

Von „gemischt“ zu „werksspezifisch“ Modelltraining für die Texterkennung von historischen Drucken und Handschriften

Christian Reul

Zentrum für Philologie und Digitalität „Kallimachos“ (ZPD)
Universität Würzburg



03.02.2022



Gliederung

1. Motivation und Grundlagen

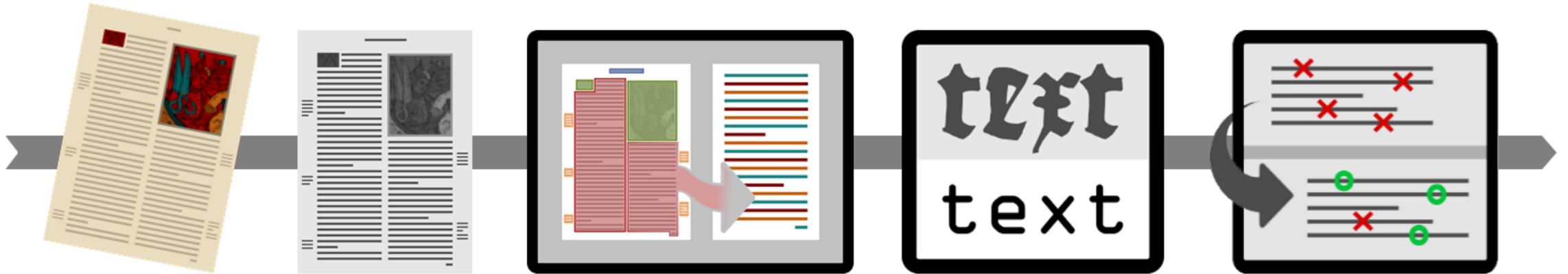
2. Methoden

3. Anwendungsbeispiel Drucke

4. Anwendungsbeispiel Handschriften

5. Diskussion und Ausblick

Grundlagen – Workflow und Modelle



- Hauptkomponenten: Vorverarbeitung, Segmentierung, OCR, Nachkorrektur
- Moderne OCR-Ansätze arbeiten auf Zeilen- und nicht mehr auf Zeichenbasis
- Sog. Modelle extrahieren Text aus Textzeilenbildern
 - Modelle müssen trainiert werden (Neuronale Netze)
 - Trainingsdaten bestehen aus Zeilenbild-Text-Paaren

Er wird eifrig gefammelt.
Er wird eifrig gefammelt.

Grundlagen – Gemischte und werksspezifische Modelle

- Gemischte Modelle:
 - Im Normalfall auf einer Vielzahl von Quellen trainiert
 - Vorteil: Out-of-the-Box Anwendung, kein werksspezifisches Training
 - Nachteil: ungenauer als werksspezifische Modelle
- Werksspezifische Modelle:
 - Im Normalfall exklusiv für die Erkennung genau einer Quelle trainiert
 - Vorteil: meist deutlich genauer als gemischte Modelle
 - Nachteil: benötigt Ground Truth (GT), also Zeilen und deren Transkriptionen
- Kombination: werksspezifisches Training ausgehend von gemischten Modellen („Finetuning“)

Motivation

- Unterschiedliche Strategien für verschiedene Anwendungsszenarien und Bedarfe
 - Massenvolltexterkennung: Out-of-the-Box Anwendung gemischter Modelle
 - Textproduktion für digitale Edition: Umfangreiches werksspezifisches Training, um bestmögliche Erkennung als Grundlage für die manuelle Nachkorrektur zu erzeugen
 - Fließender Übergang
- Idee/Ziel: Breit aufgestellte Modelle trainieren, die ...
 - out-of-the-box auf eine Vielzahl (Alter, Sprache, Schriftart) von Materialien angewendet werden können
 - als Ausgangspunkt für weiteres Finetuning dienen können
- Pragmatischer Ansatz!

