

OCR4all –

Eine (semi)automatische Open Source
Software für die OCR historischer Drucke

Christian Reul

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



05.05.2021



Gliederung

1. Einleitung

2. Submodule

3. Workflow

4. Live Demo

5. Evaluation

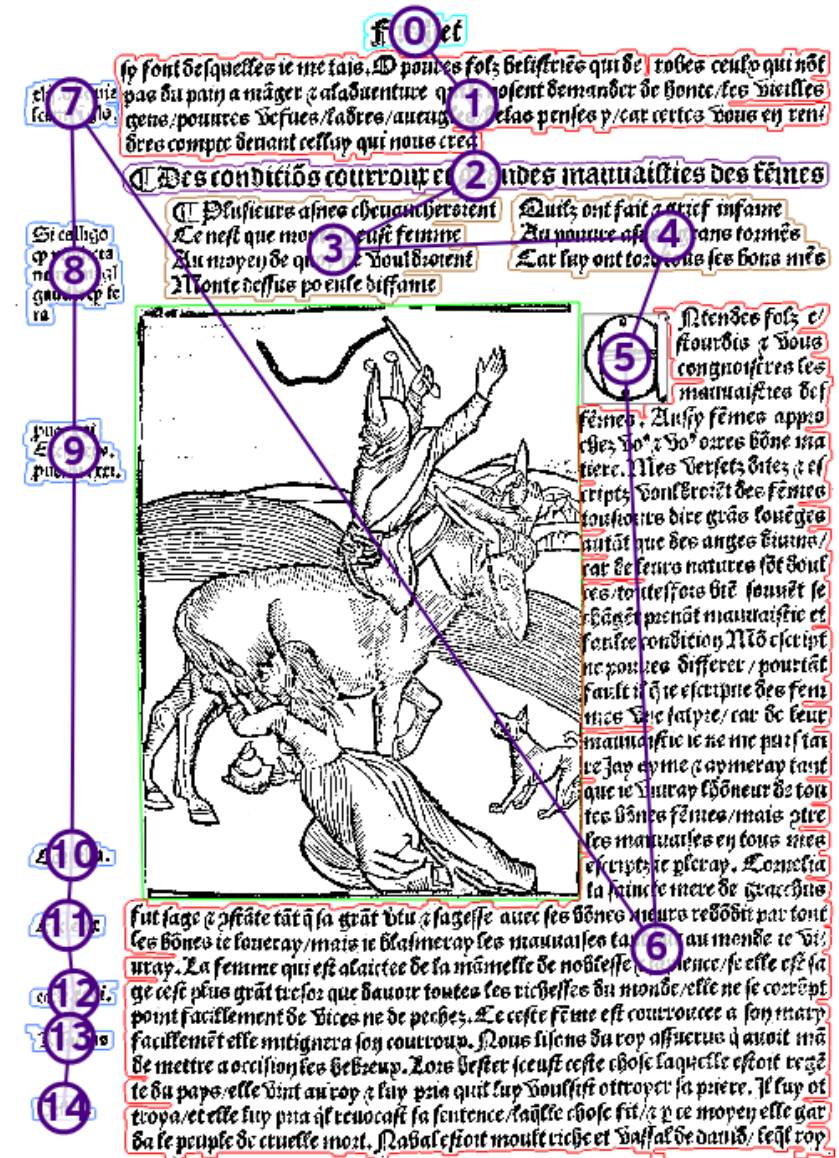
6. Diskussion und Ausblick

OCR4all – Motivation und Überblick

- Bestehende Open Source OCR Software ist mächtig, kann nicht-technische Nutzer jedoch schnell überfordern:
 - Installation nicht unbedingt trivial
 - Keine grafische Benutzeroberfläche, sondern ungewohnte Kommandozeile
 - ...
- Die Idee hinter OCR4all:
 - Verständlich und anwendbar auch für nicht-technische Nutzer
 - Nutzung von Docker/VirtualBox → Einfache Installation, Unabhängigkeit vom Betriebssystem
 - Basiert auf unterschiedlichen Open Source OCR Tools

OCR4all – Historie und Aktueller Stand

- Ursprünglich entwickelt für die OCR sehr (sehr!) alter Drucke (Projekt [Narragonien digital](#)):
 - Sehr komplizierte Layouttypisierung
 - Training werkspezifischer Modelle unerlässlich aufgrund der höchst varianten Typographie
 - Anspruch häufig: 100% Erkennungsgenauigkeit für Layout und OCR
 - Nutzer nehmen dafür einen gewissen manuellen Korrekturaufwand in Kauf
- Mittlerweile erfolgreich auf großer Bandbreite von Drucken (15. bis 21. Jh.) eingesetzt
- Hauptziel: Erhöhung des Automatisierungsgrads und der Robustheit
- Work in progress!





www.ocr4all.de

Gliederung

1. Einleitung

2. Submodule

3. Workflow

4. Live Demo

5. Evaluation

6. Diskussion und Ausblick

Aktuell integriert

- OCRopus
 - Vorverarbeitung
 - Automatische (Zeilen)Segmentierung
- Calamari
 - Texterkennung
 - Modelltraining
 - Evaluation
- LAREX
 - Ursprünglich zur interaktiven Regionensegmentierung inklusive semantischer Auszeichnung entwickelt
 - Mittlerweile umfangreiches Korrekturtool, u. a. für
 - Regionen- und Zeilenkoordinaten
 - Semantische Typisierung und Lesereihenfolge
 - Text (Ground Truth Erstellung für Training!)

Zwischenfazit und -ausblick

- Bislang begrenzte Auswahl
 - Ursprünglich für konkreten Anwendungsfall ([Narragonien digital](#)) entwickelt
 - Damaliges Ziel: erstmal funktionierenden Workflow schaffen
 - Fokus auf ausgewählte „beste“ Lösungen
- Anbieten weiterer Lösungen unbedingt intendiert
 - Problemlose Anbindung, solange klar definierte Schnittstellen eingehalten werden
 - Workflow und Lösungen für einzelne Schritte frei kombinierbar
 - Baukastenprinzip!

- Klingt nach OCR-D?!

Später mehr...

„Nimm die Binarisierung von OCROPUS, die Segmentierung von Tesseract und die Texterkennung von Calamari“

(Konstantin Baierer, vor 20min)

