



# Volltext für Wissenschaft und Lehre auf e-rara.ch Erfahrungen und Probleme mit Fraktur-OCR

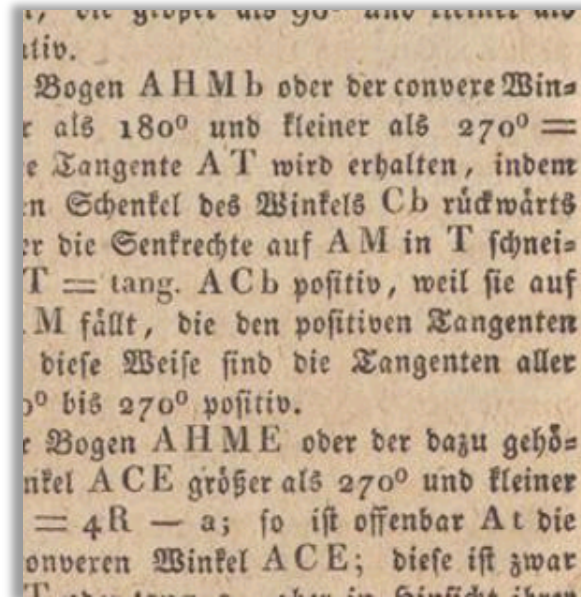
7. Bibliothekskongress Leipzig

18. März 2019

Oliver Ammann

# Agenda

- Die Plattform e-rara.ch
- Ausgangslage «e-rara.ch: Volltext» (2016-2017)
- Erweiterung Volltexterkennung: Erfahrungen, Probleme, Ergebnisse
- Fazit
- Perspektive



stumpfen Winkel, die grösser als  $90^\circ$  und  $180^\circ$  sind, negativ.

Es sey der Bogen  $ARMb$  oder der convex  
 kleiner als  $270^\circ$  —

$aR - f$  a; seine Tangente  $AT$  wird erhalte  
 man den zweiten Schenkel des Winkels  $C$   
 verlängert, bis er die Senkrechte auf  $AM$   
 positiv, weil sie auf

die Seite von  $AN$  fällt, die den positiven  
 entspricht; auf diese Weise sind die Tang  
 Winkel von  $180^\circ$  bis  $270^\circ$  positiv.

Es sey der Bogen  $AHME$  oder der dazu g  
 $270^\circ$  und kleiner

als  $360^\circ$ , oder  $-4B - a$ ; so ist offenbar  
 Tangente vom convexen Winkel  $ACB$ ; di

# Die Plattform e-rara.ch

- Plattform für digitalisierte Drucke aus Schweizer Bibliotheken
- 2007 als Projekt gestartet
- Seit 2010 online
- 5 Träger, 19 Institutionen
- 71.000 Titel online, davon 17.000 Titel der ETH-Bibliothek
- Weiterentwicklungen
  - Kollektion Privatbibliotheken
  - Volltext

The screenshot shows the e-rara.ch search results page. The search bar at the top contains 'Suche in e-rara...'. Below the search bar, there are navigation options: 'HOME', 'Alle Kollektionen', and 'Alle Titel'. The main content area displays 17198 titles, sorted by 'Titel' (Title) in ascending order. The first three results are visible, each with a thumbnail image of the book cover and a 'VOLLTEXT DURCHSUCHBAR' (Full text searchable) indicator.

