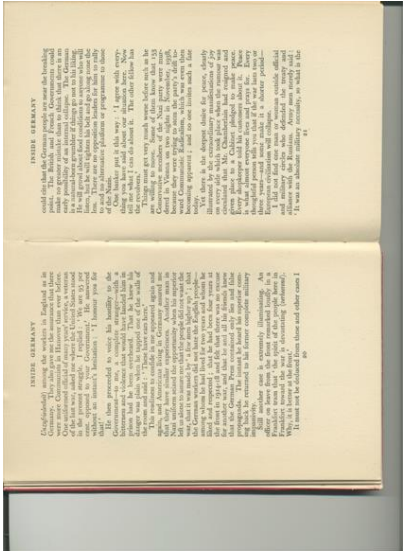
# C03: Eine Begriffsgeschichte des Vergleichens

- Doppelseitiger Farbscan
- Einheitlicher Bildausschnitt ermöglicht:
  - automatisches Drehen
  - automatisches Zuschneiden



# C03: Eine Begriffsgeschichte des Vergleichens

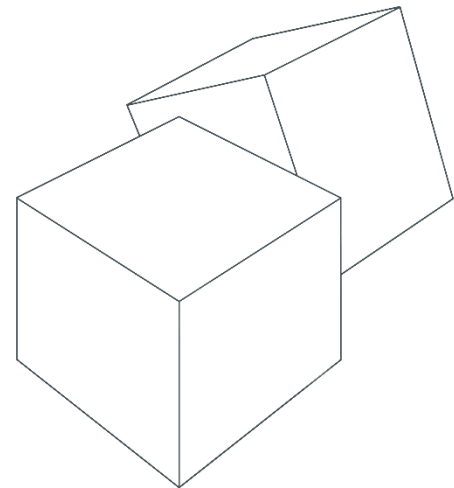
- Doppelseitiger Farbscan
- Einheitlicher Bildausschnitt ermöglicht:
  - automatisches Drehen
  - automatisches Zuschneiden





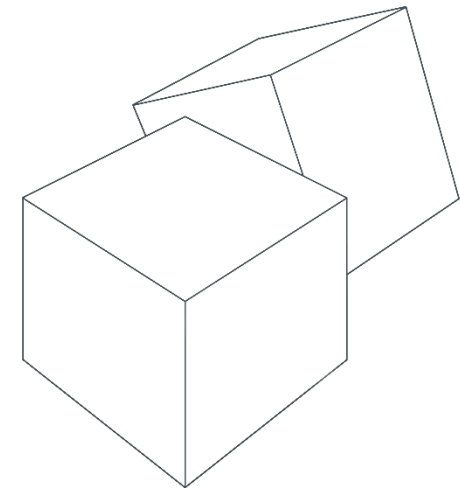
# A04: Kulturvergleiche durch Zeitschriften

- Microfichefotografie
  - Kontrastarm
  - Bildkontamination
- Zweispaltiges Seitenlayout
- „schiefe“ Bilder
- Frakturschrift



