

Digitalisierung und Erschließung historischer Schulbücher mit Schwerpunkt religiöse Bildung – zwischen Masse und gezielten manuellen Eingriffen

Christian Reul

Zentrum für Philologie und Digitalität (ZPD)
Universität Würzburg



26.06.2024



Agenda

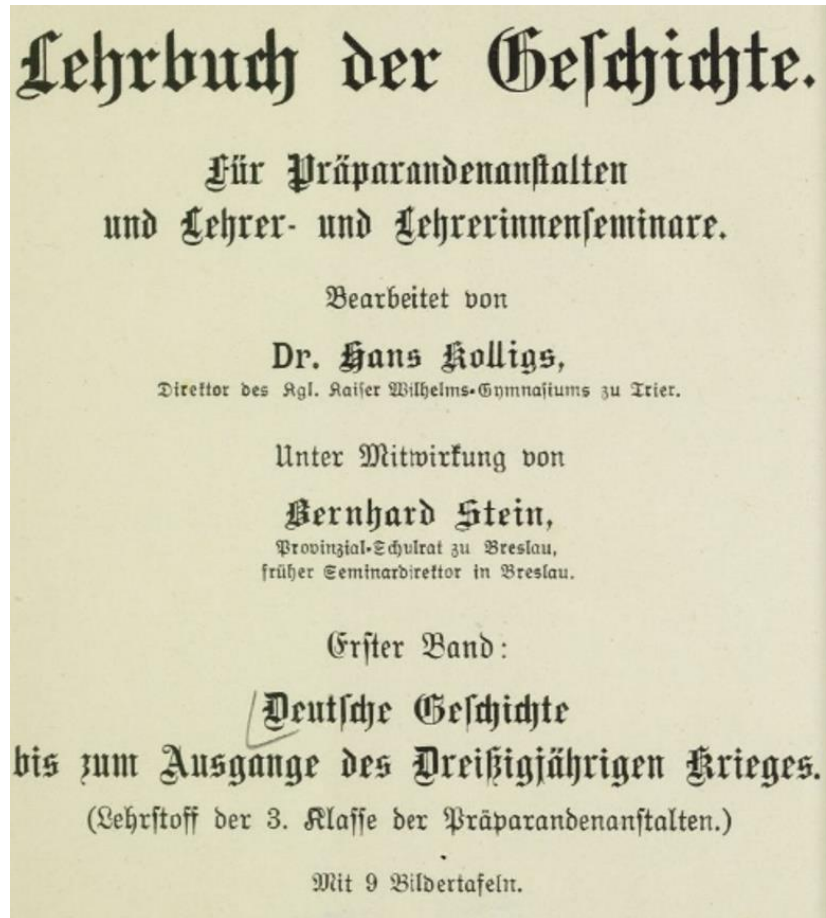
1. Allgemeine Projektinfos
2. Genereller (OCR) Workflow
3. Layoutanalyse
4. OCR-Qualitätssicherung und -Optimierung

Allgemeine Projektinfos

Digitalisierung und Erschließung historischer Schulbücher mit Schwerpunkt religiöse Bildung und Aufbau einer Wissensbasis für die bildungshistorische Forschung

- [DFG-Projekt](#), Förderlinie „Digitalisierung und Erschließung“
- Kooperation mit Dr. Anke Hertling (Leibniz-Institut für Bildungsmedien | GEI)
- Laufzeit drei Jahre, Start 07/24
- Ziele:
 - Erfassung von deutschspr. Religionsschulbücher und schulischen (Erst-)Lesebüchern
 - Knapp 2.000 Bände mit fast 500.000 Seiten
 - Zeitraum 1618 bis 1870
 - Systematische Erfassung von Akteuren aus der historischen Schulbuchproduktion unter Einsatz von Named Entity Recognition (NER) und Named Entity Linking (NEL)

Beispiel NER-Ergebnis auf fehlerhafter OCR



Lehrbuch der Geschichte MISC .
Für Präparandenanstalten
und Lehrer- und Lehrerinnenseminare.
Bearbeitet von
Dr. Hans Kolligs PER .
Direktor des Kal ORG . Kaiser Wilhelms-Gymnasiums ORG zu Trier LOC .
Unter Mitwirkung von
Vernhard Stein PER .
Provinzial. Schulrat zu Breslau LOC ,
früher Seminarlehrer in Breslau LOC .
Erster Band:
Deutsche Geschichte
bis zum Ausgange des Dreißigjährigen Krieges MISC .
(Lehrstoff der 3. Klasse der Präparandenanstalten.)
Mit 9 Bildertafeln.