**Search Results:**

- Result 1:** On the free motion of points, and on universal gravitation, including the principal propositions of books I. and III. of the principia. Whewell, William [1794-1866]. In: A treatise on dynamics. Cambridge : Deighton ; Whittaker ; London, M.DCCC.XXXII. [1832].
- Result 2:** On the motion of points constrained and resisted, and on the motion of a rigid body. Whewell, William [1794-1866]. In: A treatise on dynamics. Cambridge : Deighton ; Whittaker ; London, M.DCCC.XXXIV. [1834].
- Result 3:** A treatise on forest-trees containing not only the best methods of their culture hitherto practised, but a variety of new and useful discoveries, the result of many repeated experiments : as also, plain directions for removing most of the valuable kinds of forest-trees, to the height of thirty feet and upwards, with certain success; and, on the same principles, (with as certain success) for transplanting hedges of sundry kinds, which will at once resist cattle : to which are added, directions for the disposition, planting, and culture of hedges, by observing which, they will be handsomer and stronger fences in five years, than they now usually are in ten. Boutcher, William [18. Jht]. Edinburgh : printed by R. Fleming, and sold by the author, by J. Murray, MDCCLXXV [1775].
- Result 4:** A Treatise on the Management of Bees wherein is contained the Natural History of those Insects : with the various methods of cultivating them, both Antient and Modern and the improved Treatment of them : to which are added the Natural History of Wasps and Hornets, and the Means of destroying them.

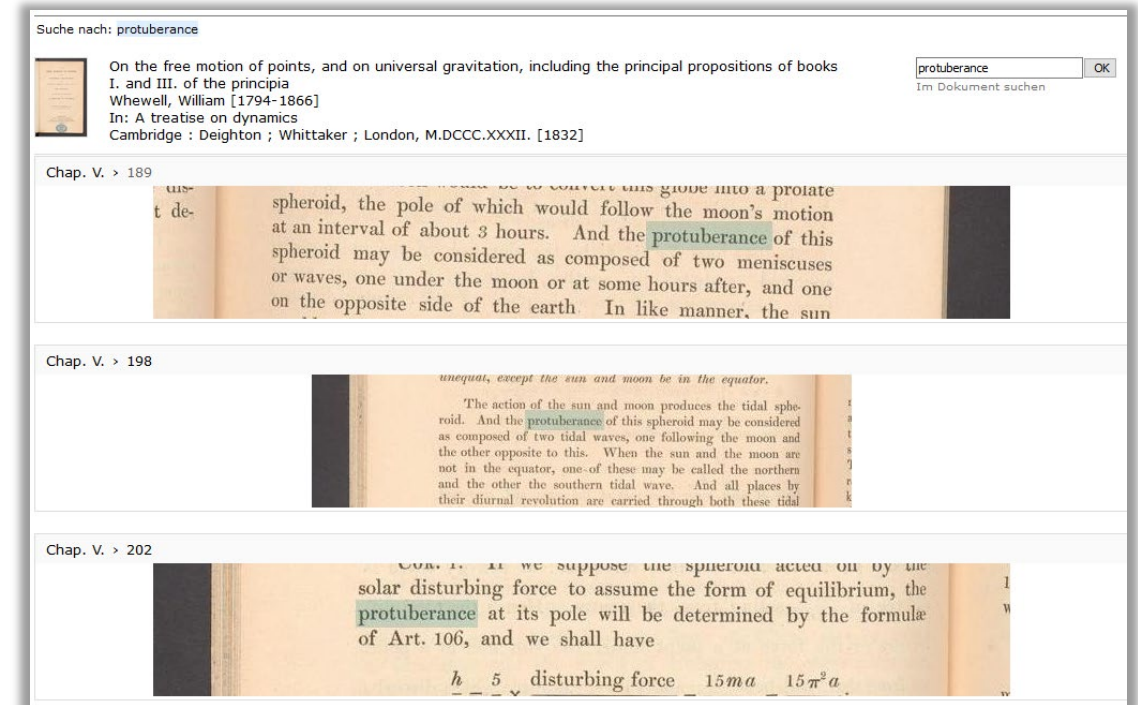
**Filtering and Statistics:**

- Sortieren nach:** Titel
- Max. Trefferanzahl:** 10
- Reihenfolge:** aufsteigend
- Blättern:** 11 - 20
- Treffer eingrenzen:**
  - Volltext:** Volltext durchsuchbar 8354
  - Autoren, Beteiligte:** Euler, Leonhard 87; Frobenius, Georg 83; Keller, Heinrich 65; Linné, Carl von 49; Martius, Carl Friedrich Philipp 44.
  - Dokumententypen:** Buch 15889; Karte 1308; Musikdruck 1.
  - Zeiträume:** 1401-1500 13; 1501-1600 621; 1601-1700 1655; 1701-1800 4922; 1801-1900 9783; 1901-2000 46.
  - Sprachen:** Deutsch 8389; Französisch 3973; Latein 2428; Englisch 1305; Italienisch 1158.
  - Druckorte:** Paris 2357; Leipzig 1128.