# A04: Kulturvergleiche durch Zeitschriften

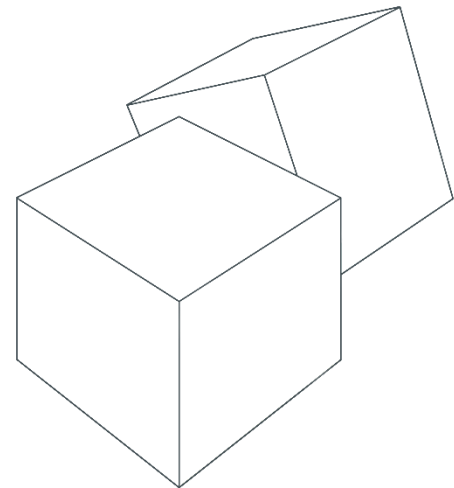
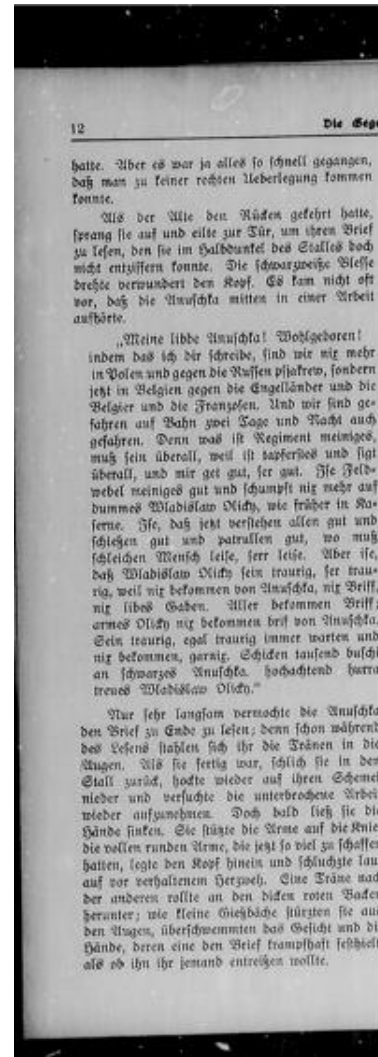
- Microfichefotografie
  - Kontrastarm
