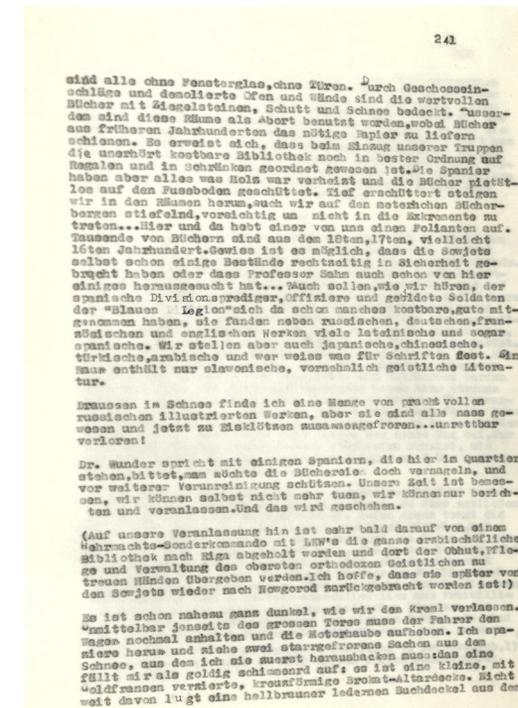
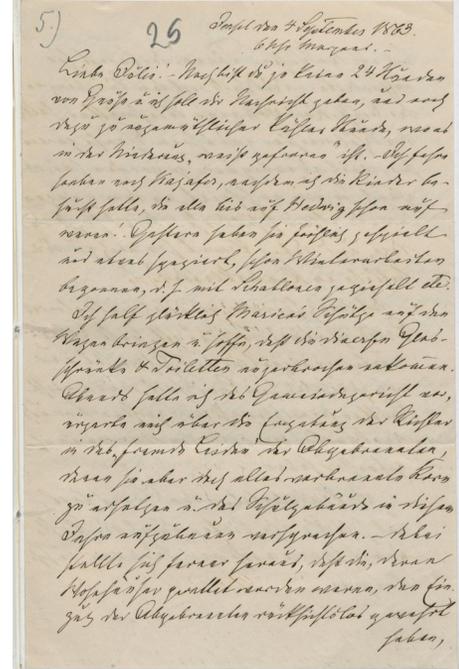


# **More Data - More Problems? Die vielfältigen Herausforderungen für Automatisierte Texterkennung auf Dokumenten (nicht nur) des Herder-Instituts**

Ole Meiners  
Arbeitsbereich Forschungsdatenmanagement  
Abteilung Digitale Geschichte und Informationssysteme  
Herder-Institut für historische Ostmitteleuropaforschung

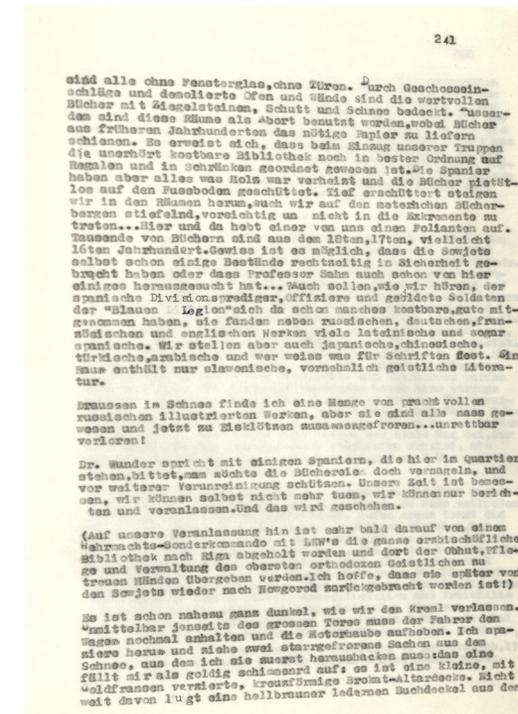
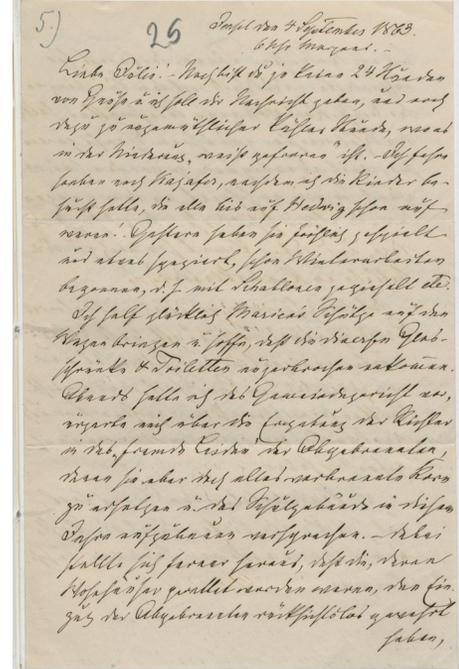
# DIE VIELFÄLTIGEN HERAUSFORDERUNGEN FÜR AUTOMATISIERTE TEXTERKENNUNG AUF DOKUMENTEN DES HERDER-INSTITUTS