Digitalisate der ETH-Bibliothek Zürich auf e-rara.ch [Screenshot]

# Ausgangslage

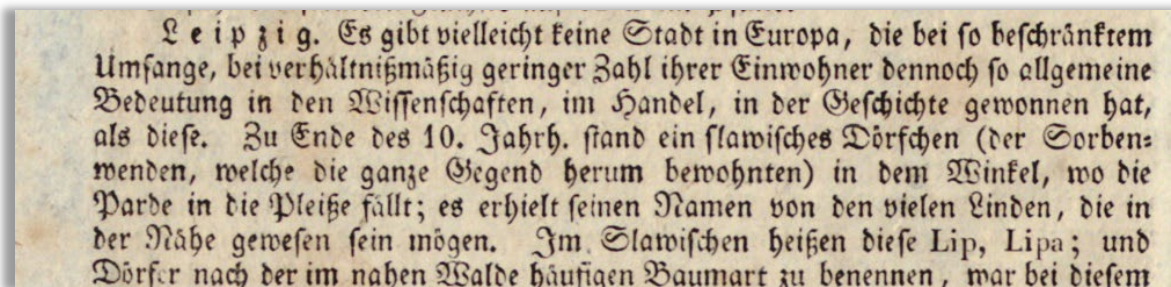
- «e-rara.ch: Volltext» (2016 – 2017)
- OCR für 3.700 Titel ab Druckjahr 1830
- nur Antiqua-Schriften
- ABBYY FineReader
- Volltextsuche online und in PDFs
- Treffer in Snippets hervorgehoben



Suchbegriff hervorgehoben in Snippets. [Screenshot e-rara.ch]

# Übersicht Erweiterung Volltexterkennung (2018)

- Volltext für 6.000 Titel (~ 3 Mio. Seiten)
- ab Druckjahr **1801**
- Antiqua- und **Frakturschrift**
- ABBYY FineReader



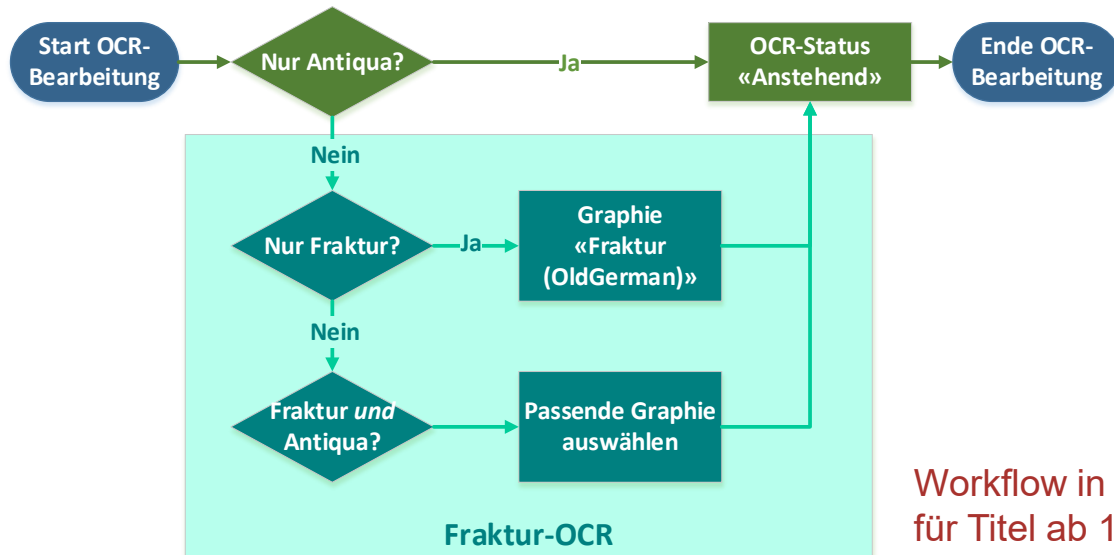
Eintrag «Leipzig» in Frakturschrift aus: *Allgemeine deutsche Real-Encyclopädie [...]*, Reutlingen 1830-1831, Seite 531. ETH-Bibliothek Zürich, Rar 3421, <http://doi.org/10.3931/e-rara-17880>.

