

## Optimierter Einsatz von OCR-Verfahren – Tesseract als Komponente im OCR-D-Workflow

# Neue Frakturmodelle für Tesseract



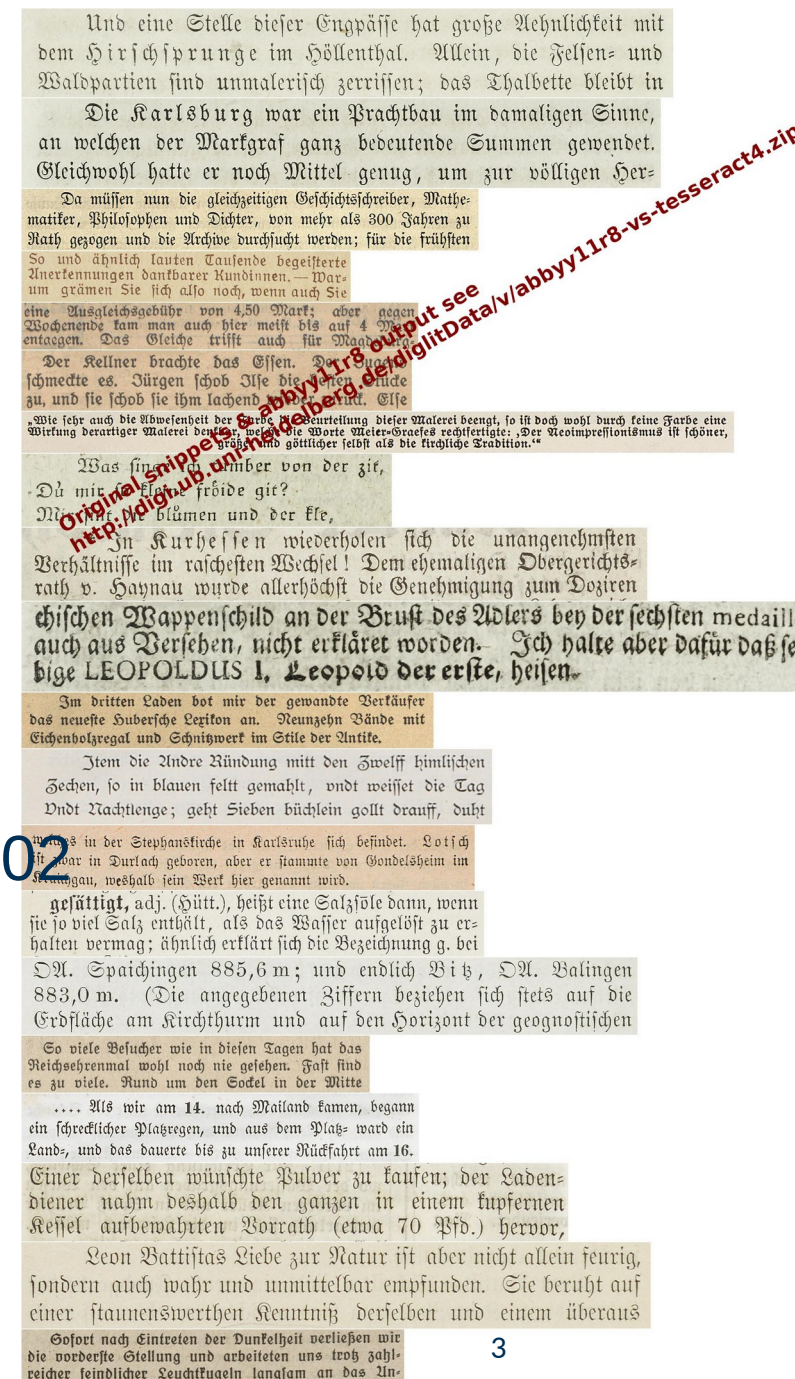
Stefan Weil  
Universitätsbibliothek Mannheim