Gliederung

1. Einleitung

2. Submodule

3. Workflow

4. Live Demo

5. Evaluation

6. Diskussion und Ausblick

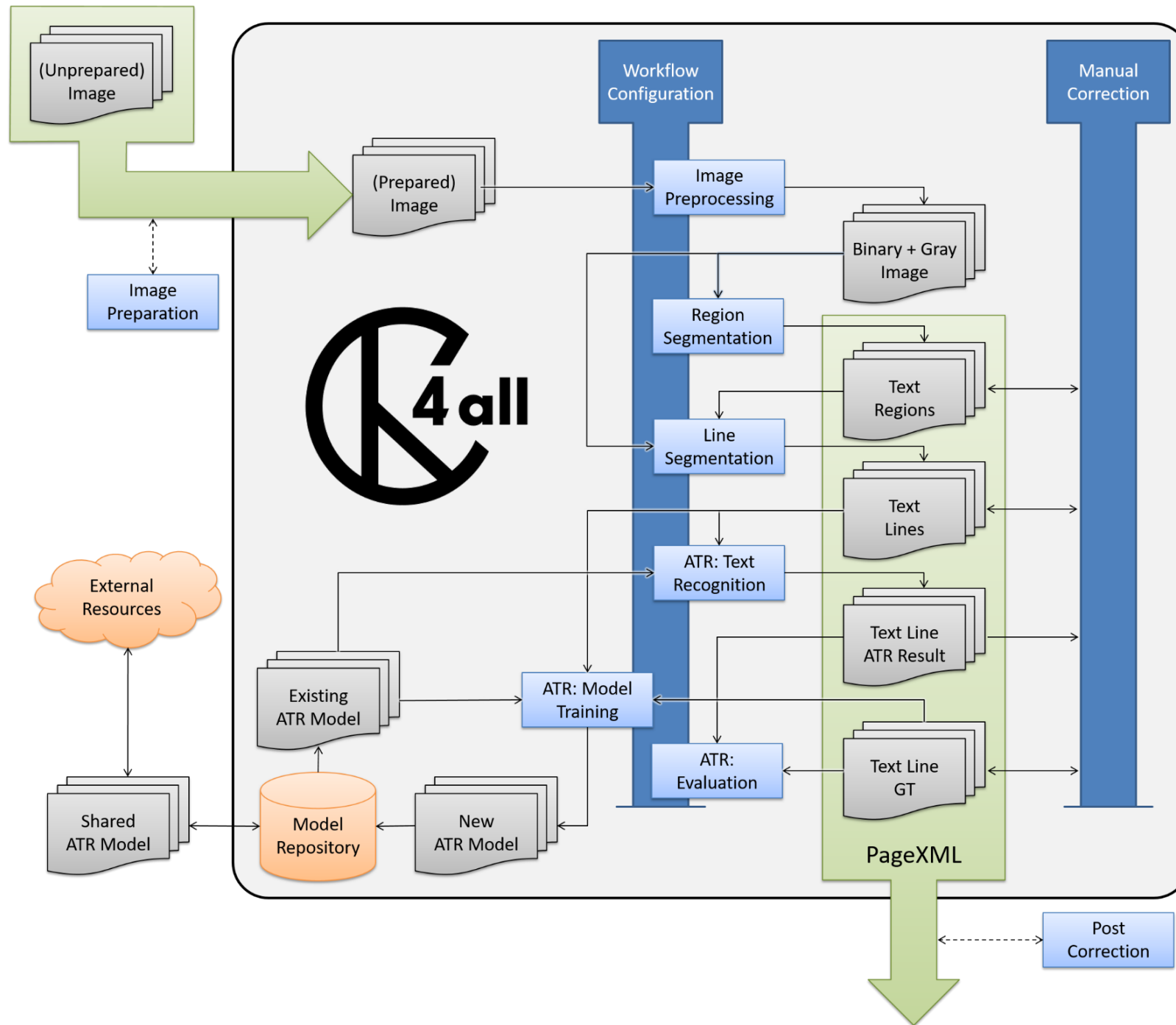


Image Preparation

- **Input:** unverarbeitete Bilder
- **Output:** vorbereitete Bilder
- Notwendige Vorbereitungen variieren von Werk zu Werk:
 - Teilung von Doppelseiten
 - Rotation
 - Entfernung von störenden Bildrändern
 - ...
- Bisher nicht in OCR4all integriert
→ Open Source Tool [ScanTailor](#)

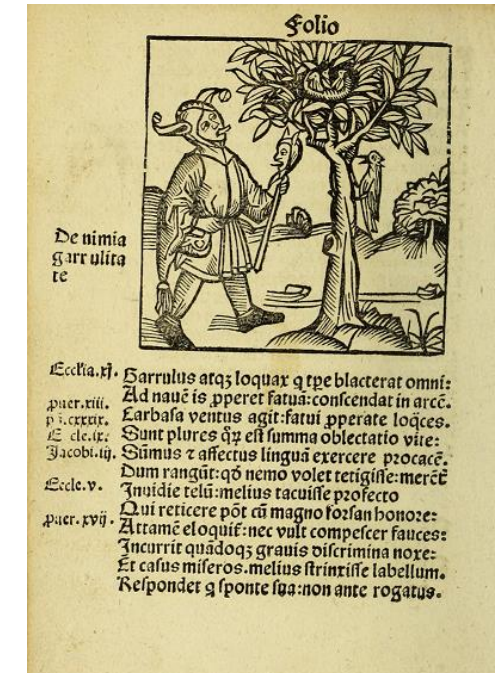
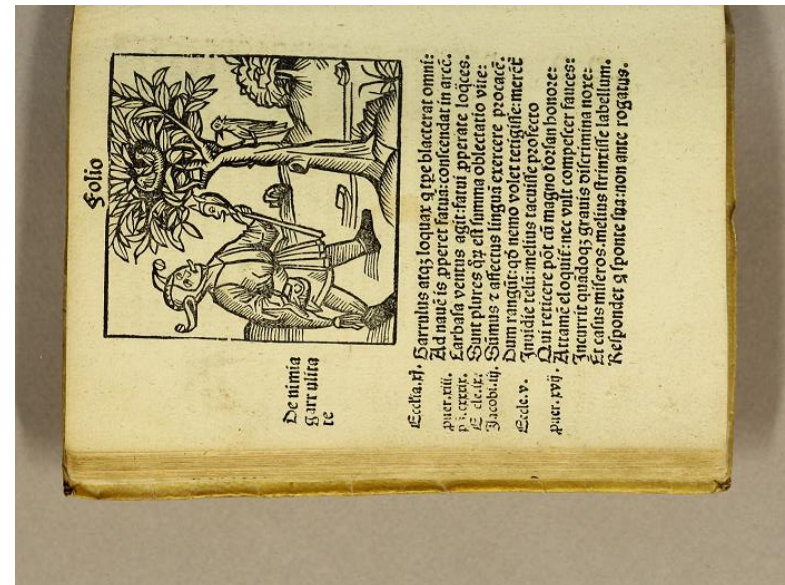
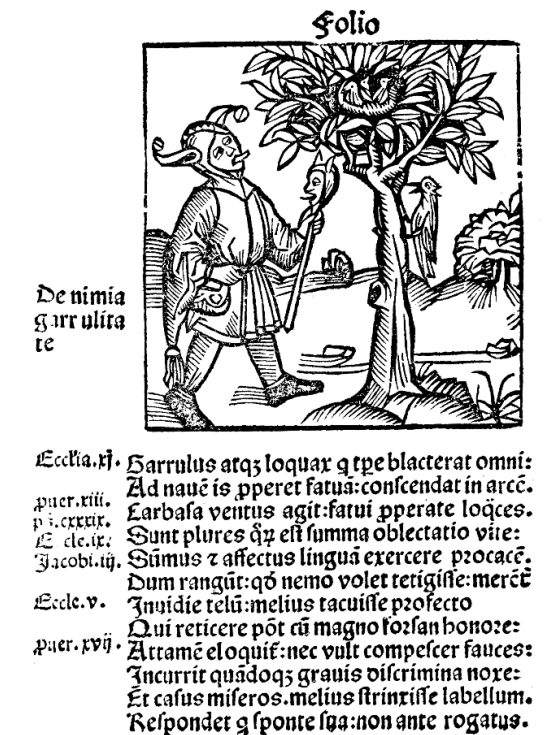
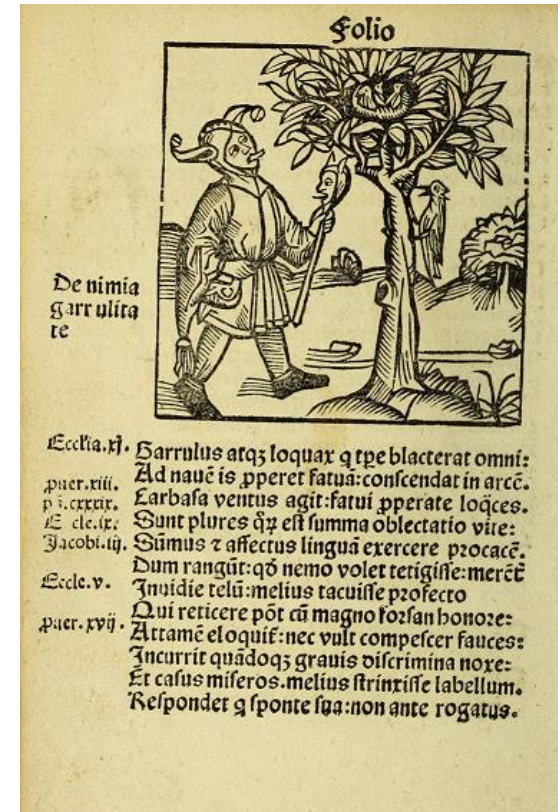