	Antiqua	Fraktur
1900	OCR ✓	Erweiterung
1830	Erweiterung	Erweiterung
1801	keine OCR	keine OCR
1450	keine OCR	keine OCR

OCR-Bearbeitung der Teilbestände auf e-rara.ch

# Umsetzung & Workflow

- Drei Phasen
  1. 3.000 Titel in Frakturschrift nach 1830
  2. 3.000 Titel nach 1801 und vor 1830: Sichtung und Feststellung der Graphie
  3. Workflow für neu hinzukommende Titel



Workflow in VLM für Titel ab 1801

	Antiqua	Fraktur
1900	OCR ✓	1. Phase
1830	2. Phase	2. Phase
1801	keine OCR	keine OCR
1450	keine OCR	keine OCR

OCR-Bearbeitung der Teilbestände auf e-rara.ch: Vorgehen

# Metadaten und Graphie\*

Katalog, MARC-Felder 008, 041: **Sprachcode**

008 990902s1818--gw-----00--ger-d



VLM: **Graphie-Einstellung**  
(Kombination aus Graphie und Sprache)

OCR

Graphie	Fraktur (OldGerman)
Status	Antiqua (RussianOldSpelling)
Bemerkung	Antiqua (Spanisch) Antiqua (Swedish) Fraktur (Danish)
Info	Fraktur (Norwegian) <b>Fraktur (OldGerman)</b> Fraktur+Antiqua (English, French, Italian, Latin, OldGerman)
Nutzungsbedingur	Fraktur+Antiqua (English, French, OldGerman)
Filterlizenz	Public Domain Mark