- 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 -> hoher Arbeitsaufwand GT-Erstellung etc. für ‚maßgeschneiderte‘ Modelle in Relation zum Output
  - nur in Teilen und nicht systematisch digitalisiert



# DIE VIELFÄLTIGEN HERAUSFORDERUNGEN FÜR AUTOMATISIERTE TEXTERKENNUNG AUF DOKUMENTEN (NICHT NUR) DES HERDER-INSTITUTS

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



# 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, eScriptorium, ocr4all, recognAlze, nopaque, **PERO** und Abbyy FineReader
- (Nach-)Training/Finetuning von Modellen zur Layouterkennung für Handschrift in Transkribus
- (Nach-)Training von OCR-Modellen für Typoskript über eScriptorium

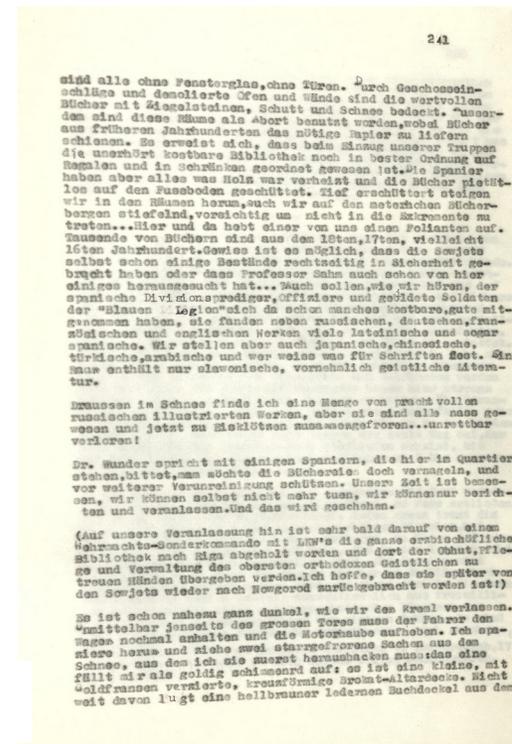
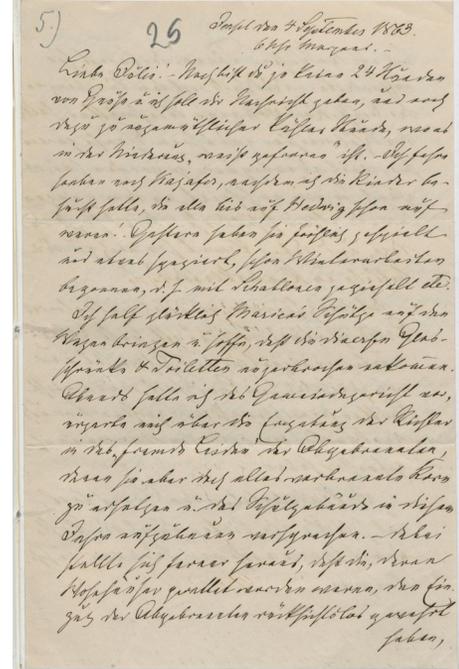
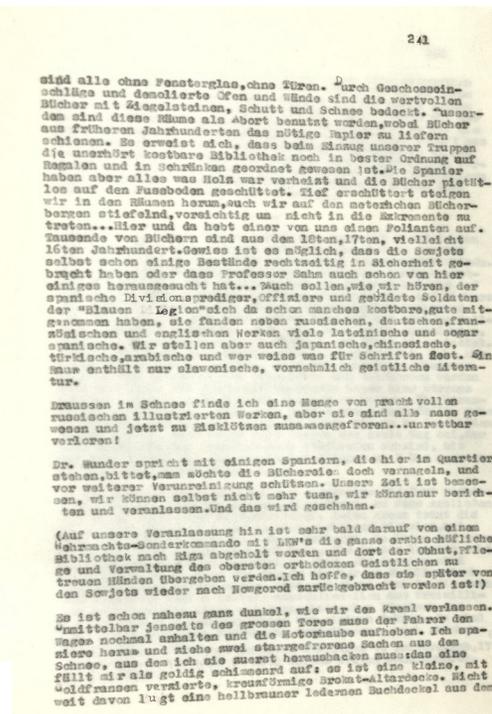
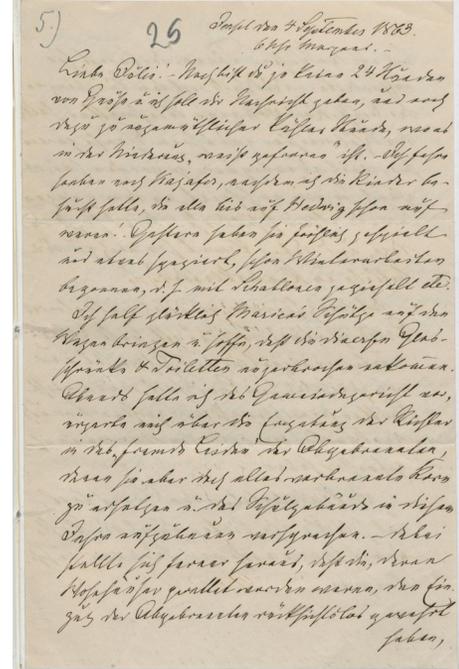