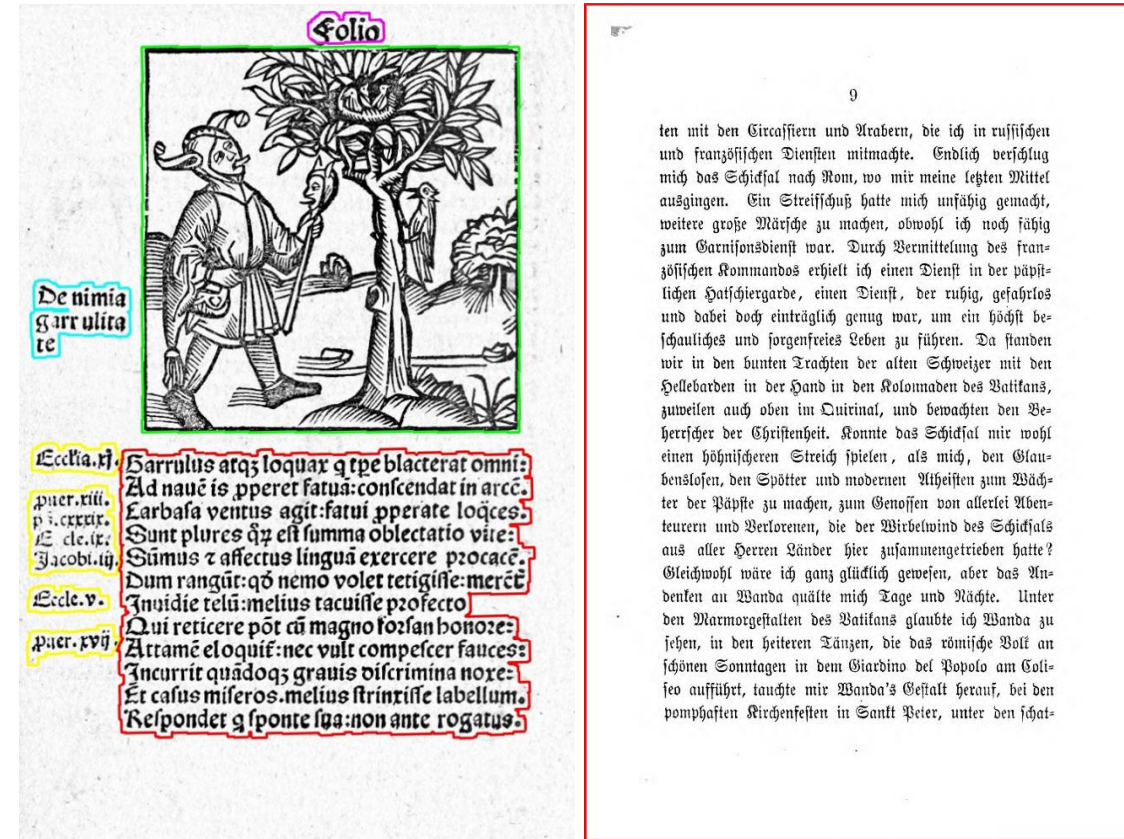
Image Preprocessing

- **Input:** Originalbilder (Farb- oder Graustufen- oder Binärbilder)
- **Output:** geradegestellte Binärbilder
- Zwei Teilschritte, die die folgenden Arbeitsschritte erleichtern:
 - Binarisierung
 - Geradestellen
- Derzeit umgesetzt durch *ocropus-nlbin*



Region Segmentation

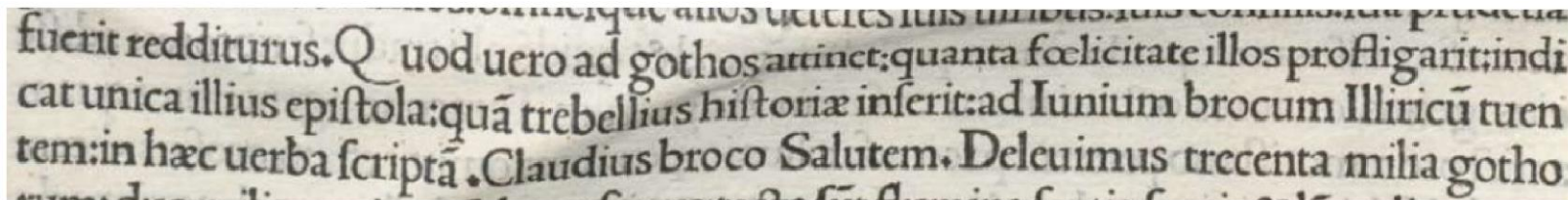
- **Input:** vorverarbeitete Bilder
- **Output:** Informationen über vorhandene Layoutelemente und deren Lesereihenfolge
- **Tools/Methoden:**
 - LAREX
 - Exakte Segmentierung inklusive semantischer Auszeichnung
 - Semiautomatisch, Intuitiv, adaptierbar und nachvollziehbar
 - Dummy Segmentation
 - Ganze Seite als ein Textsegment, Rest erledigt Zeilensegmentierung
 - Vollautomatisch und sehr schnell
 - Oft völlig ausreichend für moderate Layouts (z. B. typische Fraktur Romane aus dem 19. Jh.)



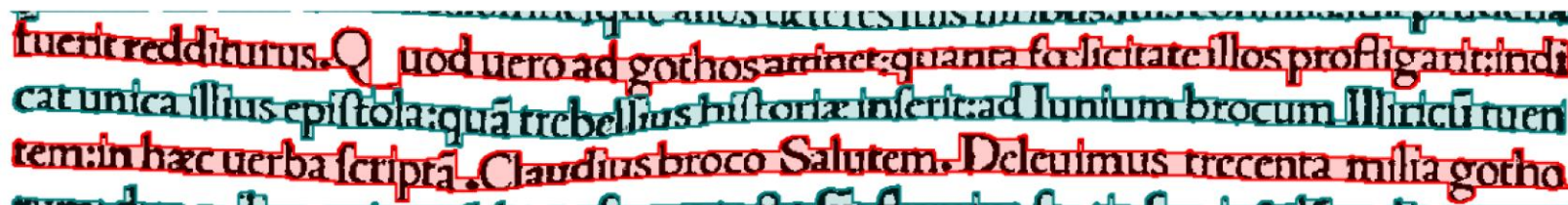
ten mit den Circassiern und Arabern, die ich in russischen und französischen Diensten mitmachte. Endlich verschlug mich das Schicksal nach Rom, wo mit meine letzten Mittel ausgingen. Ein Streifschuß hatte mich unfähig gemacht, weitere große Märsche zu machen, obwohl ich noch fähig zum Garnisonsdienst war. Durch Vermittelung des französischen Kommandos erhielt ich einen Dienst in der päpstlichen Hatfiergegarde, einen Dienst, der ruhig, gefahrlos und dabei doch eintätiglich genug war, um ein höchst beschauliches und sorgenfreies Leben zu führen. Da standen wir in den bunten Trachten der alten Schweizer mit den Hellebarben in der Hand in den Kolonnaden des Vatikans, zuweilen auch oben im Quirinal, und bewachten den Herrscher der Christenheit. Konnte das Schicksal mir wohl einen höhnischeren Streich spielen, als mich, den Glaubenslosen, den Spötter und modernen Abfasser zum Wächter der Päpste zu machen, zum Genossen von allerlei Abenteurern und Verlorenen, die der Wirbelwind des Schicksals aus aller Herren Länder hier zusammengetrieben hatte? Gleichwohl wäre ich ganz glücklich gewesen, aber das Andenken an Wanda quälte mich Tage und Nächte. Unter den Marmorgestalten des Vatikans glaubte ich Wanda zu sehen, in den heiteren Tänzen, die das römische Volk an schönen Sonntagen in dem Giardino del Popolo am Coliseo aufführt, tauchte mir Wanda's Gestalt heraus, bei den pomphaften Kirchenfesten in Sankt Peter, unter den schat-