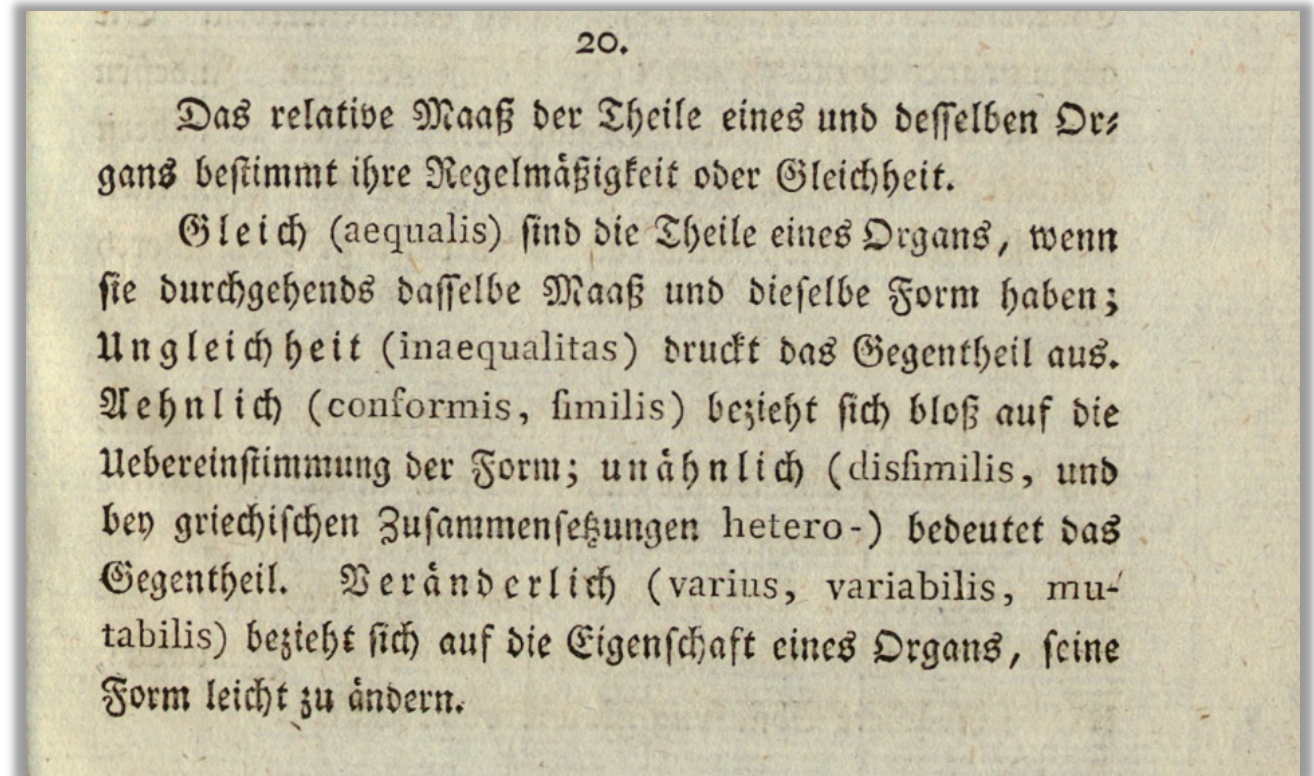
- keine Graphie-Metadaten

- manuelle Auswahl der passenden Graphie

\* Einstellung für Schriftart und Sprache bei der OCR-Bearbeitung mit ABBYY FineReader und in Visual Library Manager

# Erwartete Probleme mit OCR für Fraktur+Antiqua 1801+

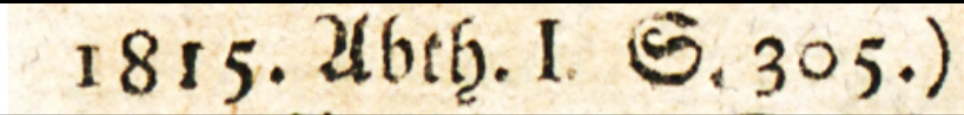
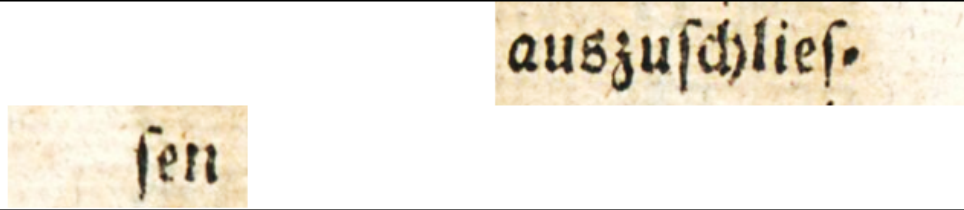
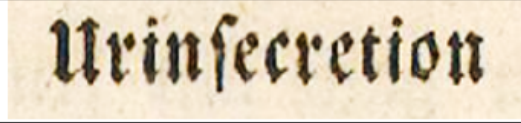
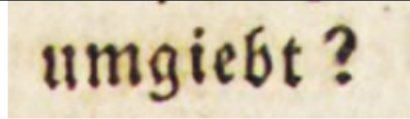
- keine Graphie-Metadaten
- schlechte Druckqualität
- (unterschiedliche) Frakturschriften
- Antiqua-Fraktur-Mischtexte



Mischtext in: *A. P. de Candolle's und K. Sprengel's Grundzüge der wissenschaftlichen Pflanzenkunde [...]*. Leipzig 1820, Seite 15.  
ETH-Bibliothek Zürich, Rar 3722, <http://doi.org/10.3931/e-rara-18699>.



# Ergebnisse I

Fehlertyp	Beispiel	OCR-Ergebnis
Ziffern		i8rz. Abth. I S. zc>;.)
Worttrennung		auszuschlies. sen
c ↔ e		Urinseeretion
u ↔ n		nmgiebt?

