

OCR-BW – Kompetenzzentrum OCR der Universitätsbibliotheken Mannheim und Tübingen







OCR-BW



Projektziele

"Aufbau eines Kompetenzzentrums für Volltexterschließung von handschriftlichen und gedruckten Werken."

Das Projekt OCR-BW unterstützt Archive, wissenschaftliche Bibliotheken und andere Institutionen in Baden-Württemberg bei der Anwendung von automatischer Texterkennungs- und Transkriptionssoftware.

UB Tübingen: Transkription und Volltexterschließung von

Autographen, Handschriften und Inkunabeln

UB Mannheim: Volltexterkennung (OCR) von Druckwerken

https://ocr-bw.bib.uni-mannheim.de/

OCR-BW



Arbeitsbeispiele für Transkription und Texterkennung

Handschrift aus dem Bestand der WLB Stuttgart (16. Jhd.)



 Historische Zeitungen (19.–20. Jhd.) für das Stadtarchiv Ladenburg, das MARCHIVUM und die BLB Karlsruhe



OCR-BW



Texterkennung von historischen Drucken mit OCR-D und Tesseract



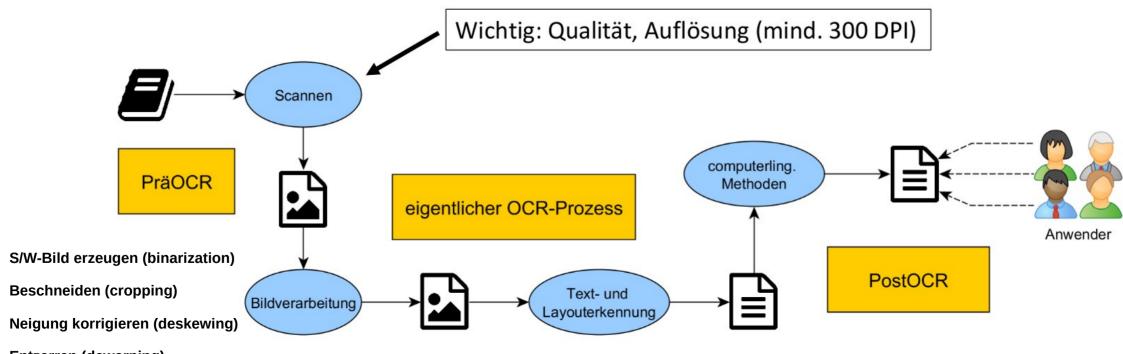
Stefan Weil, Jan Kamlah Universitätsbibliothek Mannheim





Vom Werk zum Digitalisat mit Volltext





Entzerren (dewarping)

Flecken entfernen (despeckling)

aus: Baierer, Zumstein. Verbesserung der OCR in digitalen Sammlungen von Bibliotheken

16.10.2020 5

Tesseract in Kürze https://tesseract-ocr.github.io/



- Open Source
- Komplettlösung "All-in-1"
- Kommandozeilenprogramm (CLI)
- Dient als Kernkomponente in vielen OCR-Applikationen
- Mehr als 100 Sprachen / mehr als 30 Schriften
- Liest Bilder in allen gängigen Formaten (nicht PDF!)
- Erzeugt Text, PDF, hOCR, ALTO, TSV
- Große, weltweite Anwender-Community
- Technologisch aktuell (Texterkennung mit neuronalem Netz)







OCR-D in Kürze https://ocr-d.github.io/de/



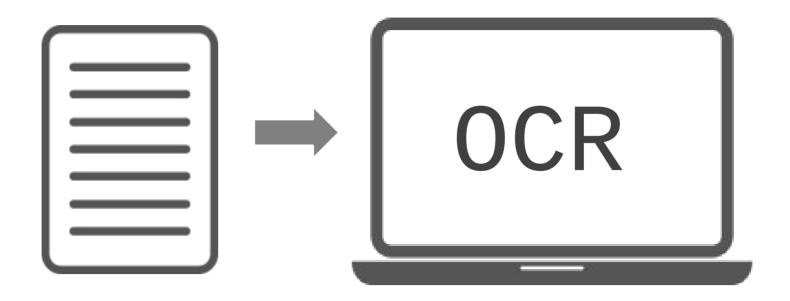
- Open Source
- Modularer Ansatz "Prozessoren" für alle Arbeitsschritte
- Kommandozeilenprogramm (CLI)
- METS als Metadatenformat
- PAGE XML für Seiteninhalte, optional auch Text, PDF, hOCR, ALTO XML, TSV, ...
- Definition von Ground-Truth-Erfassungsregeln und Bereitstellung von GT-Korpora
- Langzeitarchivierungsmöglichkeiten
- Fokus auf Drucke des 16.–18. Jahrhunderts (VD16, VD17, VD18)