  - Bildkontamination
- Zweispaltiges Seitenlayout
- „schiefe“ Bilder
- Frakturschrift





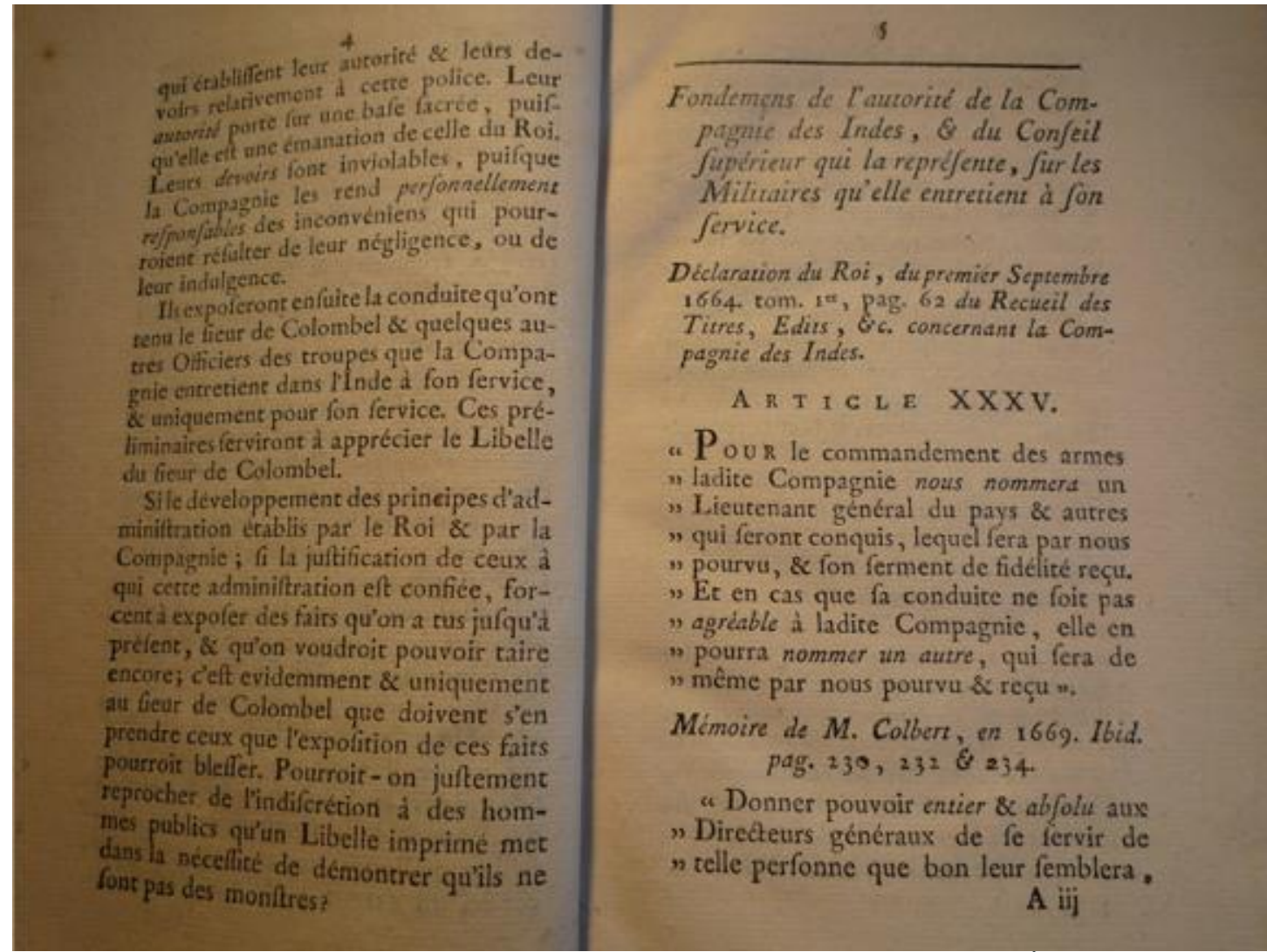
# A04: Kulturvergleiche durch Zeitschriften

