

DIE VIELFÄLTIGEN HERAUSFORDERUNGEN FÜR AUTOMATISIERTE TEXTERKENNUNG AUF DOKUMENTEN DES HERDER-INSTITUTS

- Ziel: Automatisierte Erstellung von Volltexten zu Beständen des Herder-Instituts
- ... als Grundlage für weitere NLP- und DH-Verfahren (NER, Topic Modelling etc.)
- Bestände sind...
 - sehr heterogen (zeitlich, Medientypen, Inhalte, ...)
 - vergleichsweise kleinteilig -> Arbeitsaufwand GT-Erstellung etc. für ‚maßgeschneiderte‘ Modelle hoch in Relation zum Output
 - nur in Teilen und nicht systematisch digitalisiert
- Geringe Personalausstattung

WAS BISHER GESCHAH...

- Auswahl von zwei exemplarischen Dokumenten für HTR/OCR:
 - Briefwechsel Julie und Eduard von Oettingen, DSHI 190 Livland 33, 11c, ~3000 Seiten (Handschrift 2. Hälfte 19. Jhd., siehe Abb. 1)
 - Georg v. Krusenstjern: Einsatz. Tagebuchblätter, Briefe und Notizen aus dem zweiten Weltkrieg 1941-1945, DSHI 190 Krusenstjern 2002, ~400 Seiten (Typoskript Mitte 20. Jhd., siehe Abb. 2).
- Erstellung von Evaluationskorpora (50seitige Samples, händisch transkribiert)
- Erprobung von off-the-shelf HTR/OCR-Modellen u.a. mit Transkribus, ocr4all, recogAlze, nopaque und Abby FineReader
- Training/Finetuning von Modellen zur Layouterkennung für Handschrift in Transkribus
- Training/Finetuning von OCR-Modellen für Typoskript in eScriptorium

Abb. 1: Brief Eduard v. O. an Julie v. O., 4. September 1863

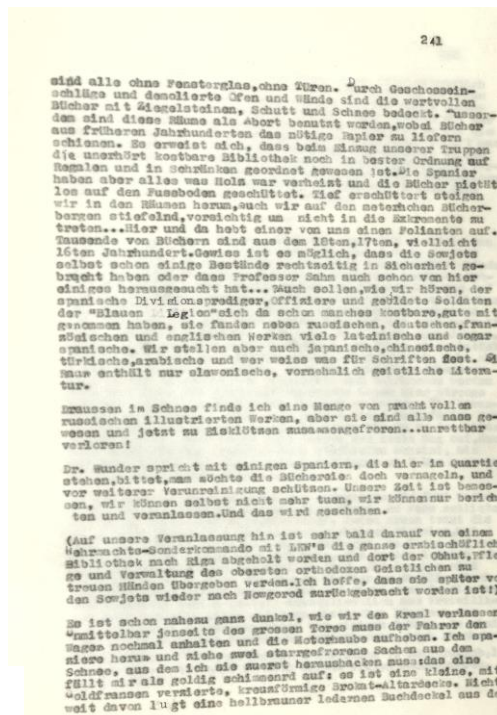
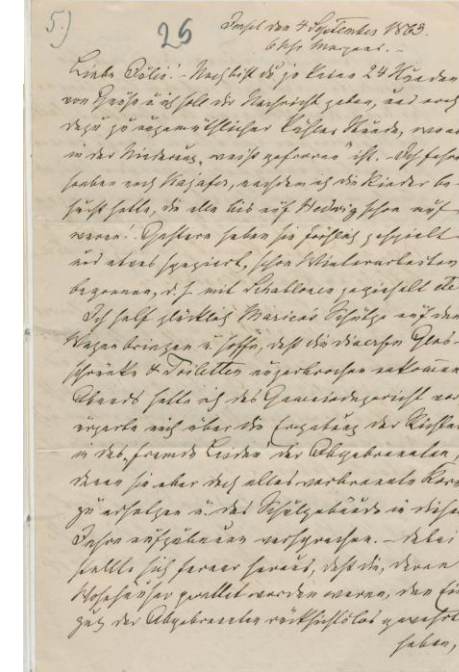


Abb. 2: Krusenstjern, Einsatz, S. 241

ERSTE ERGEBNISSE UND FRAGEN

- Erste Ergebnisse:
 - off-the-shelf Modelle mit sehr unterschiedlichen Ergebnissen

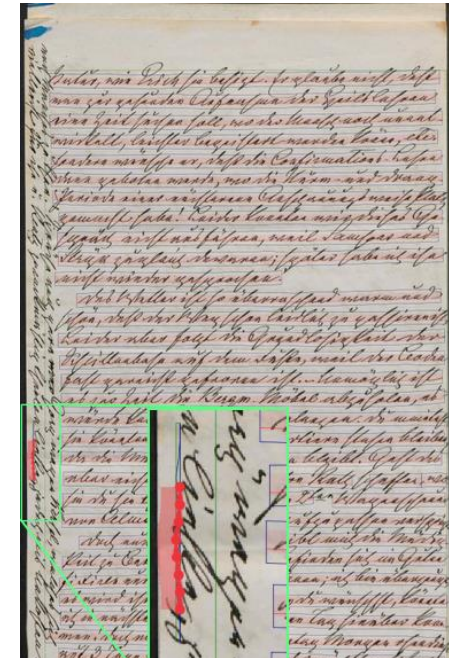
	Transkribus	eScriptorium	Abby	Tesseract	Nopaque
CER	1.96%	9.79%	17.91%	17.51%	23.07%
WER	14.99%	44.73%	49.54%	50.57%	59.83%

Ergebnisse Evaluation OCR auf 25 Seiten-Sample Bericht Krusenstjern, jeweils bestes Model

- Ergebnisse für Typoskript mit minimalen Trainingsdaten (10 Seiten) deutlich optimierbar (mittels eScriptorium: CER **2.58%**, WER: **14.14%**)
- Herausforderungen: Layout-Erkennung, v. a. für Handschriften; Evaluierung von Layout Recognition
- Für ‚kleine‘ Sprachen fehlen tw. Modelle v. a. für historische Sprach- und Schriftstufen (u. a. Polnisch)
- Bildqualität vermutlich ein Faktor, aber schwer zu beziffern

Fragen

- Layout-Erkennung ein Problem, v. a. bei Handschriften, aber auch bei schlechten Scans. Wie hiermit umgehen? Frage der Trainingsdatenmenge oder grundsätzlicheres Problem?
- Wie Layout-Erkennung evaluieren?
- Trainings-Daten-Erstellung: Standards für Transkriptionen (und Layouts) um Kompatibilität mit bestehenden Trainingsdaten zu gewährleisten?
- Ressourcen für osteuropäische Sprachen (Modelle, Trainingsdaten, ...)?



Screenshot Transkribus, Layout Analysis mit selbst trainiertem Model, aufgenommen 20.5.2024