Auswahl typischer Fehler bei Fraktur-OCR auf e-rara.ch

# Ergebnisse II

Johann Friedrich Lacomus: *Umfang und Eintheilung der Prospektive*. Königsberg 1804, Seite 19. Zentralbibliothek Zürich, NE 1859,3, <http://doi.org/10.3931/e-rara-52262>

dem Künstler zu Hülfe. Die Freyhand-Zeichnungskunst ist also nur Mitgehülfin, aber nicht schlechterdings notwendige Mitarbeiterin der Prospektive. Die Freyhand-Zeichnungskunst taugt eigentlich nur zur Vorstellung irregulärformiger Gegenstände; sie hüte sich aber ja mit ihrer Sicherheit prahlen zu wollen, welche nur die Gefährtin der Prospektive ist. Die Freyhand-Zeichnungskunst schöpft ihre Gewißheit aus unsichern Quellen, nämlich dem Augenmaße und der Fertigkeit der Hand; die Prospektive hingegen aus der Mathematik, d. h. aus einer reinen, sichern Quelle der Wahrheit.

Da die Produkte der Linien-Prospektive, ohne Aenderung der Sehwinkel, entweder kleiner, gleich groß (wenn der Gegenstand eine Fläche ist) oder größer ausfallen können als der abgebildete Gegenstand, je nachdem die Bildfläche angenommen wird, so zerfällt die Linien-Prospektive auf Flächen, nach der Entfernung der Bildfläche in Ansehung der des Gegenstandes und des Orts des Auges, in Linien-Perspektive auf Flächen, wenn sich die Tafel zwischen Aug und Gegenstand befindet; Linien-Planospek-

dem Künstler zu Hülfe. Die Freyhand-Zeichnungskunst ist also nur Mitgehülfin, aber nicht schlechterdings notwendige Mitarbeiterin der Prospektive. Die Freyhand-Zeichnungskunst taugt eigentlich nur zur Vorstellung irregulärformiger Gegenstände; sie hüte sich aber ja mit ihrer Sicherheit prahlen zu wollen, welche nur die Gefährtin der Prospektive ist. Die Freyhand-Zeichnungskunst schöpft ihre Gewißheit aus unsichern Quellen, nämlich dem Augenmaße und der Fertigkeit der Hand; die Prospektive hingegen aus der Mathematik, d. h. aus einer reinen, sichern Quelle der Wahrheit.

Da die Produkte der Linien-Prospektive, ohne Aenderung der Sehwinkel, entweder kleiner, gleich groß (wenn der Gegenstand eine Fläche ist) oder größer ausfallen können als der abgebildete Gegenstand, je nachdem die Bildfläche angenommen wird, so zerfällt die Linien-Prospektive auf Flächen, nach der Entfernung der Bildfläche in Ansehung der des Gegenstandes und des Orts des Auges, in Linien-Perspektive auf Flächen, wenn sich die Tafel zwischen Aug und Gegenstand befindet; Linien-Planospek-

# Ergebnisse III: Frakturschrift

Testseite	WG	ZG
<a href="https://www.e-rara.ch/zut/content/pageview/11419715">https://www.e-rara.ch/zut/content/pageview/11419715</a> (1817)	86,6	97,2
<a href="https://www.e-rara.ch/zut/content/pageview/4223294">https://www.e-rara.ch/zut/content/pageview/4223294</a> (1826)	93,9	98,7
<a href="https://www.e-rara.ch/sikjm/content/pageview/2993807">https://www.e-rara.ch/sikjm/content/pageview/2993807</a> (1802)	96,1	99,1
<a href="https://www.e-rara.ch/zut/content/pageview/19329070">https://www.e-rara.ch/zut/content/pageview/19329070</a> (1819)	96,2	99,1
<a href="https://www.e-rara.ch/zuz/content/pageview/14339618">https://www.e-rara.ch/zuz/content/pageview/14339618</a> (1804)	93,9	99,2
<a href="https://www.e-rara.ch/zut/content/pageview/4597935">https://www.e-rara.ch/zut/content/pageview/4597935</a> (1802)	97,5	99,2
<a href="https://www.e-rara.ch/zuz/content/pageview/8784910">https://www.e-rara.ch/zuz/content/pageview/8784910</a> (1826)	96,1	99,3
<a href="https://www.e-rara.ch/zuz/content/pageview/13231297">https://www.e-rara.ch/zuz/content/pageview/13231297</a> (1823)	96,5	99,5
<a href="https://www.e-rara.ch/zut/content/pageview/16335222">https://www.e-rara.ch/zut/content/pageview/16335222</a> (1818)	98,9	99,6
<b>Mittelwert 9 Testseiten</b>	<b>95,1</b>	<b>99,0</b>