- Microfichefotografie
  - Kontrastarm
  - Bildkontamination
- Zweispaltiges Seitenlayout
- „schiefe“ Bilder
- Frakturschrift



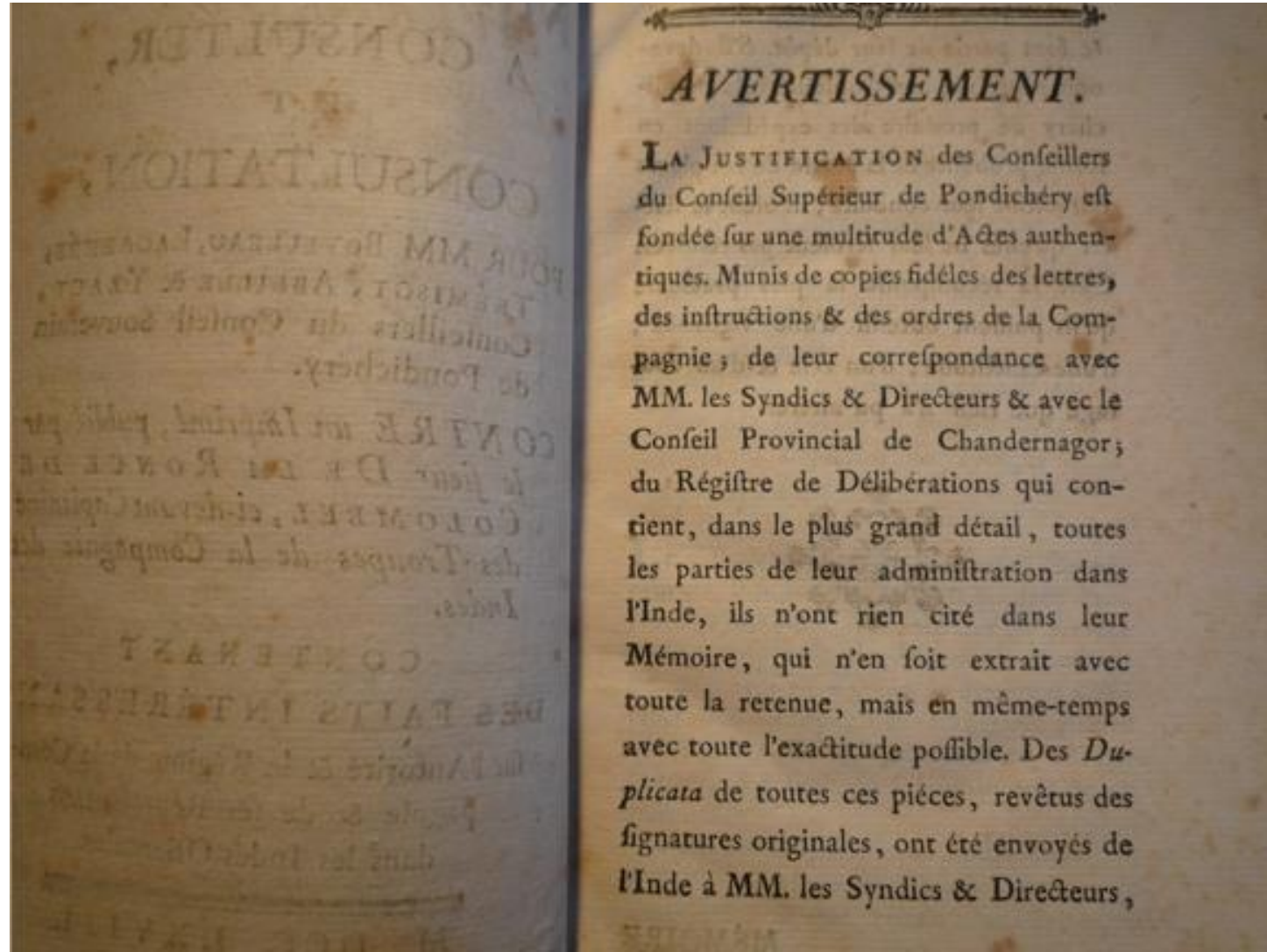
# B01: Welches Recht gilt für wen?