# Tesseract OCR

- Open Source
- Komplettlösung „All-in-1“
- 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)
- Aktive Weiterentwicklung u. a. im DFG-Projekt OCR-D
- OCR-D-Module für Binarisierung, Segmentierung, Texterkennung und weitere

# Tesseract Fraktur OCR 2018

- Drei Modelle zur Auswahl: deu\_frak, frk, Fraktur
- 14. Juni 2018: Tesseract schneidet im Frakturtest von Jochen Barth (UB Heidelberg) im Vergleich zu ABBYY FineReader schlecht ab:  
<https://github.com/tesseract-ocr/tessdata/issues/102>
- Ligaturen wie ch und ck werden häufig als Kleiner- bzw. Größerzeichen erkannt und weitere systematische Erkennungsfehler



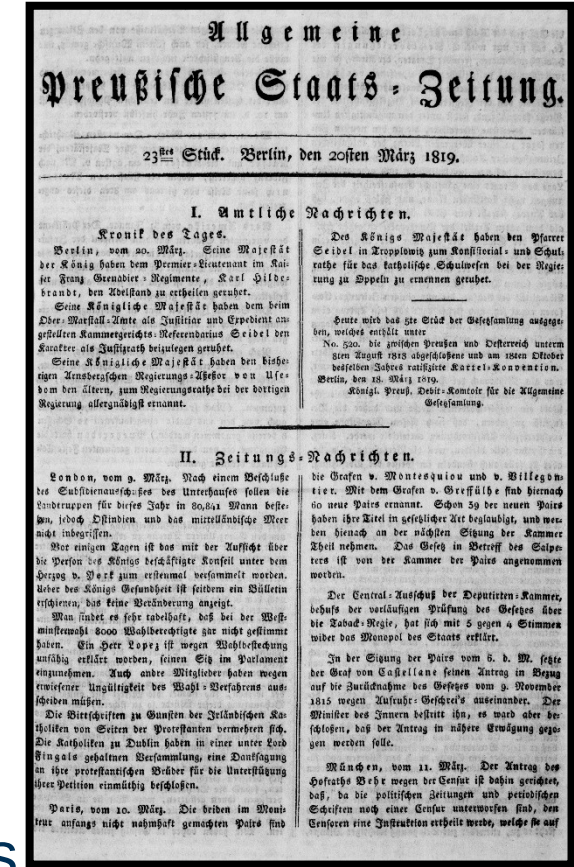
18.11.2019

- == manual/badenia1860.1862 - 354.txt ==  
UndNied eine Stelle dieser Engasse hat große Aehnlichkeit mit  
den Hirschsprünge Hirschsprünge im Hölenthal. Allein, Allein die Felsen: Felsene und  
Waldpartien sind unmalersch zerrissen; das Thalbette bleibt in
- Abbyy12r3-OldGerman vs Fraktur\_5000000\_0.466 (Stefan Weil)
- == manual/badenia1860.1862 - 491.txt ==  
DieDie Karlsruhe war ein Prachtbau in damaligen Sinne,  
an welchen der Markgraf ganz bedeutende Summen gewendet.  
Gleichwohl hatte er noch Mittel genug, um zur völligen Her- Her-
- == manual/boisseree1862bdl - 175.txt ==  
Daba müssen nun die gleichzeitigen Geschichtsschreiber, Mathe: Mathe-  
matiker, Philosophen und Dichter, von mehr als 300 Jahren zu  
Rath gezogen und die Archive durchsucht werden; für die frühesten
- == manual/badenwartel1935.1936 - 0700b.txt ==  
SoSo und ähnlich lauten Tausende begeisterte  
Anerkennungen dankbarer Aundinnen.--war: Kundinnen -- War-  
um grünen Sie sich also noch, wenn auch Sie
- == manual/hnn1936 - 057.11.txt ==  
eineeine Ausgleichsgebühr von 4,50 Mark; aber gegen  
Wochenende kan man auch hier meist bis auf 4 Mark  
entgegen. Mark  
entgegen. Das Gleiche trifft auch für Magdeburg: Magdeburg-
- == manual/hnn1936 - 069.Feierstunde.12.03.txt ==  
DerDer Kellner brächte brachte das Essen. Der Jugend  
schmeckte es. Jürgen Dürgen schob Ilse die besten Stücke  
zu, und sie schob sie ihm Lachend wieder zurück. Elfe Elfe
- == manual/kunstwart.kulturwart26.2 - 152b.txt ==  
WieWie sehr auch die Abwesenheit der Farbe die Beurteilung dieser Malerei Malerei beeengt, so ist doch wohl durch keine Farbe eine  
Wirkung derartiger Malerei denkbar, welche die Worte Meier-Graefes Meier-Graefes rechtfertigte: „Der Neopressionismus Der Neopressionismus ist schöner,  
größer und göttlicher selbst als die kirchliche Tradition.“ Tradition.“
- == manual/malk1809 - 429.txt ==  
Was FingWas singe ich Rumber tumber von der Zik, zit,  
Du mir so kleine sroude froude git?  
Nkir fint  
Mir sint die blümen blümen und der Kle Kle,
- == manual/manheimer.anzeiger1858 - 040.01.txt ==  
\* ImIm Kurhessen wiederholen sich die unangenehmsten  
Verhältnisse in raschesten Wechsel! Den ehemaligen Obergerichts- Obergerichts-  