Gliederung

1. Motivation und Grundlagen

2. Methoden

3. Anwendungsbeispiel Drucke

4. Anwendungsbeispiel Handschriften

5. Diskussion und Ausblick

Ausbalancieren der Daten

Problem:

- Daten sehr unausgeglichen: Anzahl GT Zeilen pro Quelle schwankt zwischen <50 und 10k+
- Modelle optimieren sich in Richtung der im Korpus dominanten Daten

Lösung:

- Ausbalancieren durch Definition ausgewählter Seiten (etwa 50-150 Zeilen pro Quelle, je nach Aufkommen im Korpus)
- Zweistufiger Ansatz:
 1. Training auf allen verfügbaren Daten
 2. „Refinement“ des resultierenden Modells nur auf den ausgewählten Seiten

Nutzung unterschiedlicher Binarisierungen

Problem:

- Art der Binarisierung kann Performanz gemischter Modelle erheblich beeinflussen
- Keine fixe Universallösung, sondern Einsatz bei den Endnutzern abhängig ...
 - vom Material
 - deren persönlicher Präferenz

Lösung:

- Training der gemischten Modelle auf unterschiedlichen Binarisierungen
 - Größere Robustheit gewährleistet breitere Anwendung
 - Künstlich aufgeblähte Trainingsmenge erhöht Performanz

CLARISSIMO

AC PRUDENTIS-

fimo Senatui ciuitatis Noriber-
gg, Dominis suis colendissimis,

Sebaldus Heyden

S. D,



Voties de nostrarum
Scholarum disciplina
cogito, Viri Clarissi-
mi, nonnihil admirari
soleo, cur studij literarij
loca olim Græcis ab o-
tio, Scholæ: Latinis ue-

ro Ludij dicta sint. Nam non ausim puta-
re illos artium ingenuarum. ac Philoso-
phiae studia, ita uilia ac leuia habuisse, ut
ea nō nisi per otium, & tanq̃ res ludicras
sectanda esse censerent. Presertim quum
ob ea ipsa studia, cæteras omnes gentes,
barbaras præ se uocarent, ac tanq̃ omnis
humanitatis expertes despicerent. Cu-
ius elationis causas certe non otiosas,

A 2 multo

CLARISSIMO

AC PRUDENTIS-

fimo Senatui ciuitatis Noriber-
gg, Dominis suis colendissimis,

Sebaldus Heyden

S. D,



Voties de nostrarum
Scholarum disciplina
cogito, Viri Clarissi-
mi, nonnihil admirari
soleo, cur studij literarij
loca olim Græcis ab o-
tio, Scholæ: Latinis ue-

ro Ludij dicta sint. Nam non ausim puta-
re illos artium ingenuarum. ac Philoso-
phiae studia, ita uilia ac leuia habuisse, ut
ea nō nisi per otium, & tanq̃ res ludicras
sectanda esse censerent. Presertim quum
ob ea ipsa studia, cæteras omnes gentes,
barbaras præ se uocarent, ac tanq̃ omnis
humanitatis expertes despicerent. Cu-
ius elationis causas certe non otiosas,

A 2 multo

CLARISSIMO

AC PRUDENTIS-

fimo Senatui ciuitatis Noriber-
gg, Dominis suis colendissimis,

Sebaldus Heyden

S. D,



Voties de nostrarum
Scholarum disciplina
cogito, Viri Clarissi-
mi, nonnihil admirari
soleo, cur studij literarij
loca olim Græcis ab o-
tio, Scholæ: Latinis ue-

ro Ludij dicta sint. Nam non ausim puta-
re illos artium ingenuarum. ac Philoso-
phiae studia, ita uilia ac leuia habuisse, ut
ea nō nisi per otium, & tanq̃ res ludicras
sectanda esse censerent. Presertim quum
ob ea ipsa studia, cæteras omnes gentes,
barbaras præ se uocarent, ac tanq̃ omnis
humanitatis expertes despicerent. Cu-
ius elationis causas certe non otiosas,

A 2 multo

Normalisiertes Graustufenbild

Binarisierung 1

Binarisierung 2

Gliederung

1. Motivation und Grundlagen
2. Methoden
- 3. Anwendungsbeispiel Drucke**
4. Anwendungsbeispiel Handschriften
5. Diskussion und Ausblick

Daten

- Selbst erstellt oder aus frei verfügbaren Quellen zusammengesucht
- Training
 - 21.5k Seiten aus 642 Werken
 - ca. 450 Jahre Druckgeschichte (ca. 1450 - 1900) → große Variabilität an Drucktypen
 - Verschiedene Sprachen: Deutsch, Latein, Französisch, Niederländisch
- Evaluation
 - 29 Werke zwischen 1506 und 1849
 - Repräsentative Abdeckung hinsichtlich Alter, Druckqualität, Schriftarten, etc.
 - Auswahl, Segmentierung und Transkription von vier bis fünf Seiten pro Buch

Dem hoeghen marschalck sent quirtjn
dienen. Da was der Keyser gar fro dz sein syn
en vnd materien erlengert vnd sch
Sonder beghinsel Hofmæct in Vreughden.
Stultoz ꝛc. *Qm̄ in quā stultoz vt
complexus fit: ut sub æquatore est*
Und angehängter besondern Tabell

Kallimachos Subkorpus, größtenteils
Narrenschiffe (15./16. Jh.)