OCR im Einsatz



One model to rule them all?





Neue generische Modelle sind sprachenunabhängig und decken auch viele Schriftarten ab, ebenso wie **fett** und *kursiv* gedruckte Variationen. Erleichtert Volltexterkennung für typische Druckwerke.

Erstes Modell für historische Schriftarten:

GT4HistOCR

https://ub-backup.bib.uni-mannheim.de/~stweil/ocrd-train/data/Fraktur_5000000/tessdata_fast/

https://typo-info.de/entwicklung-der-schrift/

Beispiel 1 – BLB Karlsruhe



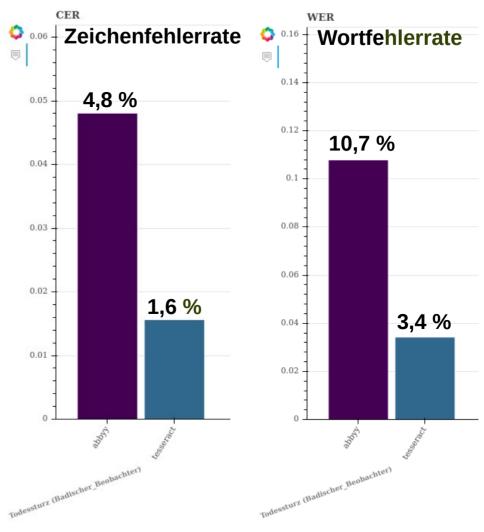
Der Zodesiturz aus der Stratosphäre

Mus ben Radiogrammen, die ber , gen feien erichwert, weil bie Inftrumente perungliidte Stratoiphären-Ballon .. Sirius" bor seinem Untergang auf die Erde herabgefunkt hatte, stellt die ruffische Breffe noch eine Sulle bemerfen & merter Gingelbeiten ausammen. Gin Teil davon ift wissenschaftlich umso wertvoller, als beim Aufschlagen ber Gondel auf den Erdboden alle wiffenschaftlichen Instrumente und Aufschreibungen durch eine Explosion vernichtet wurde. Es steht beute einwand. frei feit. daß icon in ber Luft, aber bereits außerhalb ber Stratofphare, amei Explofionen gehört murden, in beren Gefolge fich bie Gondel vom Ballon trennte. Borber mar bon amei Angenzengen beobachtet worben. baß bie Bemühungen ber Forider, ben Ballon niebergeben gu machen, auf irgend einen geheimnisvollen Wiberftanb geftoßen

Die Radiogramme ergeben bor allem, daß atmofphärifde Schwierigfeiten nicht borhanben waren. Es berrichte fast Windftille. Rach etwas mehr als 2 Stunden befand fich ber "Sirius" in 20 000 Meter Bobe in einer Temperatur, die zwischen 36 und 45 Grad Ralte ichwantte. Es waren "außerordentlich intereffante" Beobachtungen über "Art und Stärke ber fogenannten tosmifden Strablen" möglich, wie Kommandant Fedoffenko funkte. In einem der letten Explosion begünftigt, eine sol Radiogramme hieß es, weitere Beobachtun- zu befürchten gewesen wäre.

ihren Dienft nicht mehr tun. Go mor es den Luftschiffern zuletzt auch nicht mehr möglich die Position zu funken, wo der Ballon fich befand. Wie weit die Intensität der fosmischen Strahlen dabei eine Rolle ipielte, fann nur bermutet merden