Vorrede.

fer, oder einer Geschichte, verschiedene Dinge zu merken, so sind zwar solche in der Ordnung nach einander hergesetzt; Dem Lernenden aber kann die Sache erleichtert werden, wann man durch kleinere Fragen Dieselbe annoch zergliedert, und die Antwort solcher gestalt von ihm heraus locket; 3. E. p.60. werden Caligula's Laster und sonderlich seine Verschwendung beschrieben, weil diese nun, nach einander zu erzehlen, einem Anfänger in der Historie schwer fallen möchten, so kann ad §. 1. gefraget werden: Wie viel Geld hat Caligula durchgebracht? ad §. 2. Wie ging er mit seinen Schwestern um? ad §. 3. & 4. Wie tractirte er sein Pferd? und so an andern Orten mehr.

Was die gemengte oder umgekehrte Fragen pag. 22, 29, 36, 45, 2c. denen Lernenden für einen Vorteil und Impression geben, wird man bei dem ersten Anblick derselben alsobald urtheilen,

Thineus XXX.

22 MONARCHIÆ PRIMÆ
8. Aeneas succediret seinem Schwieger-Vater Latino im Regiment.

9. Er ist der erste König der Latiner von denen/die nach der Zerstörung Trojæ in Italien kommen.

10. Von dem Todt Aeneas bis auff die Erbauung der Stadt Rom werden 426. Jahr gezehlet.

CENTURIÆ TERTIÆ
Decas I.

1. Thineus regiret 30. Jahr / von An. M. 2831.

2. Die Labe des Bundes wird von den Philistern weggenommen. 1. Sam. 4. v. 11.

3. Samuel succediret dem Eli An. M. 2850. und ist so wohl ein Prophet als ein Politischer Richter gewest.

4. Dem Ascanio Aeneas Sohn folget im Regiment nach Sylvius, von welchem hernach die Successores Sylvii genennet werden / weil er in einem Walde gebohren.

5. Die Heraclidæ kommen in den Pelopponesum, und nehmen Spartam, Argos und Messenen ein/welches geschehen 80. Jahr nach der Verwüstung Trojæ.

6. Dercylus regiret 40. Jahr von A. M. 2861.

Dercylus XXXI.

7. In

DEUTSCHE KAISER.
1. Maximilian I. Kaiser...
2. Maximilian II. Kaiser...
3. Rudolph II. Kaiser...
4. Matthias Kaiser...
5. Ferdinand I. Kaiser...
6. Maximilian II. Kaiser...
7. Rudolph II. Kaiser...
8. Matthias Kaiser...
9. Ferdinand I. Kaiser...
10. Maximilian II. Kaiser...

DES DEUTSCHEN REICHS STATEN.
I. OBEROESTERREICHISCHE KREIS.
1. Fürst von Carinthien...
2. Fürst von Tyrol...
3. Fürst von Steyermark...
II. NIEDEROESTERREICHISCHE KREIS.
1. Fürst von Österreich...
2. Fürst von Kärnten...
3. Fürst von Friaul...
III. BÖHMISCHE KREIS.
1. Kaiser...
2. Fürst von Böhmen...
3. Fürst von Mähren...

IV. NIEDERRHODISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...
V. OBERRHODISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...
VI. SCHWABISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...
VII. SAARLÄNDISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...
VIII. WESTFÄLISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...
IX. OBERWESTFÄLISCHE KREIS.
1. Fürst von Brandenburg...
2. Fürst von Preußen...
3. Fürst von Ansbach...

Anmerkungen.

Die Statthalter dieser Kreise sind...
1. In der Rheinischen...
2. In der Westfälischen...
3. In der Ober- und Nieder-Rheinischen...

Genereller (OCR) Workflow

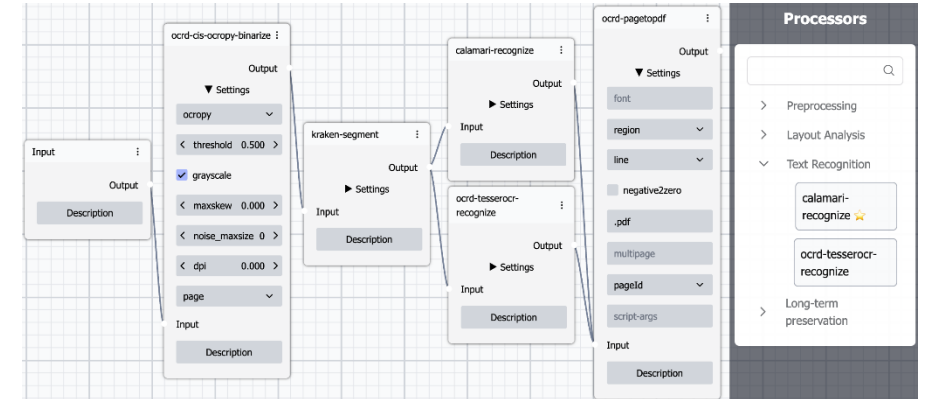
- Bearbeitung mittels OCR4all und LAREX
- Binarisierung wohl mit ocrd-sbb-binarize Gibt's was neueres/besseres?
- Layoutanalyse → nächste Folie
- Erkennung voraussichtlich mit Calamari, vorab allerdings noch Tests; Neue gemischte Modelle (hoffentlich) bis dahin fertig
- OCR-Qualitätssicherung und -Optimierung → überübernächste Folie