rath v. Haynau wurde allerhöchst die Genehmigung zum Doziren
- == manual/medaillen1737 - 00.429.txt ==  
chischenschen Wappenschild ander an der Brust des Adlers bey der sechsten  
auch aus Versehen, nicht erklärt werden. Ich halte aber dafür daß se-  
hige 1-LOKOb.vila 1. Leopold medaille,  
-e.,  
hige LEOPOLDIS I. Leopold der erste, heisere heisen.
- == manual/meggendorfer100 - 075.txt ==  
Jede dritte Läden bot mir der gewante Verkäufer  
das neueste Lubersche Lubersche Lexikon an. Neunzehn Bände mit  
Eichenholzregal und Schnitzwerk in Stile der Antike.
- == manual/mittelhochsl - 20.11.txt ==  
ZehnZehn ... in blauen feltt gemahlt, vndt weissert die Tag  
vndt Nachtlinge;  
Vndt Nachtlinge; geht Sieben bücklein gollt drauff, duht
- == manual/mone1889bdl18 - 217.txt ==  
welcheswelches in der Stephanskirche in Karlsruhe sich befindet. Lutsch  
ist zwar in Durlach geboren, aber er stammte von Gondelsheim in  
Kraichgau, weshalb sein Werk hier genannt wird.
- == manual/mothes1882bd2 - 428.txt ==  
gesättigt, aas. (Hütt.), gesättigt, adj. (Hütt.), heißt eine Salzsölze Salzsölze dann, wenn  
sie so viel Salz enthält, enthält, als das Wasser aufgelöst zu er-  
halten vermag; ähnlich erklärt sich die Bezeichnung g. bei
- == manual/oab.balingen1880 - 302.txt ==  
OA.OA Spaichingen 885,6 m; 885, 6m; und endlich Bitz, OA. Balingen  
883,0 m. Biz o. Balingen  
883,0s. (Die angegebenen Ziffern beziehen sich stets auf die  
Erdofläche am Kirchthurm und auf den Horizont der geognostischen
- == manual/olympia.zeitung1936 - 066.txt ==  
SoSo viele Besucher wie in diesen Tagen hat das  
Reichsheymal wohl noch nie gesehen. Fast sind  
es zu viele. Rund um den Sockel in der Mitte
- == manual/rj/bkg1840 - 080.txt ==  
Als wir am 14. nach Mailand kamen, begann  
ein schrecklicher Platzregen, und aus dem Platz- ward ein  
Land- und das dauerte bis zu unserer Rückfahrt am 16.
- == manual/schweizer.wochenblatt1863 - 021.3.txt ==  
EinerEiner derselben wünschte Pulver zu kaufen; der Läden- Läden-  
diener nahm deshalb den ganzen in einem kupfernen  
Kessel aufbewahrten Vorrath (etwa 70 Pfd.) hervor,
- == manual/springer1886bd1 - 000.274.txt ==  
LeonLeon Battistas Liebe zur Natur ist aber nicht allein feurig,  
sondern sehr wohl und unmittelbar empfunden, empfunden. Sie beruht auf  
einer stannswürthen staunswürthen Kenntniß derselben und einem überaus
- == manual/wehrmacht1936.1937 - 102.12.txt ==  
SofortGofort nach Eintreten der Dunkelheit verließen mir verließen wir  
die vorderste Stellung und arbeiteten uns trotz zahl: zahl:  
einer feindlicher Leuchtschein Leuchtschein an den
- == manual/wehrmacht1936.1937 - 102.12.txt ==  
SofortGofort nach Eintreten der Dunkelheit verließen mir verließen wir  
die vorderste Stellung und arbeiteten uns trotz zahl: zahl:  
einer feindlicher Leuchtschein Leuchtschein an den