Line Segmentation

- **Input:** vorverarbeitete Bilder und Segmentierungsinformationen
- **Output:** Zeilenpolygone
- Wichtiger Vorbereitungsschritt der eigentlichen Texterkennung
- Vorab Deskewing einzelner Regionen
- Zeilensegmentierung auch bei problematischer Ausgangslage möglich
- Derzeit umgesetzt durch angepasstes *ocropus-gpageseg*



fuerit redditurus. Quod uero ad gothos attinet; quanta felicitate illos profligari; indicat unica illius epistola; quam trebellius historiae inserit; ad Iunium brocum Illiricum tuentem; in haec uerba scripta. Claudius broco Salutem. Deleuimus trecenta milia gotho



fuerit redditurus. Quod uero ad gothos attinet; quanta felicitate illos profligari; indicat unica illius epistola; quam trebellius historiae inserit; ad Iunium brocum Illiricum tuentem; in haec uerba scripta. Claudius broco Salutem. Deleuimus trecenta milia gotho

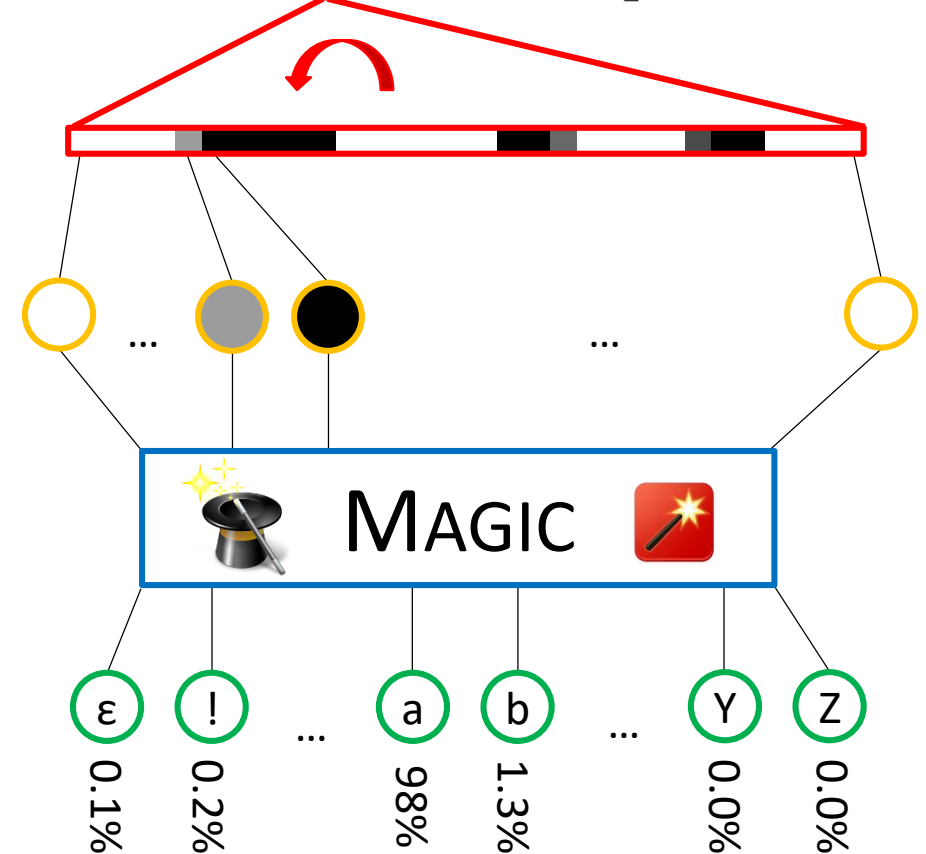
Character Recognition

- **Input:** Textzeilenbilder und OCR-Modelle
Output: erkannte Textzeilen
- Derzeit umgesetzt durch *calamari-predict*

Earrulus atq3 loquax q tpe blacterat omni:
 Ad nauē is pperet fatuā:conscenda in arcē.
 Earbasa ventus agit:fatui pperate loqces.
 Sunt plures qz est summa oblectatio vite.
 Sūmus 2 affectus lingua exercere p2ocacē.
 Dum rangūt:qđ nemo volet tetigisse:mercē
 Inuidie telū:melius tacuisse p2ofecto

Earrulus atq3 loquax q tpe blacterat omni:
 Ad nauē is pperet fatuā:conscenda in arcē
 Sunt plures qu est summa oblectatio vite:
 Sūmus 2 affectus lingua exercere p2ocacē.
 Dum rangūt:qđ nemo volet tetigisse:mercē
 Inuidie telū:melius tacuisse p2ofecto

melius tacuisse p2ofecto



Post Correction – Ground Truth Production

- **Input:** Textzeilenbilder und Erkennungsergebnisse
Output: korrigierter Text (= Ground Truth, GT)
- GT wird für das Training von OCR-Modellen benötigt
- Anpassbares Virtuelles Keyboard ermöglicht die Verwendung spezifischer Sonderzeichen

Von dem Cirurgicus

fon dem iirurgicus

IX

zc

Von dem Cirurgicus

Von dem Cirurgicus

IX

IX

Post Correction – LAREX

Fueillet

elij. di. que
lcam. i glo.

 sy font desquelles ie me tais. O poures folz belistrées qui de robes ceulz qui n'ont
 pas du pain a mager & a l'adventure quilz nosent demander de honte/les vieilles
 gens/poures desues/ladres/aveugles/helas pense y/car certes vous en rend
 dres compte deuant celluy qui nous crea

Si colligo
 q̄ vindicta
 nemo mag
 gaudet q̄ se
 ra

Des conditiōs courronpet grandes mauuaities des fēmes

**Plusieurs asnes cheuacheroient
 Le nest que monte y eust femme
 Au moyē de quoy ne voudroient
 Monte dessus po eule diffame**

**Quilz ont fait & grief infame
 Au pouure asne & grans tormēs
 Car luy ont toz tous ses bons mēs
 Car luy ont toz touo fes bons mēs**

puer. eti
 Eccle. xxx.
 puerbi. xxi.



**Ntendes folz e/
 stourdis & vous
 congnoistres les
 mauuaities des
 fēmes. Aussi fēmes appo
 ches vo' & vo' oves bone ma
 tiere. Mes versetz ditez & es
 cripts Vouldroiet des fēmes**

Quilz ont fait & grief infame

Ouilz ont fait 7 grief infame

Au pouure asne & grans tormēs

Au pouure asne 7 grans tormēs

Car luy ont toz tous ses bons mēs

Car luy ont toz touo fes bons mēs

E
 C

Ntendes folz e/

Ntendes folz e=

LOAD
SAVE

ā	æ	þ	ḃ	ç	ç	d	ð		
ē	ē	ff	ft	fl	ḡ	ḡ	h	h	ſ
ī	ī	ll	m	n	ō	œ			
ṑ	ṑ	p	p	p	p	p			
q	q	q	q	q	q	q			
ṙ	ṙ	z	s	f	th	ff	fl	ft	ſ
t	ḡ	ū	'	ḡ	ḡ				
v	v	v	w	x	y	z	&	7	
o	o	z	=	'	ll	§	¶		
+									

Evaluation

- **Input:** OCR-Ergebnisse und Ground Truth für bestimmte Textzeilen
Output: Zeichenfehlerrate (CER) und -statistiken
- Zeigt die Art und Häufigkeit der auftretenden Fehler an
- Auf Grundlage der Evaluation lässt sich erkennen, ob Fehler (noch) systematisch sind oder nicht
- Derzeit umgesetzt durch *calamari-eval*

GT	PRED	COUNT	PERCENT
{₂}	{r}	33	17.65%
{ }	{ }	11	5.88%
{ī}	{ī}	11	5.88%
{ff}	{f}	7	3.74%
{ff}	{ff}	5	5.35%
{t}	{i}	3	1.60%
{q}	{ }	2	1.07%
{ }	{ }	2	1.07%
{ }	{ }	2	1.07%
{x}	{ }	2	1.07%

Training

Umfangreiche Unterstützung ...