Mehr als die fosmifden Strahlen fteht ber Caueritoffmangel als Urfache ber Rataftrophe im Berdacht. Einige der Radio-gramme machen auf diesen Sauerstoffmangel aufmerksam. Da andererseits genug Sauerstoffporrat sich im Ballon befand, liegt die Bermutung nabe, daß der Mangel auf ein Berfagen bes Cauerftoffapparates gurudguführen mar. In diesem Falle murden die Luftidiffer den Erftidungstob ichon in ber Stratviphare gefunden haben. Das äußerft langjame Sinken des Ballons gibt Berechtigung zu der Bermutung, daß die Reißbahn überhaupt nicht gezogen war. Es ist auch anguneh men, daß ichlieklich in einer bestimmten Sphare außerhalb der Strato fich ber Cauerftoffapparat wieber normal einstellte, baß das Gas heftig ausströmte und burch Explosion die Gondel von der Ballonhülle abrif. Jedenfalls ift auffällig, daß die zwei Erplosionen in relativ geringer Sohe erfolgt find, mahrend eigentlich viel weiter oben, wo der berminderte Drud die Explosion begünftigt, eine solche Ratastrophe



Beispiel 2 - ULB Darmstadt (Preprocessing)

MANNHEIM

CER: 6.84 %

Diejenige Berren Liebhaber, welche fich bes neuen, in Rupfer geflochenen, anddiaft privilegirten 2Band Calenders, auf Diefes neue Sahr 1768. bedienen mollen, merden hiermit eraebenft benachrichtiget, daß nur noch, wegen dem ftarfen Abgang, Den Derfelbe feiner Bequemlichfeit haiber, guswartig gefunden bat, eine geringe Angabl guter Erempfarien ben mir gu haben find, bas Stud foftet, mie hekannt . 12 Areuter.

Ge mirb biermit befannt gemacht, bag burch mich Berrg Wolfgang Steeg. mener, ale Bader und Chiruraus, bas neue Bad anwiederum beiest worden ill. man berfpricht auch qualeich biermit in allen Studen gute Bedienungen. Es ift auch zu bekommen um einen civilen Preis, als ertra gutes Rosenwaffer die Maas por 24 Rreuger, fodann fdwarg Riridenmaffer , und weiß Liffenmaffer 2c. 2c.

Bu verlebnen werben angetragen:

Den bem Caffetier Becter Dabier, find zu benen bevorffebenben Mafauens Balls, feibene Bauernfleidungen und feidene Domino, vor jeden Ball befonbers a i fl. ju verlehnen, bann ju verlaufen Manne. Mafquen bas Stud a 36 fr. : Doming Masquen a 30 fr. ; Krauensimmer Masquen a 24 fr., und meiffe Manns . Banbichu, bas Paar a 18 fc.

Stem find Mafqueraben. Rleiber ju befommen, ben bem Sud Umichel Lom Beffunger, webnhaft in ber Rleifch , Schien.

Rerner find ju befommen propre Mafqueraden Rleider, ben ben benben Guben Gochem Benum, wohnhaft in ber Schlofigaf, und Daniel Benum, mobnhaft in ber grofen Ochsengaß.

Angetommene frembde Gerren Passagiers.

27om 26ften Dec. 1767. bis ben 2ten Jan. 1768. herr Soller, Feuerwerfer aus Frankreich, ben 21. Dec. berfe log. im Bren Balanterieframer, mit allerhand Waar, Engel. Bert Beit, Jager von Sachfen: Botha, ben 25. Dec. , log. in der Rron. Berr Bourgon, und Berr Boffe, swey Fangolifche Raufleute, Den 28. Dec. log in ber Kron.

Berr von Morell, ein Schweigercapitatn aus Bern, ben 27. Dec., log. im Ochfen. Berr Jofeph, ein Galanteriertdmer, log. im Storch.

Extra logirenb.

herr Refident und Sofrath be Meuffville, von Kranffurt, den 30. Dec. berein pafirt, log ben Berrn Lieutengnt De Reuffville.

216 und durchgereifte Gerren Paffagiers. Dere Baron von Gemmingen, Sachfengothaifcher Gefandter ju Wettlar, ben 26, und 31ften Dec.