- Doppelseitige Farbfotografie
- Ausschnitt von Bild zu Bild variierend
- Text ist stark verzerrt



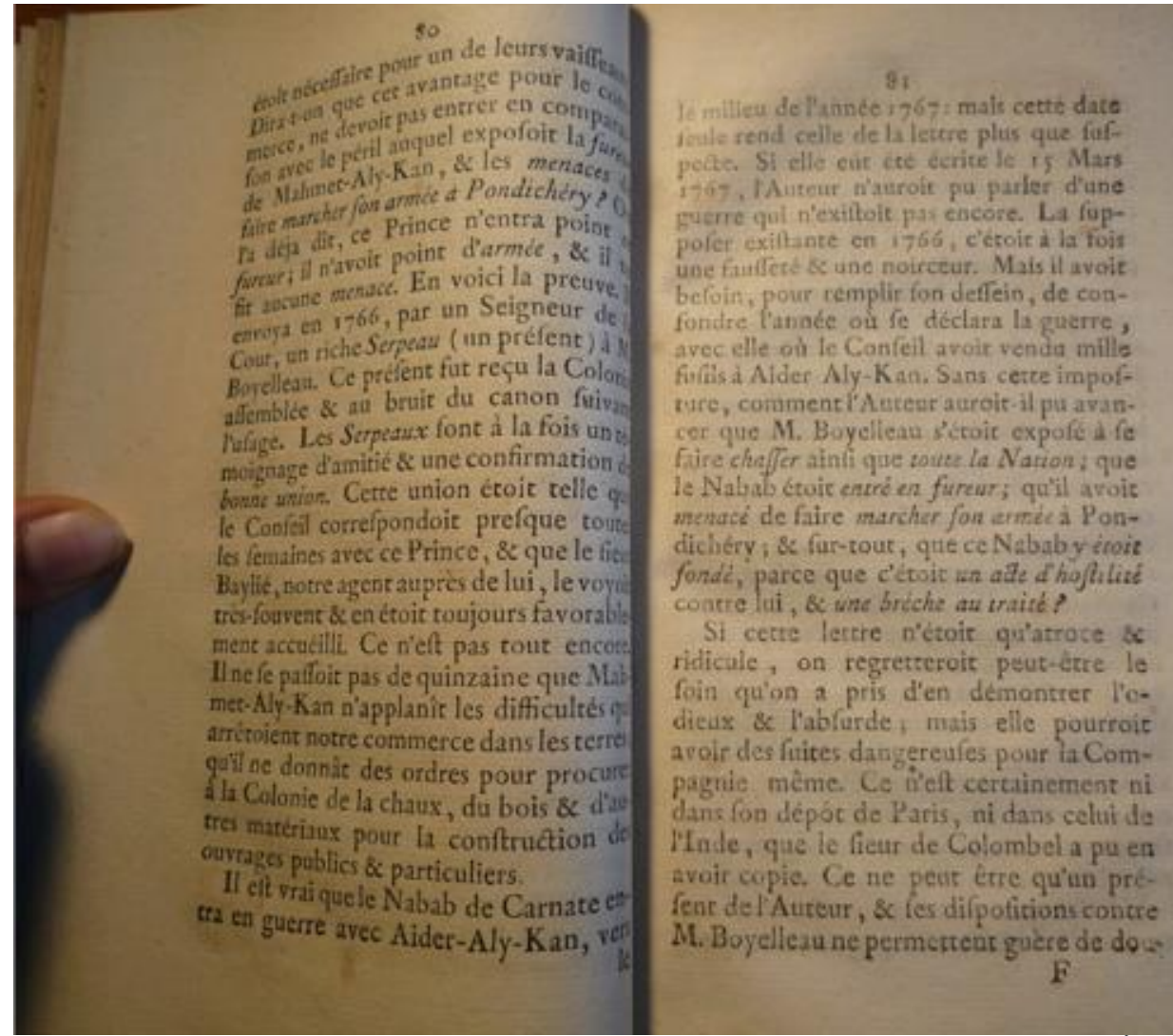
# B01: Welches Recht gilt für wen?

- Doppelseitige Farbfotografie
- Ausschnitt von Bild zu Bild variierend
- Text ist stark verzerrt



# B01: Welches Recht gilt für wen?

- Doppelseitige Farbfotografie
- Ausschnitt von Bild zu Bild variierend
- Text ist stark verzerrt



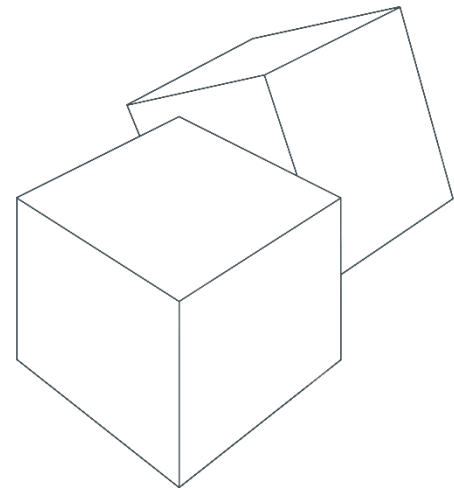
# Ergebnisse: OCR

## OCR-Pipeline: Kombination aus ocropy und tesseract

- ocropy: Binärisierung der Eingabebilder
- tesseract: Texterkennung

## Texterkennungsmodelle:

- Fraktur kein Problem
- Englisches Sprachmodell besser als deutsches (Type: Antiqua)
- Französisch problematisch
- Handschriften nicht untersucht