- ... der zahlreichen Möglichkeiten des Calamari Trainingsprozesses
 - Ensemble-Training
 - Pretraining/Finetuning
 - Datenaugmentierung
 - ...
- ... des iterativen Trainingsansatzes
 - Ständiges Durchlaufen der Schritte Erkennung, Korrektur und Training zur Effizienzsteigerung
 - Modellverwaltung

Für Details siehe Veranstaltung zum Thema „Training“ nächste Woche.

Live Demo!

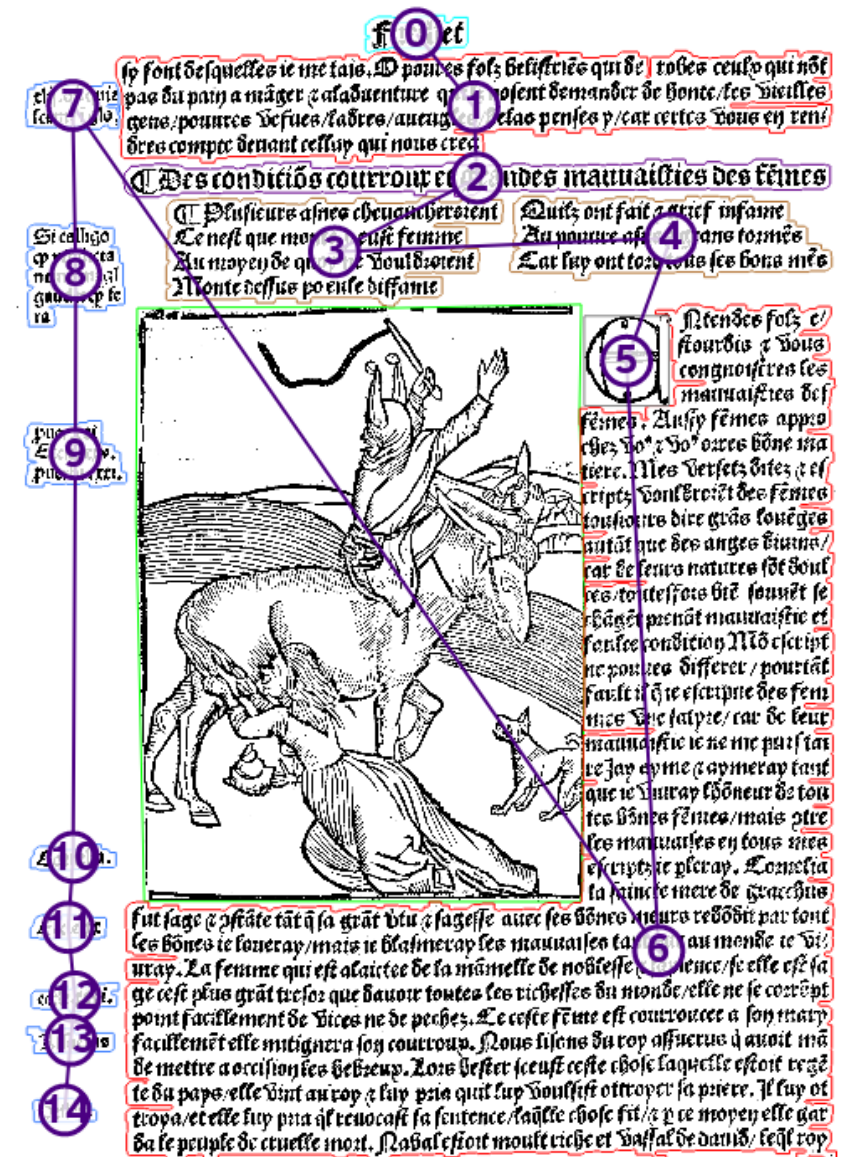


Gliederung

1. Einleitung
2. Submodule
3. Workflow
4. Live Demo
- 5. Evaluation**
6. Diskussion und Ausblick

Evaluation – Frühdrucke bis 1600

- Material
 - 5 Ausgaben des „Narrenschiffs“
 - 17 Werke des Universalgelehrten Joachim Camerarius
 - 3 gemischte, früh-neuzeitliche Drucke (Praktikum)
- Anforderungen / Ziele:
 - Fehlerfreie Segmentierung und Reading Order sowie präzise semantische Auszeichnung
 - Final: zitierfähiger Volltext (0% Fehler)
 - Vorerst: 1% CER oder besser
- Bearbeitung in zwei Gruppen:
 - Unerfahrene, nicht-technische Nutzer
 - Erfahrene Nutzer



Frühdricke bis 1600 – Beispiele I

I

HISTORIAE IESV
CHRISTI FILII DEI NATI
IN TERRA MATRE SANCTISS.
semper virgine Maria sum-
matim relata ex-
positio.

MULTI studium literarum in eo posue-
runt, vt memoria mandare scriptis
suis virorum illustrium & virtute praestan-
tium, aut varietate casuum insignium, vi-
tam, fortunas, mortem, quia vel narratio-
nem eam iucundam legentibus propter
historiae cognitionem, vel exempla propter
instructionem animorum vtilia futura esse
confidentibus existimauerunt. In qua
parte cum & nobis placuisset operae aliquid
consumere, si nihil aliud confecturis hoc
labore, saltem quietem aliquam reperituris,
& abducturis cogitationes à triflicia re-
rum & miseriae tam publicae quam priuatae,
in manus sumferamus, & aggredi ceperamus
contextum eorum quae à veteribus au-
toribus de eloquentibus Graeciae tradita &
memoriae prodita comperissemus. Cum su-
perioris temporis calamitas & communis
maeror bonorum, & sensus proprius eorum
malorum, quibus tum nostra natio obrue-
retur, aurem, vt ita dicam, velle admo-
nuit,

B 4

COMMENTARIORVM LIB. I. 19

BVIT LV MEN LITERARVM, significat de fuisse philosophiae eloquen-
tiam. Epicuri uero emuli hoc loco notantur, quos et alibi dicit minime ma-
los, sed parum acutos esse. Iam Aristotelis hoc notum est, *αὐτῶν οὐκ ἔστιν οὐδὲν ἄλλο*
ἢ τὸ ἐπιμαρτυρεῖν, quem uersum de Tragedia in illam Rhetorem interserit.
Sequitur figura sermonis inusitata uisula, IN QVAM EXER-
CITATIONEM OPERAM DEDIMVS. Simile autem est,
uel quod dicimus, Dare nomen in militiam: uel, quod magis placet, in po-
testatem esse, et in honorem, et in discrimen uersari. de quibus dubitanti ti-
bi DANIELE, aliquando copiose exposuimus sententiam nostram. Me-
mini, inquit: et quidem nuper, cum te in eadem sermonem deduxissem, pre-
sentibus aliquot amicis nostris, iucundissima fuit mentio quaedam habita à te, de
toto genere Ciceronianae imitationis. Sed quia subito tui inde quasi res illebas,
de industria an forte uito dicam nescio, sit ut adhuc permulta quidem requirā.
Ommino de industria inquam, neque enim ferenda erat executio harum rerū
illo tempore: et uitabam ijs de rebus, ijsque presentibus, longam disputationē,
ne, ut nunc sunt homines, in maleuolentia suspicionem incurrerem. Quin
igitur nunc, inquit, exple desiderium meum, conueniente et in argumento,
et loco libero. Quod nam illud: inquam. Cognoscenū silicet, quid sta-
tus de hac tota ratione sermonis antiqui, quatenus parum ille cum hodierno
usu congruit, et quatenus huius auctores, quosque ipsorum potissimum sequen-
dos existimes. Aperte inquam: in quas difficultates inducis orationem no-
stram, ex hac plana uia, quam ingressi sumus, nam uidere omnino aliquid uel
le audire nouum. Nisi enim uellem, inquit, habebam quos consulere.
Nos uero nihil habemus, inquam, noui. Quare ad illos abeas iuadeo. Sed
uitos? inquit. Scimus enim duas contra se partes constituisse: et notum disti-
dium est. Hoc uero difficillimum, inquam, decernere est. Si tamen omnino de
iudicio meo reddi certior cupis, scito me plane cum ijs facere, quos consuet
post diligentiorē curam, et studium ardentius, et magis assiduū usum
uere Latinitatis demonstrasse, et alijs directum iter, quo ipsi ad nonnullam
laudem orationis bonae et elegantis in hac lingua peruenissent. Ex quorum
numero primus tribuo P. BEMBO, cuius hac de re extat epistola ualde lucu-
lenta. Accessit enim proximè ad ueteres, quorum orationi illius oratio simil-
lima est, quo certius de re nota perspicuas sibi differere cum putandum.
Sed haec, inquit ille, et legi, et relegam saepius. Nunc uero te audire cupio.
primūque, hoc uelim ediscas, quid sit quod placere tibi tantopere in BEM-
BO ais similitudinem ueterum. Cupio enim, ut ait eam, te audire de hoc po-
tissimum