CER: 3.76 % Diejenige herren Liebhaber, welche fich bes neuen, in Rupfer geflochenen, anddigft privilegirten BBand Calenders, auf Diefes neue Sahr 1768. bedienen

mollen, merden hiermit ergebenft benadrichtiget, bag nur noch, megen bem flats fen Abagna, ben berfelbe feiner Bequemlichfeit haiber, auswärtig gefunden bat, eine geringe Ingabl guter Eremplarien beb mir ju baben find, Das Stud toftet, 5. M. Boffer. wie befannt . 12 Rreuter.

Es wird hiermit befannt gemacht, bag burch mich Beorg Bolfgang Steege mener, ale Bader und Chiruraus, Das neue Bad anwiederum befett worden ift, man perforicht auch tualeich biermit in allen Stucken aute Bedienungen. Es ift auch zu befommen um einen civilen Breis, als ertra gutes Rofenwaffer Die Maas por 24 Rreuger, fodann fdmara Rirldenmaffer, und weiß Eilienmaffer 2c. 2c.

Bu verlebnen werden angetragen:

Ben bem Caffetier Becker Dabier, find ju benen bevorftebenben Mafquens Bolle, feibene Bauernfleidungen und feibene Domino, vor ieden Ball beson. berg a ift, su perlebnen, bann ju verfaufen Manne, Mafauen bas Stuck a 36 fr.; Domino, Mafquen a 30 fr.; Rrauensimmer, Mafquen a 24 fr., und meiffe Manns , Sanbichu, bas Vaar a 18 fr.

Stem find Mafqueraben Rleiber zu befommen, ben bem Jud 2Imfchel gom Beffunger , wohnhaft in ber Bleifch , Schirn.

Rerner find zu befommen propre Mafqueraden : Rleiber, ben ben benben Studen Jochem Depum, wohnhaft in der Schlofgaß, und Daniel Depum, mobnhaft in ber grofen Ochsengaß.

Ungefommene frembde Gerren Passagiers.

Mom 26ften Dec. 1767. bis ben 2ten Jan. 1768.

Berr Coller, Feuerwerfer aus Granfreich, ben 21. Dec. Diefe log. im Berr Beegmann, Orgelmacher von Granffurt, ben 31. Dec. Diefe log. im Drep Balanterieframer, mit allerhand 2Baar,

herr Beis, Jager von Sachfen, Botha, ben 25. Dec., log. in ber Rron. Berr Bourgon, und Derr Boffe, swen Sangofifche Raufleute, Den 28. Dec. loa in ber Rron.

Berr pon Morell, ein Schweigercapitain aus Bern, ben 27. Dec., log. im Deblen. Berr Gofeph, ein Balanterieframer, log. im Storch.

Extra logirend. Berr Refibent und Sofrath be Reuffville, von Frankfurt, Den 30. Dec. berein pagirt, log ben Berrn Lieutenant De Deuffville.

21b = und durchgereiste Gerren Passagiers.

Berr Baron von Gemmingen, Sachfengothaifcher Gefandter ju Wehlar, ben 26, und 31ften Dec.

Anwendungs-Beispiel



ond leiblich als ein wahrer Bott und auch warer Menfch befigt/und burch Ehriffi willen die Residens seines Statthalters erhöcht: Sein beilig Suangelium und newes Besas / der gangen Welt offenbahrer/ und also diese nemen Ehristlichen Reichs begriff/weit / weit größer und gewaltiger worden /dann allevorherzegangene vier Monarchien gewest. Also hat der Tert in bemeltem Daniele erstillt: daß nemlich vollet von einem Stein / sohne Menschen Hand / voin Berg abgerissen aller in den großer Berg/so die gange Welt erstüllen solten / daß ist Eecles. Drauß werden missen Willen unt in Reich von orient bis zu oecident, durch seine Aposteln und Jünger / vermittels der verfündigung deß H. Euangelii, nicht nur in Aus unnd Africat sondern in Europa und Indischen Inssullen (welchs noch heutigs tags concinnictwirdt) heilsamlich propagiet worden.

Dasz Capitel.