*des siècles. Ces couronnes que j'ai si
contrasté dans les physionomies de ceux
monstrant l'intégrité de mon ame, vous
savant transporté de joie, que nous
les genitoires tant feconds, & les Vto-
Va donc pour Lifette; je n'en serai pas*

Französisch (17.-19. Jh.)

bat ihn wegen seines eigenen Glütes, sich nicht in
alle Zweige der gemischten Mathematik, wie
wins ragt Sigurds Schreckhorn hell hinaus
Borsaal, und ihre schönen Töne zogen alle
nur eine Kluft, die uns von dem gewünschten
So träumte er; ein Schwarm aufflatternder Eulen
in unsern Tagen von hundert Personen neun und
und entdeckte zuletzt einen schmalen hölzernen Steg;
lassen, daß das schwächliche Kind des reichen Hart-
er dadurch verrathen würde, daß er wirklich eine Neigung
Geschäftig eilten die Diener herbei. Der Herzog und
mich in Harnisch. Es handelt sich also nicht um einen
sind der Erste, dem ich von meinem Leben erzähle, damit

Frakturschriften 19. Jh.

**cto procederet, omnium, quorum et
nem notare, videtur initium illius iubere observari,
in parte, mentio. Differui autem ista
illius fixæ temperatū quiddam ualent ē
ferre possim gratiæ, pro summis & immortalibus**

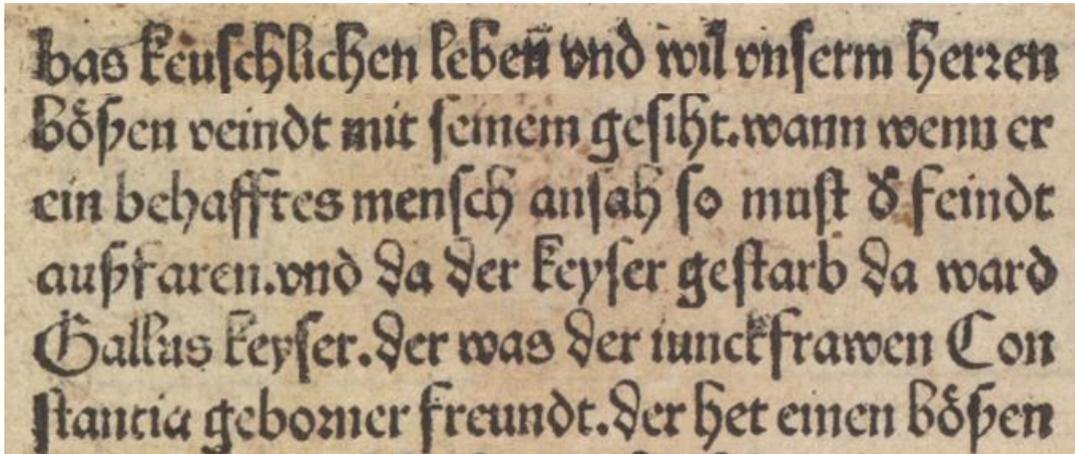
Camerarius (16. Jh.)

Evaluation

- Vergleich von Stufe 1 (alle Daten), Stufe 2 (nur ausgewählte Daten), deren Kombination sowie dem weit verbreitetem GT4HistOCR Standardmodell
- Bestimmung der Zeichenfehlerrate (Character Error Rate, CER)
- Stufe 1 und 2 vergleichbar (Balance > Masse), 1→2 am besten
- Deutlich (ca. 40%) besser als GT4HistOCR
- Mit Abstand größte verbleibende Fehlerquelle (ca. 50%): Einfügungen und Löschungen von Leerzeichen
- Durch technische Optimierungen (Datenaugmentierung, tiefere Neuronale Netze etc.) weitere Senkung auf ca. 1,6%

alle Daten	ausgew. Daten	CER in %
x		1,92
	x	1,95
x	x	1,80
GT4HistOCR		2,84

Beispielerggebnis mit ca. 1,6% CER



bas keuschlichen leben vnd wil vnserm herren
bößen veindt mit seinem gesiht. wann wenn er
ein behafftes mensch anseh so must d̄ feindt
außfaren. vnd da der keyser gestarb da ward
Gallus keyser. der was der iunckfrawen Con
stantia geborner freundt. der het einen bößen

bas ksuschlichen lcben vnd wil vnserm herren
bößen veindt mit seinem gefiht. wann wenn er
ein behafftes mensch anfah so must d̄ feindt
außtaren. vnd da der keyser gestarb da ward
tGallus keyser. der was der iunckfrawen Con
stantia geboꝛner freundt. der het einen bößen