Abb. 2: Krusenstjern, Einsatz, S. 241

# WAS SEITDEM GESCHAH...

- Texterkennung für Typoskript mit nachtrainiertem Model 'einsatzfähig' (CER: 1,82% WER: 8,87% auf Evaluationsdatensatz)
  - ... Rechtslage leider unklar daher (noch) nicht publikationsfähig
- HTR auf Briefkorporus wegen Problemen bei Layout-Erkennung zunächst auf Eis gelegt (multi-direction layout recognition bislang nicht umsetzbar)



# WAS SEITDEM GESCHAH...

- Texterkennung für Typoskript mit nachtrainiertem Model ,einsatzfähig' (CER: 1,82% WER: 8,87% auf Evaluationsdatensatz)
  - ... Rechtslage leider unklar daher (noch) nicht publikationsfähig
- HTR auf Briefkorpus wegen Problemen bei Layout-Erkennung zunächst auf Eis gelegt (multi-direction layout recognition bislang nicht umsetzbar)
- Drittes Einsatzszenario: (Vertriebenen-)Zeitungen
  - u. a. „Ostpreußenblatt“ (1950-2003), „Preußische Allgemeine Zeitung“ (2003-2023)
  - „Sudetendeutsche Zeitung“ (derzeit: 1951-1955)
  - ... und viele andere mehr (noch nicht digitalisiert)
  - Aktuell vorliegend: > 3850 Ausgaben, > 83.000 Seiten

Abb. 1: Titelseite „Das Ostpreußenblatt“, Organ der Landsmannschaft Ostpreußen, 16.01.1960



Abb. 2: Krusenstjern, Einsatz, S. 241

sind alle ohne Fensterglas, ohne Türen. Durch Gasseweinschläne und zerlegte Öfen und Hünde sind die wertvollen Bücher mit Schmutz, Schmutz und Schnee bedeckt. Unserer aus früheren Jahrhunderten das nötige Papier zu liefern die unerschöpfliche bester Biblisch noch in bester Ordnung auf Regalen und in schranken geordnet gewissen Intellektuellen haben aber alles was Holz war verheißt und die Bücher platzen auf den Fußboden geschüttet. Auf erschütterter steigen wir in den Räumen herum, such wir auf den zerbrochenen Bücherbergen stiefeln, vorsichtig um nicht in die Scherente zu treten... Hier und da hebt einer von uns einen Folianten auf. Tausende von Büchern sind aus dem 16ten, 17ten, vielleicht 18ten Jahrhundert... Gewiss ist es möglich, dass die Sowjets selbst schon einige Bestände rechtsseitig in Sicherheit gebracht haben oder dass Professor Saha auch schon von hier ein neues herausgebracht hat... Auch sollen wir hier hören, der spanische Intellektuelle, Cristóbal und geistige Soldaten der „blauen Legion“ sich da schon manches kostbare, gute mitgenommen haben, sie fanden neben russischen, deutschen, französischen und englischen Werken viele lateinische und sogar spanische. Wir stellen aber auch japanische, chinesische, türkische, arabische und was weises was für Schriften fest. Ein paar enthält nur slavonische, vornehmlich geistliche Literatur.