Augleich aber der Gottliche Prophet/die dren in die vierdte Monarchia/nemlich auffdie zwo Sulfen geset/dem menschlichen Dilde vergliechen. Also hat eskeinem Monttro oder ohn verlecken i sondern einer vollkosteren Treaturen/sonderlich einem Menschen sollen istimitet werden/der dan nur ein Naupt/zwo Hand inzellen geset habet dan einen Naupt/zwo Hand inzellen geset gestellt der Rutter gestellt der Rutter gestellt der Rutter gestellt der Rutter gestellt der gesorbin mer prainer verblieben /es in guter und von glückschlichte gesorbin mehr oder verliger Glieder/durch Gewalt oder Buschoffam bet mehr oder weniger Glieder/durch Gewalt oder Ungehorfam befomen/alsein Monttrum vind aller Beltz gut spott worden ist.

Demnach nundas Jundament gelegt/daß Römische Reich/eis nem Menschen sowiel die imilieudo dißsals zulasset zuwergleichen. Der Mensch aber an Seel vond deib bestehet/vond ohn denschieben seine lebendigen Geist erzeigen fan. Denen aber zuerhalten / megnetlich und gesundt/daß ist seinst such aber zuerhalten / megnetlich wond gesundt/daß ist seinst den wurden die Menschen mit denen eerportalischen Speisen allein/denen deuts vond vondernümstrigen Thieren vergliechen. Also muß das köstliche/zu unterscheid des Hendnischen/daß H. Reich/zu dessen des Speisen seinschen des Hendnischen Hendlich und Hendnischen des Hendnischen Hendnisc



Tesseract Beispiel



Installation: https://tesseract-ocr.github.io/tessdoc/Home.html Workflow:

```
# Aufruf-Schema
tesseract -1 model input output
# Ausgabe aller zur Verfügung stehender Modelle
# Speicherort der Modelle: Tessdata (Default: /usr/share/tesseract/tessdata/)
tesseract --list-langs
# Einsatz von GT4HistOCR-Modell mit PDF-Ausgabe
# Andere Ausgabemöglichkeiten alto, hocr, tsv, txt (Default: txt)
# Um die Ausgabe auf den Terminal auszugeben: output → stdout
tesseract -1 GT4HistOCR input output pdf
# Ausgabe der erweiterten Einstellungsmöglichkeiten
tesseract --help-extra
# Erweiterung des Aufrufs um den PageSegmentationMode (PSM)
# Das Bild ist einspaltig (PSM \rightarrow 4)
tesseract --psm 4 -1 GT4HistOCR input output
# Verarbeitung aller jpg-Dateien in einem Ordner
# Zur Parallelisierung empfohlen: GNUparallel (dabei darauf achten, dass OpenMP nicht aktiv ist)
for i in *.jpq; do echo ${i/.jpq/}; tesseract -1 GT4HistOCR $i ${i/.jpq/};done
# Alternative mit "fd-find"
fd -e jpg -x tesseract -l GT4HistOCR {} {.}
```

OCR-D Beispiel



Installation: https://github.com/OCR-D/ocrd_all Workflow für eine Einzelseite ohne bestehende METS-Datei:

```
# Step 0: Create Workspace & METS file
mkdir -p ~/projects/OCR-D/workshop/2020 02 19/
cd ~/projects/OCR-D/workshop/2020 02 19/
# Create workspace including METS file
ocrd workspace init OCRbw workspace && cd OCRbw workspace
# Step 1: Download jpg image
mkdir ./OCR-D-IMG && \
wget -0 ./OCR-D-IMG/Fraktur 1621.jpg \
https://digi.bib.uni-mannheim.de/~jkamlah/OCRbw-Workshop-2020-02-19/example/ocrd/MAX/Fraktur 1621.jpg
# Step 2: Add image to METS
# Be aware, that the ID and the GROUPID must be identical if the referenced image represents the original image
ocrd workspace add --file-grp OCR-D-IMG --file-id Fraktur 1621 --page-id OCR-D-IMG 0001 \
  --mimetype image/jpeg OCR-D-IMG/Fraktur 1621.jpg
# Install OCR model into Tesseract data path
apt-get install tesseract-ocr-script-frak
# Step 3: Run the workflow (processors)
ocrd process \
  'tesserocr-segment-region -I OCR-D-IMG -O OCR-D-SEG-BLOCK' \
  'tesserocr-segment-line -I OCR-D-SEG-BLOCK -O OCR-D-SEG-LINE' \
  'tesserocr-recognize -I OCR-D-SEG-LINE -O OCR-D-OCR-TESSEROCR -P model frk -P textequiv level word'
```