Gliederung

1. Motivation und Grundlagen
2. Methoden
3. Anwendungsbeispiel Drucke
- 4. Anwendungsbeispiel Handschriften**
5. Diskussion und Ausblick

Material: Mittelalterliche Handschriften

- [Kooperation](#) mit Dr. Stefan Tomasek (Lehrstuhl für deutsche Philologie, ältere Abteilung) anhand des Projekts „Konrad von Fußesbrunnen: Kindheit Jesu“
- Vielen Dank an die Projektgruppe und die fleißigen Hilfskräfte!
 - Dr. Stefan Tomasek, Florian Langhanki, Maximilian Wehner
 - Susanne Bremer, Lisa Gugl, Sebastian Hammer, Kiara Hart, Ursula Heß, Leonie Kampmann, Annika Müller, Anne Schmidt



Wā er erkennet ain wib
Wā er erkennet ain wib
Die nit verlagen kan Irn lib
Die nit verlagen kan Irn lib
Da ylet er vil balde hin
Da ylet er vil balde hin
Vnd wirbet waft daz ist sin sin
Vnd wirbet waft daz ist sin sin

Trainingsdaten

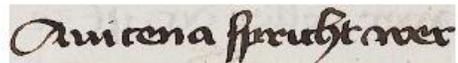
- Selbst erstellt oder aus frei verfügbaren Quellen zusammengesucht
- Frühes 13. bis spätes 15. Jh.: Kindheit Jesu, Marienleben, Parzival, ...
- 35 Werke, knapp 300 Seiten, ca. 12,5k Zeilen → zwei gemischte Modelle



dem hovbz vñ ovgen.



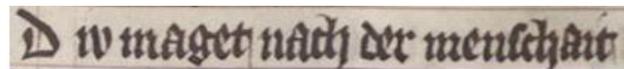
¶ Serpilleum zertribē



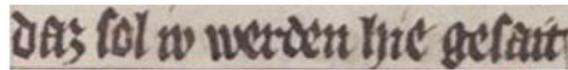
Auicena spricht wer



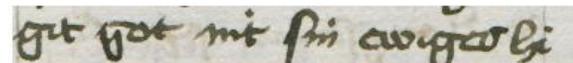
nüchter baden will



D iv maget nach der menschait



daz sol iv werden hie gefait



git got nit sin ewiges hi



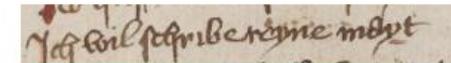
melrich Es spricht hug



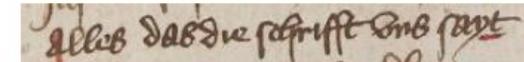
got amen Doīca t̄cia



waz ich dim. So horen



Ich wil schriben reyne mayt



Alles daz die scharfft vns sayt

Evaluationsdaten

- Vier weitere Handschriften (2x Gotische Buchschrift, 2x Bastarda)
- Kindheit Jesu, Marienleben (x2), Der Welsche Gast (x2)
- Jede Hs. noch einmal unterteilt in Trainings- (32) und Evaluationsseiten (8-18)

Daz ih des niht vol bringen chan.

Daz ih des niht vol bringen chan.

mir chom zehelfe dar an.

mir chom zehelfe dar an.