Draußen in Schnee finde ich eine Menge von prachtvollen russischen illustrierten Werken, aber sie sind alle nass geworden und jetzt zu Nichts als massenweise... unrettbar verloren!

Dr. Wunder spricht mit einigen Spaniern, die hier in Quartier stehen, hätte, man sollte die Bücherwerke doch vernageln, und vor weiterer Verunstaltung schützen. Unsere Zeit ist besonnen, wir können selbst nicht mehr tun, wir können nur berichten und vernachlässigen. Und das wird geschehen.

(Auf unsere Veranlassung hin ist sehr bald darauf von einem Gubernats-Besonderkommando mit Hilfe die ganze evangelische Bibliothek nach Hinz abgeholt worden und unter der Obhut, die treuen Händen übergeben werden. Ich hoffe, dass sie später von den Sowjets wieder nach Nowgorod zurückgebracht werden ist!)

Es ist schon nahezu ganz dunkel, wie wir den Kessel verlassen. Und toller jenseit des großen Torps muss der Fahrer den Wagen nochmal abhalten und die Notwendigkeit machen ich zu stehen herzu und nicht zu auszuhalten herzuhaben muss: das eine Schnee, aus dem ich sie zuerst herausziehen muss: das eine fällt mir als goldig schimmernd auf: es ist ein kleines, mit Goldfäden verziertes, kranzförmiges Broschett. Nicht weit davon liegt eine hellbraune Ledernen Buchdeckel aus dem

# ZEITUNGEN

- Layout wie erwartet schlecht zu erkennen
- ... beste Ergebnisse bislang mit ABBYY FineReader und PERO OCR
- ... auch hier allerdings Reading Order mitunter stark fehlerhaft

Text- und Layout Erkennung Abbyy FineReader, Sudetendeutsche Zeitung, Jg. 1 Folge 01, 07.04.1951



Text- und Layout Erkennung PERO OCR, Sudetendeutsche Zeitung, Jg. 1 Folge 01, 07.04.1951

# ZEITUNGEN

- Layout wie erwartet schlecht zu erkennen
- ... beste Ergebnisse bislang mit ABBYY FineReader und PERO OCR
- ... auch hier allerdings Reading Order mitunter stark fehlerhaft
- Herausforderungen:
  - Erstellung Ground Truth sehr arbeitsaufwändig
  - Layout-Erkennungsmodelle von ABBYY und PERO ermöglichen kein Finetuning/Adaption?
  - Gleichzeitig hohe Anforderungen an Erkennungsqualität (v.a. Text, mittelbar Layout) durch Forschungsprojekt

Text- und Layout Erkennung Abbyy  
FineReader, Sudetendeutsche  
Zeitung, Jg. 1 Folge 01, 07.04.1951



Text- und Layout Erkennung  
PERO OCR, Sudetendeutsche  
Zeitung, Jg. 1 Folge 01,  
07.04.1951

# FRAGEN

- Problem Layout-Erkennung; bei multi-direktionalen (Hand-)schriften, v. a. aber bei komplexen Zeitungs-Layouts.
- (Wie Layout-Erkennung evaluieren?)
- Tools für Layout-GT-Erstellung (eScriptorium zu umständlich für komplexe Layouts; LAREX? Weitere?)
- Publikation Evaluations-Datensätze? Verwendung als Trainingsmaterial nicht intendiert
- ... und was ist überhaupt realistisch zu erreichen/wo lohnt der eigene Arbeitsaufwand? → v. a. für Layout-Erkennung

**Vielen Dank für die  
Aufmerksamkeit!**

Ole Meiners  
Arbeitsbereich Forschungsdatenmanagement  
Abteilung Digitale Geschichte und Informationssysteme  
Herder-Institut für historische Ostmitteleuropaforschung