# Tesseract Fraktur-Modelle der UB Mannheim

- UB Mannheim benötigt Fraktur OCR u. a. für die historische Zeitung *Deutscher Reichsanzeiger und Preussischer Staatsanzeiger* (ursprünglich *Allgemeine Preussische Staatszeitung*)
- Training neuer Fraktur-Modelle aus publizierten Ground-Truth-Daten seit September 2019:
  - GT4HistOCR (Springmann et al.)
  - Austrian Newspapers (Mühlberger et al.)
- Zeichenerkennung besser als 97,5 %, Worterkennung bei 93 % für Austrian Newspapers



# Tesseract Fraktur-Training

Beispiele für Zeilen aus GT4HistOCR

*lúxχ , lúxχ ! Wie blißen ihre großen Augen! Noch  
auch nit fürchten / das in  
Den turcken flyessen·die wurden von stünd an erschauet·*

Beispiele für Zeilen aus Austrian Newspapers

**Wagner'schen Univ.-Buchhandlung in Innsbruck**  
*mußte man zahlen, wenn*  
**Neue Freie Presse.**

*weil ein Stück von der Wagschale ausgebrochen war, dasselbe*

# Literatur

- Springmann, U., Reul, Ch., Dipper, S., & Baiter, J. (2018). GT4HistOCR: Ground Truth for training OCR engines on historical documents in German Fraktur and Early Modern Latin (Version 1.0) [Data set]. Zenodo. <http://doi.org/10.5281/zenodo.1344132>
- Mühlberger, G. & Hackl, G. (2019). NewsEye / READ OCR training dataset from Austrian Newspapers (19th C.) [Data set]. Zenodo. <http://doi.org/10.5281/zenodo.3387369>
- Weil, S. (2019). Training Fraktur. GitHub.  
<https://github.com/tesseract-ocr/tesstrain/wiki>
- Metzger, N. & Weil, S. (2019). Optimierter Einsatz von OCR-Verfahren - Tesseract als Komponente im OCR-D-Workflow.  
<https://nbn-resolving.org/urn:nbn:de:bsz:180-madoc-522132>
- 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>

# Bildquellen

- Titelseite:  
<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/>
- Abbildungen von Jochen Barth, UB Heidelberg:  
[https://digi.ub.uni-heidelberg.de/diglitData/v/abby12r3-OldGerman-vs-Fraktur5000000\\_0.466\\_II.png](https://digi.ub.uni-heidelberg.de/diglitData/v/abby12r3-OldGerman-vs-Fraktur5000000_0.466_II.png)  
<https://camo.githubusercontent.com/47ac160cf86bb8f69fe98677112e5597a46da3ed/687474703a2f2f646967692e75622e756e692d68656964656c626572672e64652f6469676c6974446174612f762f6162627979313172382d76732d746573736572616374342e6a7067>
- DFG-Logo: <https://www.dfg.de/>
- GitHub Logos: <https://github.com/logos>
- OCR-D Logo: <http://www.ocr-d.de/>