man tûn fol zu allen zÿten. vñ

man tûn fol zu allen zÿten. vñ

wâ von man nit treg fol fin vnd

wâ von man nit treg fol fin vnd

Das buch bræhte her zewege. daz sie iz alle

daz buch bræhte her zewege. daz sie iz alle

mufen lesen. die gotes kint wellent we

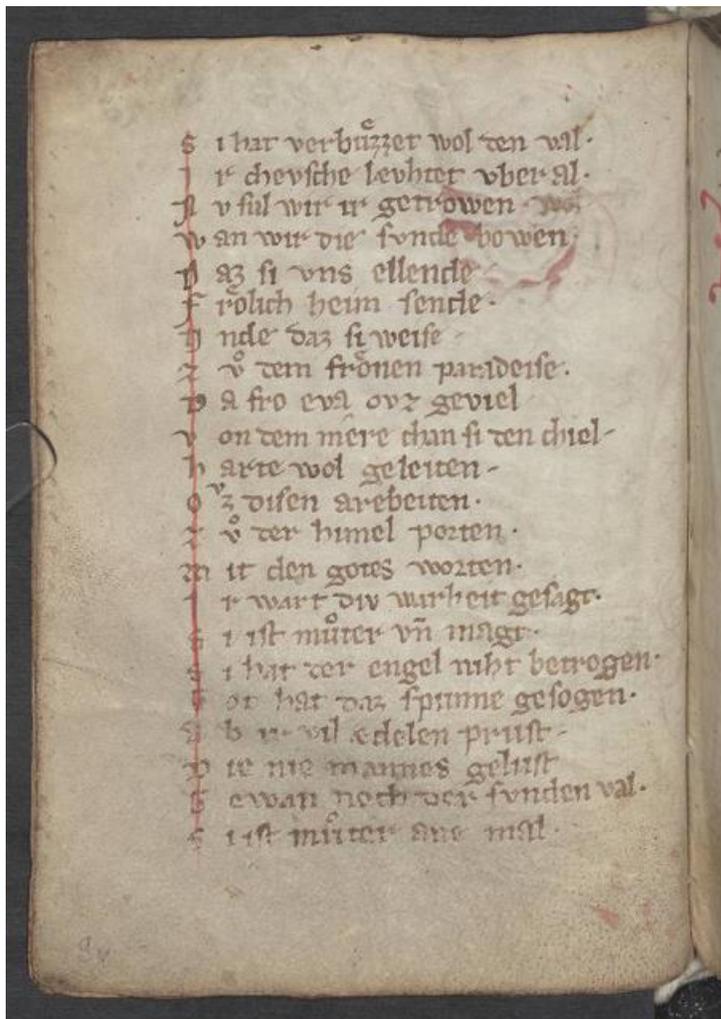
mufen lesen. die gotes kint wellent we

Das er begang mit gûter dat

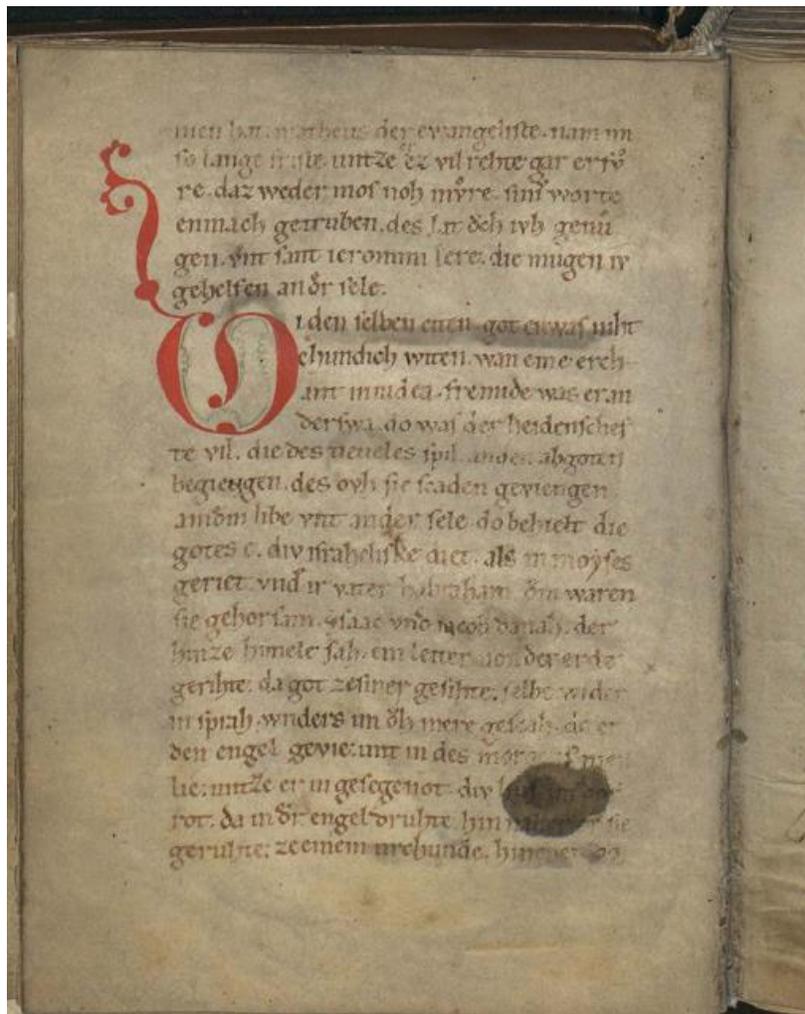
Daz er begang mit gûter dat

waz er gûtes gelesen hat

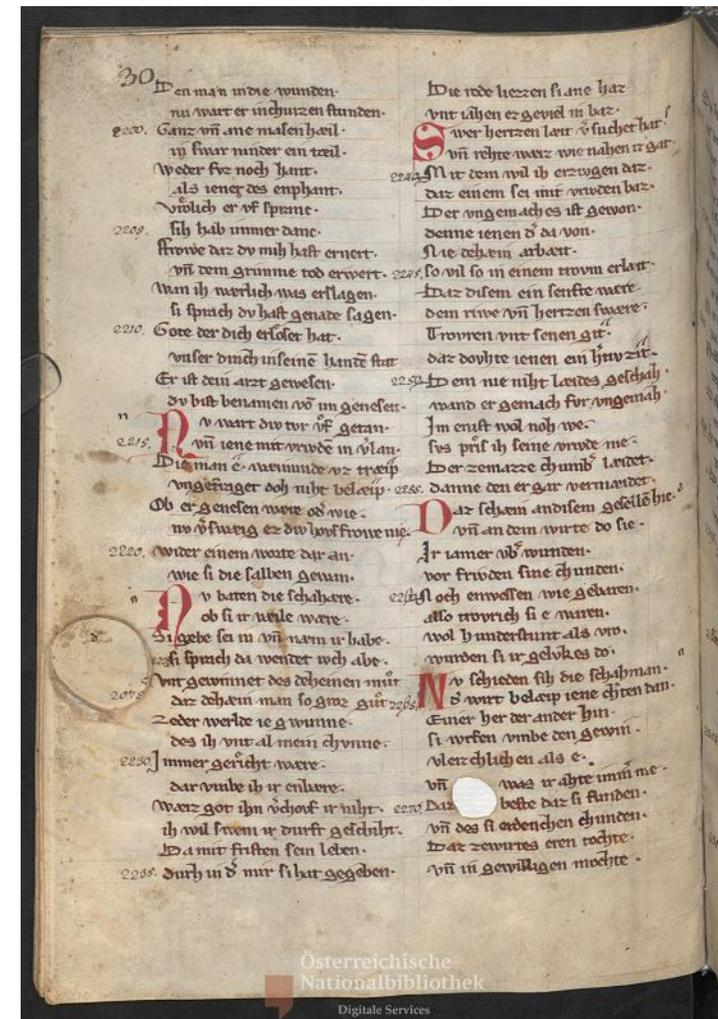
waz er gûtes gelesen hat



Wernher-Wien



Wernher-Krakau



Handschrift-B

Osterreichische
Nationalbibliothek
Digitale Services

Simulation einer schrittweisen Transkription

- Iterative Verdopplung der Trainingsmenge
 - 2, 4, 8, 16, 32 Seiten GT
 - 0 Seiten: out-of-the-box Ergebnis des gemischten Modells
- Evaluation auf fixen Evaluationsset und Berechnung der CER (%); Durchschnitt über alle Werke
- Angesichts des Materials/Modells **gute ootb CER**