PETRVS
BEMBO.

6 2

109981/3

Vita salusque homini est, is tota mente capestat,
Cogitet, ediscat, meditetur, corde uoluet
Includatque pio, hac se consoletur, eadem
Pectus ad hostiles armans communiat ictus.
QVOD patris e gremio celo descenderit alto
Filius aeterni, simul ipse aeternus, et intra
Mortalem celeste genus conclusit ortum,
Quem non ista capit totius machina mundi:
Factus homo in terris, Rerum non indigus ille
Nostrarum, sed cura fuit reparare salutem
Amisam et vitam nobis, ab origine prima
Quos peccatorum duro sub pondere pressos
Obruit ira Dei iusto commota furore,
Haec causa in terras celo hunc detraxit, ut esset
Saluator miserorum hominum, Deus vnicus et
Rex,
Σωτήρ, Λυτρωτής, Θανατηφόρος, Ἐρανα-
νοικτής.

SEQVUNTUR ALII
quidam religiosi argumenti
versus compositi a Io-
achimmo Came-
rario.

K 4

PRE-

Vita salusque homini est, is tota mente capestat,
Cogitet, ediscat, meditetur, corde uoluet
Includatque pio, hac se consoletur, eadem
Pectus ad hostiles armans communiat ictus.
QVOD patris e gremio celo descenderit alto
Filius aeterni, simul ipse aeternus, et intra
Mortalem celeste genus conclusit ortum,
Quem non ista capit totius machina mundi:
Factus homo in terris, Rerum non indigus ille
Nostrarum, sed cura fuit reparare salutem
Amisam et vitam nobis, ab origine prima
Quos peccatorum duro sub pondere pressos
Obruit ira Dei iusto commota furore,
Haec causa in terras celo hunc detraxit, ut esset
Saluator miserorum hominum, Deus vnicus et
Rex,
Σωτήρ, Λυτρωτής, Θανατηφόρος, Ἐρανα-
νοικτής.

SEQVUNTUR ALII
quidam religiosi argumenti
versus compositi a Io-
achimmo Came-
rario.


K 4

PRE-

Frühdrucke bis 1600 – Beispiele II


wolt er sy mit gewalte nemen. die kna-
ben müßten sein aigen sein. vñ die maid
wolt er dabey in das offen fra w hauf
thun das sy im pfenning gewinnen solt
ten. hört wie ein schantliche vnbeschay-
dene botschafft das vñ einem kung was.
Vnd der er sich billich er geschampft het
zu gedencken. denn das er es überlawte
hieß außsprüffen.

Albensteur. Wie Tristrant ritter gema-
cht ward. vñnd sich verwilliget mit
Morcholtzen zu wechten.



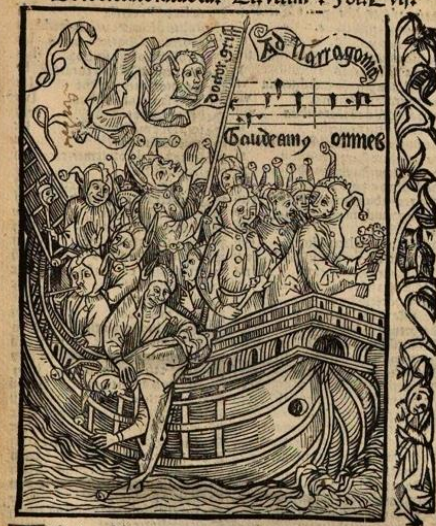
wenn ainer ain schönen billen über kam / so wolt
der annder noch ain hübschere habenn / es kostere
was es wolt / das triben sy bey ainem halben tar.
Do begund es nachnen das sy nit vil bar gelt mer
hetten. doch het ainer mer onworden dann der an-
der.

Wie fortunatus zu böser gesel-
schafft kam / mit denen / vñnd mit
leichten frawen / als sein gelt ver-
ther / vñnd sich darnach vil armüt
leiden müße.



Fortunatus der hett am minsten / der ward
auch am ersten gerecht. Er het sine klainat
vñnd als onworden. desgleichen die andern


De societate fatuorum. Li. vltim^o. Fol. Lvij.



Sopia multū ugi deducit gaudia stulti
Gaudentes socii ⁊ multitudinē cetus
Non minus arduū q̄ si tū ardeat vnus.

Sopia ꝛc. Scribit in stultos se sua multitudine defendentes. Nā
interdū multitudinē p̄cedunt est. Verū qđ Augustinus ait ve in
diceret̄ habet̄ nō min⁹ ardebit̄ q̄ cū multis ardebit̄: qđ q̄to cōbu-
stibile copiosius: t̄to ē calor uehemētior. Nec multitudo errātū uenit̄ meo
ref: qđ sicut bonū quāto cōmuni⁹ t̄to est meli⁹: ita malū tanto peius. Nec
sodoma ⁊ gomorra perissent̄ si plusculos innoxios habuissent. Quia so stul-
torū infini⁹ est̄ numer⁹ ⁊ infini⁹te species ne qđ se p̄missum q̄rat: idē in-
struimus hoc cōmune receptaculū. S; receptū canēdū est: nauis vrget.
f̄nis bulus operis. Sequitur index

De societate fatuorum. Li. vltim^o. Fol. Lvij.



Sopia multū ugi deducit gaudia stulti
Gaudentes socii ⁊ multitudinē cetus
Non minus arduū q̄ si tū ardeat vnus.


Sopia ꝛc. Scribit in stultos se sua multitudine defendentes. Nā
interdū multitudinē p̄cedunt est. Verū qđ Augustinus ait ve in
diceret̄ habet̄ nō min⁹ ardebit̄ q̄ cū multis ardebit̄: qđ q̄to cōbu-
stibile copiosius: t̄to ē calor uehemētior. Nec multitudo errātū uenit̄ meo
ref: qđ sicut bonū quāto cōmuni⁹ t̄to est meli⁹: ita malū tanto peius. Nec
sodoma ⁊ gomorra perissent̄ si plusculos innoxios habuissent. Quia so stul-
torū infini⁹ est̄ numer⁹ ⁊ infini⁹te species ne qđ se p̄missum q̄rat: idē in-
struimus hoc cōmune receptaculū. S; receptū canēdū est: nauis vrget.
f̄nis bulus operis. Sequitur index

Frühdrucke bis 1600 – Beispiele III

Dann keiner mag mehr reich tumb han
Dann wer ganz nichts begeret darvon.
Der geytzig nimmer füller sich
Wen wol benigete der ist vast reich.æ.

Von newen Fünden.

Wer vil new fünd macht durch die landt
Der gibt vil ergernuß vnd schandt/
Beseh er sich im spiegel wol
Balt wist er was er lassen soll.