- NER mit (angepassten) Standardlösungen
- Unterstütztes (GND-API + eigene GUI) NEL

Tester gesucht!

Layoutanalyse

- Ursprünglicher Plan: OCR-D Workflow Suche
 - Erstellung via NodeFlow
 - Ergebnisabgleich via LAREX
 - Gezielter Einsatz via OCR4all-Tagging-Funktionalität



- Aber: lieber direkt voll auf trainierbare Methoden setzen
 - Besser und nachhaltiger
 - Paralleles [Projekt zur Layoutanalyse](#) (mit Stabi Berlin und SLUB Dresden) → Synergieeffekte
 - Fokus auf Kraken und (evtl. später) Eynollah
 - gezielt Trainingsdaten erstellen, Active Learning

Gerne mitmachen! (mit eigenen Daten)

OCR-Qualitätssicherung und -Optimierung

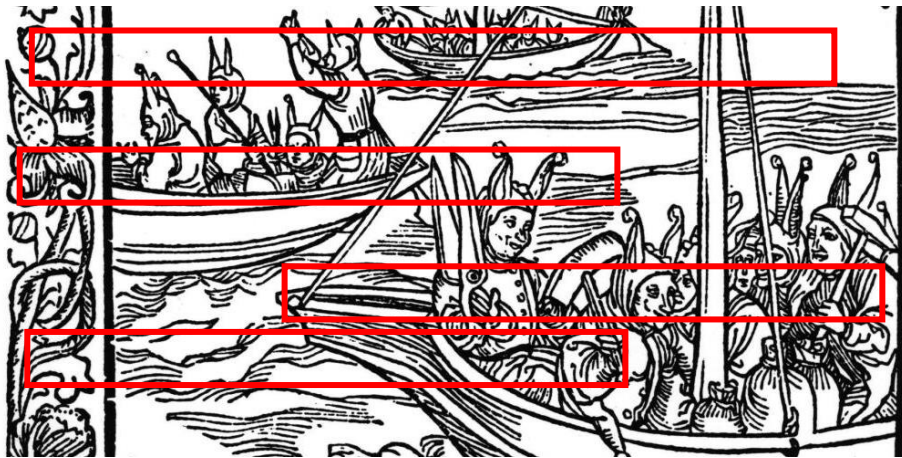
- Eingehende manuelle Prüfung aller Werke/Seiten nicht gangbar
- Stattdessen gezielte Prüfung, basierend auf Qualitätsabschätzung
 - Auf Werk-/Seitenebene → jeweils sortiert nach (vermuteter) „Qualität“
 - Flexibel und experimentell!
- Art des Eingriffs dann je nachdem
 - Anpassung von Parametern, Anwendung eines anderen Modells, ...
 - Erstellung von Trainingsdaten und gezieltes Nachtraining (auf Werk-/Clusterebene)
 - ...?
- Frage: Wie Qualität abschätzen?
 - Bereits gute Erfahrungen mit OCR-Konfidenz, allerdings auch klare Schwächen
 - Nur pragmatische (und idealerweise direkt einsetzbare) Lösungen sinnvoll

Erfahrungen? Ideen?

Stärken/Schwächen OCR-Konfidenz

Text Text Text Text Text

teilsegmentierte Zeile, Zeilen in Bild
→ OCR-Konfidenz klarer Indikator



Lorem ipsum dolor sit gubergren, no sea
amet, consetetur takimata sanctus est

sadipscing elit, sed
diam nonumy eirmod
tempor invidunt ut
labore et dolore magna
aliquyam erat, sed diam
voluptua. At vero eos et
accusam et justo duo
dolores et ea rebum.
Stet clita kasd

Lorem ipsum dolor sit
amet. Lorem ipsum dolor
sit amet, consetetur
sadipscing elit, sed diam
nonumy eirmod tempor
invidunt ut labore et
dolore magna aliquyam
erat, sed diam

Übersehene Spaltentrennung
→ OCR-Konfidenz ggf. top,
Ergebnis jedoch Quatsch