16.10.2020 15

OCR-D Beispiel



Installation: https://github.com/OCR-D/ocrd_all Workflow für einen Datensatz mit METS:

```
# Create workspace including METS file
ocrd workspace clone -a METS-URL
# Run the recommended workflow (processors)
ocrd process \
  "cis-ocropy-binarize -I OCR-D-IMG -O OCR-D-BIN" \
  "anybaseocr-crop -I OCR-D-BIN -O OCR-D-CROP" \
  "skimage-binarize -I OCR-D-CROP -O OCR-D-BIN2 -P method li" \
  "skimage-denoise -I OCR-D-BIN2 -O OCR-D-BIN-DENOISE -P level-of-operation page" \
  "tesserocr-deskew -I OCR-D-BIN-DENOISE -O OCR-D-BIN-DENOISE-DESKEW -P operation level page" \
  "cis-ocropy-segment -I OCR-D-BIN-DENOISE-DESKEW -O OCR-D-SEG-REG -P level-of-operation page" \
  "tesserocr-deskew -I OCR-D-SEG-REG -O OCR-D-SEG-REG-DESKEW" \
  "cis-ocropy-clip -I OCR-D-SEG-REG-DESKEW -O OCR-D-SEG-REG-DESKEW-CLIP" \
  "cis-ocropy-segment -I OCR-D-SEG-REG-DESKEW-CLIP -O OCR-D-SEG-LINE" \
  "cis-ocropy-clip -I OCR-D-SEG-LINE -O OCR-D-SEG-CLIP-LINE -P level-of-operation line" \
  "cis-ocropy-dewarp -I OCR-D-SEG-CLIP-LINE -O OCR-D-SEG-LINE-RESEG-DEWARP" \
  "calamari-recognize -I OCR-D-SEG-LINE-RESEG-DEWARP -O OCR-D-OCR -P checkpoint /path/to/models/\*.ckpt.json"
```

Literatur



- Weil, S. (2019). Training Fraktur. GitHub. https://github.com/tesseract-ocr/tesstrain/wiki
- Weil, S. (2019). Vom Bild zum Text.
 Automatisierte Texterkennung in historischen Drucken mit der freien Software Tesseract.
 https://nbn-resolving.org/urn:nbn:de:0290-opus4-163511
- Weil, S., & Zumstein, P. (2016). Mit freier Software Text in Digitalisaten erkennen. https://speakerdeck.com/zuphilip/mit-freier-software-text-in-digitalisaten-erkennen-ocr-praxis-an-der-ub-mannheim

16.10.2020 17

Bildquellen



- Titelseite:
 - https://pixabay.com/photos/letter-handwriting-written-ink-447577/ https://pixabay.com/photos/wash-angle-hook-book-printing-705674/
- Vektorgrafiken:
 https://pixabay.com/de/vectors/flach-design-symbol-icon-www-2126884/
 https://pixabay.com/de/vectors/flach-design-symbol-icon-www-2126884/
 https://pixabay.com/de/vectors/flach-design-symbol-icon-www-2126880/
 https://pixabay.com/de/vectors/werkzeug-schraubenschl%C3%BCssel-3456474/
- Schriftarten: https://typo-info.de/entwicklung-der-schrift/
- GitHub Logos: https://github.com/logos
- Tux: https://de.wikipedia.org/wiki/Linux#/media/Datei:Tux.svg
- Win10: https://commons.wikimedia.org/wiki/File:Windows_10_Logo.svg
- Mac Logo: https://commons.wikimedia.org/wiki/File:Apple_logo_black.svg
- OCR-D Logo: https://www.ocr-d.de/
- DFG-Logo: https://www.dfg.de/service/logo_corporate_design/index.html