- Training auf lediglich zwei Seiten **halbiert die CER**
- Weitere deutliche Steigerungen (bis deutlich unter 2% CER), jedoch abnehmende Effizienz
- Verbleibende Fehler von Leerzeichen und Punkten dominiert

# Seiten	Pretrained CER in %	Verb. in %
0 (ootb)	6,22	-
2	3,27	48
4	2,58	21
8	2,17	16
16	1,94	11
32	1,65	15

Nutzen gemischter Modelle

- Vergleich der Ergebnisse ...
 - beim Starten des Trainings From Scratch
 - bei Nutzung eines vortrainierten Modells als Ausgangspunkt (Pretrained)
- Verb. PT/FS: Verbesserungsfaktor Pretrained im Vergleich zu From Scratch
- Riesiger Unterschied, speziell bei wenig Trainingsmaterial (**85% bei 2 Seiten!**)
- Effekt nimmt mit steigender Seitenzahl stetig ab, ist jedoch **selbst bei 32 Seiten noch klar erkennbar** (knapp 40%)

# Seiten	From Scratch CER in %	Pretrained CER in %	Verb. PT/FS in %
2	21,12	3,27	85
4	10,74	2,58	76
8	5,82	2,17	63
16	3,78	1,94	49
32	2,72	1,65	39