ASettwann was ein schandlich ding
Das wagt mā jetzt schlecht vñ gering/
E in eh; was ettwan tragen bārt/
Das was gar mannlich/schon vnd wert/
Do wurden man auch billich geehrt/
Jetzt handt die weybschen gench geleert
Vñnd schaben all tag jr zwilchbacken

Sic

So groß gewalt vff erd nye kam
Der nitt zū zyeten/end ouch nam
Wann im syn 37l/vnd stündlin kam



Von end des gewalttes

Noch fyndt man narren manigfalt
Die sich verkont vff iren gewalt
Als ob er ewilich solt ston
Der doch dūt/wie der schne zergon

Fueillec

qui veult mettre sijn a plusieurs choses ne scauroit estre constant. / s'ij veult a plu
ieurs gens plaire. il luy conuient estre humble & vider de douls langaiges. Et se
quelque aduersite luy suruient il faultdra quil se porte paciemment sans se esba
hir. Auz nobles doit vider de langaige esto quant ce quil dira plaise a celuy a qui
il parlera. Saluer doit humblemēt & point a nul ne doit courroucer. sil veult estre
de tous apme. car pour la grant charge quil a il se doit faire de to' apmer. Il nau
ra repos nullement pensant toujours a ses biens & a ses offices peent grant peine
en se damnant. et na sejour repos ne soulas. tellement quil ne pense a dieu ne a le
secur. tant a de pensees en la teste. Touchant de tels fols; ie me y deposite a me tai
rap pour le present. mais mieulx seroit de secur bier & loy aulement. Vng boy mai
fere et se faire de luy apmer que de se vouloir de plusieurs offrir. a en la fin offre
en mille grace d'ung cheuuy. et perdie & consumer pareillemēt soy temps en espe
rance de vouloir trop acquerre.

De trop parler.

Qui siet sa langue referenc
Et de trop parler la referaindre
Ne se voit a mal profierner
Par tristesse que loy doit craindre

Mais qui parletrop sans se faidre
De deshonore con; fait la pie
Par trop garruler quoy lespie

Deus ebetes
lāgues dragon
ques q a trop pl
et blasmer aut
truy vo' arestes toute natu
re de bestes. doiscaultz & d ser
pēs. & autres de nature hu
maine se peuent chafiter/
mais mais la lāgue d'hom
me ne se peult chafiter car el
se est pleine de mal & d veni
mortifere. elle macule tout
se corps qui garde sa bouche
garde soy ame. Pour ce lan
gues arceuesqes bades. Vo' d ce
sue doctrine. car p' vault vng
cop de lāgue q vng cop de lā
ce. celui q parle trop & vaine
mēt en to' tēps q ne vūt il a
mēt folle nef. Venez y tost
pour gouverner les voilles/a
uāces Vo' fols q trop plees et
raualtes vos langaiges. plu
sirs a q se delactēt & ne pētē



Prouer. xij
ps. cxxix.
Ecclesia. ix.
Jacob. iij.
ecclesiast. v.

Fueillec


qui veult mettre sijn a plusieurs choses ne scauroit estre constant. / s'ij veult a plu
ieurs gens plaire. il luy conuient estre humble & vider de douls langaiges. Et se
quelque aduersite luy suruient il faultdra quil se porte paciemment sans se esba
hir. Auz nobles doit vider de langaige esto quant ce quil dira plaise a celuy a qui
il parlera. Saluer doit humblemēt & point a nul ne doit courroucer. sil veult estre
de tous apme. car pour la grant charge quil a il se doit faire de to' apmer. Il nau
ra repos nullement pensant toujours a ses biens & a ses offices peent grant peine
en se damnant. et na sejour repos ne soulas. tellement quil ne pense a dieu ne a le
secur. tant a de pensees en la teste. Touchant de tels fols; ie me y deposite a me tai
rap pour le present. mais mieulx seroit de secur bier & loy aulement. Vng boy mai
fere et se faire de luy apmer que de se vouloir de plusieurs offrir. a en la fin offre
en mille grace d'ung cheuuy. et perdie & consumer pareillemēt soy temps en espe
rance de vouloir trop acquerre.

De trop parler.

Qui siet sa langue referenc
Et de trop parler la referaindre
Ne se voit a mal profierner
Par tristesse que loy doit craindre

Mais qui parletrop sans se faidre
De deshonore con; fait la pie
Par trop garruler quoy lespie

Deus ebetes
lāgues dragon
ques q a trop pl
et blasmer aut
truy vo' arestes toute natu
re de bestes. doiscaultz & d ser
pēs. & autres de nature hu
maine se peuent chafiter/
mais mais la lāgue d'hom
me ne se peult chafiter car el
se est pleine de mal & d veni
mortifere. elle macule tout
se corps qui garde sa bouche
garde soy ame. Pour ce lan
gues arceuesqes bades. Vo' d ce
sue doctrine. car p' vault vng
cop de lāgue q vng cop de lā
ce. celui q parle trop & vaine
mēt en to' tēps q ne vūt il a
mēt folle nef. Venez y tost
pour gouverner les voilles/a
uāces Vo' fols q trop plees et
raualtes vos langaiges. plu
sirs a q se delactēt & ne pētē



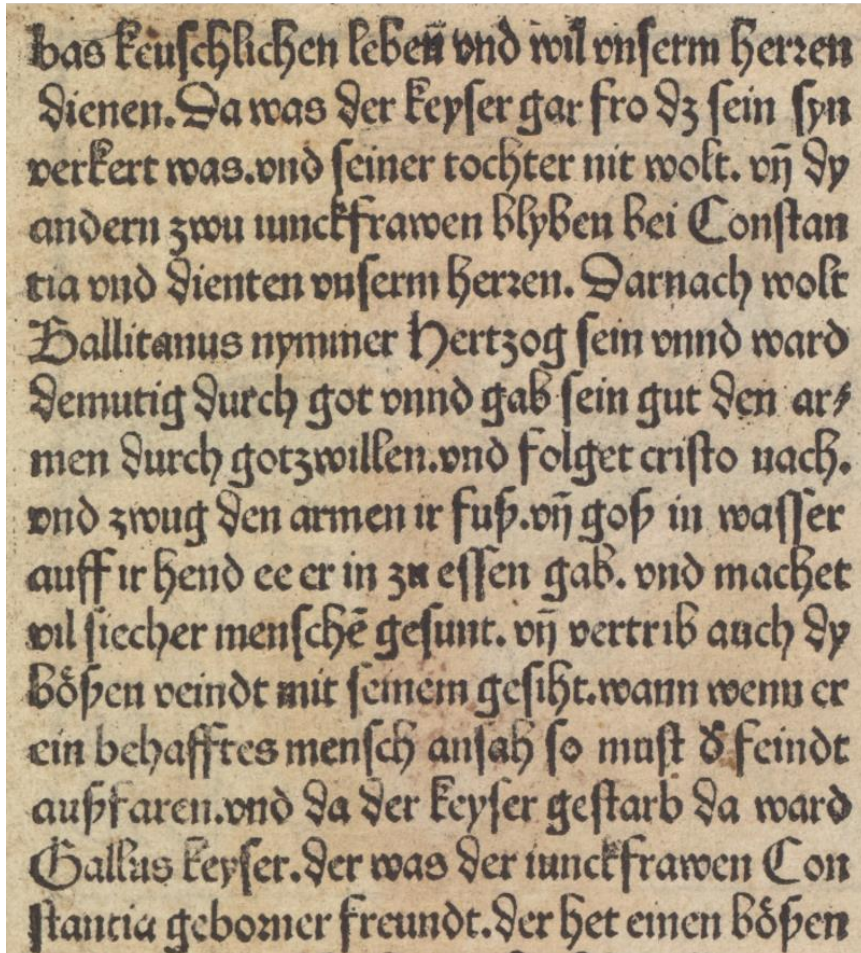
Prouer. xij
ps. cxxix.
Ecclesia. ix.
Jacob. iij.
ecclesiast. v.