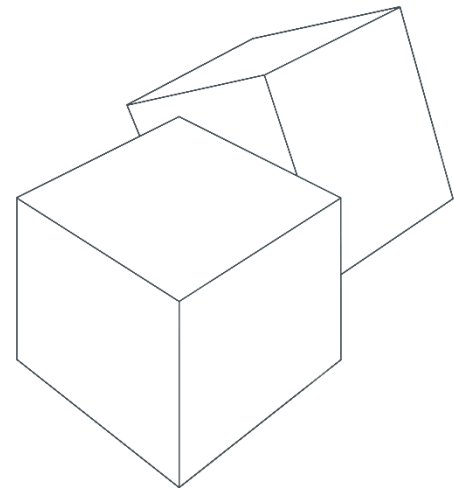
# Ergebnisse: Postkorrektur und Markup

## PoCoTo fehlerbehaftet

- Support wurde eingestellt
- Ersatz wird entwickelt (webbasiert)
- Ausgabeformate fehlerhaft

## Steile Lernkurve für Markup

- Im weiteren Verlauf sehr flüssige Umsetzung
- Tabellen und Listen weiterhin aufwendig umzusetzen



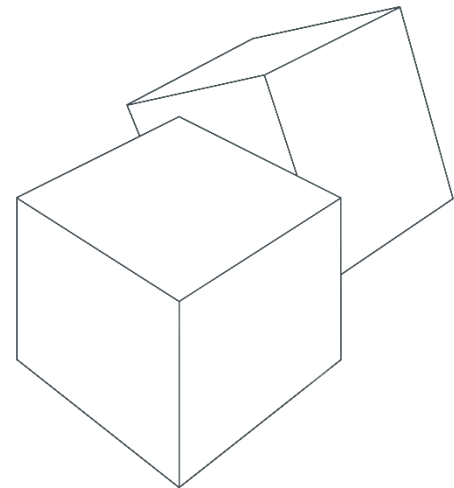
# Fazit: Eingabeformate

Bildformate: TIFF, PNG, JPEG (möglichst verlustfrei komprimiert)  
oder  
PDF mit Bildern

Auflösung nicht aussagekräftig, sondern Punktdichte

- 300ppi
- Große Zeichen (wenig Inhalt pro Seite): 150ppi
- Viel Bildinhalt (A04): 450ppi/ 600ppi

Möglichst gleichbleibender Bildaufbau ohne große Verzerrungen



# Fazit: Workflow

## Postkorrektur

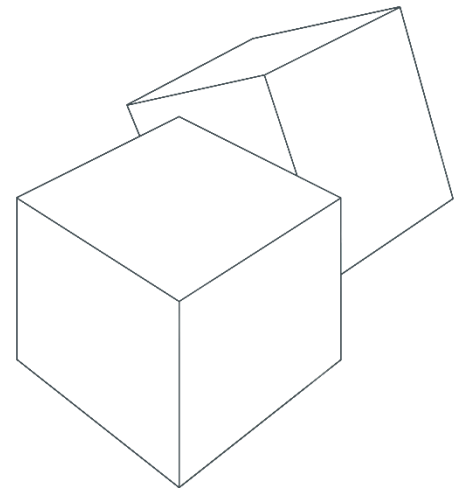
- Bei Fehlertoleranz: Unnötig für rein quantitative Analyse
- Direkt im reinen vorausgezeichneten Textformat