Gliederung

1. Motivation und Grundlagen
2. Methoden
3. Anwendungsbeispiel Drucke
4. Anwendungsbeispiel Handschriften
- 5. Diskussion und Ausblick**

Zusammenfassung

- Gemischte Modelle trainiert für ...
 - Drucke in lateinischer Schrift zwischen 1450 und 1900
 - Mittelalterliche Handschriften (Gotische Buchschriften und Bastarden)
- Hochperformant ...
 - bei der out-of-the-box Anwendung
 - als Ausgangspunkt für Finetuning
- Balance > Masse → lieber viele Quellen mit jeweils wenig Zeilen, als umgekehrt
- Selbst mit wenigen Seiten große Verbesserungen durch Finetuning möglich
- Weitere Erkenntnisse zu Einfluss von Netzstrukturen und Binarisierungen sowie in-domain und out-of-domain Pretraining; für Details s.
 - Reul, Wick, Nöth, Büttner, Wehner, Springmann: [Mixed Model OCR Training on Historical Latin Script for Out-of-the-Box Recognition and Finetuning](#). HiP'21.
 - Reul, Tomasek, Langhanki, Springmann: [Open Source Handwritten Text Recognition on Medieval Manuscripts using Mixed Models and Document-Specific Finetuning](#). submitted to DAS22.

Aktuelle Arbeiten I

Et quant ce vint apres nonne que gauuāi

Et quant ce vint apres nonne que gauuāi

Et quant ce vint apres nonne que gauuāi

fu auques lasses et que le bras lycommenca

fu auques lasses et que le bras lycommenca

fu auques lasses et que le bras lycommenca

adouloir le morholt qui bien sen aparceuoit

adouloir le morholt qui bien sen aparceuoit

adouloir le morholt qui bien sen aparceuoit

Frz. Handschrift aus dem 15. Jh.

Fehlerrate: etwa 0,6%

Transkription/Korrektur: Prof. Burrichter

li dist Sire chl̄r Il est huymais tard et voʹ

li dist Sire chl̄r Il est huymais tard et voʹ

li dist Sire chl̄r Il est huymais tard et voʹ

estes lasses et trauailles et Je aussi Sy aly

estes lasses et trauailles et Je aussi Sy aly

estes lasses et trauailles et Je aussi Sy aly

vngs tant essaie laũt que bien nous deuons

vngs tant essaie laũt que bien nous deuons

vngs tant essaie laũt que bien nous deuons

en̄cognoistre Je ne ledy pas ne pō voʹ louer

en̄cognoistre Je ne ledy pas ne pō voʹ louer

en̄cognoistre Je ne ledy pas ne pō voʹ louer

Aktuelle Arbeiten II

S. D. Saepenumero mī in aīō et ante oculos est moeror

S.D. Saepenumero mī in aīō et ante oculos est moeror

S.D. Saepenumero mī in aīō et ante oculos est moeror

Illustr. principū, & vereor interdū ne magnitudo

illustr. principū, & vereor interdū ne magnitudo

illustr. principū, & vereor interdū ne magnitudo

Pr. Georgij, gravius affligat. Ego quidem huius

Pr. Georgij, gravius affligat. Ego quidem huius

pr. Georgij, gravius affligat. Ego quidem huius

Joachim Camerarius der Ältere Schmierfink

Fehlerrate: etwa 3,2%

Transkription/Korrektur: Dr. Schlegelmilch

Protectio & iam nūc poterit affiduitate precatōnis christi

Protectio & iam nūc poterit affiduitate precatōnis christi

protectio & iam nūc poterit affiduitate precatōnis christi

anae exorari. hinc igitur speremus caetera. His die

anae exorari. hinc igitur speremus caetera. His die

anae exorari. hinc igitur speremus caetera. His die

bus accepi lrās scriptas ad Illustr. Pr. Georgiū.

bus accepi lrās scriptas ad Illustr. Pr. Georgiū.