Frühdrucke bis 1600 – Ergebnisse

	Unerfahrene Nutzer	Erfahrene Nutzer
Erreichte CER	0,47% ± 0,22%	0,49% ± 0,30%
Transkribiertes Trainingsmaterial	988 Zeilen	927 Zeilen
Korrekturzeit pro Zeile	10s ± 5,2s	5,5s ± 2,4s
Segmentierungszeit pro Seite	1,1min ± 0,5min	0,6min ± 0,2min

*Jeweils Durchschnittswerte, ggf. mit Standardabweichung

Beispielerggebnis mit ca. 0,6% CER



Das k^uschlichen leben vnd wil vnserm herren
dienen. Da was der keyser gar fro dz sein syn
verkert was. vnd seiner tochter nit wolt. vñ dy
andern zwu iunckfrawen blyben bei Constan
tia vnd dienten vnserm herren. Darnach wolt
Ballitanus nymmer Hertzog sein vnnd ward
demutig durch got vnnd gab sein gut den ar
men durch gotzwillen. vnd folget cristo nach.
vnd zwug den armen ir fuß. vñ goß in wasser
auff ir hend ee er in zu essen gab. vnd machet
vil siecher mensche gesunt. vñ vertrib auch dy
bößen veindt mit seinem gesiht. wann wenn er
ein behafftes mensch ansah so must ð feindt
außtaren. vnd da der keyser gestarb da ward
Gallus keyser. der was der iunckfrawen Con
stantia geborner freundt. der het einen bößen

bas k^uschlichen l^oben vnd wil vnserm herren
dienen. Da was der keyser gar fro dz sein syn
verkert was. vnd seiner tochter nit wolt. vñ dy
andern zwu iunckfrawen blyben bei Constan
tia vnd dienten vnserm herren. Darnach wolt
Ballitanus nymmer Hertzog sein vnnd ward
demutig durch got vnnd gab sein gut den ar
men durch gotzwillen. vnd folget cristo nach.
vnd zwug den armen ir fuß. vñ goß in wasser
auff ir hend ee er in zu essen gab. vnd machet
vil siecher mensche gesunt. vñ vertrib auch dy
bößen veindt mit seinem gesiht. wann wenn er
ein behafftes mensch ansah so must ð feindt
außtaren. vnd da der keyser gestarb da ward
tGallus keyser. der was der iunckfrawen Con
stantia geborner freundt. der het einen bößen

Gliederung

1. Einleitung
2. Submodule
3. Workflow
4. Live Demo
5. Evaluation

6. Diskussion und Ausblick

Zusammenfassung

- OCR4all vielseitig einsetzbar
 - Erfolgreiche Verarbeitung von (historischen) Drucken des 15. bis 21. Jahrhunderts
 - Hauptanwendung: lokale Installation beim Nutzer
 - Nutzung als Serveranwendung prinzipiell bereits jetzt möglich
- Ergebnisse und dafür notwendiger Aufwand stark abhängig ...
 - ... vom Material
 - Automatische Segmentierung möglich?
 - Passendes gemischtes Modell vorhanden?
 - ...
 - ... von den Nutzeranforderungen
 - Grad der semantischen Auszeichnung?
 - Ansprüche hinsichtlich Genauigkeit?
 - ...

- Kooperationsvereinbarung Sommer 2020
 - Umsetzung von OCR-D Spezifikationen und Schnittstellen in OCR4all zum beiderseitigen Vorteil:
 - für OCR-D: bei Bedarf vereinfachter Zugang für die Nutzer, größere Reichweite
 - für OCR4all: erweiterte Auswahl an Werkzeugen, Flexibilität
 - Fortlaufender Austausch (Schnittstellen, Skalierbarkeit, kommende OCR Entwicklungen, GT, ...)
- Erfolgreicher Projektantrag *OCR4all-libraries* in dritter OCR-D Phase:
 - Kooperation zwischen GEI Braunschweig und Uni Würzburg (HCI, ZPD)
 - Volle Unterstützung der OCR-D Lösungen sowie deren Steuerung und Konfiguration über die GUI
 - Ermöglichung einer Massenverarbeitung von Werken bzw. Werkclustern
 - Ausbau von LAREX als visuelle Erklärungskomponente (Fehleranalyse, Vergleich von Workflows, ...)
 - Optimierung der Usability
 - ...

Verarbeitung von Handschriften

- Erfassung von Drucken und Handschriften konzeptionell sehr ähnlich
- OCR4all bereits jetzt vielseitig im Handschrifteneinsatz
- Systematische Evaluation in Kooperation mit Dr. Stefan Tomasek (Lehrstuhl für Ältere Deutsche Literatur, Uni Würzburg) anhand des Projekts „Konrad von Fußesbrunnen: Kindheit Jesu“
- Größtes Desiderat: robuste automatische Segmentierung
- Zeitnahe Anbindung vielversprechender, Baseline-basierter, trainierbarer (!) Lösungen:
 - [Kraken](#) ([Kiessling: A Modular Region and Text Line Layout Analysis System](#))
 - Entwicklungen am Würzburger Lehrstuhl für Künstliche Intelligenz (Fischer, Hartelt, Gehrke, Puppe)
 - ...

Aktuelle Entwicklungen und Planungen

- Zeitnah: weitreichend überarbeitete Version mit zahlreichen Verbesserungen hinsichtlich Stabilität und Flexibilität
- Mittelfristig: größere Flexibilität hinsichtlich Material und Anwendungsszenarien
 - Handschriften
 - Massenvolltextdigitalisierung
- Weiterer DFG Antrag in Vorbereitung, u. a. Fokus auf
 - Optimierung des kollaborativen Arbeitens (Projekt-, Nutzer- und Taskverwaltung, Ressourcenmanagement, Backup und Versionierung, ...)
 - Bei Bedarf Steuerung einzelner Prozesse und Workflows aus LAREX heraus
 - Konfidenzbasierte Qualitätsanalyse und interaktive Nachkorrektur
 - Usability (Nutzerstudien, Überarbeitung des Schulungskonzepts und der Anleitungen, ...)

Interesse an OCR4all?

- Probieren Sie OCR4all auf Ihrem eigenem Rechner mit Ihrem Material aus
 - www.ocr4all.de
 - https://github.com/OCR4all/getting_started
- Wir helfen gerne bei Fragen und Problemen: ocr4all@uni-wuerzburg.de
 - Installation und Nutzung
 - Bug Fixes und Feature Requests
 - Projekt-spezifisches Consulting
 - Server Support, falls die eigene Hardware nicht ausreicht
 - ...
- Wir sind auf Ihr (unverblühtes) Feedback angewiesen!
- ... aber bedenken Sie bitte immer: OCR4all ist umsonst und Work in Progress 😊

Vielen Dank für Ihre Aufmerksamkeit!

Acknowledgements:

- **OCR4all Web App:** Dr. Herbert Baier-Saip, Maximilian Nöth, Dennis Christ, Alexander Hartelt, Nico Balbach, Kevin Chadbourne
- **LAREX Web App:** Maximilian Nöth, Kevin Chadbourne, Nico Balbach
- **Calamari:** Dr. Christoph Wick, Andreas Büttner
- **Tests, Anleitungen, Nutzer Support und Artwork:** Maximilian Wehner, Raphaelle Jung, ...
- **Distribution via Docker und VirtualBox:** Björn Eyeselein, Yannik Herbst
- **Ideen und Feedback:** Dr. Uwe Springmann, Maximilian Wehner, Prof. Dr. Frank Puppe, Christine Grundig, Prof. Dr. Brigitte Burrichter, Prof. Dr. Joachim Hamm, ...
- **Förderung:** DFG Förderinitiative „OCR-D“ Phase 2 und 3, BMBF Projekt „Kallimachos“, Lehrstuhl für Künstliche Intelligenz (Prof. Dr. Frank Puppe) und Zentrum für Philologie und Digitalität der Uni Würzburg