## Notwendigkeit für Markup je nach Anwendungsfall

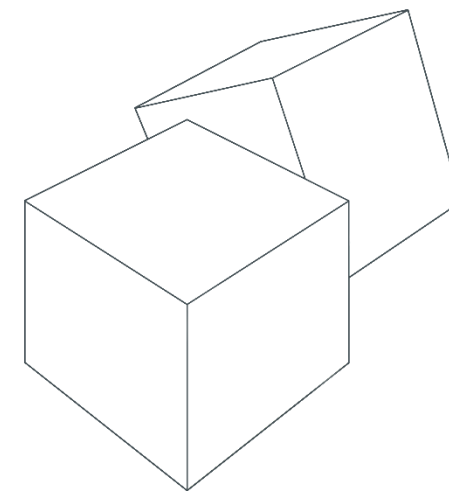
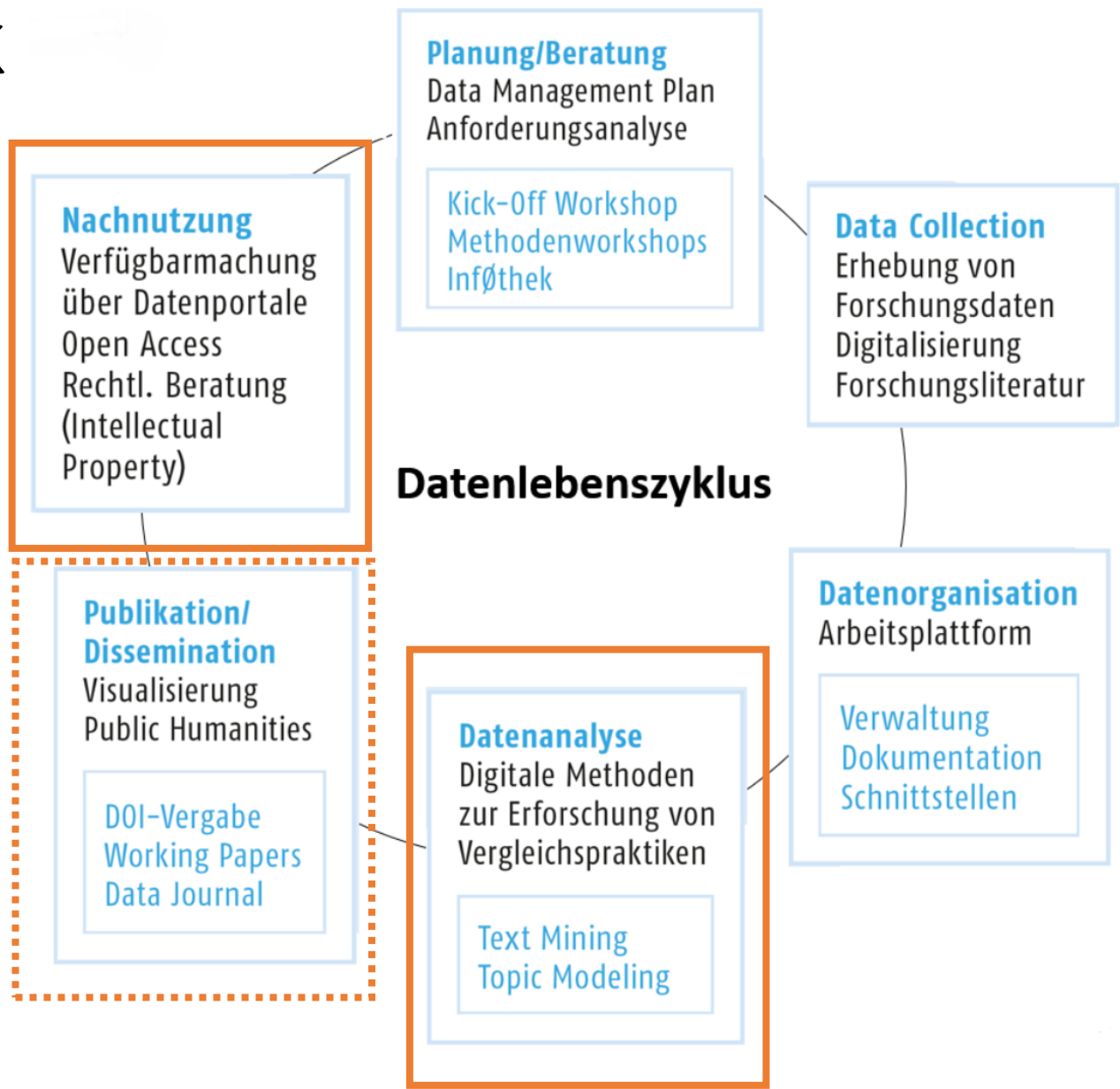
- Korrektur (vorausgezeichnetes Textformat) und Markup in einem Arbeitsschritt möglich
- Kosten-Nutzen-Frage