theils der sogenante blaue Quarz von der Oberpfalz, theils endlich das blaue birien, was, mit Feldspath gemengt, ben gefunden worden ist. Davon ist aber sen der blaue wirkliche Quarz von Ubo in Finland, welcher auch Saphirgen heilic genant worden ist (darüber Leonhard 1815. Abth. I. S. 305.) und der blaue von Golling in Salzburg, der auch

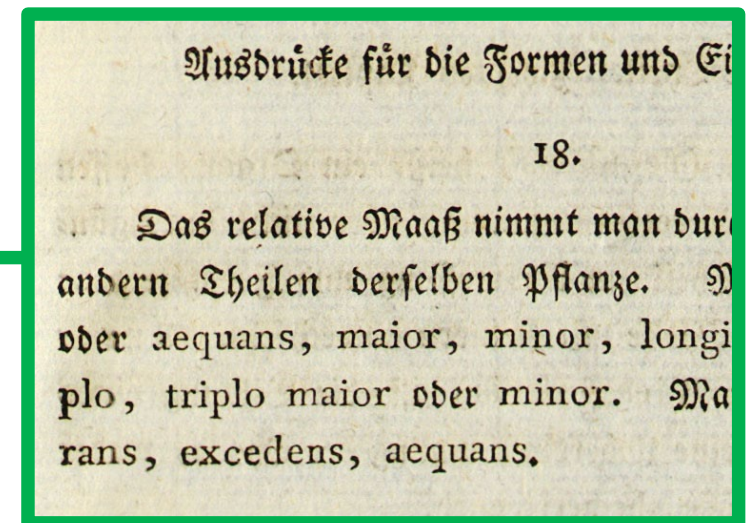
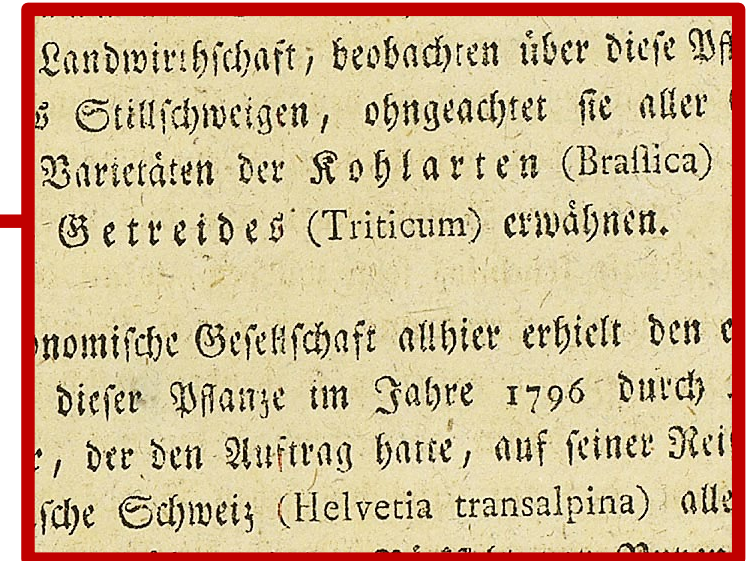
Hingegen bei dem Ausmessen auf man sich weniger, weil man dieselbe mehrmal nacheinander messen kann. Am besten ist, wenn man beide Me der verbindet. Daß man nemlich die auf dem Felde gemessen hat, nachher wieder nachmift. Man sieht dann, ob getragen, und sich an den Ruthen nicht verschrieben hat.