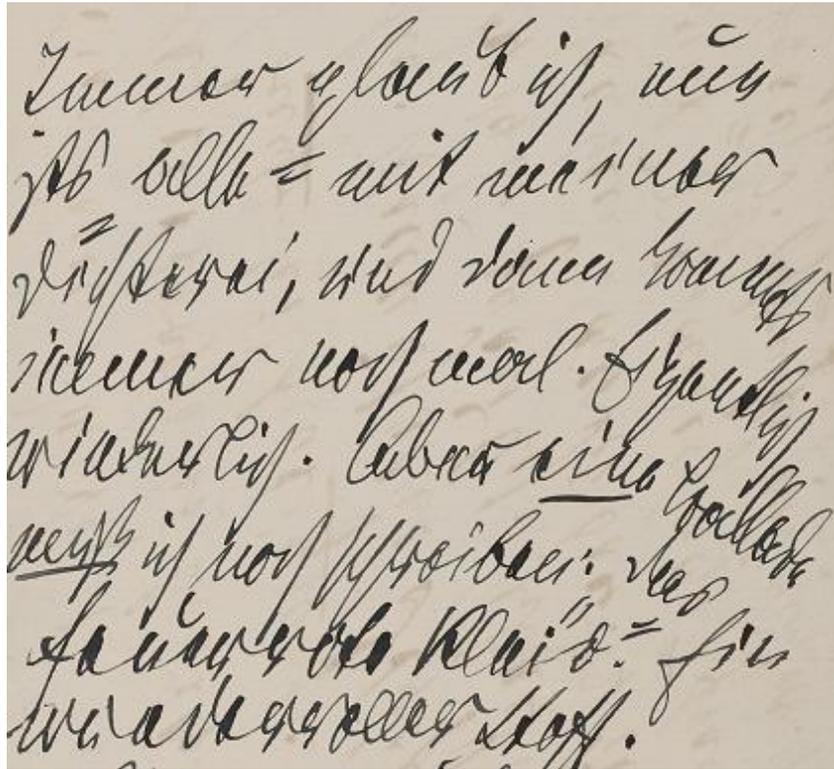
bus accepi lrās scriptas ad Illustr. Pr. Georgiū.

Cum autē nulla spes mī fuerit me brevi reperturū

Cum autē nulla spes mī fuerit me brevi reperturū

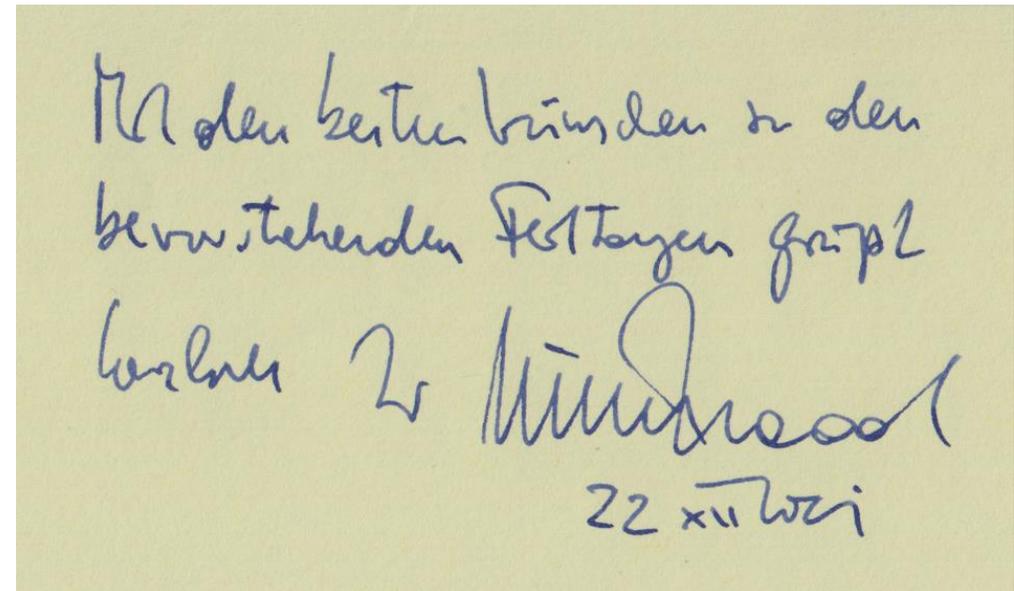
cum autē nulla spes mī fuerit me brevi reperturū

Ausblick: Schlimmer geht immer



Zuversicht glaub ich, nicht
ich solle - mit mir was
drückerai, was davon kann
sich aus was auch. Glaubt
nicht. Aber eine
weil ich was probieren: was
für ein was. Was für
was das alles soll.

Ein Herr Liliencron (?); Mitte 19. Jh. (?)
Antragskooperation mit Univ. Hamburg



Allen besten Wünschen zu den
bevorstehenden Festtagen ganz
besonders zu Weihnachten
22. X. 2021

Würzburger Musikprofessor, der anonym
bleiben möchte; Dezember 2021