## Arbeitsaufwand (Normseite):

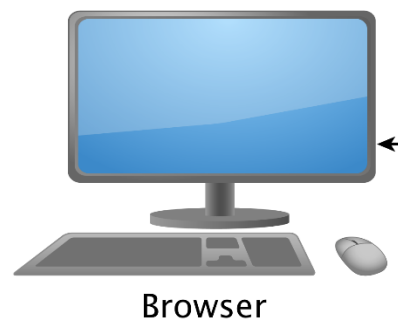
- Postkorrektur: 10-15 min
- Markup: 5 min



# Ausblick

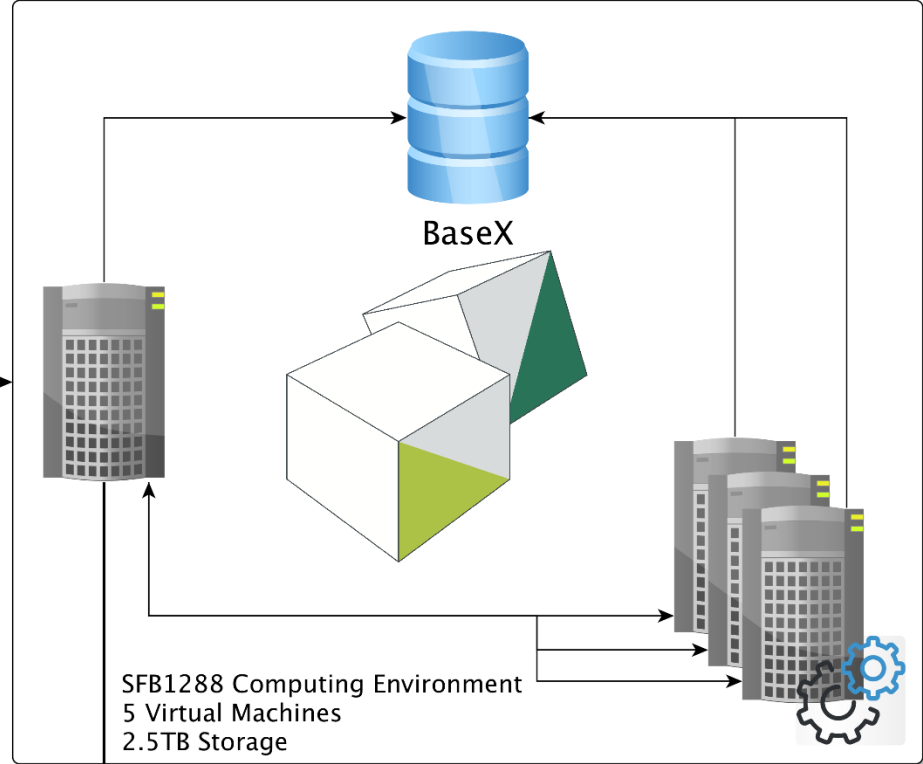


# Ausblick

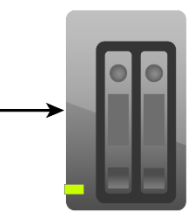


HTTPS

HTTPS

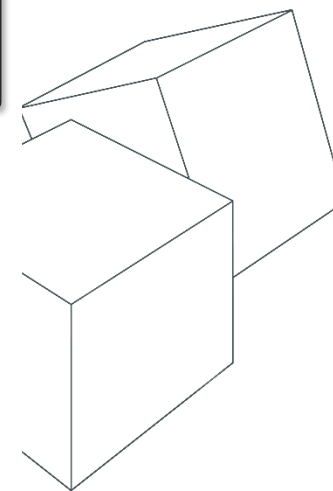


GIT via pack



Versionierung

- Dienstplattform für browserbasierte DH-Tools
- OCR-Pipeline für Digitalisate erste Anwendung auf Dienstplattform



**VIELEN DANK!**