Auszählung Erkennungsfehler in Testseiten aus Testsample «Fraktur 1801-1830»;  
WG: Wortgenauigkeit in %, ZG: Zeichengenauigkeit in %. Groß- und Kleinschreibung,  
Sonderzeichen und Ziffern sind nicht berücksichtigt.

# Ergebnisse IV: Mischschrift

Testseite	WG	ZG
<a href="https://www.e-rara.ch/zuz/content/pageview/8266791">https://www.e-rara.ch/zuz/content/pageview/8266791</a> (1804)	91,8	97,6
<a href="https://www.e-rara.ch/zut/content/pageview/5253211">https://www.e-rara.ch/zut/content/pageview/5253211</a> (1810)	90,7	97,8
<a href="https://www.e-rara.ch/zut/content/pageview/18577029">https://www.e-rara.ch/zut/content/pageview/18577029</a> (1821)	90,6	97,9
<a href="https://www.e-rara.ch/zuz/content/pageview/10777711">https://www.e-rara.ch/zuz/content/pageview/10777711</a> (1803)	90,8	98,0
<a href="https://www.e-rara.ch/zut/content/pageview/13602392">https://www.e-rara.ch/zut/content/pageview/13602392</a> (1819)	94,6	98,8
<a href="https://www.e-rara.ch/zut/content/pageview/12378413">https://www.e-rara.ch/zut/content/pageview/12378413</a> (1824-1825)	94,2	98,9
<a href="https://www.e-rara.ch/zuz/content/pageview/8759470">https://www.e-rara.ch/zuz/content/pageview/8759470</a> (1805)	95,1	99,0
<a href="https://www.e-rara.ch/bau_1/content/pageview/9770818">https://www.e-rara.ch/bau_1/content/pageview/9770818</a> (1828)	98,1	99,6
<a href="https://www.e-rara.ch/zut/content/pageview/5892271">https://www.e-rara.ch/zut/content/pageview/5892271</a> (1820)	99,0	99,7
<b>Mittelwert 9 Testseiten</b>	<b>93,9</b>	<b>98,6</b>

Auszählung Erkennungsfehler in Testseiten aus Testsample «Fraktur+Antiqua 1801-1830»; WG: Wortgenauigkeit in %, ZG: Zeichengenauigkeit in %. Groß- und Kleinschreibung, Sonderzeichen und Ziffern sind nicht berücksichtigt.



# Fazit

- 11.700 von 71.000 Titeln mit Volltext
- Relativ aufwendige Feststellung der Graphie
  - Bessere Ergebnisse rechtfertigen den Aufwand
  - Kosten für eine kombinierte Lizenz  
Fraktur+Antiqua rund doppelt so teuer wie einfache Antiqua-Lizenz
- NutzerInnen schätzen die guten Volltexte
- Visits e-rara.ch ↑ 33%

	Antiqua	Fraktur
1900	OCR ✓	OCR ✓
1830	OCR ✓	OCR ✓
1801	keine OCR	keine OCR
1450	keine OCR	keine OCR

OCR-Status der Teilbestände auf e-rara.ch  
nach der Erweiterung 2018

# Perspektiven: Aktuelle Weiterentwicklungen e-rara.ch

- IIF-Manifest und neuer Viewer
- JPEG-Download
- Social Media-Share Buttons
- **Volltext-Download (Textdatei)**
- Export Metadaten für Literaturverwaltung im Format .ris  
*(verfügbar später im Frühjahr 2019)*



# Vielen Dank für Ihre Aufmerksamkeit

7. Bibliothekskongress Leipzig | 18. März 2019

ETH Zürich, ETH-Bibliothek, Oliver Ammann, Alte und Seltene Drucke, Rämistrasse 101, 8092 Zürich  
Tel. +41 44 632 49 05, [oliver.ammann@library.ethz.ch](mailto:oliver.ammann@library.ethz.ch), [www.library.ethz.ch](http://www.library.ethz.ch)



Creative Commons - CC BY. Von dieser Lizenz nicht erfasst sind das ETH-Logo sowie die Fotografien auf den Seiten 1 und 15.